1  **Harmonization of resting-state functional MRI data across multiple imaging sites**

2  **via the separation of site differences into sampling bias and measurement bias**

3

4  Ayumu Yamashita[1,*], Noriaki Yahata[1,2,3], Takashi Itahashi[4], Giuseppe Lisi[1], Takashi Yamada[1,4], Naho
5  Ichikawa[5], Masahiro Takamura[5], Yujiro Yoshihara[6], Akira Kunimatsu[7,8], Naohiro Okada[2,9], Hirotaka
6  Yamagata[10], Koji Matsuo[10,11], Ryuichiro Hashimoto[1,4,12], Go Okada[5], Yuki Sakai[1,13], Jun Morimoto[1],
7  Jin Narumoto[1,13], Yasuhiro Shimada[14], Kiyoto Kasai[1,2,9], Nobumasa Kato[1,4], Hidehiko Takahashi[6],
8  Yasumasa Okamoto[5], Saori C Tanaka[1], Mitsuo Kawato[1], Okito Yamashita[1,15,*], and Hiroshi
9  Imamizu[1,16,*]
10
11  [1] Brain Information Communication Research Laboratory Group, Advanced Telecommunications
12    Research Institutes International, Kyoto, Japan
13  [2] Department of Neuropsychiatry, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan
14  [3] Department of Molecular Imaging and Theranostics, National Institute of Radiological Sciences,
15    National Institutes for Quantum and Radiological Science and Technology, Chiba, Japan
16  [4] Medical Institute of Developmental Disabilities Research, Showa University, Tokyo, Japan
17  [5] Department of Psychiatry and Neurosciences, Hiroshima University, Hiroshima, Japan
18  [6] Department of Psychiatry, Kyoto University Graduate School of Medicine, Kyoto, Japan
19  [7] Department of Radiology, IMSUT Hospital, Institute of Medical Science, The University of Tokyo,
20    Tokyo, Japan
21  [8] Department of Radiology, Graduate School of Medicine, The University of Tokyo, Tokyo, Japan
22  [9] The International Research Center for Neurointelligence (WPI-IRCN) at the University of Tokyo
23    Institutes for Advanced Study (UTIAS), Tokyo, Japan
24  [10] Division of Neuropsychiatry, Department of Neuroscience, Yamaguchi University Graduate School of
25    Medicine, Yamaguchi, Japan
26  [11] Department of Psychiatry, Faculty of Medicine, Saitama Medical University, Saitama, Japan
27  [12] Department of Language Sciences, Tokyo Metropolitan University, Tokyo, Japan
28  [13] Department of Psychiatry, Graduate School of Medical Science, Kyoto Prefectural University of
29    Medicine, Kyoto, Japan
30  [14] Brain Activity Imaging Center, ATR-Promotions Inc., Kyoto, Japan
31  [15] Center for Advanced Intelligence Project, RIKEN, Tokyo, Japan
32  [16] Department of Psychology, Graduate School of Humanities and Sociology, The University of Tokyo,
33    Tokyo, Japan
34
35  **\* Correspondence:** imamizu@gmail.com (H.I.), oyamashi@atr.jp (O.Y.), or ayumu722@gmail.com
36    (A.Y.)

37

38  **Abstract**

39  When collecting large neuroimaging data associated with psychiatric disorders, images
40  must be acquired from multiple sites because of the limited capacity of a single site.
41  However, site differences represent the greatest barrier when acquiring multi-site
42  neuroimaging data. We utilized a traveling-subject dataset in conjunction with a multi-
43  site, multi-disorder dataset to demonstrate that site differences are composed of biological
44  sampling bias and engineering measurement bias. Effects on resting-state functional MRI
45  connectivity because of both bias types were greater than or equal to those because of

46  psychiatric disorders. Furthermore, our findings indicated that each site can sample only

47  from among a subpopulation of participants. This result suggests that it is essential to

48  collect large neuroimaging data from as many sites as possible to appropriately estimate

49  the distribution of the grand population. Finally, we developed a novel harmonization

50  method that removed only the measurement bias by using traveling-subject dataset and

51  achieved the reduction of the measurement bias by 29% and the improvement of the

52  signal to noise ratios by 40%.

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

## Introduction

Acquiring and sharing large neuroimaging data have recently become critical for bridging the gap between basic neuroscience research and clinical applications such as the diagnosis and treatment of psychiatric disorders (Human Connectome Project (HCP) [1], [http://www.humanconnectomeproject.org/]; Human Brain Project [https://www.humanbrainproject.eu/en/]; UK Biobank [http://www.ukbiobank.ac.uk/]; and Strategic Research Program for Brain Sciences (SRPBS) [2] [https://www.amed.go.jp/program/list/01/04/001_nopro.html]) [3-5]. When collecting large data associated with psychiatric disorders, it is necessary to acquire images from multiple sites because it is nearly impossible for a single site to collect large neuroimaging data (Connectomes Related to Human Disease (CRHD), [https://www.humanconnectome.org/disease-studies]; Autism Brain Imaging Data Exchange (ABIDE); and SRPBS) [2, 6-8]. In 2013, the Japan Agency for Medical Research and Development (AMED) organized the Decoded Neurofeedback (DecNef) Project. The project determined the unified imaging protocol on 28th February 2014 (http://www.cns.atr.jp/rs-fmri-protocol-2) and have collected multisite resting-state functional magnetic resonance imaging (rs-fMRI) data using twelve scanners across eight research institutes for recent five years. The collected dataset encompasses 2,239 samples and five disorders and is publicly shared through the SRPBS multisite multi-disorder database (https://bicr-resource.atr.jp/decnefpro/). This project has enabled the identification of resting-state functional connectivity (rs-fcMRI)-based biomarkers of several psychiatric disorders that can be generalized to completely independent cohorts [2, 8-10]. However, multisite dataset with multiple disorders raises difficult problems never included in a single-site based dataset of healthy population (e.g., HCP and UK Biobank). That is, our experience in the SRPBS database demonstrated difficulty in fully control of scanner type, imaging protocol, patient demographics [10-13] even if the unified protocol is determined. Moreover, there often exists unpredictable difference in participant population among sites. Therefore, researchers must work with heterogeneous neuroimaging data. In particular, site differences represent the greatest barrier when extracting disease factors by applying machine-learning techniques to such heterogeneous data [14] because disease factors tend to be confounded with site factors [2, 8, 10-13, 15]. This confounding occurs because a single site (or hospital) is apt to

110    sample only a few types of psychiatric disorders (e.g., primarily schizophrenia from site

111    A and primarily autism spectrum disorder from site B). To properly manage such

112    heterogeneous data, it is necessary to harmonize the data among the sites [16-19].

113    Moreover, a deeper understanding of these site differences is essential for efficient

114    harmonization of the data.

115        Site differences essentially consist of two types of biases: engineering bias (i.e.,

116    measurement bias) and biological bias (i.e., sampling bias). Measurement bias includes

117    differences in the properties of MRI scanners such as imaging parameters, field strength,

118    MRI manufacturers, and scanner models, whereas sampling bias refers to differences in

119    participant groups among sites. Previous studies have investigated the effect of

120    measurement bias on resting-state functional connectivity by using a traveling-subject

121    design [20] wherein multiple participants travel to multiple sites for the assessment of

122    measurement bias [7]. By contrast, researchers to date have only speculated with regard

123    to sampling bias. For example, differences in the clinical characteristics of patients

124    examined at different sites are presumed to underlie the stagnant accuracy of certain

125    biomarkers, even after combining the data from multiple sites [12]. Furthermore, to our

126    knowledge, no study has mathematically defined sampling bias or conducted quantitative

127    analyses of its effect size, which is likely because the decomposition of site differences

128    into measurement bias and sampling bias is a complex process. To achieve this aim, we

129    combined a separate traveling-subject rs-fMRI dataset with the SRPBS multi-disorder

130    dataset. Simultaneous analysis of the datasets enabled us to divide site differences into

131    measurement bias and sampling bias and to quantitatively compare their effect sizes on

132    resting-state functional connectivity with those of psychiatric disorders.

133        Furthermore, our detailed analysis of measurement and sampling biases enabled

134    us to investigate the origin of each bias in multisite datasets for the first time. For

135    measurement bias, we quantitatively compared the magnitude of the effects among the

136    different imaging parameters, fMRI manufacturers, and number of coils in each fMRI

137    scanner. We further examined two alternative hypotheses regarding the mechanisms

138    underlying sampling bias: one hypothesis assumes that each site samples subjects from a

139    common population. In this situation, sampling bias occurs because of the random

140    sampling of subjects, which results in incidental differences in the patients' characteristics

141    among the sites. The second hypothesis assumes that each site samples subjects from

142    different subpopulations. In this situation, sampling bias occurs because of sampling from

143    subpopulations with different characteristics. For example, assume multiple sites plan to

144    collect data from the same population of patients with major depressive disorder.

145    Subtypes of major depressive disorder exist within the population such as atypical

146    depression and melancholic depression [21, 22]; therefore, one subpopulation may

147    contain a large proportion of patients with atypical depression, whereas another

148    subpopulation may contain a large proportion of patients with melancholic depression.

149    Therefore, in some instances, atypical depression may be more frequent among patients

150    at site A, whereas melancholic depression may be more frequent among patients at site

151    B. The basic protocol for collecting large-scale datasets differ between these two

152    hypotheses; thus, it is necessary to determine the hypothesis that most appropriately

153    reflects the characteristics of the SRPBS dataset. In the former situation, one would

154    simply need to collect data from a large number of subjects, even with a small number of

155    sites. In the latter situation, a larger number of sites would be required to obtain truly

156    representative data.

157         To overcome these limitations associated with site differences, we developed a

158    novel harmonization method that enabled us to subtract only the measurement bias by

159    using a traveling-subject dataset. We investigated that how much our proposed method

160    could reduce the measurement bias and could improve the signal to noise ratio. We

161    compared its performance to those of other commonly used harmonization methods. All

162    data utilized in this study can be downloaded publicly from the DecNef Project Brain

163    Data Repository at https://bicr-resource.atr.jp/decnefpro/.

164

165    **Results**

166    **Datasets**

167    We used two rs-fMRI datasets: the (1) SRPBS multi-disorder dataset, (2) a traveling-

168    subject dataset.

169

170    *SRPBS multi-disorder dataset*

171    This dataset included patients with five different disorders and healthy controls (HCs)

172    who were examined at nine sites belonging to eight research institutions. A total of 805

173    participants were included: 482 HCs from nine sites, 161 patients with major depressive

174    disorder (MDD) from five sites, 49 patients with autism spectrum disorder (ASD) from

175    one site, 65 patients with obsessive-compulsive disorder (OCD) from one site, and 48

176    patients with schizophrenia (SCZ) from three sites (Supplementary Table 1). The rs-fMRI

177    data were acquired using a unified imaging protocol at all but three sites (Supplementary

178    Table 2; http://www.cns.atr.jp/rs-fmri-protocol-2/). Site differences in this dataset

179    included both measurement and sampling biases (Fig. 1a). For bias estimation, we only

180    used data obtained using the unified protocol. (Patients with OCD were not scanned using

181    this unified protocol; therefore, the disorder factor could not be estimated for OCD.)

182

183    *Traveling-subject dataset*

184    We acquired a traveling-subject dataset to estimate measurement bias across sites in the

185    SRPBS dataset. Nine healthy participants (all men; age range: 24–32 years; mean age:

186    27±2.6 years) were scanned at each of 12 sites, which included the nine sites in the SRPBS

187    dataset, and produced a total of 411 scan sessions (see "Participants" in the Methods

188    section). Although we had attempted to acquire this dataset using the same imaging

189    protocol as that in the SRPBS multi-disorder dataset, there were some differences in the

190    imaging protocol across sites because of limitations in parameter settings or the scanning

191    conventions of each site (Supplementary Table 3). There were two phase-encoding

192    directions (P→A and A→P), three MRI manufacturers (Siemens, GE, and Philips), four

193    numbers of coil channels (8, 12, 24, and 32), and seven scanner types (TimTrio, Verio,

194    Skyra, Spectra, MR750W, SignaHDxt, and Achieva). Site differences in this dataset

195    included measurement bias only as the same nine participants were scanned across the 12
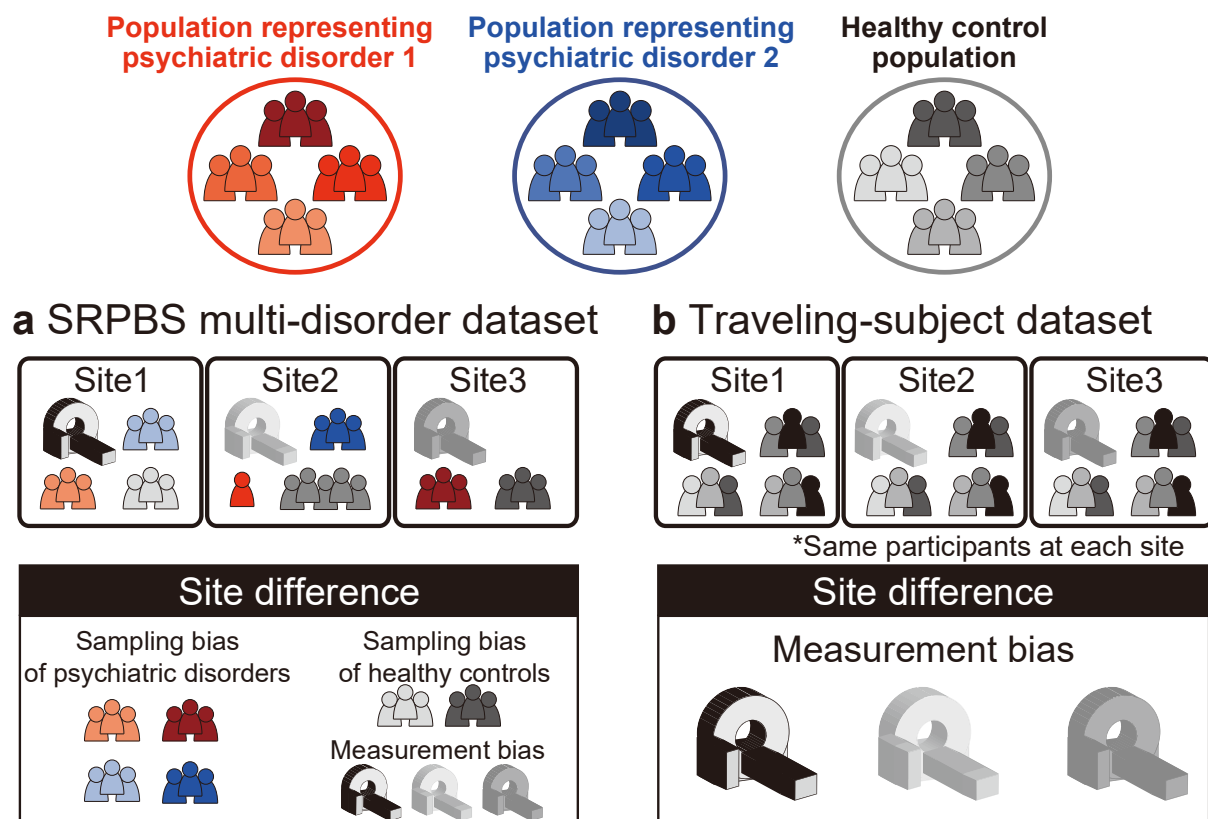
196    sites (Fig. 1b).

197

**Figure 1: Schematic examples illustrating the two main datasets.**

(a) The SRPBS multi-disorder dataset includes patients with psychiatric disorders and healthy controls. The number of patients and scanner types differed among sites. Thus, site differences consist of sampling bias and measurement bias. (b) The traveling-subject dataset includes only healthy controls, and the participants were the same across all sites. Thus, site differences consist of measurement bias only. SRPBS: Strategic Research Program for Brain Sciences.

215

**Visualization of site differences and disorder effects**

We first visualized the site differences and disorder effects in the SRPBS multi-disorder rs-fcMRI dataset while maintaining its quantitative properties by using a principal component analysis (PCA)—an unsupervised dimension reduction method. Functional connectivity was calculated as the temporal correlation of rs-fMRI blood-oxygen-level dependent (BOLD) signals between two brain regions for each participant. There are some candidates for the measure of functional connectivity such as the tangent method and partial correlation [11, 23]; however, we used Pearson's correlation coefficients because they have been the most commonly used values in previous studies. Functional connectivity was defined based on a functional brain atlas consisting of 268 nodes (i.e., regions) covering the whole brain, which has been widely utilized in previous studies [20, 24-26]. The Fisher's $z$-transformed Pearson's correlation coefficients between the preprocessed BOLD signal time courses of each possible pair of nodes were calculated and used to construct $268 \times 268$ symmetrical connectivity matrices in which each element represents a connection strength, or edge, between two nodes. We used 35,778 connectivity values [i.e., $(268 \times 267)/2$] of the lower triangular matrix of the connectivity matrix. All participant data in the SRPBS multi-disorder dataset were plotted on two axes consisting of the first two principal components (Fig. 2, small, light-colored symbols). The averages of the HCs within individual sites and the averages of individual psychiatric or developmental disorders are presented as dark-colored symbols in Fig. 2. There was a clear separation of the Hiroshima University Hospital (HUH) site for principal component 1, which explained most of the variance in the data. Furthermore, there were no differences between the differences of the sites and the disorder factors. Patients with ASD were only scanned at the Showa University (SWA) site; therefore, the averages for patients with ASD (▲) and HCs (blue ●) scanned at this site were projected to nearly identical positions (Fig. 2).
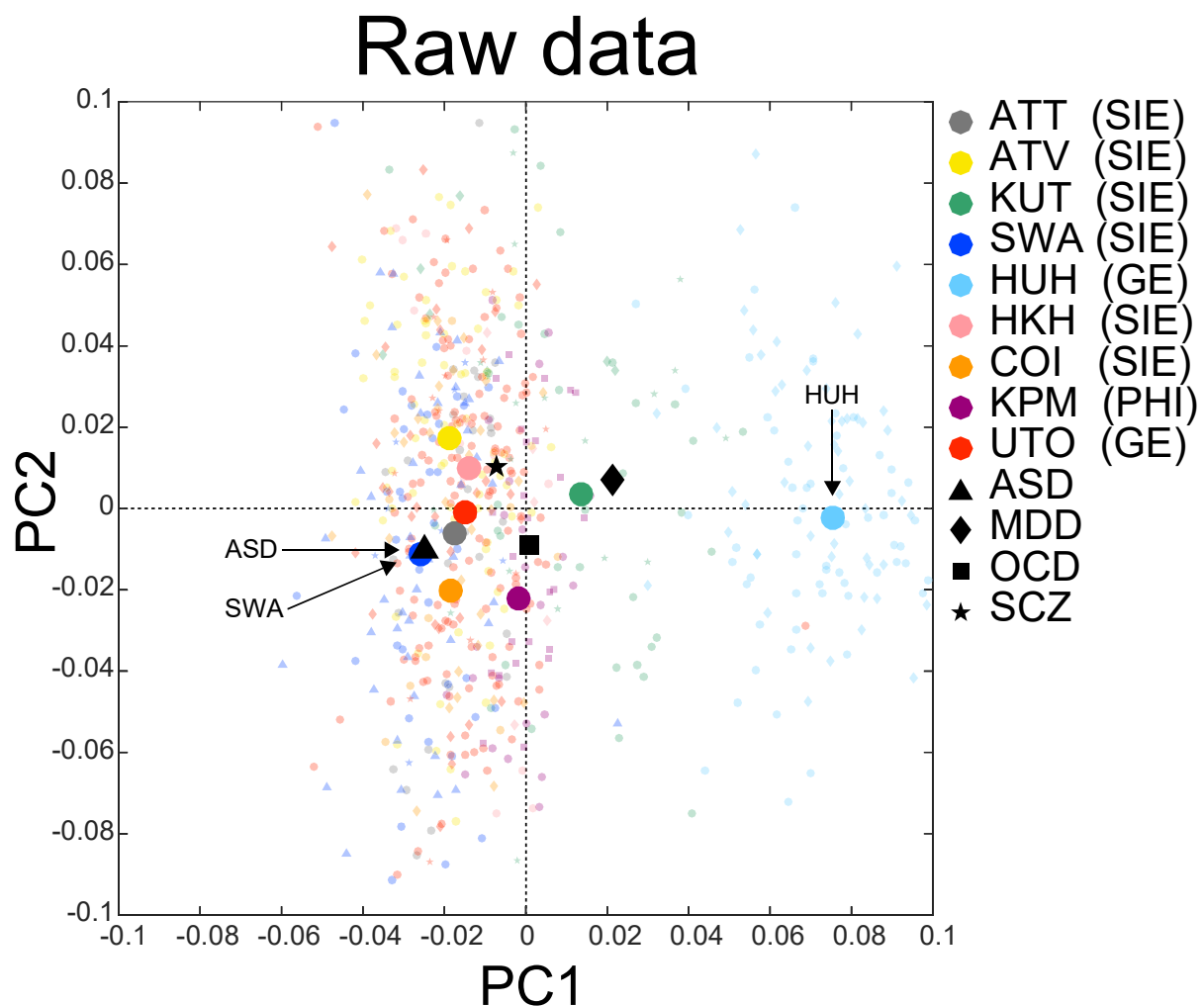
# Raw data



**Figure 2: PCA dimension reduction in the SRPBS multi-disorder dataset.**

All participants in the SRPBS multi-disorder dataset projected into the first two principal components (PCs), as indicated by small, light-colored markers. The average across all healthy controls in each site and the average within each psychiatric disorder are depicted as dark-colored makers. The color of the marker represents the site, while the shape represents the psychiatric disorder. PCA: principal component analysis; SRPBS: Strategic Research Program for Brain Sciences; ATT: Siemens TimTrio scanner at Advanced Telecommunications Research Institute International; ATV: Siemens Verio scanner at Advanced Telecommunications Research Institute International; KUT: Siemens TimTrio scanner at Kyoto University; SWA: Showa University; HUH: Hiroshima University Hospital; HKH: Hiroshima Kajikawa Hospital; COI: Center of Innovation in Hiroshima University; KPM: Kyoto Prefectural University of Medicine; UTO: University of Tokyo; ASD: Autism Spectrum Disorder. MDD: Major Depressive Disorder. OCD: Obsessive Compulsive Disorder. SCZ: Schizophrenia. SIE: Siemens fMRI. GE: GE fMRI. PHI: Philips fMRI.

257 **Bias estimation**

258 To quantitatively investigate the site differences in the rs-fcMRI data, we identified

259 measurement biases, sampling biases, and disorder factors. We defined measurement bias for

260 each site as a deviation of the correlation value for each functional connection from its average

261 across all sites. We assumed that the sampling biases of the HCs and patients with psychiatric

262 disorders differed from one another. Therefore, we calculated the sampling biases for each site

263 separately for HCs and patients with each disorder. Disorder factors were defined as deviations

264 from the HC values. Sampling biases were estimated for patients with MDD and SCZ because

265 only these patients were sampled at multiple sites. Disorder factors were estimated for MDD,

266 SCZ, and ASD because patients with OCD were not scanned using the unified protocol.

267 It is difficult to separate site differences into measurement bias and sampling bias

268 using only the SRPBS multi-disorder dataset because the two types of bias covaried across

269 sites. Different samples (i.e., participants) were scanned using different parameters (i.e.,

270 scanners and imaging protocols). In contrast, the traveling-subject dataset included only

271 measurement bias because the participants were fixed. By combining the traveling-subject

272 dataset with the SRPBS multi-disorder dataset, we simultaneously estimated measurement bias

273 and sampling bias as different factors affected by different sites. We utilized a linear mixed-

274 effects model to assess the effects of both types of bias and disorder factors on functional

275 connectivity, as follows.

276

277 *Linear mixed-effects model for the SRPBS multi-disorder dataset*

278 In this model, the connectivity values of each participant in the SRPBS multi-disorder dataset

279 were composed of fixed and random effects. Fixed effects included the sum of the average

280 correlation values across all participants and all sites at baseline, the measurement bias, the

281 sampling bias, and the disorder factors. The combined effect of participant factors (i.e.,

282 individual difference) and scan-to-scan variations was regarded as the random effect (see

283    "Estimation of biases and factors" in the Methods section).

284

285    *Linear mixed-effects model for the traveling-subject dataset*

286    In this model, the connectivity values of each participant for a specific scan in the traveling-

287    subject dataset were composed of fixed and random effects. Fixed effects included the sum of

288    the average correlation values across all participants and all sites, participant factors, and

289    measurement bias. Scan-to-scan variation was regarded as the random effect. For each

290    participant, we defined the participant factor as the deviation of connectivity values from the

291    average across all participants.

292        We estimated all biases and factors by simultaneously fitting the aforementioned two

293    regression models to the functional connectivity values of the two different datasets. For this

294    regression analysis, we used data from participants scanned using a unified imaging protocol

295    in the SRPBS multi-disorder dataset and from all participants in the traveling-subject dataset.

296    In summary, each bias or each factor was estimated as a vector that included a dimension

297    reflecting the number of connectivity values (i.e., 35,778). Vectors included in our further

298    analyses are those for measurement bias at 12 sites, sampling bias of HCs at six sites, sampling

299    bias for patients with MDD at three sites, sampling bias for patients with SCZ at three sites,

300    participant factors of nine traveling-subjects, and disorder factors for MDD, SCZ, and ASD.

301

302    **Quantification of site differences**

303    To quantitatively evaluate the effect of measurement and sampling biases on functional

304    connectivity, we compared the magnitudes of both types of bias with the magnitudes of

305    psychiatric disorders and participant factors. For this purpose, we investigated the magnitude

306    distribution of both biases, as well as the effects of psychiatric disorders and participant factors

307    on functional connectivity overall 35,778 elements in a 35,778-dimensional vector to see how

308    many functional connectivities were largely affected (Fig. 3a: the *x*-axis shows the magnitude

309    as Fisher's *z*-transformed Pearson's correlation coefficients, while the *y*-axis shows the density

310    of the number of connectivities). Figure 3b shows the same data, except the *y*-axis represents

311    the log-transformed number of connectivities for better visualization of small values. These

312    distributions show that, on average, connectivity was unaffected by either type of bias or by

313    each factor because the averages of each distribution were nearly 0. However, there were

314    significant differences among biases and factors for larger magnitudes near the tails of their

315    distributions. For example, the number of connectivities, which was largely affected (i.e., a

316    magnitude larger than 0.2), was more than 100 for the participant factor, approximately 100

317    for measurement bias, and nearly 0 for all sampling biases, as well as all disorder factors.

318           To quantitatively summarize the effect of each factor, we calculated the first, second,

319    and third statistical moments of each histogram (Fig. 3c). Based on the mean values and the

320    cube roots of the third moments, all distributions could be approximated as bilaterally

321    symmetric with a mean of zero. Thus, distributions with larger squared roots of the second

322    moments (standard deviations) affect more connectivities with larger effect sizes. The value

323    of the standard deviation was largest for the participant factor (0.0662), followed by these

324    values for the measurement bias (0.0411), the SCZ factor (0.0377), the MDD factor (0.0328),

325    the ASD factor (0.0297), the sampling bias for HCs (0.0267), sampling bias for patients with

326    SCZ (0.0217), and sampling bias for patients with MDD (0.0214). To compare the sizes of

327    the standard deviation between participant factors and measurement bias, we evaluated the

328    variance of each distribution. All pairs of variances were analyzed using Ansari–Bradley tests.

329    Our findings indicated that all variances of the participant factors were significantly larger

330    than all variances of the measurement biases (nine participant factors × 12 measurement

331    biases = 108 pairs; $W^*$: mean = -59.80, max = -116.81, min = -3.69; *p* value after Bonferroni

332    correction: max = 0.011, min = 0, *n* = 35,778). In addition, the variances of 10 of 12

333    measurement biases were significantly larger than the variance of the MDD factor, the

334    variances of seven of 12 measurement biases were significantly larger than the variance of

335    the SCZ factor, and the variances of all measurement biases were significantly larger than the

336    variance of the MDD factor (Supplementary Table 8). Furthermore, we plotted fractions of

337    the data variance determined using the aforementioned factors (i.e., contribution size) in our

338    linear mixed-effects model (Fig. 3d; see "Analysis of contribution size" in the Methods

339    section). The results were consistent with the analysis of the standard deviation (Fig. 3c,

340    middle). These results indicate that the effect size of measurement bias on functional

341    connectivity is smaller than that of the participant factor but is mostly larger than those of the

342    disorder factors, which suggest that measurement bias represents a serious limitation in

343    research regarding psychiatric disorders. The largest variance in sampling bias was

344    significantly larger than the variance of the MDD factor (Supplementary Table 9), whereas

345    the smallest variance in sampling bias was one-half the size of the variance for disorder factors.

346    These findings indicate that sampling bias also represents a major limitation in psychiatric

347    research.

348         The standard deviation of the participant factor was approximately twice that for SCZ,

349    MDD, and ASD; therefore, individual variability within the healthy population was much

350    greater than that among patients with SCZ, MDD, or ASD when all functional connections

351    were considered. Furthermore, the standard deviations of the measurement biases were mostly

352    larger than those of the disorder factors, while the standard deviations of the sampling biases

353    were comparable with those of the disorder factors. Such relationships make the development

354    of rs-fcMRI-based classifiers of psychiatric or developmental disorders very challenging.

355    Only when a small number of disorder-specific and site-independent abnormal functional

356    connections can be selected from among a vast number does it become feasible to develop

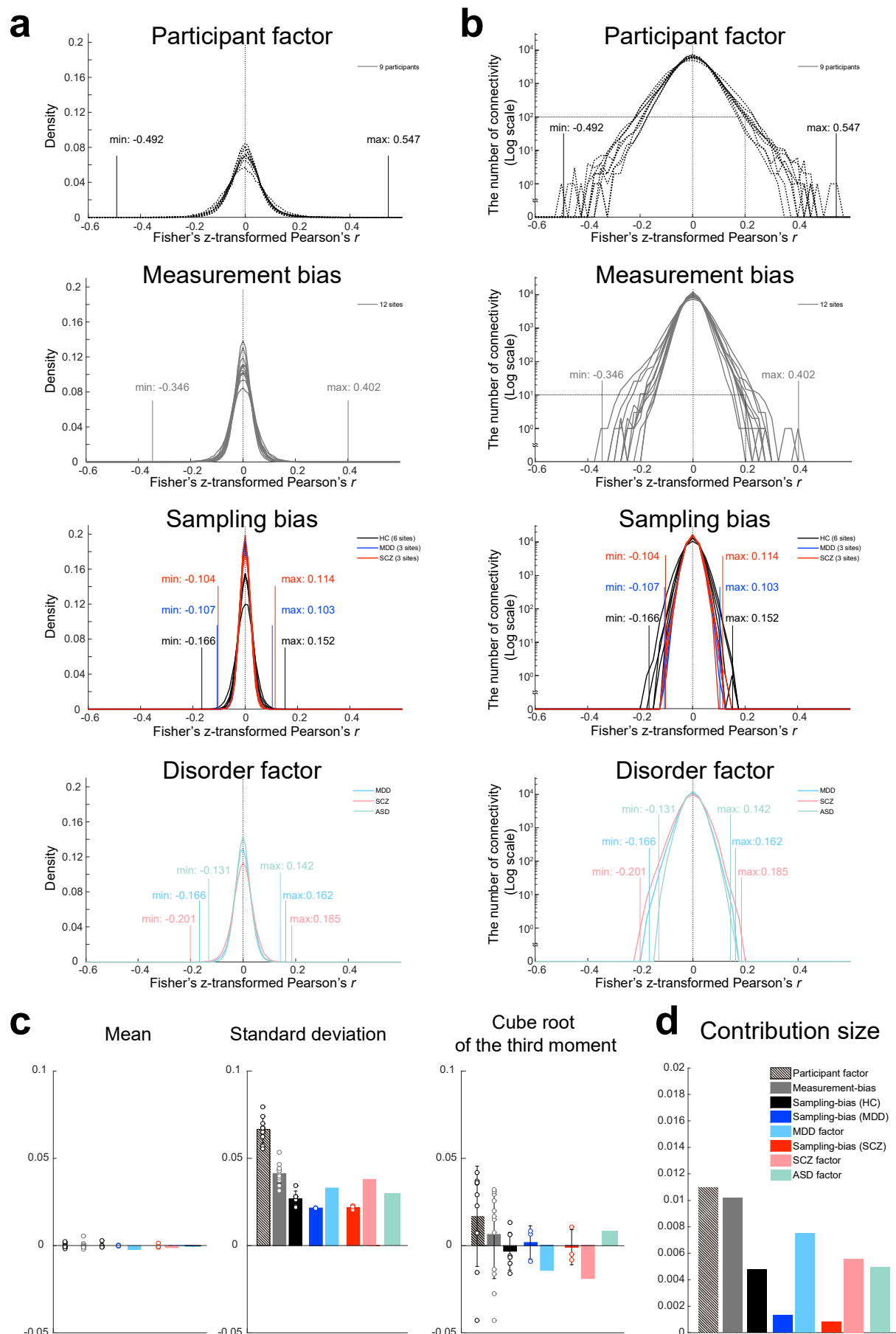357    robust and generalizable classifiers across multiple sites [2, 8-10, 15].

**Figure 3: Distributions and statistics for each type of bias and each factor.**

360  (a, b) The distribution of the effects of each bias and each factor on functional connectivity

361  vectors. Functional connectivity was measured based on Fisher's z-transformed Pearson's

362  correlation coefficients. The *x*-axis represents the effect size of the Fisher's z-transformed

363  Pearson's correlation coefficients. In (a) and (b), the *y*-axis represents the density of

364  connectivity and the log-transformed the number of connections, respectively. Each line

365  represents one participant or one site. (c) The means, standard deviations, and third moments

366  standardized to the same scale on the vertical axis (i.e., cube root) for each type of bias and

367  each factor. Bars represent the average value, while the error bars represent the standard

368  deviation across sites or participants. Each data point represents one participant or one site. (d)

369  Contribution size of each bias and each factor. HC: healthy controls; SCZ: schizophrenia;

370  MDD: major depressive disorder; ASD: autism spectrum disorder.

371

372  **Brain regions contributing most to biases and associated factors**

373  To evaluate the spatial distribution of the two types of bias and all factors in the whole brain,

374  we utilized a previously described visualization method [27] to project connectivity

375  information to anatomical regions of interest (ROIs). We first quantified the effect of a bias or

376  a factor on each functional connectivity as the median of its absolute values across sites or

377  across participants. Thus, we obtained 35,778 values, each of which was associated with one

378  connectivity and represented the effect of a bias or factor on the connectivity. We then

379  summarized these effects on connectivity for each ROI by averaging the values of all

380  connectivities connected with the ROI (see "Spatial characteristics of measurement bias,

381  sampling bias, and each factor in the brain" in the Methods section). The average value

382  represents the extent the ROI contributes to the effect of a bias or factor. By repeating this

383  procedure for each ROI and coding the averaged value based on the color of an ROI, we were

384  able to visualize the relative contribution of individual ROIs to each bias or factor in the whole

385  brain (Fig. 4). Consistent with the findings of previous studies, the effect of the participant

386  factor was large for several ROIs in the cerebral cortex, especially in the prefrontal cortex, but

387  small in the cerebellum and visual cortex [24]. The effect of measurement bias was large in

388  inferior brain regions where functional images are differentially distorted depending on the

389  phase-encoding direction [28, 29]. Connections involving the medial dorsal nucleus of the

390 thalamus were also heavily affected by both MDD, SCZ and ASD. Effects of the MDD factor

391 were observed in the dorsomedial prefrontal cortex and the superior temporal gyrus in which

392 abnormalities have also been reported in previous studies [22, 30, 31]. Effects of the SCZ factor

393 were observed in the left inferior parietal lobule, bilateral anterior cingulate cortices, and left

394 middle frontal gyrus in which abnormalities have been reported in previous studies [32-34].

395 Effects of the ASD factor were observed in the putamen, the medial prefrontal cortex, and the

396 right middle temporal gyrus in which abnormalities have also been reported in previous studies

397 [10, 11, 35]. The effect of sampling bias for HCs was large in the inferior parietal lobule and

398 the precuneus, both of which are involved in the default mode network and the middle frontal

399 gyrus. Sampling bias for disorders was large in the medial dorsal nucleus of the thalamus, left

400 dorsolateral prefrontal cortex, dorsomedial prefrontal cortex, and cerebellum for MDD [22];

401 and in the prefrontal cortex, cuneus, and cerebellum for SCZ [33].

402

403

404

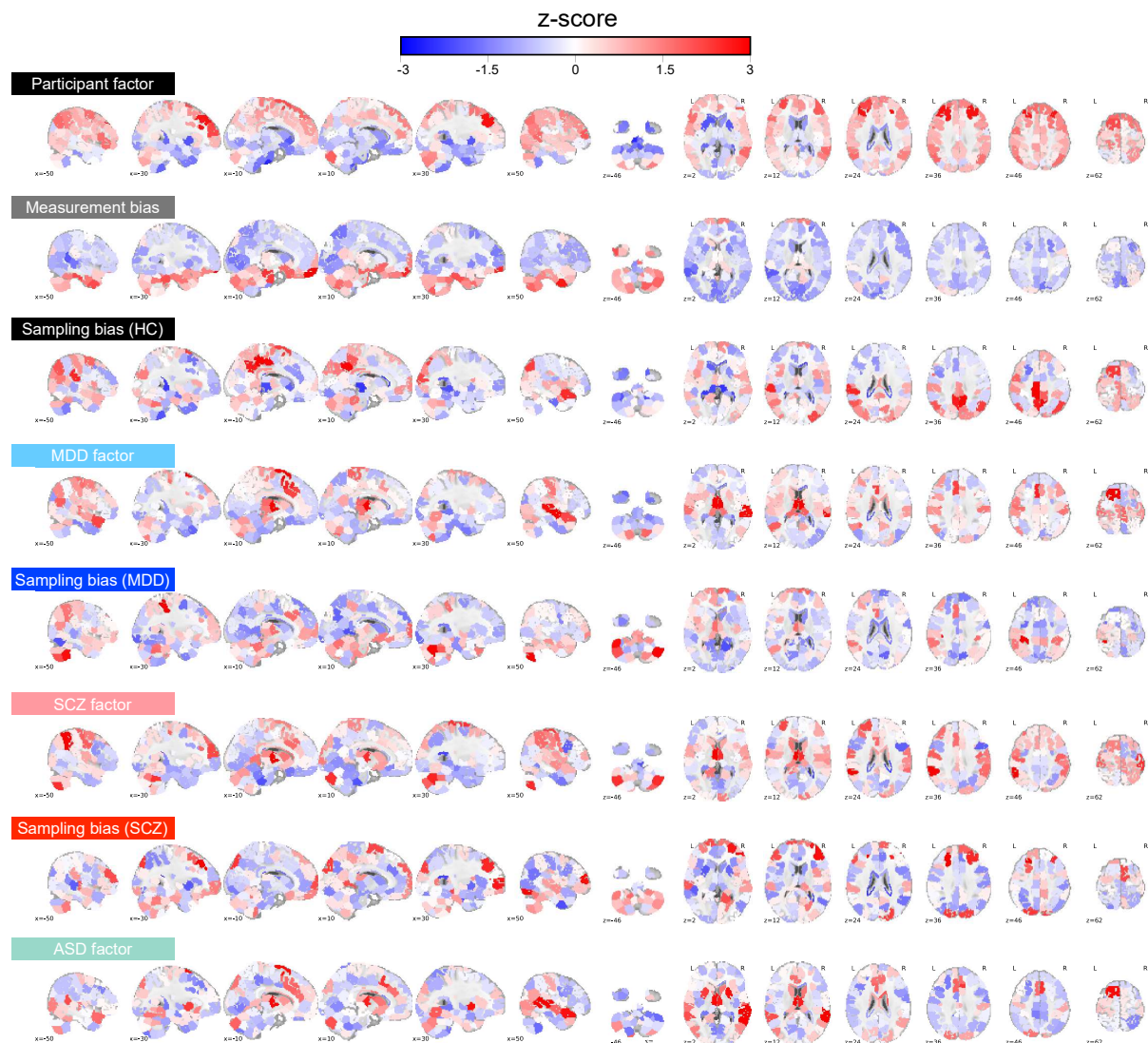405

406

407

408

409

410

411

412

**Figure 4: Spatial distribution of each type of bias and each factor in various brain regions.** Mean effects of connectivity for all 268 ROIs. For each ROI, the mean effects of all functional connections associated with that ROI were calculated for each bias and each factor. Warmer (red) and cooler (blue) colors correspond to large and small effects, respectively. The magnitudes of the effects are normalized within each bias or each factor (*z*-score). ROI: region of interest; HC: healthy control; SCZ: schizophrenia; MDD: major depressive disorder; ASD: autism spectrum disorder.

427 **Characteristics of measurement bias**

428 We next investigated the characteristics of measurement bias. We first examined whether

429 similarities among the estimated measurement bias vectors for the 12 included sites reflect

430 certain properties of MRI scanners such as phase-encoding direction, MRI manufacturer, coil

431 type, and scanner type. We used hierarchical clustering analysis to discover clusters of similar

432 patterns for measurement bias. This method has previously been used to distinguish subtypes

433 of MDD, based on rs-fcMRI data [22]. As a result, the measurement biases of the 12 sites were

434 divided into phase-encoding direction clusters at the first level (Fig. 5a). They were divided

435 into fMRI manufacturer clusters at the second level, and further divided into coil type clusters,

436 followed by scanner model clusters. Furthermore, we quantitatively verified the magnitude

437 relationship among factors by using the same model to assess the contribution of each factor

438 (Fig. 5b; see "Analysis of contribution size" in the Methods section). The contribution size was

439 largest for the phase-encoding direction (0.0391), followed by the contribution sized for fMRI

440 manufacturer (0.0318), coil type (0.0239), and scanner model (0.0152). These findings indicate

441 that the main factor influencing measurement bias is the difference in the phase-encoding

442 direction, followed by fMRI manufacturer, coil type, and scanner model, respectively.
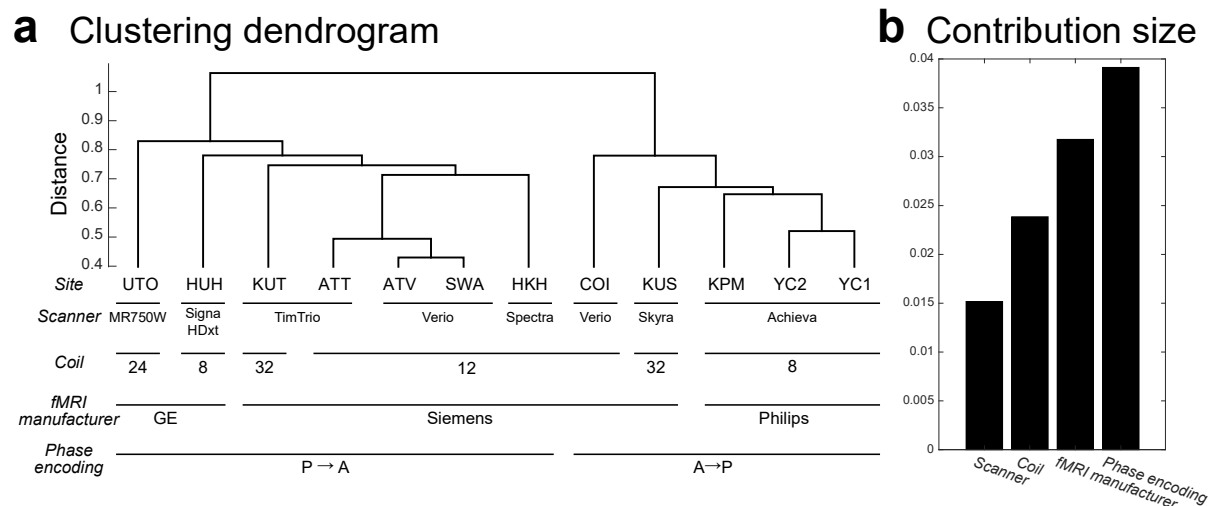
443

444

445

446

**Figure 5: Clustering dendrogram for measurement bias.**

(a) The height of each linkage in the dendrogram represents the dissimilarity (1 - $r$) between the clusters joined by that link. (b) Contribution size of each factor. UTO: University of Tokyo; HUH: Hiroshima University Hospital; KUT: Siemens TimTrio scanner at Kyoto University; ATT: Siemens TimTrio scanner at Advanced Telecommunications Research Institute International; ATV: Siemens Verio scanner at Advanced Telecommunications Research Institute International; SWA: Showa University; HKH: Hiroshima Kajikawa Hospital; COI: Center of Innovation in Hiroshima University; KUS: Siemens Skyra scanner at Kyoto University; KPM: Kyoto Prefectural University of Medicine; YC1: Yaesu Clinic 1; YC2: Yaesu Clinic 2.

470 **Sampling bias is because of sampling from among a subpopulation**

471 We investigated two alternative models for the mechanisms underlying sampling bias. In the

472 "single-population model", which assumes that participants are sampled from a common

473 population (Fig. 6a), the functional connectivity values of each participant were generated from

474 a Gaussian distribution (see "Comparison of models for sampling bias" in the Methods section).

475 In the "different-subpopulation model," which assumes that sampling bias occurs partly

476 because participants are sampled from among a different subpopulation at each site (Fig. 6b),

477 we assumed that the average of the subpopulation differed among sites and was generated from

478 a Gaussian distribution. In addition, the functional connectivity values of each participant were

479 generated from a Gaussian distribution, based on the average of the subpopulation at each site.

480 It is necessary to determine which model is more suitable for collecting big data across multiple

481 sites: If the former model is correct, then the data can be used to represent a population by

482 increasing the number of participants, even if the number of sites is small. If the latter model

483 is correct, data should be collected from many sites, as a single site does not represent the true

484 grand population distribution, even with a very large sample size.

485 For each model, we first investigated how the number of participants at each site

486 determined the effect of sampling bias on functional connectivity. We measured the magnitude

487 of the effect, based on the variance values for sampling bias across functional connectivity (see

488 the "Quantification of site differences" section). We used variance instead of the standard

489 deviation to simplify the statistical analysis, although there is essentially no difference based

490 on which value is used. We theorized that each model represents a different relationship

491 between the number of participants and the variance of sampling bias. Therefore, we

492 investigated which model best represents the actual relationships in our data by comparing the

493 corrected Akaike information criterion (AICc) [36, 37] and Bayesian information criterion

494 (BIC). Moreover, we performed leave-one-site-out cross-validation evaluations of predictive

495 performance in which all but one site was used to construct the model and the variance of the

496   sampling bias was predicted for the remaining site. We then compared the predictive

497   performances between the two models. Our results indicated that the different-subpopulation

498   model provided a better fit for our data than the single-population model (Fig. 6c; different-

499   subpopulation model: AICc = -108.80 and BIC = -113.22; single-population model: AICc = -

500   96.71 and BIC = -97.92). Furthermore, the predictive performance was significantly higher for

501   the different-subpopulation model than for the single-population model (one-tailed Wilcoxon

502   signed-rank test applied to absolute errors: $Z = 1.67$, $p = .0469$, $n = 6$; Figs. 6d and 6e). This

503   result indicates that sampling bias is not only caused by random sampling from a single grand

504   population, depending on the number of participants among sites, but also by sampling from

505   among different subpopulations. Sampling biases thus represent a major limitation in

506   attempting to estimate a true single distribution of HC or patient data based on measurements

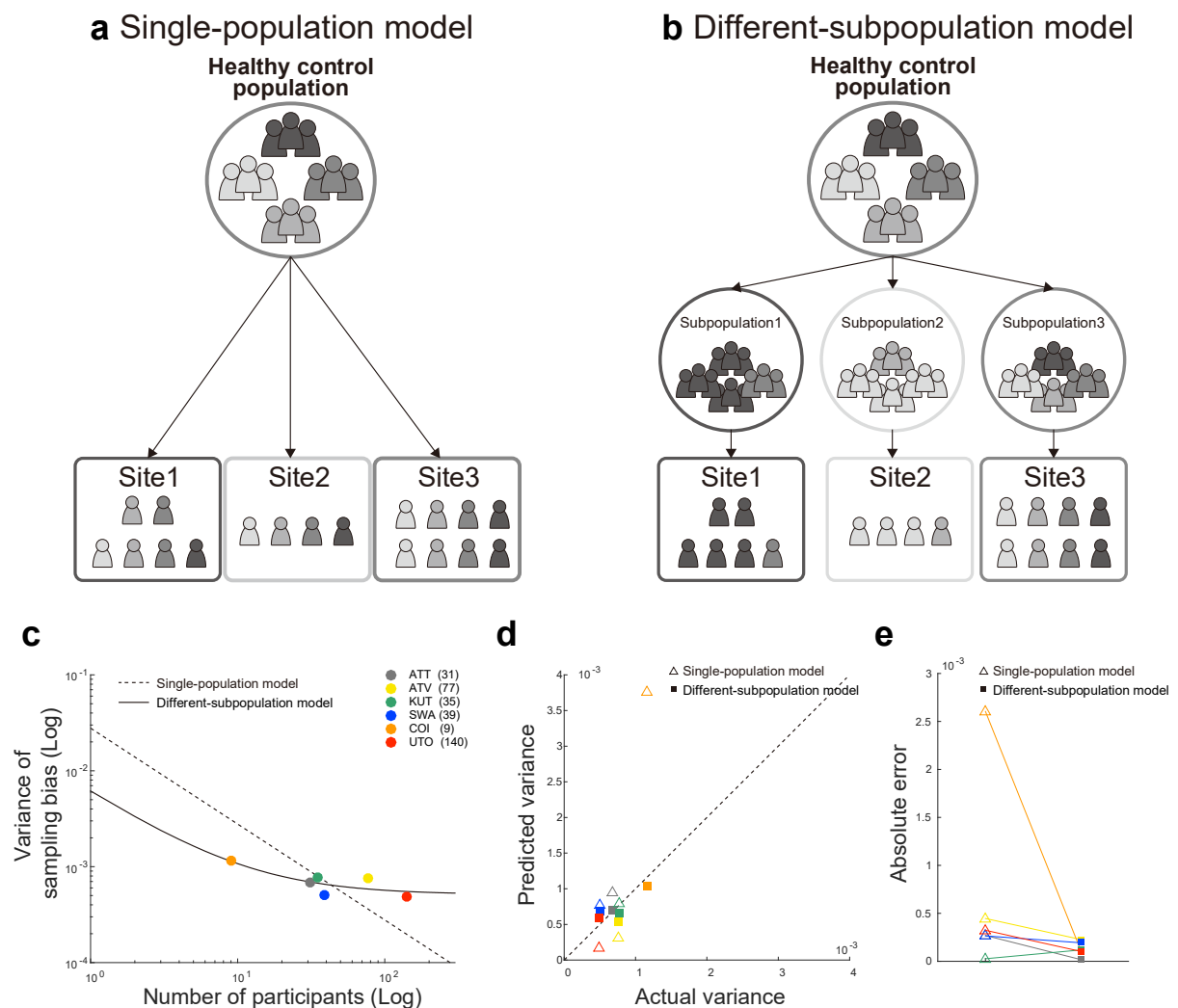507   obtained from a finite number of sites and participants.

**Figure 6: Comparison of the two models of sampling bias.**

Schematic examples illustrating the single-population (a) and different-subpopulation models (b) and the results of model fitting (c). The *x*-axis represents the number of participants on a logarithmic scale, while the *y*-axis represents the variance of sampling bias on a logarithmic scale. The broken line represents the prediction of the single-population model, while the solid line represents the prediction of the different-subpopulation model. Each data point represents one site. (d) Results of the predictions determined by using each model. The *x*-axis represents the actual variance, while the *y*-axis represents the predicted variance. Open triangles correspond to the single-population model, while filled squares correspond to the different-subpopulation model. (e) Performance of prediction using the two models, based on the absolute error between the actual and predicted variance. UTO: University of Tokyo; COI: Center of Innovation in Hiroshima University; SWA: Showa University; KUT: Siemens TimTrio scanner at Kyoto University; ATT: Siemens TimTrio scanner at Advanced Telecommunications Research Institute International; ATV: Siemens Verio scanner at Advanced Telecommunications Research Institute International.

525 **Visualization of the effect of harmonization**

526 We next developed a novel harmonization method that enabled us to subtract only the

527 measurement bias using the traveling-subject dataset. Using a linear mixed-effects model, we

528 estimated the measurement bias separately from sampling bias (see the "Bias estimation" in

529 the Methods section). Thus, we could remove the measurement bias from the SRPBS multi-

530 disorder dataset (i.e., traveling-subject method, see "Traveling-subject harmonization" in the

531 Methods section). To visualize the effects of the harmonization process, we plotted the data

532 after subtracting only the measurement bias from the SRPBS multi-disorder dataset as

533 described in the "Visualization of site differences and disorder effects" section (Fig. 7).

534 Relative to the data reported in Fig. 2, which reflects the data before harmonization, the HUH

535 site moved much closer to the origin (i.e., grand average) and showed no marked separation

536 from the other sites. This result indicates that the separation of the HUH site observed in Fig.

537 2 was caused by measurement bias, which was removed following harmonization. Furthermore,

538 harmonization was effective in distinguishing patients and HCs scanned at the same site. Since

539 patients with ASD were only scanned at the Showa University (SWA) site, the averages for

540 patients with ASD (▲) and HCs (blue ●) scanned at this site were projected to nearly identical

541 positions (Fig. 2). However, the two symbols are clearly separated from one another in Fig. 7.

542 The effect of a psychiatric disorder (ASD) could not be observed in the first two PCs without

543 harmonization but became detectable following the removal of measurement bias.
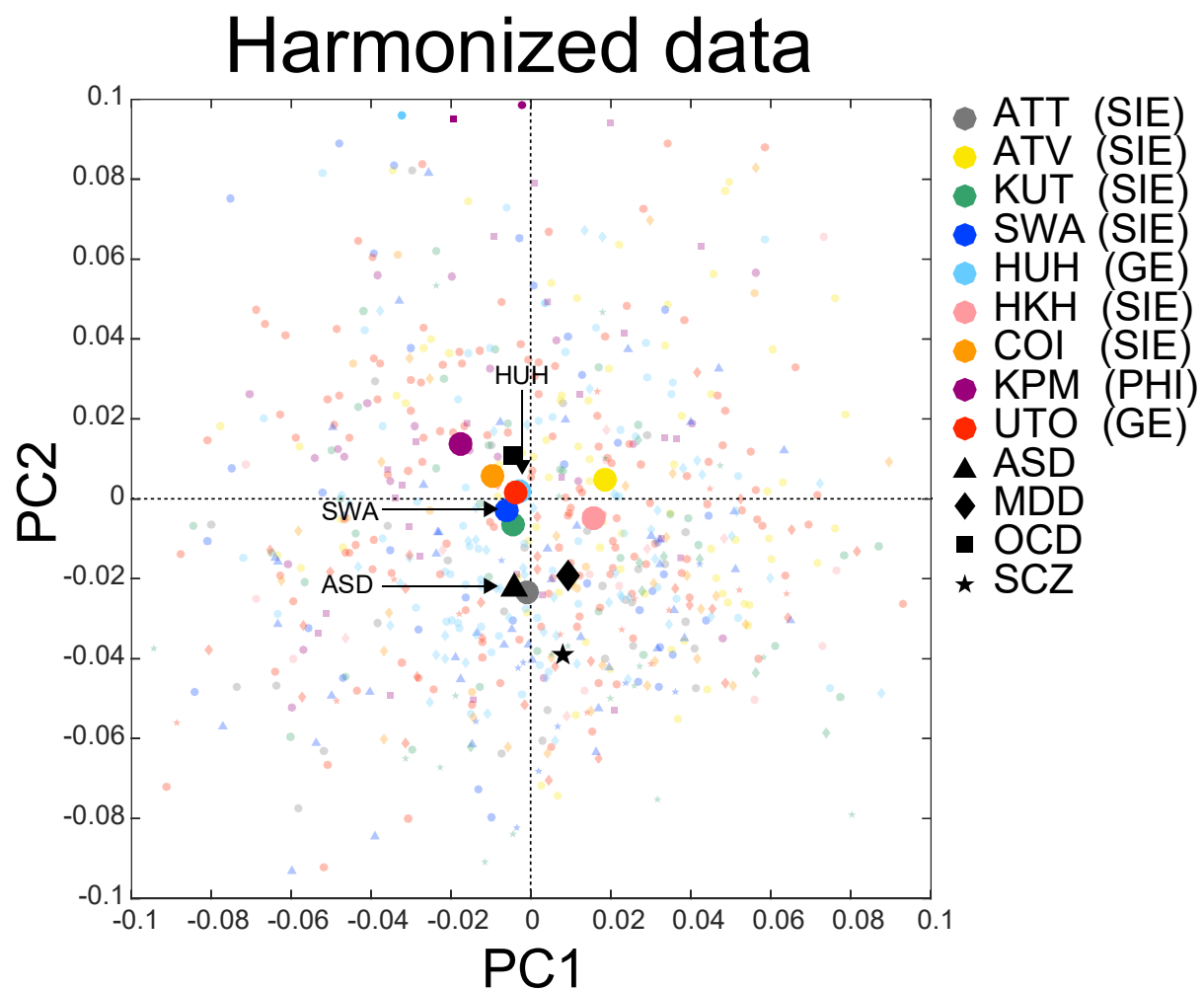
**Figure 7: PCA dimension reduction in the SRPBS multi-disorder dataset after harmonization.**

All participants in the SRPBS multi-disorder dataset after harmonization projected into the first two principal components (PCs), as indicated by small, light-colored markers. The average across all healthy controls in each site and the average within each psychiatric disorder are depicted as dark-colored makers. The color of the marker represents the site, while the shape represents the psychiatric disorder. PCA: principal component analysis; SRPBS: Strategic Research Program for Brain Sciences; ATT: Siemens TimTrio scanner at Advanced Telecommunications Research Institute International; ATV: Siemens Verio scanner at Advanced Telecommunications Research Institute International; KUT: Siemens TimTrio scanner at Kyoto University; SWA: Showa University; HUH: Hiroshima University Hospital; HKH: Hiroshima Kajikawa Hospital; COI: Center of Innovation in Hiroshima University; KPM: Kyoto Prefectural University of Medicine; UTO: University of Tokyo; ASD: Autism Spectrum Disorder. MDD: Major Depressive Disorder. OCD: Obsessive Compulsive Disorder. SCZ: Schizophrenia. SIE: Siemens fMRI. GE: GE fMRI. PHI: Philips fMRI.

**561** **Quantification of the effect of traveling-subject harmonization**

**562** To correct difference among sites there are three commonly used harmonization methods: (1)

**563** a ComBat method [16, 17, 19, 38], a batch-effect correction tool commonly used in genomics,

**564** site difference was modeled and removed; (2) a generalized linear model (GLM) method, site

**565** difference was estimated without adjusting for biological covariates (e.g., diagnosis) [16, 18,

**566** 22]; and (3) an adjusted GLM method, site difference was estimated while adjusting for

**567** biological covariates [16, 18] (see the "Harmonization procedures" in the Methods section).

**568** However, all these methods estimate the site difference without separating site difference into

**569** the measurement bias and the sampling bias and subtract the site difference from data.

**570** Therefore, existing harmonization methods might have pitfall to eliminate not only biologically

**571** meaningless measurement bias but also eliminate biologically meaningful sampling bias. Here,

**572** we tested whether the traveling-subject harmonization method indeed removes only the

**573** measurement bias and whether the existing harmonization methods simultaneously remove the

**574** measurement and sampling biases. Specifically, we performed 2-fold cross-validation

**575** evaluations in which the SRPBS multi-disorder dataset was partitioned into two equal-size

**576** subsamples (fold1 data and fold2 data) with the same proportions of sites. Between these two

**577** subsamples, the measurement bias is common, but the sampling bias is different (because the

**578** scanners are common, and participants are different). We estimated the measurement bias (or

**579** site difference including the measurement bias and the sampling bias for the existing methods)

**580** by applying the harmonization methods to the fold1 data and subtracted the measurement bias

**581** or site difference from the fold2 data. We then estimated the measurement bias in the fold2

**582** data. For the existing harmonization methods, if the site difference estimated by using fold1

**583** contains only the measurement bias, the measurement bias estimated in fold2 data after

**584** subtracting the site difference should be smaller than that of without subtracting the site

**585** difference (Raw). To separately estimate measurement bias and sampling bias in both

**586** subsamples while avoiding information leak, we also divided the traveling-subject dataset into

587    two equal-size subsamples with the same proportions of sites and subjects. We concatenated

588    one subsample of traveling-subject dataset to fold1 data to estimate the measurement bias for

589    traveling-subject method (estimating dataset) and concatenated the other subsample of

590    traveling-subject dataset to fold2 data for testing (testing dataset). That is, in the traveling-

591    subject harmonization method, we estimated the measurement bias using the estimating dataset

592    and removed the measurement bias from the testing dataset. By contrast, in the other

593    harmonization methods, we estimated the site difference using the fold1 data (not including the

594    subsample of traveling-subject dataset) and removed the site difference from the testing dataset.

595    We then estimated the measurement bias using the testing dataset and evaluated the standard

596    deviation of the magnitude distribution of measurement bias calculated in the same way as

597    described in "Quantification of site differences" section. To verify whether important

598    information such as participant factors and disorder factors are kept in the testing dataset, we

599    also estimated the disorder factors and participant factors and calculated the ratio of the

600    standard deviation of measurement bias to the standard deviation of participant factor and

601    disorder factor as signal to noise ratios. This procedure was done again by exchanging the

602    estimating dataset and the testing dataset (see the "2-fold cross-validation evaluation procedure"

603    in the Methods section).

604         Fig. 8 shows that the standard deviation of measurement bias and the ratio of the

605    standard deviation of measurement bias to the standard deviation of participant factor and

606    disorder factor in the both fold data for the four harmonization methods and without

607    harmonization (Raw). Our result shows that the reduction of the standard deviation of

608    measurement bias from the Raw was highest in the traveling-subject method among all

609    methods (29% reduction compared to 3% in the second highest value for ComBat method).

610    Moreover, improvement in the signal to noise ratios were also highest in our method for

611    participant factor (41% improvement compared to 3% in the second highest value for ComBat

612    method) and for disorder factor (39% improvement compared to 3% in the second highest value

613 for ComBat method). These results indicate that the traveling-subject harmonization method

614 indeed removed the measurement bias and improved the signal to noise ratios.

615



617 **Figure 8: Reduction of the measurement bias and improvement of signal to noise ratios**
618 **for different harmonization methods.**
619 (a) Standard deviation of the measurement bias. (b) Ratio of standard deviation of the
620 measurement bias to standard deviation of the participant factor. (c) Ratio of standard deviation
621 of the measurement bias to standard deviation of the disorder factor. Different colored columns
622 show the results from different harmonization method. Two columns of the same color show
623 the results of the two folds. GLM: generalized linear model.

624

625

626

627

628

629

630

631

632

633

**Discussion**

In the present study, by acquiring a separate traveling-subject dataset and the SRPBS multi-disorder dataset, we separately estimated measurement and sampling biases for multiple sites, which we then compared with the magnitude of disorder factors. Furthermore, we investigated the origin of each bias in multi-site datasets. Finally, to overcome the problem of site difference, we developed a novel harmonization method that enabled us to subtract the measurement bias by using a traveling-subject dataset and achieved the reduction of the measurement bias by 29% and the improvement of the signal to noise ratios by 40%.

We assessed the effect sizes of measurement and sampling biases in comparison with the effects of psychiatric disorders on resting-state functional connectivity. Our findings indicated that measurement bias exerted significantly greater effects than disorder factors, whereas sampling bias was comparable to (or even larger than) the disorder effects (Fig. 3). However, we did not control for variations in disease stage and treatment in our dataset. Although controlling for such heterogeneity may increase the effect size of disorder factors, such control is not feasible when collecting big data from multiple sites. Therefore, it is important to appropriately remove measurement bias from heterogeneous patient data to identify relatively small disorder effects. This issue is essential for investigating the relationships among different psychiatric disorders because disease factors are often confounded by site differences. As previously mentioned, it is common for a single site to sample only a few types of psychiatric disorders (e.g., SCZ from site A and ASD from site B). In this situation, it is critical to dissociate disease factors from site differences. This dissociation can be accomplished by subtracting only the measurement bias which is estimated from traveling subject dataset.

Our results indicated that measurement bias is primarily influenced by differences in the phase-encoding direction, followed by differences in fMRI manufacturer, coil type, and scanner model (Fig. 5). These results are consistent with our finding of large measurement

660     biases in the inferior brain regions (Fig. 4), the functional imaging of which is known to be

661     influenced by the phase-encoding direction [28, 29]. Previous studies have reported that the

662     effect because of the difference in the phase-encoding direction can be corrected using the field

663     map obtained at the time of imaging [28, 39-41]. The field map was acquired in parts of the

664     traveling-subject dataset; therefore, we investigated the effectiveness of field map correction

665     by comparing the effect size of the measurement bias and the participant factor between

666     functional images with and without field map correction. Our prediction was as follows: if field

667     map correction is effective, the effect of measurement bias will decrease, while that of the

668     participant factor will increase following field map correction. Field map correction using

669     SPM12 (http://www.fil.ion.ucl.ac.uk/spm/software/spm12) reduced the effect of measurement

670     bias in the inferior brain regions (whole brain: 3% reduction in the standard deviation of

671     measurement bias) and increased the effect of the participant factor in the whole brain (3%

672     increase in the standard deviation of the participant factor; Supplementary Figures 2a and 2b).

673     However, the effect of measurement bias remained large in inferior brain regions

674     (Supplementary Figure 2a), and hierarchical clustering analysis revealed that the clusters of the

675     phase-encoding direction remained dominant (Supplementary Figure 2c). These results

676     indicate that, even with field map correction, it is largely impossible to remove the influence

677     of differences in phase-encoding direction on functional connectivity. Thus, harmonization

678     methods are still necessary to remove the effect of these differences and other measurement-

679     related factors. However, some distortion correction methods have been developed (e.g., top-

680     up method and symmetric normalization) [42, 43], and further studies are required to verify the

681     efficacy of these methods.

682        Our data supported the different-subpopulation model rather than the single-

683     population model (Fig. 6), which indicates that sampling bias is caused by sampling from

684     among different subpopulations. Furthermore, these findings suggest that, during big data

685     collection, it is better to sample participants from several imaging sites than to sample many

686    participants from a few imaging sites. These results were obtained only by combining the

687    SRPBS multi-disorder database with a traveling-subject dataset

688    (http://www.cns.atr.jp/decnefpro/). To the best of our knowledge, the present study is the first

689    to demonstrate the presence of sampling bias in rs-fcMRI data, the mechanisms underlying this

690    sampling bias, and the effect size of sampling bias on resting-state functional connectivity,

691    which was comparable to that of psychiatric disorders. We analyzed sampling bias among HCs

692    only, because the number of sites was too small to conduct an analysis of patients with

693    psychiatric diseases.

694        We developed a novel harmonization method using a traveling-subject dataset (i.e.,

695    traveling-subject method), which was then compared with existing harmonization methods.

696    Our results demonstrated that the traveling-subject method outperformed other conventional

697    GLM-based harmonization methods and ComBat method. The traveling-subject method

698    achieved reduction of the measurement bias by 29% compared to 3% in the second highest

699    value for ComBat method and improvement of the signal to noise ratios by 40% compared to

700    3% in the second highest value for ComBat method. This result indicates that the traveling-

701    subject dataset helps to properly estimate the measurement bias and also helps to harmonize

702    the rs-fMRI data across imaging sites. To further quantitatively evaluate the harmonization

703    method, we constructed biomarkers for psychiatric disorders based on rs-fcMRI data, which

704    distinguishes between HCs and patients, and a regression model to predict participants' age

705    based on rs-fcMRI data using SRPBS multi-disorder dataset (see "Classifiers for MDD and

706    SCZ, based on the four harmonization methods" and "Regression models of participant age

707    based on the four harmonization methods" in Supplementary Information). We evaluated the

708    generalization performance to independent validation dataset, which was not included in

709    SRPBS multi-disorder dataset. The traveling-subject harmonization method improved the

710    generalization performance of all these prediction models as compared with the case where

711    harmonization was not performed. These results indicate that the traveling-subject dataset also

712    helps the constructing a prediction model based on multi-site rs-fMRI data.

713        The present study possesses some limitations of note. The accuracy of measurement

714    bias estimation may be improved by further expanding the traveling-subject dataset. This can

715    be achieved by increasing the number of traveling participants or sessions per site. However,

716    as mentioned in a previous traveling-subject study [20], it is costly and time-consuming to

717    ensure that numerous participants travel to every site involved in big database projects. Thus,

718    the cost-performance tradeoff must be evaluated in practical settings. The numbers of traveling

719    participants and MRI sites used in this study (nine and 12, respectively) were larger than those

720    used in a previous study (eight and eight, respectively) [20], and the number of total sessions

721    in this study (411) was more than three times larger than that used in the previous study (128)

722    [20]. Furthermore, although we estimated the measurement bias for each connectivity,

723    hierarchical models of the brain (e.g., ComBat) may be more appropriate for improving the

724    estimates of measurement bias.

725        In summary, by acquiring a separate traveling-subject dataset and the SRPBS multi-

726    disorder database, we revealed that site differences were composed of biological sampling bias

727    and engineering measurement bias. The effect sizes of these biases on functional connectivity

728    were greater than or equal to the effect sizes of psychiatric disorders, highlighting the

729    importance of controlling for site differences when investigating psychiatric disorders.

730    Furthermore, using the traveling-subject dataset, we developed a novel traveling-subject

731    method that harmonizes the measurement bias only by separating sampling bias from site

732    differences. Our findings verified that the traveling-subject method outperformed conventional

733    GLM-based harmonization methods and ComBat method. These results suggest that a

734    traveling-subject dataset can help to harmonize the rs-fMRI data across imaging sites.

735

736

737

738 # Methods

739 ## Participants

740 We used two resting-state functional MRI datasets for all analyses: (1) the SRPBS multi-disorder dataset, which encompasses

741 multiple psychiatric disorders; (2) a traveling-subject dataset. The SRPBS multi-disorder dataset contains data for 805

742 participants (482 HCs from nine sites, 161 patients with MDD from five sites, 49 patients with ASD from one site, 65 patients

743 with OCD from one site, and 48 patients with SCZ from three sites (Supplementary Table 1). Data were acquired using a

744 Siemens TimTrio scanner at Advanced Telecommunications Research Institute International (ATT), a Siemens Verio scanner

745 at Advanced Telecommunications Research Institute International (ATV), a Siemens Verio at the Center of Innovation in

746 Hiroshima University (COI), a GE Signa HDxt scanner at HUH, a Siemens Spectra scanner at Hiroshima Kajikawa Hospital

747 (HKH), a Philips Achieva scanner at Kyoto Prefectural University of Medicine (KPM), a Siemens Verio scanner at SWA, a

748 Siemens TimTrio scanner at Kyoto University (KUT), and a GE MR750W scanner at the University of Tokyo (UTO). Each

749 participant underwent a single rs-fMRI session for 5–10 min. The rs-fMRI data were acquired using a unified imaging protocol

750 at all but three sites (Supplementary Table 2; http://www.cns.atr.jp/rs-fmri-protocol-2/). During the rs-fMRI scans, participants

751 were instructed as follows, except at one site: "Please relax. Don't sleep. Fixate on the central crosshair mark and do not think

752 about specific things." At the remaining site, participants were instructed to close their eyes rather than fixate on a central

753 crosshair.

754 In the traveling-subject dataset, nine healthy participants (all male participants; age range, 24–32 years; mean age,

755 $27\pm2.6$ years) were scanned at each of 12 sites in the SRPBS consortium, producing a total of 411 scan sessions. Data were

756 acquired at the sites included in the SRPBS multi-disorder database (i.e., ATT, ATV, COI, HUH, HKH, KPM, SWA, KUT,

757 and UTO) and three additional sites: Kyoto University (KUS; Siemens Skyra) and Yaesu Clinic 1 and 2 (YC1 and YC2; Philips

758 Achieva) (Supplementary Table 3). Each participant underwent three rs-fMRI sessions of 10 min each at nine sites, two

759 sessions of 10 min each at two sites (HKH & HUH), and five cycles (morning, afternoon, next day, next week, next month)

760 consisting of three 10-minute sessions each at a single site (ATT). In the latter situation, one participant underwent four rather

761 than five sessions at the ATT site because of a poor physical condition. Thus, a total of 411 sessions were conducted [8

762 participants $\times (3\times9+2\times2+5\times3\times1) + 1$ participant $\times (3\times9+2\times2+4\times3\times1)$]. During each rs-fMRI session, participants were

763 instructed to maintain a focus on a fixation point at the center of a screen, remain still and awake, and to think about nothing

764 in particular. For sites that could not use a screen in conjunction with fMRI (HKH & KUS), a seal indicating the fixation point

765 was placed on the inside wall of the MRI gantry. Although we attempted to ensure imaging was performed using the same

766 parameters at all sites, there were two phase-encoding directions (P→A and A→P), three MRI manufacturers (Siemens, GE,

767 and Philips), four different numbers of coils (8, 12, 24, 32), and seven scanner types (TimTrio, Verio, Skyra, Spectra, MR750W,

768 SignaHDxt, Achieva) (Supplementary Table 3).

769 All participants in all datasets provided written informed consent, and all recruitment procedures and experimental

770 protocols were approved by the Institutional Review Boards of the principal investigators' respective institutions (Advanced

771 Telecommunications Research Institute International (ATR), Hiroshima University, Kyoto Prefectural University of Medicine,

772 Showa University, The University of Tokyo).

773

774 ## Preprocessing and calculation of the resting-state functional connectivity matrix

775 The rs-fMRI data were preprocessed using SPM8 implemented in MATLAB. The first 10 s of data were discarded to allow

776    for T1 equilibration. Preprocessing steps included slice-timing correction, realignment, co-registration, segmentation of T1-

777    weighted structural images, normalization to Montreal Neurological Institute (MNI) space, and spatial smoothing with an

778    isotropic Gaussian kernel of 6 mm full-width at half-maximum. For the analysis of connectivity matrices, ROIs were delineated

779    according to a 268-node gray matter atlas developed to cluster maximally similar voxels [26]. The BOLD signal time courses

780    were extracted from these 268 ROIs. To remove several sources of spurious variance, we used linear regression with 36

781    regression parameters [44] such as six motion parameters, average signals over the whole brain, white matter, and cerebrospinal

782    fluid. Derivatives and quadratic terms were also included for all parameters. A temporal band-pass filter was applied to the

783    time series using a first-order Butterworth filter with a pass band between 0.01 Hz and 0.08 Hz to restrict the analysis to low-

784    frequency fluctuations, which are characteristic of rs-fMRI BOLD activity [44]. Furthermore, to reduce spurious changes in

785    functional connectivity because of head motion, we calculated frame-wise displacement (FD) and removed volumes with FD

786    > 0.5 mm, as proposed in a previous study [45]. The FD represents head motion between two consecutive volumes as a scalar

787    quantity (i.e., the summation of absolute displacements in translation and rotation). Using the aforementioned threshold, 5.4%

788    ± 10.6% volumes (i.e., the average [approximately 13 volumes] ± 1 SD) were removed per 10 min of rs-fMRI scanning (240

789    volumes) in the traveling-subject dataset, 6.2% ± 13.2% volumes were removed per rs-fMRI session in the SRPBS multi-

790    disorder dataset. If the number of volumes removed after scrubbing exceeded the average of –3 SD across participants in each

791    dataset, the participants or sessions were excluded from the analysis. As a result, 14 sessions were removed from the traveling-

792    subject dataset, 20 participants were removed from the SRPBS multi-disorder dataset. Furthermore, we excluded participants

793    for whom we could not calculate functional connectivity at all 35,778 connections, primarily because of the lack of BOLD

794    signals within an ROI. As a result, 99 participants were further removed from the SRPBS multi-disorder dataset.

795

796    **Principal component analysis**

797    We developed bivariate scatter plots of the first two principal components based on a PCA of functional connectivity values

798    in the SRPBS multi-disorder dataset (Fig. 2). To visualize whether most of the variation in the SRPBS multi-disorder dataset

799    was still associated with imaging site after harmonization, we performed a PCA of functional connectivity values in the

800    harmonized SRPBS multi-disorder dataset (Fig. 7). We used the traveling-subject method for harmonization, as described in

801    the following section.

802

803    **Estimation of biases and factors**

804    The participant factor ($p$), measurement bias ($m$), sampling biases ($s_{hc}, s_{mdd}, s_{scz}$), and psychiatric disorder factor ($d$) were

805    estimated by fitting the regression model to the functional connectivity values of all participants from the SRPBS multi-

806    disorder dataset and the traveling-subject dataset. In this instance, vectors are denoted by lowercase bold letters (e.g., $m$) and

807    all vectors are assumed to be column vectors. Components of vectors are denoted by subscripts such as $m_k$. To represent

808    participant characteristics, we used a 1-of-K binary coding scheme in which the target vector (e.g., $\mathbf{x}_m$) for a measurement

809    bias $m$ belonging to site $k$ is a binary vector with all elements equal to zero—except for element $k$, which equals 1. If a

810    participant does not belong to any class, the target vector is a vector with all elements equal to zero. A superscript T denotes

811    the transposition of a matrix or vector, such that $\mathbf{x}^T$ represents a row vector. For each connectivity, the regression model can

812    be written as follows:

813

814 $$Connectivity = \mathbf{x}_m{}^T m + \mathbf{x}_{s_{hc}}{}^T s_{hc} + \mathbf{x}_{s_{mdd}}{}^T s_{mdd} + \mathbf{x}_{s_{scz}}{}^T s_{scz} + \mathbf{x}_d{}^T d + \mathbf{x}_p{}^T p + const + e,$$

815  such that $\sum_{j}^{9} p_j = 0, \sum_{k}^{12} m_k = 0, \sum_{k}^{6} s_{hc_k} = 0, \sum_{k}^{3} s_{mdd_k} = 0, \sum_{k}^{3} s_{scz_k} = 0, d_1(\text{HC}) = 0,$

816  in which $\boldsymbol{m}$ represents the measurement bias (12 sites $\times$ 1), $\boldsymbol{s_{hc}}$ represents the sampling bias of HCs (six sites $\times$ 1),

817  $\boldsymbol{s_{mdd}}$ represents the sampling bias of patients with MDD (three sites $\times$ 1), $\boldsymbol{s_{scz}}$ represents the sampling bias of patients

818  with SCZ (three sites $\times$ 1), $\boldsymbol{d}$ represents the disorder factor (3 $\times$ 1), $\boldsymbol{p}$ represents the participant factor

819  (nine traveling subjects $\times$ 1), $const$ represents the average functional connectivity value across all participants from all

820  sites, and $e \sim \mathcal{N}(0, \gamma^{-1})$ represents noise. For each functional connectivity value, we estimated the respective parameters

821  using regular ordinary least squares regression with L2 regularization, as the design matrix of the regression model is rank-

822  deficient. When regularization was not applied, we observed spurious anticorrelation between the measurement bias and the

823  sampling bias for HCs, and spurious correlation between the sampling bias for HCs and the sampling bias for patients with

824  psychiatric disorders (Supplementary Figure 3a, left). These spurious correlations were also observed in the permutation data

825  in which there were no associations between the site label and data (Supplementary Figure 3a, right). This finding suggests

826  that the spurious correlations were caused by the rank-deficient property of the design matrix. We tuned the hyper-parameter

827  lambda to minimize the absolute mean of these spurious correlations (Supplementary Figure 3c, left).

828

829  **Analysis of contribution size**

830  To quantitatively verify the magnitude relationship among factors, we calculated and compared the contribution size to

831  determine the extent to which each bias type and factor explain the variance of the data in our linear mixed-effects model (Fig.

832  3d). After fitting the model, the $b$-th connectivity from subject $a$ can be written, as follows:

833

834  $$Connectivity_{a,b} = \mathbf{x}_m^{a\,\mathrm{T}} \boldsymbol{m^b} + \mathbf{x}_{s_{hc}}^{a\,\mathrm{T}} \boldsymbol{s_{hc}^b} + \mathbf{x}_{s_{mdd}}^{a\,\mathrm{T}} \boldsymbol{s_{mdd}^b} + \mathbf{x}_{s_{scz}}^{a\,\mathrm{T}} \boldsymbol{s_{scz}^b} + \mathbf{x}_d^{a\,\mathrm{T}} \boldsymbol{d^b} + \mathbf{x}_p^{a\,\mathrm{T}} \boldsymbol{p^b} + const + e,$$

835

836  For example, the contribution size of measurement bias (i.e., the first term) in this model was calculated as

837

838  $Contribution\ size_m$

839  $$= \frac{1}{N_m} \frac{1}{N_s * N} \sum_{a=1}^{N_s} \sum_{b=1}^{N} \frac{\left(\mathbf{x}_m^{a\,\mathrm{T}} \boldsymbol{m^b}\right)^2}{\left(\mathbf{x}_m^{a\,\mathrm{T}} \boldsymbol{m^b}\right)^2 + \left(\mathbf{x}_{s_{hc}}^{a\,\mathrm{T}} \boldsymbol{s_{hc}^b}\right)^2 + \left(\mathbf{x}_{s_{mdd}}^{a\,\mathrm{T}} \boldsymbol{s_{mdd}^b}\right)^2 + \left(\mathbf{x}_{s_{scz}}^{a\,\mathrm{T}} \boldsymbol{s_{scz}^b}\right)^2 + \left(\mathbf{x}_d^{a\,\mathrm{T}} \boldsymbol{d^b}\right)^2 + \left(\mathbf{x}_p^{a\,\mathrm{T}} \boldsymbol{p^b}\right)^2 + e^2},$$

840

841  in which $N_m$ represents the number of components for each factor, $N$ represents the number of connectivities, $N_s$

842  represents the number of subjects, and $Contribution\ size_m$ represents the magnitude of the contribution size of

843  measurement bias. These formulas were used to assess the contribution sizes of individual factors related to measurement bias

844  (e.g., phase-encoding direction, scanner, coil, and fMRI manufacturer: Fig. 5b). We decomposed the measurement bias into

845  these factors, after which the relevant parameters were estimated. Other parameters were fixed at the same values as previously

846  estimated.

847

848  **Spatial characteristics of measurement bias, sampling bias, and each factor in the brain**

849  To evaluate the spatial characteristics of each type of bias and each factor in the brain, we calculated the magnitude of the

850  effect on each ROI. First, we calculated the median absolute value of the effect on each functional connection among sites or

851  participants for each bias and participant factor. We then calculated the absolute value of each connection for each disorder

852　factor. The uppercase bold letters (e.g., $\boldsymbol{M}$) and subscript vectors (e.g., $\boldsymbol{m}_k$) represent the vectors for the number of functional

853　connections:

854

855
$$\boldsymbol{M} = \underset{k}{\text{median}}(|\boldsymbol{m}_k|), \boldsymbol{S}_{hc} = \underset{k}{\text{median}}(|\boldsymbol{s}_{hc_k}|), \boldsymbol{S}_{mdd} = \underset{k}{\text{median}}(|\boldsymbol{s}_{mdd_k}|), \boldsymbol{S}_{scz} = \underset{k}{\text{median}}(|\boldsymbol{s}_{scz_k}|), \boldsymbol{D}_2 = |\boldsymbol{d}_2|, \boldsymbol{D}_3 = |\boldsymbol{d}_3|, \boldsymbol{P} = \underset{j}{\text{median}}(|\boldsymbol{p}_j|)$$

856

857　We next calculated the magnitude of the effect on ROIs as the average connectivity value between all ROIs, except for

858　themselves.

859

860
$$Effect\_on\_ROI_n = \frac{1}{N_{ROI} - 1} \sum_{v \neq n}^{N_{ROI}} Effect\_on\_connectivity_{n,v},$$

861

862　in which $N_{ROI}$ represents the number of ROIs, $Effect\_on\_ROI_n$ represents the magnitude of the effect on the *n*-th ROI, and

863　$Effect\_on\_connectivity_{n,v}$ represents the magnitude of the effect on connectivity between the *n*-th ROI and *v*-th ROI.

864

865　**Hierarchical clustering analysis for measurement bias**

866　We calculated the Pearson's correlation coefficients among measurement biases $\boldsymbol{m}_k$ ($N \times$

867　$1$, where $N$ is the number of functional connections) for each site *k*, and performed a hierarchical clustering analysis

868　based on the correlation coefficients across measurement biases. To visualize the dendrogram (Fig. 5), we used the

869　"*dendrogram*", "*linkage*", and "*optimalleaforder*" functions in MATLAB (R2015a, Mathworks, USA).

870

871

872　**Comparison of models for sampling bias**

873　We investigated whether sampling bias is caused by the differences in the number of participants among imaging sites, or by

874　sampling from different populations among imaging sites. We constructed two models and investigated which model provides

875　the best explanation of sampling bias. In the single-population model, we assumed that participants were sampled from a single

876　population across imaging sites. In the different-population model, we assumed that participants were sampled from different

877　populations among imaging sites. We first theorized how the number of participants at each site affects the variance of

878　sampling biases across connectivity values, as follows:

879　　In the *single-population model*, we assumed that the functional connectivity values of each participant were

880　generated from an independent Gaussian distribution, with a mean of 0 and a variance of $\xi^2$ for each connectivity value.

881　Then, the functional connectivity vector for participant *j* at site *k* can be described as

882

883
$$\boldsymbol{c}_j^k \sim \mathcal{N}(\boldsymbol{0}, \xi^2 \boldsymbol{I}).$$

884

885　Let $\boldsymbol{c}_k$ be the vector of functional connectivity at site *k* averaged across participants. In this model, $\boldsymbol{c}_k$ represents the

886　sampling bias and can be described as

887
$$\boldsymbol{c}_k = \frac{1}{N_k} \sum_{j=1}^{N_k} \boldsymbol{c}_j^k \sim \mathcal{N}\left(\boldsymbol{0}, \frac{\xi^2}{N_k} \boldsymbol{I}\right),$$

888　in which $N_k$ represents the number of participants at site *k*. The variance across functional connectivity values for $\boldsymbol{c}_k$ is

889    described as

890

$$V_k = \frac{1}{N}\sum_{i=1}^{N}(c_{ki} - \overline{c_k})^2 = \frac{1}{N}\boldsymbol{c_k}^{\mathrm{T}}\left(\mathbf{I} - \frac{1}{N}\mathbf{11'}\right)^{\mathrm{T}}\left(\mathbf{I} - \frac{1}{N}\mathbf{11'}\right)\boldsymbol{c_k} \approx \frac{1}{N}\boldsymbol{c_k}^{\mathrm{T}}\boldsymbol{c_k},$$

891    in which $\mathbf{1}$ represents the $N \times 1$ vector of ones and $\mathbf{I}$ represents the $N \times N$ identity matrix. Since $N$ equals 35,778 and

892    $\frac{1}{35778}$ is sufficiently smaller than 1, we can approximate

893

$$\mathbf{I} - \frac{1}{N}\mathbf{11'} \approx \mathbf{I}.$$

894    Then, the expected value of variance across functional connectivity values for sampling-bias can be described as

895

$$\mathbb{E}[V_k] \approx \frac{1}{N}\mathbb{E}[\boldsymbol{c_k}^{\mathrm{T}}\boldsymbol{c_k}] = \frac{1}{N}Tr\left(\frac{\xi^2}{N_k}\mathbf{I}\right) = \frac{\xi^2}{N_k}.$$

896

897        In the different-population model, we assumed that the functional connectivity values of each participant were

898    generated from a different independent Gaussian distribution, with an average of $\boldsymbol{\beta_k}$ and a variance of $\xi^2$ depending on the

899    population of each site. In this situation, the functional connectivity vector for participant $j$ at site $k$ can be described as

900

$$\boldsymbol{c}_j^k \sim \mathcal{N}(\boldsymbol{\beta_k}, \xi^2\mathbf{I}).$$

901    Here, we assume that the average of the population $\boldsymbol{\beta_k}$ is sampled from an independent Gaussian distribution with an average

902    of 0 and a variance of $\sigma^2$. That is, $\boldsymbol{\beta_k}$ is expressed as

903

904

$$\boldsymbol{\beta_k} \sim \mathcal{N}(\mathbf{0}, \sigma^2\mathbf{I}).$$

905

906    The vector of functional connectivity for site $k$ averaged across participants can then be described as

907

$$\boldsymbol{c_k} \sim \mathcal{N}\left(\mathbf{0}, \left(\frac{\xi^2}{N_k} + \sigma^2\right)\mathbf{I}\right).$$

908    The variance across functional connectivity values for $\boldsymbol{c_k}$ can be described as

909

$$\mathbb{E}[V_k] \approx \frac{\xi^2}{N_k} + \sigma^2.$$

910

911        In summary, the variance of sampling bias across functional connectivity values in each model is expressed by the

912    number of participants at a given site, as follows:

913

914

$$\text{single-population model: } y_k = -x_k + 2\log_{10}\xi$$

915

916

$$\text{different-population model: } y_k = -\log_{10}(\xi^2 10^{-x_k} + \sigma^2),$$

917

918    in which $y_k = \log_{10}(v_k)$, $v_k$ represents the variance across functional connectivity values for $\boldsymbol{s_{hc_k}}$, $\boldsymbol{s_{hc_k}}$ represents the

919    sampling bias of HCs at site $k$ ($N \times 1$: $N$ is the number of functional connectivity), $x_k = \log_{10}(N_k)$, and $N_k$

920    represents the number of participants at site $k$. We estimated the parameters $\xi$ and $\sigma$ using the MATLAB (R2015a,

921    Mathworks, USA) optimization function "*fminunc*". To simplify statistical analyses, sampling bias was estimated based on

922    functional connectivity in which the average across all participants was set to zero.

923    We aimed to determine which model provided the best explanation of sampling bias in our data by calculating the

924    corrected Akaike information criterion (AICc; under the assumption of a Gaussian distribution) for small-sample data [36, 37],

925    as well as BIC:

926

927
$$\text{AICc} = \sum_{k=1}^{6} \ln \varphi_k{}^2 + 2q + \frac{2q(q+1)}{(6-q-1)},$$

928
$$\text{BIC} = \sum_{k=1}^{6} \ln \varphi_k{}^2 + q * \log(6),$$

929

930    in which $\varphi_k = v_k - \widehat{v_k}$, $\widehat{v_k}$ represents the estimated variance, and $q$ represents the number of parameters in each model (1

931    or 2).

932    To investigate prediction performance, we used leave-one-site-out-cross-validation in which we estimated the

933    parameters $\xi$ and $\sigma$ using data from five sites. The variance of sampling bias was predicted based on the number of

934    participants at the remaining site. This procedure was repeated to predict variance values for sampling bias at all six sites. We

935    then calculated the absolute errors between predicted and actual variances for all sites.

936

937    **Harmonization procedures**

938    We compared four different harmonization methods for the removal of site differences, as described in the main text.

939

940    *Traveling-subject harmonization*

941    Measurement biases were estimated by fitting the regression model to the combined SRPBS multi-disorder and traveling-

942    subject datasets in the same way in "Estimation of biases and factors" section. For each connectivity, the regression model can

943    be written as follows:

944

945
$$\textit{Connectivity} = \mathbf{x}_m{}^{\mathbf{T}} \boldsymbol{m} + \mathbf{x}_{s_{hc}}{}^{\mathbf{T}} \boldsymbol{s}_{hc} + \mathbf{x}_{s_{mdd}}{}^{\mathbf{T}} \boldsymbol{s}_{mdd} + \mathbf{x}_{s_{scz}}{}^{\mathbf{T}} \boldsymbol{s}_{scz} + \mathbf{x}_d{}^{\mathbf{T}} \boldsymbol{d} + \mathbf{x}_p{}^{\mathbf{T}} \boldsymbol{p} + \boldsymbol{const} + \boldsymbol{e}. \tag{1}$$

946

947    Measurement bias were removed by subtracting the estimated measurement biases. Thus, the harmonized functional

948    connectivity values were set, as follows:

949

950
$$\textit{Connectivity}^{Traveling-subje} = \textit{Connectivity} - \mathbf{x}_m{}^{\mathbf{T}} \widehat{\boldsymbol{m}},$$

951

952    in which $\widehat{\boldsymbol{m}}$ represents the estimated measurement bias.

953

954    *GLM harmonization*

955    The GLM harmonization method adjusts the functional connectivity value for site difference using GLM. Site differences were

956    estimated by fitting the regression model, which included site label only, to the SRPBS multi-disorder dataset only. The

957    regression model can be written as

958

959
$$\textit{Connectivity} = \boldsymbol{const} + \mathbf{x}_s{}^{\mathbf{T}} \boldsymbol{s}^{GLM} + \boldsymbol{e}, \tag{2}$$

960

961   in which $s^{GLM}$ represents the site difference (nine sites $\times$ 1). For each functional connectivity value, we estimated the

962   parameters using regular ordinary least squares regression. Site differences were removed by subtracting the estimated site

963   differences. Thus, the harmonized functional connectivity values were set, as follows:

964

965   $$Connectivity^{GLM} = Connectivity - \mathbf{x}_s{}^T \widehat{s^{GLM}},$$

966

967   in which $\widehat{s^{GLM}}$ represents the estimated site difference.

968

969   *Adjusted GLM harmonization*

970   Site differences were estimated by fitting the regression model, which included site label and diagnosis label, to the SRPBS

971   multi-disorder dataset. The regression model can be written as

972

973   $$\boldsymbol{Connectivity} = \boldsymbol{const} + \mathbf{x}_s{}^T \boldsymbol{s}^{Adj} + \mathbf{x}_d{}^T \boldsymbol{d}^{Adj} + \boldsymbol{e}, \qquad (3)$$

974

975   In which $s^{Adj}$ represents the site difference (nine sites $\times$ 1). For each functional connectivity value, we estimated the

976   parameters via regular ordinary least squares regression. Site differences were removed by subtracting the estimated site

977   difference only. Thus, the harmonized functional connectivity values were set, as follows:

978

979   $$Connectivity^{Adj} = Connectivity - \mathbf{x}_s{}^T \widehat{s^{Adj}},$$

980

981   in which $\widehat{s^{Adj}}$ represents the estimated site difference.

982

983   *ComBat harmonization*

984   The ComBat harmonization model [16, 17, 19, 38] extends the adjusted GLM harmonization method in two ways: (1) it models

985   site-specific scaling factors and (2) it uses empirical Bayesian criteria to improve the estimation of site parameters for small

986   sample sizes. The model assumes that the expected connectivity value can be modeled as a linear combination of the biological

987   variables and the site differences in which the error term is modulated by additional site-specific scaling factors.

988

989   $$\boldsymbol{Connectivity} = \boldsymbol{const} + \mathbf{x}_s{}^T \boldsymbol{s}^{ComBat} + \mathbf{x}_d{}^T \boldsymbol{d}^{ComBat} + \delta_k \boldsymbol{e}, \qquad (4)$$

990

991   in which $s^{ComBat}$ represents the site difference (nine sites $\times$ 1), and $\delta_k$ represents the scale parameter for site differences at

992   site *k* for the respective connectivity value. The harmonized functional connectivity values were set, as follows:

993

995   $$Connectivity^{ComBat} = \frac{Connectivity - const - \mathbf{x}_s{}^T \widehat{s^{ComBat}} - \mathbf{x}_d{}^T \widehat{d^{ComBat}}}{\widehat{\delta_k}} + const + \mathbf{x}_d{}^T \widehat{d^{ComBat}},$$

994

996   in which $\widehat{\delta_k}$, $\widehat{d^{ComBat}}$, and $\widehat{s^{ComBat}}$ are the empirical Bayes estimates of $\delta_k$, $d^{ComBat}$, and $s^{ComBat}$, respectively using

997   "combat" function in https://github.com/Jfortin1/ComBatHarmonization. Thus, ComBat simultaneously models and estimates

998   biological and nonbiological terms and algebraically removes the estimated additive and multiplicative site differences. Of

999   note, in the ComBat model, we included diagnosis as covariates to preserve important biological trends in the data and avoid

1000  overcorrection.

1001

1002  **2-fold cross-validation evaluation procedure**

1003  We compared four different harmonization methods for the removal of site difference or measurement bias by 2-fold cross-

1004  validation, as described in the main text. In the traveling-subject harmonization method, we estimated the measurement bias

1005  by applying the regression model written in equation (1) in "Harmonization procedures" section to the estimating dataset. Thus,

1006  the harmonized functional connectivity values in testing dataset were set, as follows:

1007  $$connectivity^{Traveling-subje}_{testing\ dataset} = Connectivity_{testing\ dataset} - \mathbf{x_m}^\mathrm{T}\hat{\boldsymbol{m}}_{estimating\ dataset},$$

1008  in which $\hat{\boldsymbol{m}}_{estimating\ dataset}$ represents the estimated measurement bias using the estimating dataset.

1009  By contrast, in the other harmonization methods, we estimated the site differences by applying the regression models written

1010  in equations (2)–(4) in "Harmonization procedures" section to the estimating dataset (fold1 data). Thus, the harmonized

1011  functional connectivity values in testing dataset were set, as follows:

1012  $$connectivity^{GLM}_{testing\ dataset} = Connectivity_{testing\ dataset} - \mathbf{x_s}^\mathrm{T}\widehat{\boldsymbol{s^{GLM}}}_{fold1},$$

1013  $$connectivity^{Adj}_{testing\ dataset} = Connectivity_{testing\ dataset} - \mathbf{x_s}^\mathrm{T}\widehat{\boldsymbol{s^{Adj}}}_{fold1},$$

1014  $$connectivity^{ComBat}_{testing\ dataset} = Connectivity_{testing\ dataset} - \mathbf{x_s}^\mathrm{T}\boldsymbol{s^{\widehat{ComBat}}}_{fold1},$$

1015  in which $\widehat{\boldsymbol{s^{GLM}}}_{fold}$ , $\widehat{\boldsymbol{s^{Adj}}}_{fold1}$, $\boldsymbol{s^{\widehat{ComBat}}}_{fo}$   represents the estimated site differences using fold1 data.

1016  We then estimated the measurement bias, participant factor, and disorder factors by applying the regression model written in

1017  equation (1) to the harmonized functional connectivity values in the testing dataset. Finally, we evaluated the standard

1018  deviation of the magnitude distribution of measurement bias calculated in the same way as described in "Quantification of site

1019  differences" section among the harmonization methods. This procedure was done again by exchanging the estimating dataset

1020  and the testing dataset.

1021

1022  **Code availability:** All codes used for the analyses are available from the authors on request.

1023  **Data availability:** All relevant data are available from the authors on request. All data can be downloaded publicly from the

1024  following site: https://bicr-resource.atr.jp/decnefpro/.

1030  **Author contributions:** A.Y., N.Y., and H.I. designed the study. N.Y., T.I., T.Y., N.I., M.T., Y.Y., A.K., N.O., T.Y., K.M.,

1031  R.H., G.O., Y.S., J.N., Y.S., K.K., N.K., H.T., Y.O. and S.T. recruited participants of the study, collected their clinical and

1032  imaging data and constructed the database. A.Y. performed data preprocessing and data analysis under the super vision of

1033  G.L., J.M., O.Y., M.K., and H.I., and A.Y., O.Y., M.K. and H.I. primarily wrote the manuscript.

1034  **Competing financial interests:** The authors declare no competing financial interests.

1035

1036

## Reference

1. Glasser MF, Smith SM, Marcus DS, Andersson JL, Auerbach EJ, Behrens TE, et al. The Human Connectome Project's neuroimaging approach. Nat Neurosci. 2016;19(9):1175-87. doi: 10.1038/nn.4361. PubMed PMID: 27571196.

2. Yamada T, Hashimoto RI, Yahata N, Ichikawa N, Yoshihara Y, Okamoto Y, et al. Resting-State Functional Connectivity-Based Biomarkers and Functional MRI-Based Neurofeedback for Psychiatric Disorders: A Challenge for Developing Theranostic Biomarkers. Int J Neuropsychopharmacol. 2017;20(10):769-81. doi: 10.1093/ijnp/pyx059. PubMed PMID: 28977523; PubMed Central PMCID: PMCPMC5632305.

3. Biswal BB, Mennes M, Zuo XN, Gohel S, Kelly C, Smith SM, et al. Toward discovery science of human brain function. Proc Natl Acad Sci U S A. 2010;107(10):4734-9. doi: 10.1073/pnas.0911855107. PubMed PMID: 20176931; PubMed Central PMCID: PMCPMC2842060.

4. Woo CW, Chang LJ, Lindquist MA, Wager TD. Building better biomarkers: brain models in translational neuroimaging. Nat Neurosci. 2017;20(3):365-77. doi: 10.1038/nn.4478. PubMed PMID: 28230847.

5. Xia M, He Y. Functional connectomics from a "big data" perspective. NeuroImage. 2017;160:152-67. doi: 10.1016/j.neuroimage.2017.02.031. PubMed PMID: 28232122.

6. Di Martino A, Yan CG, Li Q, Denio E, Castellanos FX, Alaerts K, et al. The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain architecture in autism. Mol Psychiatry. 2014;19(6):659-67. doi: 10.1038/mp.2013.78. PubMed PMID: 23774715; PubMed Central PMCID: PMCPMC4162310.

7. Pearlson G. Multisite collaborations and large databases in psychiatric neuroimaging: advantages, problems, and challenges. Schizophr Bull. 2009;35(1):1-2. doi: 10.1093/schbul/sbn166. PubMed PMID: 19023121; PubMed Central PMCID: PMCPMC2643967.

8. Yahata N, Kasai K, Kawato M. Computational neuroscience approach to biomarkers and treatments for mental disorders. Psychiatry Clin Neurosci. 2017;71(4):215-37. doi: 10.1111/pcn.12502. PubMed PMID: 28032396.

9. Takagi Y, Sakai Y, Lisi G, Yahata N, Abe Y, Nishida S, et al. A Neural Marker of Obsessive-Compulsive Disorder from Whole-Brain Functional Connectivity. Sci Rep. 2017;7(1):7538. Epub 2017/08/10. doi: 10.1038/s41598-017-07792-7. PubMed PMID: 28790433; PubMed Central PMCID: PMCPMC5548868.

10. Yahata N, Morimoto J, Hashimoto R, Lisi G, Shibata K, Kawakubo Y, et al. A small number of abnormal brain connections predicts adult autism spectrum disorder. Nat Commun. 2016;7:11254. Epub 2016/04/15. doi: 10.1038/ncomms11254. PubMed PMID: 27075704; PubMed Central PMCID: PMCPMC4834637.

11. Abraham A, Milham MP, Di Martino A, Craddock RC, Samaras D, Thirion B, et al. Deriving reproducible biomarkers from multi-site resting-state data: An Autism-based example. NeuroImage. 2017;147:736-45. doi: 10.1016/j.neuroimage.2016.10.045. PubMed PMID: 27865923.

12. Nieuwenhuis M, Schnack HG, van Haren NE, Lappin J, Morgan C, Reinders AA, et al. Multi-center MRI prediction models: Predicting sex and illness course in first episode psychosis patients. NeuroImage. 2017;145(Pt B):246-53. doi: 10.1016/j.neuroimage.2016.07.027. PubMed PMID: 27421184; PubMed Central PMCID: PMCPMC5193177.

13. Orban P, Dansereau C, Desbois L, Mongeau-Perusse V, Giguere CE, Nguyen H, et al. Multisite

1075    generalizability of schizophrenia diagnosis classification based on functional brain connectivity. Schizophrenia

1076    research. 2018;192:167-71. Epub 2017/06/12. doi: 10.1016/j.schres.2017.05.027. PubMed PMID: 28601499.

1077    14. Dansereau C, Benhajali Y, Risterucci C, Pich EM, Orban P, Arnold D, et al. Statistical power and prediction

1078    accuracy in multisite resting-state fMRI connectivity. NeuroImage. 2017;149:220-32. doi:

1079    10.1016/j.neuroimage.2017.01.072. PubMed PMID: 28161310.

1080    15. Watanabe T, Sasaki Y, Shibata K, Kawato M. Advances in fMRI Real-Time Neurofeedback. Trends Cogn

1081    Sci. 2017;21(12):997-1010. doi: 10.1016/j.tics.2017.09.010. PubMed PMID: 29031663; PubMed Central

1082    PMCID: PMCPMC5694350.

1083    16. Fortin JP, Cullen N, Sheline YI, Taylor WD, Aselcioglu I, Cook PA, et al. Harmonization of cortical thickness

1084    measurements across scanners and sites. NeuroImage. 2017;167:104-20. doi:

1085    10.1016/j.neuroimage.2017.11.024. PubMed PMID: 29155184.

1086    17. Fortin JP, Parker D, Tunc B, Watanabe T, Elliott MA, Ruparel K, et al. Harmonization of multi-site diffusion

1087    tensor imaging data. NeuroImage. 2017;161:149-70. doi: 10.1016/j.neuroimage.2017.08.047. PubMed PMID:

1088    28826946; PubMed Central PMCID: PMCPMC5736019.

1089    18. Rao A, Monteiro JM, Mourao-Miranda J, Alzheimer's Disease I. Predictive modelling using neuroimaging

1090    data in the presence of confounds. NeuroImage. 2017;150:23-49. doi: 10.1016/j.neuroimage.2017.01.066.

1091    PubMed PMID: 28143776; PubMed Central PMCID: PMCPMC5391990.

1092    19. Yu M, Linn KA, Cook PA, Phillips ML, McInnis M, Fava M, et al. Statistical harmonization corrects site

1093    effects in functional connectivity measurements from multi-site fMRI data. Hum Brain Mapp. 2018. doi:

1094    10.1002/hbm.24241. PubMed PMID: 29962049.

1095    20. Noble S, Scheinost D, Finn ES, Shen X, Papademetris X, McEwen SC, et al. Multisite reliability of MR-based

1096    functional connectivity. NeuroImage. 2017;146:959-70. doi: 10.1016/j.neuroimage.2016.10.020. PubMed

1097    PMID: 27746386; PubMed Central PMCID: PMCPMC5322153.

1098    21. Benazzi F. Various forms of depression. Dialogues in Clinical Neuroscience. 2006;8(2):151-61. PubMed

1099    PMID: PMC3181770.

1100    22. Drysdale AT, Grosenick L, Downar J, Dunlop K, Mansouri F, Meng Y, et al. Resting-state connectivity

1101    biomarkers define neurophysiological subtypes of depression. Nat Med. 2017;23(1):28-38. doi:

1102    10.1038/nm.4246. PubMed PMID: 27918562; PubMed Central PMCID: PMCPMC5624035.

1103    23. Ng B, Dressler M, Varoquaux G, Poline JB, Greicius M, Thirion B. Transport on Riemannian manifold for

1104    functional connectivity-based classification. Medical image computing and computer-assisted intervention :

1105    MICCAI    International Conference on Medical Image Computing and Computer-Assisted Intervention.

1106    2014;17(Pt 2):405-12. Epub 2014/12/09. PubMed PMID: 25485405.

1107    24. Finn ES, Shen X, Scheinost D, Rosenberg MD, Huang J, Chun MM, et al. Functional connectome

1108    fingerprinting: identifying individuals using patterns of brain connectivity. Nat Neurosci. 2015;18(11):1664-

1109    71. doi: 10.1038/nn.4135. PubMed PMID: 26457551; PubMed Central PMCID: PMCPMC5008686.

1110    25. Rosenberg MD, Finn ES, Scheinost D, Papademetris X, Shen X, Constable RT, et al. A neuromarker of

1111    sustained attention from whole-brain functional connectivity. Nat Neurosci. 2016;19(1):165-71. doi:

1112    10.1038/nn.4179. PubMed PMID: 26595653; PubMed Central PMCID: PMCPMC4696892.

1113    26. Shen X, Tokoglu F, Papademetris X, Constable RT. Groupwise whole-brain parcellation from resting-state

fMRI data for network node identification. NeuroImage. 2013;82:403-15. doi: 10.1016/j.neuroimage.2013.05.081. PubMed PMID: 23747961; PubMed Central PMCID: PMCPMC3759540.

27. Noble S, Spann MN, Tokoglu F, Shen X, Constable RT, Scheinost D. Influences on the Test-Retest Reliability of Functional Connectivity MRI and its Relationship with Behavioral Utility. Cereb Cortex. 2017;27(11):5415-29. doi: 10.1093/cercor/bhx230. PubMed PMID: 28968754.

28. Jezzard P, Clare S. Sources of distortion in functional MRI data. Hum Brain Mapp. 1999;8(2-3):80-5. Epub 1999/10/19. PubMed PMID: 10524596.

29. Weiskopf N, Hutton C, Josephs O, Deichmann R. Optimal EPI parameters for reduction of susceptibility-induced BOLD sensitivity losses: a whole-brain analysis at 3 T and 1.5 T. NeuroImage. 2006;33(2):493-504. doi: 10.1016/j.neuroimage.2006.07.029. PubMed PMID: 16959495.

30. Kaiser RH, Andrews-Hanna JR, Wager TD, Pizzagalli DA. Large-Scale Network Dysfunction in Major Depressive Disorder: A Meta-analysis of Resting-State Functional Connectivity. JAMA Psychiatry. 2015;72(6):603-11. doi: 10.1001/jamapsychiatry.2015.0071. PubMed PMID: 25785575; PubMed Central PMCID: PMCPMC4456260.

31. Mulders PC, van Eijndhoven PF, Schene AH, Beckmann CF, Tendolkar I. Resting-state functional connectivity in major depressive disorder: A review. Neurosci Biobehav Rev. 2015;56:330-44. doi: 10.1016/j.neubiorev.2015.07.014. PubMed PMID: 26234819.

32. Kuhn S, Gallinat J. Resting-state brain activity in schizophrenia and major depression: a quantitative meta-analysis. Schizophr Bull. 2013;39(2):358-65. doi: 10.1093/schbul/sbr151. PubMed PMID: 22080493; PubMed Central PMCID: PMCPMC3576173.

33. Li T, Wang Q, Zhang J, Rolls ET, Yang W, Palaniyappan L, et al. Brain-Wide Analysis of Functional Connectivity in First-Episode and Chronic Stages of Schizophrenia. Schizophr Bull. 2017;43(2):436-48. doi: 10.1093/schbul/sbw099. PubMed PMID: 27445261; PubMed Central PMCID: PMCPMC5605268.

34. Minzenberg MJ, Laird AR, Thelen S, Carter CS, Glahn DC. Meta-analysis of 41 functional neuroimaging studies of executive function in schizophrenia. Archives of general psychiatry. 2009;66(8):811-22. Epub 2009/08/05. doi: 10.1001/archgenpsychiatry.2009.91. PubMed PMID: 19652121; PubMed Central PMCID: PMCPMC2888482.

35. Anderson JS, Nielsen JA, Froehlich AL, DuBray MB, Druzgal TJ, Cariello AN, et al. Functional connectivity magnetic resonance imaging classification of autism. Brain. 2011;134(Pt 12):3742-54. doi: 10.1093/brain/awr263. PubMed PMID: 22006979; PubMed Central PMCID: PMCPMC3235557.

36. Burnham KP, Anderson DR. Model selection and multimodel inference: a practical information-theoretic approach: Springer Science & Business Media; 2003.

37. Cortese A, Amano K, Koizumi A, Kawato M, Lau H. Multivoxel neurofeedback selectively modulates confidence without changing perceptual performance. Nat Commun. 2016;7:13669. doi: 10.1038/ncomms13669. PubMed PMID: 27976739; PubMed Central PMCID: PMCPMC5171844 inventor of patents related to the neurofeedback method used in this study, and the original assignee of the patents is ATR, with which M.K. is affiliated. The remaining authors declare no competing financial interests.

38. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. Biostatistics. 2007;8(1):118-27. doi: 10.1093/biostatistics/kxj037. PubMed PMID: 16632515.

1153    39. Hutton C, Bork A, Josephs O, Deichmann R, Ashburner J, Turner R. Image distortion correction in fMRI: A
1154        quantitative evaluation. NeuroImage. 2002;16(1):217-40. doi: 10.1006/nimg.2001.1054. PubMed PMID:
1155        11969330.

1156    40. Jenkinson M. Fast, automated, N-dimensional phase-unwrapping algorithm. Magn Reson Med.
1157        2003;49(1):193-7. doi: 10.1002/mrm.10354. PubMed PMID: 12509838.

1158    41. Jezzard P, Balaban RS. Correction for geometric distortion in echo planar images from B0 field variations.
1159        Magn Reson Med. 1995;34(1):65-73. Epub 1995/07/01. PubMed PMID: 7674900.

1160    42. Andersson JLR, Skare S, Ashburner J. How to correct susceptibility distortions in spin-echo echo-planar
1161        images: application to diffusion tensor imaging. NeuroImage. 2003;20(2):870-88. doi: 10.1016/s1053-
1162        8119(03)00336-7.

1163    43. Wang S, Peterson DJ, Gatenby JC, Li W, Grabowski TJ, Madhyastha TM. Evaluation of Field Map and
1164        Nonlinear Registration Methods for Correction of Susceptibility Artifacts in Diffusion MRI. Front Neuroinform.
1165        2017;11:17. doi: 10.3389/fninf.2017.00017. PubMed PMID: 28270762; PubMed Central PMCID:
1166        PMCPMC5318394.

1167    44. Ciric R, Wolf DH, Power JD, Roalf DR, Baum GL, Ruparel K, et al. Benchmarking of participant-level
1168        confound regression strategies for the control of motion artifact in studies of functional connectivity.
1169        NeuroImage. 2017;154:174-87. doi: 10.1016/j.neuroimage.2017.03.020. PubMed PMID: 28302591; PubMed
1170        Central PMCID: PMCPMC5483393.

1171    45. Power JD, Barnes KA, Snyder AZ, Schlaggar BL, Petersen SE. Spurious but systematic correlations in
1172        functional connectivity MRI networks arise from subject motion. NeuroImage. 2012;59(3):2142-54. doi:
1173        10.1016/j.neuroimage.2011.10.018. PubMed PMID: 22019881; PubMed Central PMCID: PMCPMC3254728.
1174