

¹ Statistical power and prediction accuracy in multisite
² resting-state fMRI connectivity

³ Christian Dansereau^{a,b}, Yassine Benhajali^{a,c}, Celine Risterucci^d, Emilio Merlo
⁴ Pich^{d,**}, Pierre Orban^a, Douglas Arnold^e, Pierre Bellec^{a,b,*}

⁵ ^a*Centre de Recherche de l'Institut Universitaire de Gériatrie de Montréal, Montréal, CA*

⁶ ^b*Department of Computer Science and Operations Research, University of Montreal,
Montreal, CA*

⁷ ^c*Département d'anthropologie, Université Montréal, Montréal, CA*

⁸ ^d*Clinical Imaging, pRED, F.Hoffmann-La Roche, Basel, CH*

⁹ ^e*NeuroRx inc., Montréal, CA*

¹¹ **Abstract**

Connectivity studies using resting-state functional magnetic resonance imaging are increasingly pooling data acquired at multiple sites. While this may allow investigators to speed up recruitment or increase sample size, multisite studies also potentially introduce systematic biases in connectivity measures across sites. In this work, we measure the inter-site effect in connectivity and its impact on our ability to detect individual and group differences. Our study was based on real, as opposed to simulated, multisite fMRI datasets collected in $N = 345$ young, healthy subjects across 8 scanning sites with 3T scanners and heterogeneous scanning protocols, drawn from the 1000 functional connectome project. We first empirically show that typical functional networks were reliably found at the group level in all sites, and that the amplitude of the inter-site effects was small to moderate, with a Cohen's effect size below 0.5 on average across brain connections. We then implemented a series of Monte-Carlo simulations, based on real data, to evaluate the impact of the multisite effects on detection power in statistical tests comparing two groups (with and without the effect) using a general linear model, as well as on the prediction of group labels with a support-vector machine. As a reference, we also implemented the same simulations with fMRI data collected at a single site using an identical sample size. Simulations revealed that using data from heterogeneous sites only slightly decreased our ability to detect changes compared to a monosite study with the GLM, and had a greater impact on prediction accuracy. However, the deleterious effect of multisite data pooling tended to decrease as the total sample size increased, to a point where differences between monosite and multisite simulations were small with $N = 120$ subjects. Taken together, our results support the feasibility of multisite studies in rs-fMRI provided the sample size is large enough.

¹² **Keywords:** multisite, statistical power, prediction accuracy, Monte-Carlo
¹³ simulation, sample size, resting-state, fMRI connectivity, SVM

^{*}Corresponding author: pierre.bellec@criugm.qc.ca

Preprint submitted to CNS/PAIN, Fakuda Development Centre Europe, London, UK 27, 2017

Email address: pierre.bellec@criugm.qc.ca (Pierre Bellec)

14 **Highlights**

- 15 • Small to moderate systematic site effects in fMRI connectivity.
- 16 • Small impact of site effects on the detection of group differences for sample
17 size > 100.
- 18 • Linear regression of the sites prior to multivariate prediction do not im-
19 prove prediction accuracy.

20 **1. Introduction**

21 *Main objective.* Multisite studies are becoming increasingly common in resting-
22 state functional magnetic resonance imaging (rs-fMRI). In particular, some
23 consortia have retrospectively pooled rs-fMRI data from multiple independent
24 studies comparing clinical cohorts with control groups, e.g. normal controls in
25 the 1000 functional connectome project (FCP) (Biswal et al., 2010), children
26 and adolescents suffering from attention deficit hyperactivity disorder from the
27 ADHD200 (Milham et al., 2012; Fair et al., 2012), individuals diagnosed with
28 autism spectrum disorder in ABIDE (Nielsen et al., 2013), individuals suffer-
29 ing from schizophrenia (Cheng et al., 2015), or elderly subjects suffering from
30 mild cognitive impairment (Tam et al., 2015). The rationale behind such initia-
31 tives is to dramatically increase the sample size at the cost of decreased sample
32 homogeneity. The systematic variations of connectivity measures derived using
33 different scanners, called site effects, may decrease the statistical power of group
34 comparisons, and somewhat mitigate the benefits of having a large sample size
35 (Brown et al., 2011; Jovicich et al., 2016). In this work, our main objective
36 was to quantitatively assess the impact of site effects on group comparisons in
37 rs-fMRI connectivity.

38 *Group comparison in rs-fMRI connectivity.* In this work, we focused on the most
39 common measure of individual functional connectivity, which is the Pearson's
40 correlation coefficient between the average rs-fMRI time series of two brain re-
41 gions. To compare two groups, a general linear model (GLM) is typically used
42 to establish the statistical significance of the difference in average connectiv-
43 ity between the groups. Finally a *p*-value is generated for each connection to
44 quantify the probability that the difference in average connectivity is signifi-
45 cantly different from zero (Worsley and Friston, 1995; Yan et al., 2013). If the
46 estimated *p*-value is smaller than a prescribed tolerable level of false-positive
47 findings (see for more detail Table 1), generally adjusted for the number of tests
48 performed across connections, say $\alpha = 0.001$, then the difference in connectivity
49 is deemed significant.

50 *Statistical power in group comparisons at multiple sites.* The statistical power of
51 a group comparison study is the probability of finding a significant difference,
52 when there is indeed a true difference. A careful study design involves the
53 selection of a sample size that is large enough to reach a set level of statistical

		Detected value	
		patho	no patho
Actual value	patho	True Positive	False Negative
	no patho	False Positive	True Negative

Table 1: Confusion matrix.

power, e.g. 80%. In the GLM, the statistical power actually depends on a series of parameters (Desmond and Glover, 2002; Durnez et al., 2014): (1) the sample size (the larger the better); (2) the absolute size of the group difference (the larger the better), and, (3) the intrinsic variability of measurements (the smaller the better) (4) the rejection threshold α for the null hypothesis.

Sources of variability: factors inherent to the scanning protocol. In a multisite (or multi-protocol) setting, differences in imaging or study parameters may add variance to rs-fMRI measures, e.g. the scanner make and model (Friedman et al., 2006, 2008), repetition time, flip angle, voxel resolution or acquisition volume (Friedman and Glover, 2006), experimental design such as eyes-open/eyes-closed (Yan et al., 2009), experiment duration (Van Dijk et al., 2010), and scanning environment such as sound attenuation measures (Elliott et al., 1999), or head-motion restraint techniques (Edward et al., 2000; Van Dijk et al., 2012), amongst others. These parameters can be harmonized to some extent, but differences are unavoidable in large multisite studies. The recent work of Yan et al. (2013) has indeed demonstrated the presence of significant site effects in rs-fMRI measures in the 1000 FCP. Site effects will increase the variability of measures, and thus decrease statistical power. To the best of our knowledge, it is not yet known how important this decrease in statistical power may be.

Sources of variability: within-subject. The relative importance of site effects in rs-fMRI connectivity depends on the amplitude of the many other sources of variance. First, rs-fMRI connectivity only has moderate-to-good test-retest reliability using standard 10-minute imaging protocols (Shehzad et al., 2009), even when using a single scanner and imaging session. Differences in functional connectivity across subjects are also known to correlate with a myriad of behavioural and demographic subject characteristics (Anand et al., 2007; Sheline et al., 2010; Kilpatrick et al., 2006). Taken together, these sources of variance reflect a fundamental volatility of human physiological signals.

Sources of variability: factors inherent to the site. In addition to physiology, some imaging artefacts will vary systematically from session to session, even at a single site. For example, intensity non-uniformities across the brain depend

85 on the positioning of subjects (Caramanos et al., 2010). Room temperature
86 has also been shown to impact MRI measures (Vanhoutte et al., 2006). Given
87 the good consistency of key findings in resting-state connectivity across sites,
88 such as the organization of distributed brain networks (Biswal et al., 2010),
89 it is reasonable to hypothesize that site effects will be small compared to the
90 combination of physiological and within-site imaging variance.

91 *Multivariate analysis.* Another important consideration regarding the impact of
92 site effects on group comparison in rs-fMRI connectivity is the type of method
93 used to identify differences. The concept of statistical power is very well estab-
94 lished in the GLM framework, which tests one brain connection at a time (mass
95 univariate testing). However, multivariate methods that combine several or all
96 connectivity values in a single prediction are also widely used and likely affected
97 by the site effects. A popular multivariate technique in rs-fMRI is support-
98 vector machine (SVM) (Cortes and Vapnik, 1995). In this approach, the group
99 sample is split into a training set and a test set. The SVM is trained to predict
100 group labels on the training set, and the accuracy of the prediction is evaluated
101 independently on the test set. The accuracy level of the SVM captures the
102 quality of the prediction of clinical labels from resting-state connectivity, but
103 does not explicitly tell which brain connection is critical for the prediction. The
104 accuracy score can thus be seen as a separability index between the individuals
105 of two groups in high dimensional space. Altogether, the objectives and mea-
106 sures of statistical risk for SVM and GLM are quite different. Because SVM
107 has the ability to combine measures across connections, unlike univariate GLM
108 tests, we hypothesized that the GLM and SVM will be impacted differently by
109 site effects. Even though the accuracy is expected to be lower for the multisite
110 than the monosite configuration, it has been shown that the generalizability of a
111 predictive model to unseen sites is greater for models trained on multisite than
112 monosite datasets as shown by Abraham et al. (2016).

113 *Specific objectives.* Our first objective was to characterize, using real data, the
114 amplitude of systematic site effects in rs-fMRI connectivity measures across
115 sites, as a function of within-site variance. We based our evaluation on images
116 generated from independent groups at 8 sites equipped with 3T scanners, in a
117 subset ($N = 345$) of the 1000 FCP. Our second objective was to evaluate the
118 impact of site effects on the detection power of group differences in rs-fMRI
119 connectivity. To answer this question directly, one would need to scan two dif-
120 ferent cohorts of participants at least twice, once in a multisite setting and once
121 in a monosite setting. Such an experiment may be too costly to implement for
122 addressing a purely technical objective. As a more feasible alternative, we im-
123 plemented a series of Monte Carlo simulations, adding synthetic “pathological”
124 effects in the 1000 FCP sample. One interesting feature of the “1000 FCP”
125 dataset is the presence of one large site of ~ 200 subjects and 7 small sites of
126 ~ 20 subjects per site. We were therefore able to implement realistic scenarios
127 following either a monosite or a multisite design (with 7 sites), with the same
128 total sample size. Our simulations gave us full control on critical aspects for the

Site	Magnet	Scanner	Channels	N	Nfinal	Sex	Age	TR	#Slices	#Frames	
Baltimore, USA	3T	Philips Achieva	8	23	21	8M/15F	20-40	2.5	47	123	
Berlin, DE	3T	Siemens Tim Trio	12	26	26	13M/13F	23-44	2.3	34	195	
Cambridge, USA	3T	Siemens Tim Trio	12	198	195	75M/123F	18-30	3	47	119	
Newark, USA	3T	Siemens Allegra	12	19	17	9M/10F	21-39	2	32	135	
NewYork.b, USA	3T	Siemens Allegra	1	20	18	8M/12F	18-46	2	33	175	
Oxford, UK	3T	Siemens Tim Trio	12	22	20	12M/10F	20-35	2	34	175	
Queensland, AU	3T	Bruker		1	19	17	11M/8F	20-34	2.1	36	190
SaintLouis, USA	3T	Siemens Tim Trio	12	31	31	14M/17F	21-29	2.5	32	127	

Table 2: Sites selected from the 1000 Functional Connectome Project.

129 detection of group differences, such as the amplitude of the group difference,
 130 sample size, and the balancing of groups across sites. We evaluated the ability
 131 of detecting group differences both in terms of sensitivity for a GLM and in
 132 terms of accuracy for a SVM model.

133 **2. Method**

134 *2.1. Imaging sample characteristics*

135 The full 1000 FCP sample includes 1082 subjects, with images acquired over
 136 33 sites spread across North America, Europe, Australia and China. As the
 137 1000 FCP is a retrospective study, no effort was made to harmonize popula-
 138 tion characteristics or imaging acquisition parameters (Biswal et al., 2010). A
 139 subset of sites was selected based on the following criteria: (1) 3T scanner field
 140 strength, (2) full brain coverage for the rs-fMRI scan, and, (3) a minimum of 15
 141 young or middle aged adult participants, with a mixture of males and females
 142 (4) samples drawn from a population with a predominant Caucasian ethnicity.
 143 In addition, only young and middle aged participants (18-46 years old) were
 144 included in the study, and we further excluded subjects with excessive motion
 145 (see next Section). The final sample for our study thus included 345 cognitively
 146 normal young adults (150 males, age range: 18-46 years, mean±std: 23.8 ± 5.14)
 147 with images acquired across 8 sites located in Germany, the United Kingdom,
 148 Australia and the United States of America. The total time of available rs-
 149 fMRI data for these subjects ranged between 6 and 7.5 min and only one run
 150 was available per subject. See Table 2 for more details on the demographics
 151 and imaging parameters at each site selected in the study. The experimental
 152 protocols for all datasets as well as data sharing in the 1000 FCP were approved
 153 by the respective ethics committees of each site. This secondary analysis of the
 154 1000 FCP sample was approved by the local ethics committee at CRIUGM,
 155 University of Montreal, QC, Canada.

156 *2.2. Computational environment*

157 All experiments were performed using the NeuroImaging Analysis Kit, NIAK¹
158 (Bellec et al., 2011) version 0.12.18, under CentOS version 6.3 with Octave²
159 version 3.8.1 and the Minc toolkit³ version 0.3.18. Analyses were executed in
160 parallel on the Mammouth supercomputer⁴, using the pipeline system for Oc-
161 tave and Matlab, PSOM (Bellec et al., 2012) version 1.0.2. The scripts used
162 for processing can be found on Github⁵. Prediction was performed using the
163 LibSVM library (Chang and Lin, 2011). Visualization was implemented using
164 Python 2.7.9 from the Anaconda 2.2.0⁶ distribution, along with Matplotlib⁷
165 (Hunter, 2007), Seaborn⁸ and Nilearn⁹ for brain map visualizations.

166 *2.3. Preprocessing*

167 Each fMRI dataset was corrected for slice timing; a rigid-body motion was
168 then estimated for each time frame, both within and between runs, as well as
169 between one fMRI run and the T1 scan for each subject (Collins et al., 1994).
170 The T1 scan was itself non-linearly co-registered to the Montreal Neurological
171 Institute (MNI) ICBM152 stereotaxic symmetric template (Fonov et al., 2011),
172 using the CIVET pipeline (Ad-Dab'bagh et al., 2006). The rigid-body, fMRI-
173 to-T1 and T1-to-stereotaxic transformations were all combined to re-sample the
174 fMRI in MNI space at a 3 mm isotropic resolution. To minimize artifacts due
175 to excessive motion, all time frames showing a frame displacement, as defined
176 in Power et al. (2012), greater than 0.5 mm were removed and a residual motion
177 estimated after scrubbing. A minimum of 50 unscrubbed volumes per run was
178 required for further analysis (13 subjects were rejected). The following nuisance
179 covariates were regressed out from fMRI time series: slow time drifts (basis of
180 discrete cosines with a 0.01 Hz highpass cut-off), average signals in conservative
181 masks of the white matter and the lateral ventricles (average Pearson correlation
182 across all subjects is 0.242 between gray matter and white matter signals, and
183 0.031 between gray matter and ventricles signals) as well as the first principal
184 components (accounting for 95% variance) of the six rigid-body motion parame-
185 ters and their squares (Giove et al., 2009; Lund et al., 2006). The fMRI volumes
186 were finally spatially smoothed with a 6 mm isotropic Gaussian blurring kernel.
187 A more detailed description of the pipeline can be found on the NIAK website¹⁰
188 and Github¹¹.

¹<http://simexp.github.io/niak/>

²<http://gnu.octave.org/>

³<http://www.bic.mni.mcgill.ca/ServicesSoftware/ServicesSoftwareMincToolKit>

⁴<http://www.calculquebec.ca/index.php/en/resources/compute-servers/>

^{mammouth-serie-ii}

⁵<https://github.com/SIMEXP/Projects/tree/master/multisite>

⁶<http://docs.continuum.io/anaconda/index>

⁷<http://matplotlib.org/>

⁸<http://stanford.edu/~mwaskom/software/seaborn/index.html>

⁹<http://nilearn.github.io/>

¹⁰http://niak.simexp-lab.org/pipe_preprocessing.html

¹¹<https://github.com/SIMEXP>

189 *2.4. Inter-site bias in resting-state connectivity*

190 *Functional connectomes.* We compared the functional connectivity measures
191 derived from different sites of the 1000 FCP. A functional brain parcellation
192 with 100 regions was first generated using a bootstrap analysis of stable clusters
193 (Bellec et al., 2010b), on the Cambridge cohort of the 1000 FCP ($N = 195$), as
194 described in Orban et al. (2015). For a given pair of regions, the connectivity
195 measure was defined by the Fisher transformation of the Pearson’s correlation
196 coefficient between the average temporal rs-fMRI fluctuations of the two regions.
197 For each subject, a 100×100 functional connectome matrix was thus generated,
198 featuring the connections for every possible pair of brain regions.

199 *Inter-site effects.* The inter-site effects at a particular connection were defined
200 as the absolute difference in average connectivity between two sites. In order to
201 formally test the significance of the inter-site effects, we used a GLM including
202 age, sex and residual motion as covariates (corrected to have a zero mean across
203 subjects), as well as dummy variables coding for the average connectivity at
204 each site. For each site, a “contrast” vector was coded to measure the difference
205 in average connectivity between this site and the grand average of functional
206 connectivity combining all other sites. A p -value was generated for each connec-
207 tion to quantify the probability that the observed effect using this contrast was
208 significantly different from zero (Worsley and Friston, 1995). The number of
209 false discovery was also controlled ($q = 0.05$) using a Benjamini-Hochberg false
210 discovery rate (FDR) procedure (Benjamini and Hochberg, 1995). To quantify
211 the severity of inter-site effects, we derived Cohen’s d effect size measure for each
212 connection: $|\beta_c|/\hat{\sigma}$, with β_c being the weight associated with the contrast. The
213 standard deviation from the noise $\hat{\sigma}$ was calculated as $\hat{\sigma} = \sqrt{\sum e^2/(N - K)}$,
214 e being the residuals from the GLM, N the sample size and K the number of
215 covariates in the model. As secondary analyses, t -tests were also implemented
216 in the GLM to validate that age, sex as well as residual motion made significant
217 contributions to the model.

218 *2.5. Simulations*

219 *Data generation process.* We implemented Monte-Carlo simulations to assess
220 the detection sensitivity of group differences in rs-fMRI connectivity. The sim-
221 ulations were based on the 1000 FCP sample, with 8 sites totaling 345 subjects.
222 The multisite simulations were sampled from 148 subjects, available across $S = 7$
223 sites. The monosite simulations were sampled from 195 subjects available at
224 $S = 1$ site (Cambridge). For each simulation, a subset of subjects of a given
225 size N was selected randomly and stratified by site. For each site, a ratio W
226 of the selected subjects was randomly assigned to a so-called “patient” group.
227 We focus our analysis on connections showing a fair-to-good test-retest reli-
228 ability based on a previous study reporting 11 connections likely impacted by
229 Alzheimer’s disease, see Orban et al. (2015) for details. For each connection, a
230 “pathology” effect was added to the connectivity measures of the subjects be-
231 longing to the “patient” group. This additive shift in connectivity for “patients”
232 was selected as to achieve a specified effect size, defined below.

Effect size (Cohen's d). The Cohen's d was used to quantify the effect size. For a group comparison, Cohen's d is defined as the difference μ between the means of the two groups, divided by the standard deviation of the measures within each group, here assumed to be equal. For a given connection between brain regions i and j , let $y_{i,j}$ be the functional connectivity measure for a particular subject of the 1000 FCP sample. If the subject was assigned to the "patient" group in a particular simulation, an effect was added to generate a simulated connectivity measure $y_{i,j}^*$ equal to $y_{i,j} + \mu$. For a specified effect size d , the parameter μ was set to $d \times s_{i,j}$, where $s_{i,j}$ is the standard deviation of connectivity between region i and j . The parameter $s_{i,j}$ was estimated as the standard deviation of connectivity measures across subjects in the mono-site sample (Cambridge), without any "pathological" effect simulated.

GLM tests. In order to detect changes between the simulated groups at each connection, a GLM was estimated from the simulated data, using age, sex and frame displacement as confounds (corrected to have a zero mean across subjects). To account for site-specific effects, $S - 1$ dummy variables (binary vectors coding for each site) were added to the model, with S being the total number of sites used in the study, in addition to an intercept accounting for the global average. Finally, one dummy variable coded for the "patient" group. The regression coefficients of the linear model were estimated with ordinary least squares, and a t -test, with associated p -value, was calculated for the coefficient of the "patient" variable. A significant pathology effect was detected if the p value was smaller than a prescribed α level. The α level needs to be adjusted for multiple comparisons (in our case 11 connections, but this would depend on the number of connections selected in a particular study), which can be done in an adaptive manner using FDR. When connections are pre-specified, such as in e.g. Wang et al. (2012), a more liberal threshold can be applied. In our case, since we wanted to have a constant behavior independent of the effect size, we tested different typical values for α in $\{0.001, 0.01, 0.05\}$. For each simulation sample b and each connection, we derived a p -value $p^{(b)}$, and the effect was deemed detected if $p^{(b)}$ was less than α . The sensitivity of the test for a particular connection was evaluated by the frequency of positive detections over all simulation samples.

Prediction accuracy. In addition to mass univariate GLM tests, we also investigated a linear SVM (Cortes and Vapnik, 1995) using a Monte-Carlo simulation of the prediction of clinical labels based on cross-validation. For SVM simulations, all possible connections between the 100 brain regions were used simultaneously to predict the presence of the simulated pathology in a given subject. For a participant assigned to the "patient" group, a "pathology" effect was only simulated in a set percentage of connections, which were randomly selected. The proportion of connections with a non-null effect was denoted as π_1 . For a given simulation at sample size N , the SVM model was trained on N subjects selected randomly and stratified by site. The accuracy of the model was evaluated on a separate sample consisting of the remaining subjects, unused during train-

277 ing. For example, for a multisite simulation with $N = 80$ subjects for training,
278 the model accuracy was tested on the remaining 68 subjects: 148 (available
279 subjects) minus 80 (subjects in the training set). During training, a 10-fold
280 cross-validation was used to optimize the hyper-parameters of the SVM inde-
281 pendently for each simulation. The mean and standard deviation of accuracy
282 scores across all samples were derived for each simulation scenario.

283 *Simulation experiments.* All the simulation parameters have been summarized
284 below:

- 285 • Sample size N .
286 • Patient allocation ratio W .
287 • Number of sites S .
288 • The type of detection method, either GLM or SVM.
289 • For GLM tests, the false-positive rate α .
290 • For SVM tests, the proportion of “pathological” connections π_1 .
291 • The effect size d .

292 For a given set of simulation parameters, we generated $B = 10^3$ Monte-Carlo
293 samples to estimate either the sensitivity (for GLM test) or the accuracy (for
294 SVM prediction) of the method. For all experiments, we investigated effect sizes
295 $d \in \{0, 2\}$ with a step of 0.01 and $\alpha \in \{0.001, 0.01, 0.05\}$. The number of site(s)
296 was $S = 1$ for the monosite analysis and $S = 7$ for the multisite analysis. We
297 implemented the following experiments:

- 298 • (\mathcal{E}_1) Test the impact of the sample size on GLM $N \in \{40, 80, 120\}$, with
299 a fixed allocation ratio $W = 0.5$.
300 • (\mathcal{E}_2) Test the impact of the allocation ratio on GLM $W \in \{0.5, 0.3, 0.15\}$
301 for a fixed sample size $N = 120$.
302 • (\mathcal{E}_3) Test the impact of multisite correction (regressing out the site effects
303 using dummy variables coding for each site) and affected connection vol-
304 ume (π_1) on the prediction accuracy. For the prediction scenario, we used
305 a range of $\pi_1 \in \{0.1, 1, 5\%\}$, and two sample sizes $N \in \{80, 120\}$ subjects
306 for training, with model accuracy estimated on $N = 68$ and $N = 28$,
307 respectively.

308 **3. Results**

309 **3.1. Inter-site effects in fMRI connectivity**

310 *Site effects in the default-mode network.* We first focused on the connections as-
311 sociated with a seed region located in the posterior cingulate cortex, a key node

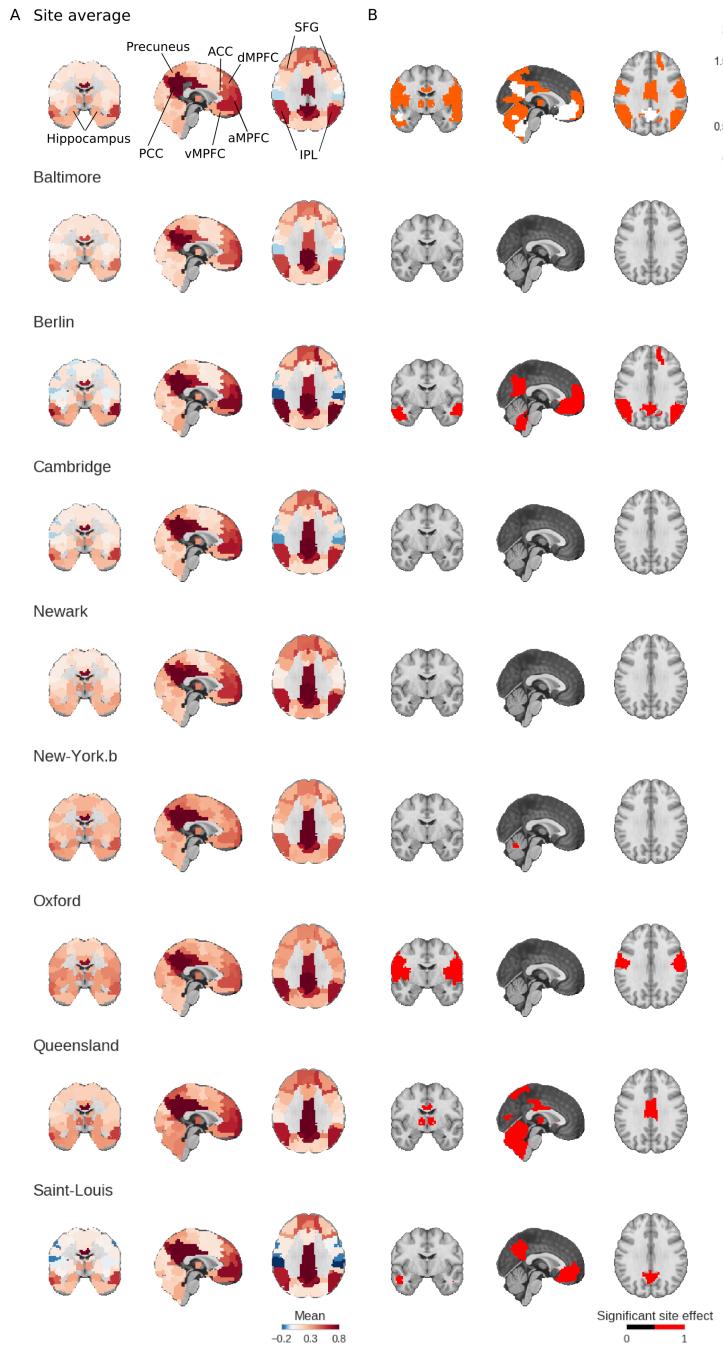


Figure 1: Panel A: map of the DMN obtained using a seed in the posterior cingulate cortex, averaging all subjects and sites together (first row) and then averaging all subjects for each of the 8 sites (subsequent rows). Panel B shows the number of sites with a significant inter-site difference for each brain region (first row) and the significant differences between the average functional connectivity maps of one site versus all the others (subsequent rows).

of the default-mode network (DMN), which is one of the most widely studied resting-state networks (Greicius et al., 2004). The connections were based on the Cambridge 100 parcellation, and were represented as a connectivity map, (Figure 1). Figure 1A shows the posterior cingulate cortex connectivity map, averaged across all subjects and all sites. The key regions of the DMN are easily identifiable, and include the posterior cingulate cortex, precuneus, inferior parietal lobule, anterior cingulate cortex, medial pre-frontal cortex (dorsal, anterior and ventral), superior frontal gyri and the medial temporal lobe (Damoiseaux et al., 2006; Dansereau et al., 2014; Yan et al., 2013). The average connectivity map of the DMN was then extracted for each site, Figure 1A. Qualitatively, the DMN maps were consistent across sites, as expected based on the literature. We then tested for the significance of the site effects (Figure 1B), i.e. the difference in average connectivity at a given site and the average connectivity at all remaining sites. The statistical maps were corrected for multiple comparisons across the brain with FDR at $q \leq 0.05$ (Benjamini and Hochberg, 1995). A significant site effect for at least one connection could be identified for every site, without exception, Figure 1B. Figure 1C shows how reproducible the significant site effects were in connectivity across the brain and sites. The identified significant connections were quite variable across sites, most of them being identified at less than three sites.

Site effects across the connectome. In order to extend these observations outside of the DMN, we derived the entire connectome using the Cambridge 100 parcellation. Figure 2A shows the average connectome, pooling all subjects and sites together. The regions have been re-ordered based on a hierarchical clustering with Ward criterion. A network structure is clearly visible as squares of high connectivity on the diagonal of the connectome (as outlined by black lines). Each diagonal square corresponds to the intra-network connectivity for a partition into 7 networks (Figure 2A). These 7 networks¹² were consistent with the major resting-state networks reported using a cluster analysis in previous works (e.g. van den Heuvel et al., 2008; Bellec et al., 2010a; Yeo et al., 2011; Power et al., 2011): the DMN, visual, sensorimotor, dorsal and ventral attentional networks, mesolimbic and cerebellar networks were identified (Figure 2B). Figure 2C shows how this large-scale connectome organization varied from site to site. The average connectivity per site as well as significant differences with the average of the remaining sites ($q \leq 0.05$) is shown in Figure 2C. Visually, consistent with our previous observations in the DMN, the organization of the average connectome into large-scale resting-state networks was preserved across all sites.

Some significant site effects were still detected in the connectivity both within each network, as well as between networks. By counting the number of sites showing a significant effect for each pair of regions, it was apparent that significant site effects were quite variable in their localization and spread across

¹²<http://neurovault.org/images/39184/>

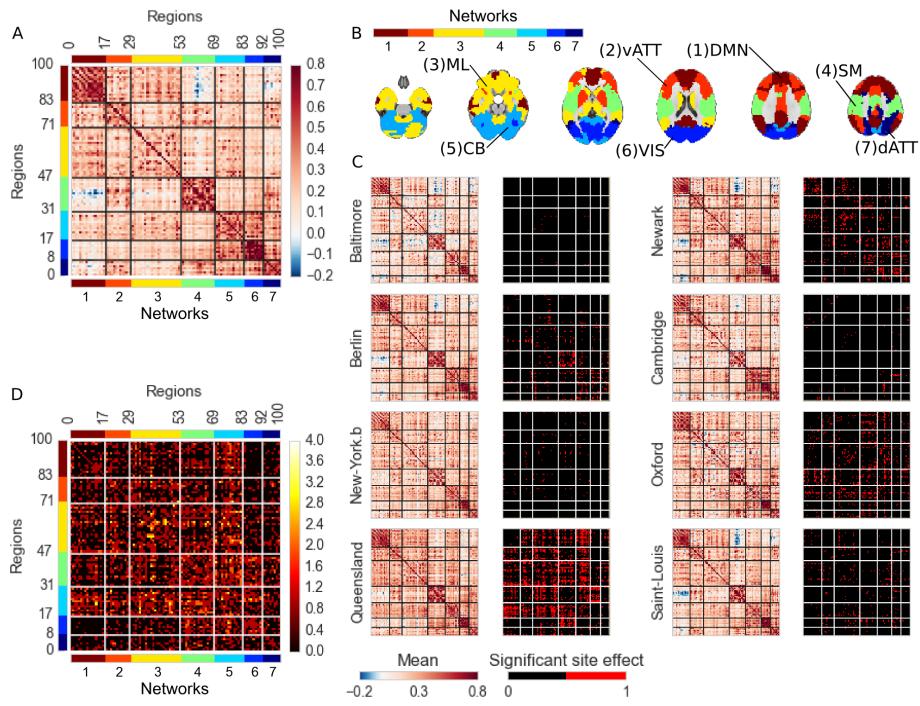


Figure 2: Panel A shows the average functional connectomes for 8 sites of the 1000 FCP. Colors next to the x and y axis correspond to different networks in a 7-cluster solution of the matrix, obtained from a hierarchical clustering (Ward criterion). Panel B presents the corresponding 7 brain networks, along with labels. Panel C shows average connectomes for individual sites, as well as connections with a significant site effect. Panel D shows the number of sites at which a given connection was detected as significant. ML: mesolimbic, CB: cerebellar, VIS: visual, vATT: ventral attentional, dATT: dorsal attentional, DMN: default mode network, SM: sensorimotor.

354 the full connectome (Figure 2D). Concerning the association with the other
 355 confounding variables in the model (sex, age and motion) many connections
 356 were found to be significantly associated with motion, see Supplementary Ma-
 357 terial Figure S5, although very few connections were found to be significantly
 358 associated with the sex and age, see Supplementary Material Figure S6 and
 359 S7. We also checked that the analysis was not predominantly driven by the
 360 larger Cambridge site. We thus ran the same analysis excluding that site (see
 361 Supplementary Material Figure S8). The number of significant pairs remained
 362 very similar, although the spatial location of half of the significant connectivity
 363 pairs changed when the large Cambridge site was removed from the analysis.
 364 Those findings do not qualitatively change our conclusion, but they influence
 365 the location of the significant connections. These differences may be due to the
 366 intrinsic variability in the statistical test, and not just the size of the Cambridge
 367 site. In summary, those findings support the inclusion of age, sex and motion
 368 parameters in a GLM in order to remove their confounding effects in addition
 369 to site effects.

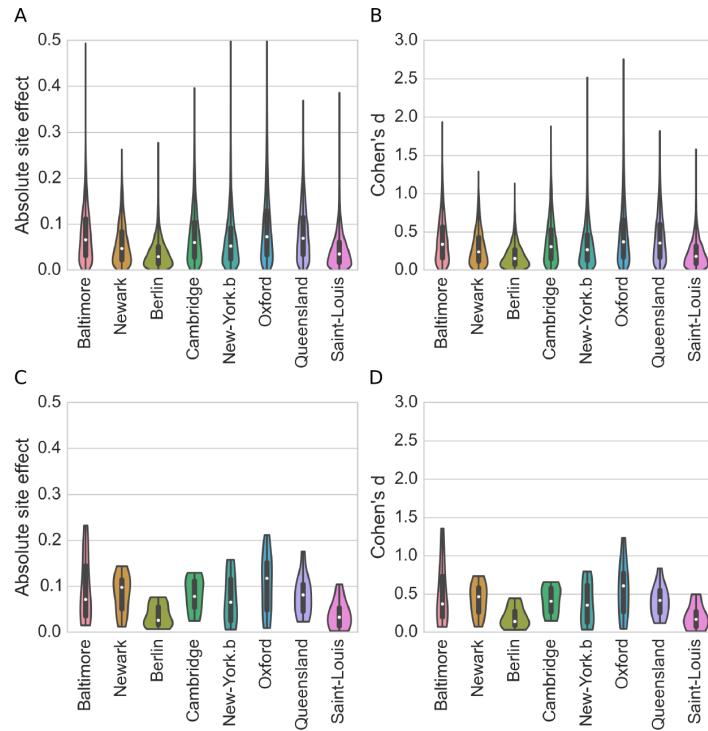


Figure 3: Effect size of the inter-site effects from a subset of 8 sites from the 1000 FCP. Panels A,C show the distribution of absolute differences in functional connectivity, while panels B,D show Cohen's d measures of inter-site effects. Panels A,B show violin plots across every connections in the BASC Cambridge 100 parcellation, while Panels C,D focus on the selected 11 functional connections used in simulations, only.

370 *Site effects vs. within-site variations across subjects.* We measured the amplitude
371 of inter-site effects, represented as violin plots across connections using
372 either the absolute difference in average connectivity (Figure 3A,C) or Cohen’s
373 d effect size measures (Figure 3B,D). The violin plots include either every con-
374 nection from the BASC Cambridge parcellation (Figure 3A,B), or only the 11
375 connections selected for Monte-Carlo simulations (Figure 3C,D). For absolute
376 differences, the distributions were mostly consistent across sites, with a me-
377 dian around 0.06, 5% percentile near 0 and 95% percentiles in the 0.08- 0.1
378 range. For Cohen’s d , the distributions were also consistent across sites, with
379 a median around 0.33, 5% percentile near 0 and 95% percentiles in the 0.4-
380 0.6 range. These effect sizes are typically deemed small-to-moderate (Cohen,
381 1992), although such a qualitative assessment needs to be refined based on each
382 application. This result thus suggests that the impact of additive inter-site ef-
383 fects on statistical tests will be limited. Similar findings were observed across
384 all possible connections, or across the 11 pairs of connections selected in the
385 simulation study.

386 *Differences in standard deviation across sites.* We also investigated the site dif-
387 ferences in standard deviation of connectivity across subjects, see Supplemen-
388 tary Figure S1 for the DMN, Supplementary Material S2 for the connectomes.
389 The standard group GLM assumes equal variance of resting-state connectiv-
390 ity across all subjects, or “homoscedasticity”. Significant differences in across-
391 subject standard deviation between sites violates the homoscedastic assump-
392 tion, and may jeopardize the validity of the false-positive rates of the model.
393 Qualitatively, we first observed that the sites showing the larger number of dif-
394 ferences were the one with the most temporal variance among connections see
395 Supplementary Figure S3. We then ran a White’s test aimed at rejecting ho-
396 moscedasticity at each connection, independently. The White’s tests resulted in
397 a family of p-values, which was corrected for multiple comparisons using FDR
398 ($q < 0.05$). The homoscedastic hypothesis was rejected in a large portion of
399 connections. This was expected due to the large overall number of subjects and
400 consequently large statistical power of White’s procedure. However, despite
401 reaching significance, the absolute difference in the average standard deviation
402 between two sites was 19% of the grand average standard deviation, on average
403 across pairs of sites. Such a small departure from homoscedasticity likely has
404 only a mild impact on the GLM, which we formally investigated using Monte-
405 Carlo simulations.

406 *3.2. Multisite Monte-Carlo simulations*

407 *Validity of the control of false positives in the GLM.* An excellent control of
408 the false positive rate was observed at all nominal levels $\alpha \in \{0.001, 0.01, 0.05\}$,
409 both in monosite simulations or in multisite simulations, when site covariates
410 were included in the GLM, see Figure 4. This means that the nominal, user-
411 specified, false positive rate matched precisely with the effective false positive
412 rate measured in the simulations. This observation held for any combination
413 of allocation ratio, $W \in \{15\%, 30\%, 50\%\}$, and sample size, $N \in \{40, 80, 120\}$.

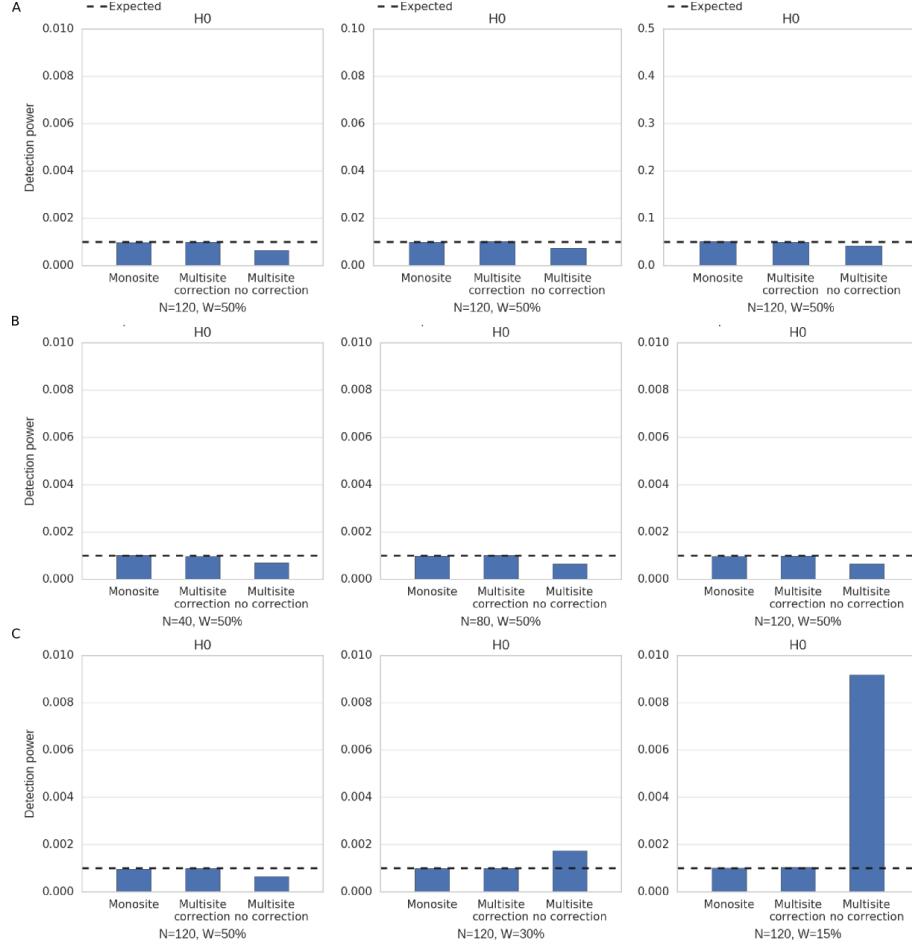


Figure 4: Monte-Carlo simulation of the false positive rate in the absence of group differences ($d = 0$), either for a monosite ($S = 1$, left), a multisite ($S = 7$) with (middle) or without (right) site covariates included in the GLM. In panel A, three different α values have been tested, $\alpha \in \{0.001, 0.01, 0.05\}$ with a fixed sample size and patient allocation ratio ($N = 120, W = 50\%$). In panel B, three different sample sizes have been tested, $N \in \{40, 80, 120\}$ with a fixed patient allocation ratio ($W = 50\%$) (Experiment (\mathcal{E}_1)). In panel C, three different patient allocation ratios have been tested, $W \in \{50\%, 30\%, 15\%\}$ with a fixed sample size ($N = 120$) (Experiment (\mathcal{E}_2))).

414 By contrast, when no site covariates were included in the GLM, the false positive
 415 rate was not controlled appropriately, sometimes by a wide margin. In
 416 the absence of site covariates, the procedure was sometimes too conservative,
 417 e.g. $W = 50\%$, and sometimes very liberal, e.g. $N = 120, W = 15\%$. This
 418 experiment showed that, despite the mild departure from homoscedasticity re-
 419 ported above, the GLM does control for false-positive rate at each connection
 420 very precisely, if and only if site covariates are included in the model.

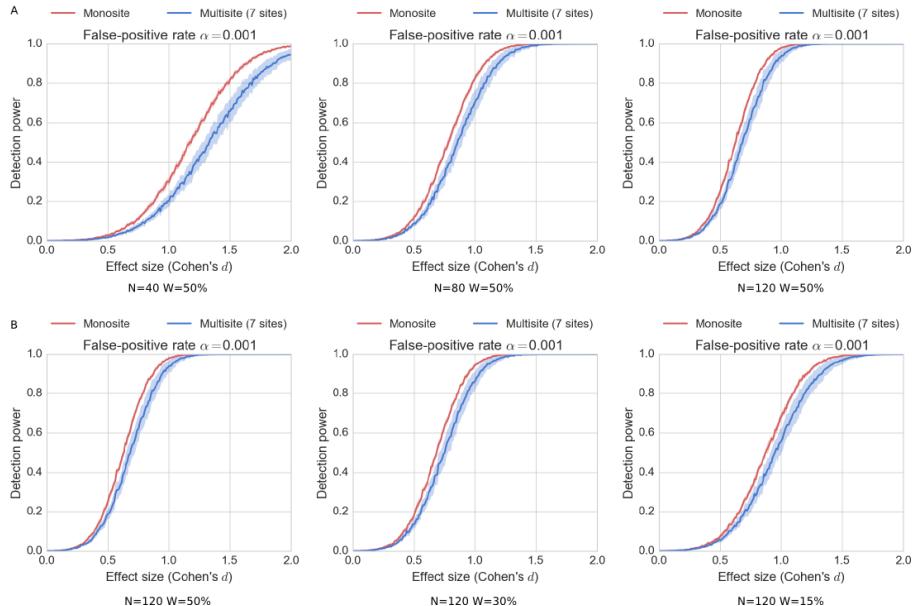


Figure 5: Monte-Carlo simulation of detection power as a function of the effect size $d \in [0, 2]$, either for a monosite ($S = 1$, in red) or a multisite ($S = 7$, in blue) sample, when testing differences between two groups with a GLM and a false-positive rate $\alpha = 0.001$. The plain curves are the average statistical power across 11 connections, and the shaded area represents ± 1 standard deviation across connections. In panel A, the patient allocation ratio is fixed ($W = 50\%$) and three different sample sizes have been tested, $N \in \{40, 80, 120\}$ (Experiment (\mathcal{E}_1)). In panel B, the sample size is fixed ($N = 120$) and three different patient allocation ratios have been tested $W \in \{15\%, 30\%, 50\%\}$ (Experiment (\mathcal{E}_2)).

421 *Statistical power and effect size.* Figure 5A shows the relationship between effect
 422 size and a GLM detection power in experiment (\mathcal{E}_1) , i.e. for a fixed allocation ra-
 423 dio ($W = 50\%$) and three different sample sizes, $N \in \{40, 80, 120\}$. The average
 424 and std of detection power was plotted across the 11 selected connections. The
 425 variations of statistical power across connections were very small for monosite
 426 simulations, as the effect size was adjusted based on the standard deviation
 427 of each connection within that sample. As expected, the sensitivity increased
 428 with sample size, quite markedly. In multisite simulations ($S = 7$), for a large
 429 effect size ($d = 1$), the detection power was 20% with 40 subjects , 80% with

430 80 subjects and 95% with 120 subjects. The sensitivity was larger with a single
 431 site than a multisite sample, yet the difference between the two decreased as
 432 sample size increased. With $N = 40$ and $d = 1$, the detection power was close
 433 to 30% for a single site sample, compared to 20% for the multisite sample. With
 434 $N = 120$ and $d = 1$, the difference in sensitivity was only of a few percent. The
 435 same trend was apparent for all tested effect sizes as well as for $\alpha \in \{0.01, 0.05\}$
 436 (not shown).

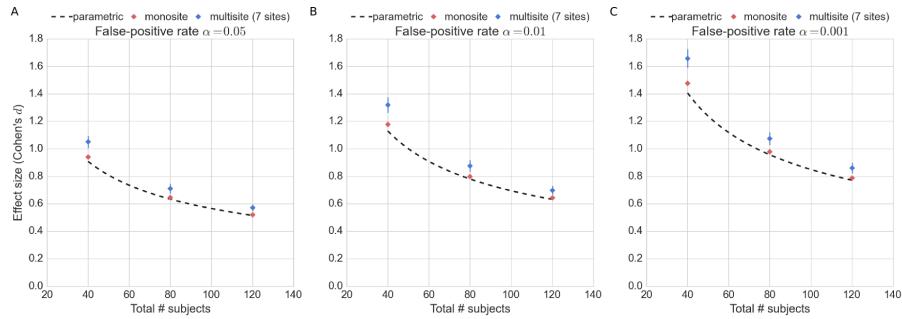


Figure 6: Effect size detectable at 80% sensitivity as a function of sample size, for different false-positive rate $\alpha \in \{0.05, 0.01, 0.001\}$ (experiment (\mathcal{E}_1)). All simulations used a balanced patient allocation ratio $W = 50\%$. The monosite performance is shown in red and the multisite in blue. The dotted black line shows the detectable effect size for a classical parametric t -test.

437 *Statistical power and group allocation ratio.* Figure 5B shows the relationship
 438 between effect size and a GLM detection power in experiment (\mathcal{E}_2), i.e. for a
 439 fixed sample size ($N = 120$) and three different patient allocation ratio, $W \in$
 440 $\{15\%, 30\%, 50\%\}$. Overall, we found that the detection power increased with
 441 W . For example, with $d = 1$, the detection power was 65% for $W = 15\%$, and
 442 increased to 90% with $W = 30\%$, and finally 95% for $W = 50\%$. The impact
 443 of W was observed in both monosite and multisite samples, with an optimal
 444 allocation ratio of $W = 50\%$ for both. This observation was also made for
 445 $\alpha \in \{0.01, 0.05\}$ (not shown).

446 *Detectable effect size, as a function of sample size.* An alternative summary of
 447 experiment (\mathcal{E}_1) is to represent the effect size that can be detected with 80%
 448 sensitivity, as a function of sample size for monosite and multisite configurations,
 449 see Figure 6. As a reference, we computed the same curve for parametric t -test
 450 comparisons, under assumptions of normality. As expected, the detectable effect
 451 size for parametric t -tests closely followed the monosite estimation. For a small
 452 sample size ($N = 40$), the detectable effect size was notably larger in multisite
 453 configurations than in a monosite configuration (difference of about 0.25 in
 454 Cohen's d for $\alpha = 0.001$). However, the difference decreased for large sample
 455 sizes to become smaller than 0.1 with $N = 120$ and $\alpha = 0.001$. The lowest
 456 detectable effect size for a sensitivity of 80% at $\alpha = 0.05$ was about $d = 0.8$,
 457 achieved in a monosite configuration with $N = 120$. At this sample size, the

458 difference between single and multisite configurations was marginal, with only
 459 a few percent's of difference in detectable effect sizes.

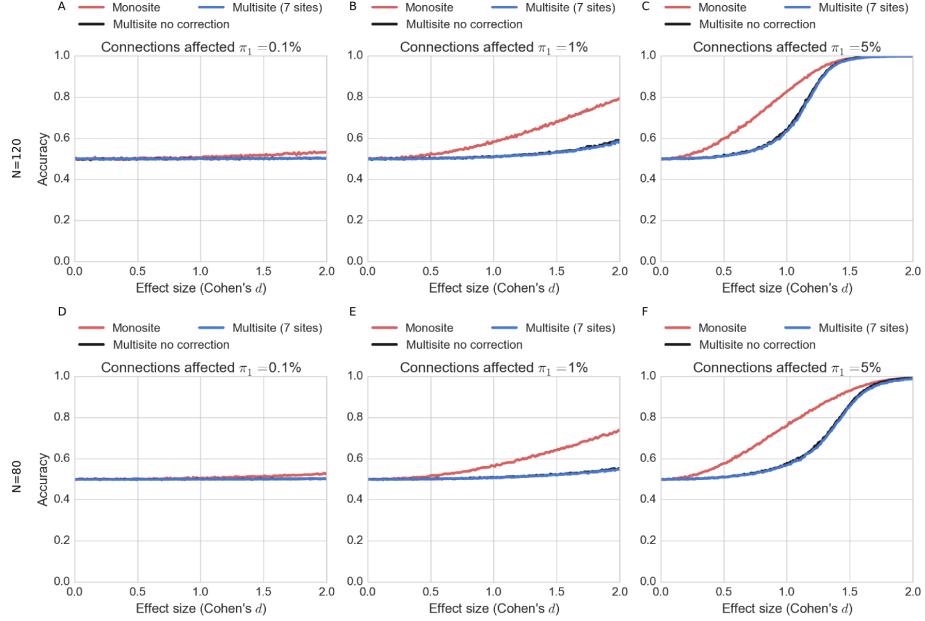


Figure 7: Prediction accuracy of patient vs. controls as a function of effect size. Three simulation settings are presented on each plot: monosite (red curve), multisite with regression of site effects ($S = 7$, blue curve), and multisite without regression of site effects ($S = 7$, black curve). Accuracy was estimated over $B = 10^3$ simulation samples with a patient allocation ratio $W = 50\%$ and 3 volumes of affected connections $\pi_1 = 0.1\%$ (left column), $\pi_1 = 1\%$ (middle column) and $\pi_1 = 5\%$ (right column). Two sample sizes were tested: $N = 120$ randomly selected subjects for training, with the remaining $N = 28$ to estimate accuracy (first row), and $N = 80$ randomly selected subjects for training, with the remaining $N = 68$ to estimate accuracy (second row).

460 *Prediction accuracy.* In experiment (\mathcal{E}_3), we examined the impact of effect size
 461 and the volume of affected connections on prediction accuracy in a SVM, see
 462 Figure 7. The volume of changes π_1 had a major impact on prediction accu-
 463 racy. At $\pi_1 = 0.1\%$ (around 5 connections) the accuracy level was at chance
 464 level across all tested effect sizes, (Figure 7A). With $\pi_1 = 1\%$, accuracy slightly
 465 increased, but effect sizes larger than $d = 2$ were still required to reach over 80%
 466 accuracy (Figure 7B). With $\pi_1 = 5\%$, 95% accuracy was achieved at the same
 467 effect size (about $d = 1.5$) for monosite and multisite simulations, although the
 468 accuracy in multisite simulations was notably lower than for monosite simu-
 469 lations across most effect sizes (Figure 7C). The relationship between effect size
 470 and accuracy followed a sigmoidal curve in both settings, yet a sharper, and
 471 later transition between very low and very high accuracy was observed in mul-
 472 tisite simulations. Interestingly, correcting for site effects by regressing out the

473 dummy variable before running the SVM classifier had no impact on accuracy
474 levels. The sample size ($N = 80$ vs $N = 120$ for training) did have a moderate
475 effect on prediction accuracy: for $\pi_1 = 5\%$ and $d = 1$ and monosite simulations,
476 accuracy was about 85% with $N = 120$ (Figure 7C) and 75% with $N = 80$
477 (Figure 7F).

478 **4. Discussion and conclusions**

479 *Inter-site effects in rs-fMRI connectivity.* Typical resting-state networks, such
480 as the DMN, the attentional, visual and sensorimotor networks, were reliably
481 found across sites. This was strongly expected given the relative consistency of
482 their distribution across individuals, studies, preprocessing approaches or even
483 methods used to extract networks (e.g. Damoiseaux et al., 2006; van den Heuvel
484 et al., 2008; Bellec et al., 2010b; Yeo et al., 2011; Power et al., 2011). We how-
485 ever found that significant differences in average connectivity existed between
486 sites, as previously reported by Yan et al. (2013). These site effects in con-
487 nectivity may undermine the generalization of the results derived at a single
488 site. The inter-subject (intra-site) standard deviation of the connections was
489 found to be more than twice as large as the inter-site absolute effect, on aver-
490 age across brain connections. This effect size measured in Cohen's d would be
491 deemed small-to-moderate, which suggests that the impact of additive inter-site
492 effects on statistical tests will be limited. This is a reassuring finding support-
493 ing the feasibility of statistical tests pooling fMRI data across multiple sites.
494 Previous studies (Sutton et al., 2008; Brown et al., 2011) had reported inter-site
495 variance up to 10 times smaller than inter-subject variability, but these studies
496 had much more homogeneous scanning environments than ours and also used
497 different fMRI outcome measures. In our case, we still investigated only 3T
498 scanners, mostly Siemens, and inter-site effects may be larger when considering
499 other manufacturers or field strengths.
500

501 *Statistical power and multisite rs-fMRI.* After accounting for site-related addi-
502 tive effects in a GLM, the multisite simulation pooling 7 sites together showed
503 detection power close to that of a monosite simulation with equivalent sample
504 size. The difference was noticeable for small sample size (total $N = 40$), and
505 became very small for a sample size $N = 120$. Another observation was that,
506 for a given detection power, the lowest effect size that we were able to detect
507 was more variable across connections for a low sample size. We demonstrated
508 that a parametric group GLM does control precisely for the rate of false positive
509 discoveries, even in multisite settings, as long as site covariates are included in
510 the model. Taken together, these observations suggest to use sample sizes larger
511 than 100 subjects for GLM multisite studies. This conclusion may depend on
512 the number of sites pooled in the study and the actual number of subjects in
513 each of those sites, which we could not test in this work due to the size of the
514 available sample.

515 *Modeling site effects as random variables.* We modeled the effect of each site
516 on the average connectivity between any given pair of regions as a fixed effect.
517 This means that the proposed GLM inference does apply only to collection of
518 sites included in a given analysis. The linear mixed-effects model (Chen et al.,
519 2013) would allow more powerful inferences: by modeling site effects as random
520 variables, following a specific distribution (e.g. Gaussian), we would be able
521 to generalize observations potentially to any collection of sites, provided our
522 assumptions are accurate. The sample of sites available for this study (7 at
523 most) is however too small in our view to correctly estimate the variability of
524 effects across sites. This work would also require to formulate and investigate
525 empirically as well as on simulations different models for the distribution of
526 inter-site variations of site effects (e.g. Gaussian distribution).

527 *Site heteroscedasticity.* We observed mild heteroscedasticity across sites. Our
528 simulations showed that this does not compromise the control of false positive
529 rate in the GLM, even under homoscedastic assumptions, with the range of
530 contrasts we investigated. Regression models more robust to heteroscedasticity
531 may be investigated in the future, e.g. weighted least squares regression or linear
532 mixed-effects modeling (Chen et al., 2013).

533 *Statistical power and sample size.* For a medium effect size, e.g. $d = 0.5$, the
534 sensitivity was low (below 20%), even for monosite simulations with $N = 120$
535 subjects. This sobering result supports the current trend in the literature to
536 pool multiple data samples to increase sample size, at the cost of decreased
537 homogeneity. We also found that resting-state studies based on 40 subjects or
538 less, even at a single site, are seriously underpowered, except for extremely large
539 effect sizes (Cohen's d greater than 1.5). Finally, unbalanced patient allocation
540 ratio in site samples greatly reduces sensitivity, even in monosite studies. Bal-
541 anced datasets, i.e. with equal numbers of patients and controls at each site,
542 should therefore be favored.

543 *Prediction.* Comparing the monosite and the multisite accuracy curves reveals
544 a substantial drop in accuracy from monosite to multisite across a broad range
545 of effect sizes. However, it should be noted that classifiers trained across multi-
546 ple data sources will likely generalize better to new observations, which is likely
547 a critical feature in most applications and reflects the true potential clinical
548 utility of this type of technique. Our conclusions are consistent with the work
549 of Nielsen et al. (2013), which compares the prediction of a clinical diagnosis
550 of autism in monosite vs. multisite settings. The authors concluded that the
551 prediction accuracy for the multisite sample was significantly smaller than for
552 the monosite sample. A somewhat surprising observation in our analysis was
553 that linear correction for site-specific effects did not improve accuracy of pre-
554 diction using SVM. The SVM model seems to learn features that are invariant
555 across sites, maybe focusing on connections with the smallest site effect, or
556 looking at differences between connections similarly impacted by a site effect.
557 Finally, an important conclusion of our simulations was that the volume of brain

558 connections affected by a disease impacts accuracy as much as the effect size
559 per connection. This suggests that feature reduction and/or selection is a very
560 important step to improve sensitivity to small effect sizes.

561 *Beyond additive site effect.* An important limitation to our study is that we only
562 investigated the impact of additive effects in brain connectivity across sites.
563 Areas of future work include interactions between site effects and pathology,
564 possibly in the form of polynomial and non-linear interactions. We hope that,
565 in the future, fMRI data acquired on clinical cohorts at tens of sites will become
566 available, which will enable researchers to test empirically the presence of such
567 interaction effects.

568 *Other types of multisite data.* Another limitation of our study is that we only
569 investigated multisite data featuring roughly equal sample sizes with fairly bal-
570 anced patient allocation ratios at each site. Multisite studies including a very
571 large number of sites with sometimes only a few subjects per site are however
572 quite common, e.g. the Alzheimer’s disease neuroimaging initiative (ADNI)
573 (Mueller et al., 2005) and many pharmaceutical clinical trials at phase II and
574 III ¹³. In this type of design, the multisite effect may play a much more pro-
575 nounced role than in our simulations as it cannot be modeled in the GLM, and
576 will become an intrinsic added source of inter-subject variance (Feaster et al.,
577 2011). Unfortunately, this type of design could not be tested with the cur-
578 rent dataset due to the limited number of sites available. This represents an
579 important avenue of future work.

580 *Underlying causes of the site effects.* Not all sites seemed to be equally impacted
581 by the site effects, with sites like Berlin or Saint-Louis showing a small number
582 of connections significantly different than the grand average connectivity matrix,
583 while sites like Baltimore, Queensland and Oxford showed many more connec-
584 tions affected by the site effects. Interestingly this can potentially be due to
585 temporal variance of the connections (see Supplementary Figure S3) partly ex-
586 plained by the scanner make since Queensland and Baltimore site used scanners
587 from different makers (namely Bruker and Philips) than the rest of the sites used
588 in this study (Siemens scanners). This may suggest that scanners SNR (signal
589 to noise ratio) may partly explain the variance of connectivity. These differences
590 may not be statistically significant, or they may reflect real differences due to
591 protocol, scanner characteristics at these sites or differences in sampling across
592 sites. Multiple causes may be interacting together to produce the site effects,
593 as reported by Yan et al. (2013), although some of these sources of variance
594 could be better controlled like the scanner parameters, paired with the use of a
595 phantom to promote more homogeneous configurations across sites (Friedman
596 et al., 2006; Friedman and Glover, 2006; Glover et al., 2012; Friedman et al.,
597 2008). Even in standardized experiments, it should be noted that differences in
598 scanner protocols remain (Brown et al., 2011). A much larger multisite sample

¹³<http://www.rochetrials.com/trialDetailsGet.action?studyNumber=BP28248>

599 with systematically varying parameters could enable a data-driven identification
600 of the critical parameters impacting site effects. The various releases made
601 by the INDI initiative may fill that gap in the literature in the future, as the
602 scanner protocols are much better described in recent releases, such as Corr
603 (Zuo et al., 2014), than they were in the initial FCP release. These findings
604 stress the need for more work to find the source of that variance rather than
605 ad-hoc procedures to correct for them.

606 **5. Acknowledgments**

607 Parts of this work were presented at the 2013 annual meeting of the Orga-
608 nization for Human Brain Mapping, as well as the 2013 Alzheimer’s Associa-
609 tion International Conference (AAIC) (Dansereau et al., 2013). The authors
610 are grateful to the members of the 1000 functional connectome consortium for
611 publicly releasing their datasets. The computational resources used to per-
612 form the data analysis were provided by ComputeCanada¹⁴ and CLUMEQ¹⁵,
613 which is funded in part by NSERC (MRS), FQRNT, and McGill University.
614 This project was funded by NSERC grant number RN000028 and the Canadian
615 Consortium on Neurodegeneration in Aging (CCNA), through a grant from the
616 Canadian Institute of Health Research and funding from several partners in-
617 cluding SANOFI-ADVENTIS R&D. PB is supported by a salary award from
618 “Fonds de recherche du Québec – Santé” and the Courtois Foundation.

619 **References**

- 620 Abraham, A., Milham, M., Di Martino, A., Craddock, R. C., Samaras, D.,
621 Thirion, B., Varoquaux, G., 2016. Deriving reproducible biomarkers from
622 multi-site resting-state data: An autism-based example. *NeuroImage*.
- 623 Ad-Dab’bagh, Y., Einarson, D., Lyttelton, O., Muehlboeck, J. S., Mok, K.,
624 Ivanov, O., Vincent, R. D., Lepage, C., Lerch, J., Fombonne, E., Evans,
625 A. C., 2006. The CIVET Image-Processing environment: A fully automated
626 comprehensive pipeline for anatomical neuroimaging research. In: Corbetta,
627 M. (Ed.), *Proceedings of the 12th Annual Meeting of the Human Brain Map-*
628 *ping Organization*. *Neuroimage*, Florence, Italy.
- 629 Anand, A., Li, Y., Wang, Y., Gardner, K., Lowe, M. J., 2007. Reciprocal ef-
630 fects of antidepressant treatment on activity and connectivity of the mood
631 regulating circuit: an fmri study. *The Journal of neuropsychiatry and clinical*
632 *neurosciences*.
- 633 Bellec, P., Carbonell, F., Perlberg, V., Evans, A. C., 2010a. A neuroimaging
634 analysis kit for Octave and Matlab.
635 URL <http://code.google.com/p/niak/>

¹⁴<https://compute-canada.org/>

¹⁵<http://www.clumeq.mcgill.ca/>

- 636 Bellec, P., Carbonell, F. M., Perlberg, V., Lepage, C., Lyttelton, O., Fonov,
 637 V., Janke, A., Tohka, J., Evans, A. C., 2011. A neuroimaging analysis kit for
 638 Matlab and Octave. In: Proceedings of the 17th International Conference on
 639 Functional Mapping of the Human Brain. pp. In Press+.
- 640 Bellec, P., Lavoie-Courchesne, S., Dickinson, P., Lerch, J. P., Zijdenbos, A. P.,
 641 Evans, A. C., 2012. The pipeline system for Octave and Matlab (PSOM): a
 642 lightweight scripting framework and execution engine for scientific workflows.
 643 Frontiers in neuroinformatics 6.
 644 URL <http://dx.doi.org/10.3389/fninf.2012.00007>
- 645 Bellec, P., Rosa-Neto, P., Lyttelton, O. C., Benali, H., Evans, A. C., Jul. 2010b.
 646 Multi-level bootstrap analysis of stable clusters in resting-state fMRI. Neu-
 647 roImage 51 (3), 1126–1139.
 648 URL <http://dx.doi.org/10.1016/j.neuroimage.2010.02.082>
- 649 Benjamini, Y., Hochberg, Y., 1995. Controlling the false-discovery rate: a prac-
 650 tical and powerful approach to multiple testing. J. Roy. Statist. Soc. Ser. B
 651 57, 289–300.
- 652 Biswal, B. B., Mennes, M., Zuo, X.-N. N., Gohel, S., Kelly, C., Smith, S. M.,
 653 Beckmann, C. F., Adelstein, J. S., Buckner, R. L., Colcombe, S., Dogonowski,
 654 A.-M. M., Ernst, M., Fair, D., Hampson, M., Hoptman, M. J., Hyde, J. S.,
 655 Kiviniemi, V. J., Kötter, R., Li, S.-J. J., Lin, C.-P. P., Lowe, M. J., Mackay,
 656 C., Madden, D. J., Madsen, K. H., Margulies, D. S., Mayberg, H. S., McMa-
 657 hon, K., Monk, C. S., Mostofsky, S. H., Nagel, B. J., Pekar, J. J., Peltier,
 658 S. J., Petersen, S. E., Riedl, V., Rombouts, S. A., Rypma, B., Schlaggar, B. L.,
 659 Schmidt, S., Seidler, R. D., Siegle, G. J., Sorg, C., Teng, G.-J. J., Veijola, J.,
 660 Villringer, A., Walter, M., Wang, L., Weng, X.-C. C., Whitfield-Gabrieli, S.,
 661 Williamson, P., Windischberger, C., Zang, Y.-F. F., Zhang, H.-Y. Y., Castel-
 662 lanos, F. X., Milham, M. P., Mar. 2010. Toward discovery science of human
 663 brain function. Proceedings of the National Academy of Sciences of the United
 664 States of America 107 (10), 4734–4739.
 665 URL <http://dx.doi.org/10.1073/pnas.0911855107>
- 666 Brown, G. G., Mathalon, D. H., Stern, H., Ford, J., Mueller, B., Greve, D. N.,
 667 McCarthy, G., Voyvodic, J., Glover, G., Diaz, M., et al., 2011. Multisite
 668 reliability of cognitive bold data. Neuroimage 54 (3), 2163–2175.
- 669 Caramanos, Z., Fonov, V. S., Francis, S. J., Narayanan, S., Pike, G. B., Collins,
 670 D. L., Arnold, D. L., 2010. Gradient distortions in mri: Characterizing
 671 and correcting for their effects on siena-generated measures of brain volume
 672 change. NeuroImage 49 (2), 1601–1611.
- 673 Chang, C.-C., Lin, C.-J., 2011. LIBSVM: A library for support vector machines.
 674 ACM Transactions on Intelligent Systems and Technology 2, 27:1–27:27, soft-
 675 ware available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.

- 676 Chen, G., Saad, Z. S., Britton, J. C., Pine, D. S., Cox, R. W., 2013. Linear
677 mixed-effects modeling approach to fmri group analysis. *Neuroimage* 73, 176–
678 190.
- 679 Cheng, W., Palaniyappan, L., Li, M., Kendrick, K. M., Zhang, J., Luo, Q., Liu,
680 Z., Yu, R., Deng, W., Wang, Q., Ma, X., Guo, W., Francis, S., Liddle, P.,
681 Mayer, A. R., Schumann, G., Li, T., Feng, J., May 2015. Voxel-based, brain-
682 wide association study of aberrant functional connectivity in schizophrenia
683 implicates thalamocortical circuitry. *Npj Schizophrenia* 1, –.
684 URL <http://dx.doi.org/10.1038/npjschz.2015.16>
- 685 Cohen, J., 1992. A power primer. *Psychological bulletin* 112 (1), 155.
- 686 Collins, D. L., Neelin, P., Peters, T. M., Evans, A. C., 1994. Automatic 3D
687 intersubject registration of MR volumetric data in standardized Talairach
688 space. *Journal of computer assisted tomography* 18 (2), 192–205.
689 URL <http://view.ncbi.nlm.nih.gov/pubmed/8126267>
- 690 Cortes, C., Vapnik, V., Sep. 1995. Support-vector networks. *Machine Learning*
691 20 (3), 273–297.
692 URL <http://dx.doi.org/10.1007/BF00994018>
- 693 Damoiseaux, J. S., Rombouts, S. A. R. B., Barkhof, F., Scheltens, P., Stam,
694 C. J., Smith, S. M., Beckmann, C. F., Sep. 2006. Consistent resting-state
695 networks across healthy subjects. *Proceedings of the National Academy of
696 Sciences* 103 (37), 13848–13853.
697 URL <http://dx.doi.org/10.1073/pnas.0601417103>
- 698 Dansereau, C., Bellec, P., Lee, K., Pittau, F., Gotman, J., Grova, C., 2014.
699 Detection of abnormal resting-state networks in individual patients suffering
700 from focal epilepsy: An initial step toward individual connectivity assess-
701 ment. *Frontiers in Neuroscience* 8 (419).
702 URL [http://www.frontiersin.org/brain_imaging_methods/10.3389/
703 fnins.2014.00419/abstract](http://www.frontiersin.org/brain_imaging_methods/10.3389/fnins.2014.00419/abstract)
- 704 Dansereau, C., Risterucci, C., Pich, E. M., Arnold, D., Bellec, P., 2013. A
705 power analysis for multisite studies in resting-state functional connectivity,
706 with an application to clinical trials in alzheimer's disease. Vol. 9. pp. P248
707 – P249, alzheimer's Association International Conference 2013 Alzheimer's
708 Association International Conference 2013.
709 URL [http://www.sciencedirect.com/science/article/pii/
710 S1552526013011461](http://www.sciencedirect.com/science/article/pii/S1552526013011461)
- 711 Desmond, J., Glover, G., Aug. 2002. Estimating sample size in functional mri
712 (fmri) neuroimaging studies: Statistical power analyses. *Journal of Neuro-
713 science Methods* 118 (2), 115–128.
714 URL [http://dx.doi.org/10.1016/s0165-0270\(02\)00121-8](http://dx.doi.org/10.1016/s0165-0270(02)00121-8)

- 715 Durnez, J., Moerkerke, B., Nichols, T. E., 2014. Post-hoc power estimation for
716 topological inference in fmri. Neuroimage 84, 45–64.
- 717 Edward, V., Windischberger, C., Cunningham, R., Erdler, M., Lanzenberger, R.,
718 Mayer, D., Endl, W., Beisteiner, R., Nov 2000. Quantification of fmri artifact
719 reduction by a novel plaster cast head holder. Hum Brain Mapp 11 (3), 207–
720 213.
- 721 Elliott, M. R., Bowtell, R. W., Morris, P. G., Jun 1999. The effect of scanner
722 sound in visual, motor, and auditory functional mri. Magn Reson Med 41 (6),
723 1230–1235.
- 724 Fair, D. A., Nigg, J. T., Iyer, S., Bathula, D., Mills, K. L., Dosenbach, N.
725 U. F., Schlaggar, B. L., Mennes, M., Gutman, D., Bangaru, S., Buitelaar,
726 J. K., Dickstein, D. P., Martino, A. D., Kennedy, D. N., Kelly, C., Luna,
727 B., Schweitzer, J. B., Velanova, K., Wang, Y.-F., Mostofsky, S., Castellanos,
728 F. X., Milham, M. P., 2012. Distinct neural signatures detected for adhd
729 subtypes after controlling for micro-movements in resting state functional
730 connectivity mri data. Front Syst Neurosci 6, 80.
731 URL <http://dx.doi.org/10.3389/fnsys.2012.00080>
- 732 Feaster, D., Mikulich-Gilbertson, S., Brincks, A., Sep. 2011. Modeling site effects
733 in the design and analysis of multi-site trials. The American journal of drug
734 and alcohol abuse 37 (5), 383–391.
735 URL <http://dx.doi.org/10.3109/00952990.2011.600386>
- 736 Fonov, V., Evans, A. C., Botteron, K., Almlí, C. R., McKinstry, R. C., Collins,
737 D. L., Brain Development Cooperative Group, Jan. 2011. Unbiased average
738 age-appropriate atlases for pediatric studies. NeuroImage 54 (1), 313–327.
739 URL <http://dx.doi.org/10.1016/j.neuroimage.2010.07.033>
- 740 Friedman, L., Glover, G., Jun. 2006. Report on a multicenter fmri quality assur-
741 ance protocol. Journal of magnetic resonance imaging : JMRI 23 (6), 827–839.
742 URL <http://dx.doi.org/10.1002/jmri.20583>
- 743 Friedman, L., Glover, G., Consortium, T. F., Nov. 2006. Reducing interscanner
744 variability of activation in a multicenter fmri study: Controlling for signal-to-
745 fluctuation-noise-ratio (sfnr) differences. NeuroImage 33 (2), 471–481.
746 URL <http://dx.doi.org/10.1016/j.neuroimage.2006.07.012>
- 747 Friedman, L., Stern, H., Brown, G. G., Mathalon, D. H., Turner, J., Glover,
748 G. H., Gollub, R. L., Lauriello, J., Lim, K. O., Cannon, T., et al., 2008.
749 Test-retest and between-site reliability in a multicenter fmri study. Human
750 brain mapping 29 (8), 958–972.
- 751 Giove, F., Gili, T., Iacobella, V., Macaluso, E., Maraviglia, B., Oct. 2009.
752 Images-based suppression of unwanted global signals in resting-state func-
753 tional connectivity studies. Magnetic resonance imaging 27 (8), 1058–1064.
754 URL <http://dx.doi.org/10.1016/j.mri.2009.06.004>

- 755 Glover, G. H., Mueller, B. A., Turner, J. A., van Erp, T. G., Liu, T. T., Greve,
756 D. N., Voyvodic, J. T., Rasmussen, J., Brown, G. G., Keator, D. B., et al.,
757 2012. Function biomedical informatics research network recommendations for
758 prospective multicenter functional mri studies. Journal of Magnetic Reso-
759 nance Imaging 36 (1), 39–54.
- 760 Greicius, M. D., Srivastava, G., Reiss, A. L., Menon, V., Mar. 2004. Default-
761 mode network activity distinguishes Alzheimer's disease from healthy aging:
762 Evidence from functional MRI. Proceedings of the National Academy of Sci-
763 ences of the United States of America 101 (13), 4637–4642.
764 URL <http://dx.doi.org/10.1073/pnas.0308627101>
- 765 Hunter, J. D., 2007. Matplotlib: A 2d graphics environment. Computing In
766 Science & Engineering 9 (3), 90–95.
- 767 Jovicich, J., Minati, L., Marizzoni, M., Marchitelli, R., Sala-Llonch, R., Bartrés-
768 Faz, D., Arnold, J., Benninghoff, J., Fiedler, U., Roccatagliata, L., et al.,
769 2016. Longitudinal reproducibility of default-mode network connectivity in
770 healthy elderly participants: A multicentric resting-state fmri study. Neu-
771 roImage 124, 442–454.
- 772 Kilpatrick, L., Zald, D., Pardo, J., Cahill, L., 2006. Sex-related differences
773 in amygdala functional connectivity during resting conditions. Neuroimage
774 30 (2), 452–461.
- 775 Lund, T. E., Madsen, K. H., Sidaros, K., Luo, W.-L., Nichols, T. E., Jan. 2006.
776 Non-white noise in fMRI: does modelling have an impact? NeuroImage 29 (1),
777 54–66.
778 URL <http://dx.doi.org/10.1016/j.neuroimage.2005.07.005>
- 779 Milham, M. P., Fair, D., Mennes, M., Mostofsky, S. H., 2012. The adhd-200
780 consortium: a model to advance the translational potential of neuroimaging
781 in clinical neuroscience. Frontiers in Systems Neuroscience 6 (62).
782 URL [http://www.frontiersin.org/systems_neuroscience/10.3389/
783 fnsys.2012.00062/full](http://www.frontiersin.org/systems_neuroscience/10.3389/fnsys.2012.00062/full)
- 784 Mueller, S. G., Weiner, M. W., Thal, L. J., Petersen, R. C., Jack, C., Jagust,
785 W., Trojanowski, J. Q., Toga, A. W., Beckett, L., Nov 2005. The alzheimer's
786 disease neuroimaging initiative. Neuroimaging Clin N Am 15 (4), 869–77, xi-
787 xii.
788 URL <http://dx.doi.org/10.1016/j.nic.2005.09.008>
- 789 Nielsen, J., Zielinski, B., Fletcher, T., Alexander, A., Lange, N., Bigler, E.,
790 Lainhart, J., Anderson, J., 2013. Multisite functional connectivity mri classi-
791 fication of autism: Abide results. Frontiers in human neuroscience 7, –.
792 URL <http://view.ncbi.nlm.nih.gov/pubmed/24093016>
- 793 Orban, P., Madjar, C., Savard, M., Dansereau, C., Tam, A., Das, S., Evans,
794 A. C., Rosa-Neto, P., Breitner, J. C., Bellec, P., oct 2015. Test-retest resting-
795 state fMRI in healthy elderly persons with a family history of alzheimer's

- 796 disease. *Scientific Data* 2, 150043.
797 URL <http://dx.doi.org/10.1038/sdata.2015.43>
- 798 Power, J. D., Barnes, K. A., Snyder, A. Z., Schlaggar, B. L., Petersen, S. E.,
799 Feb. 2012. Spurious but systematic correlations in functional connectivity
800 MRI networks arise from subject motion. *NeuroImage* 59 (3), 2142–2154.
801 URL <http://dx.doi.org/10.1016/j.neuroimage.2011.10.018>
- 802 Power, J. D., Cohen, A. L., Nelson, S. M., Wig, G. S., Barnes, K. A., Church,
803 J. A., Vogel, A. C., Laumann, T. O., Miezin, F. M., Schlaggar, B. L., Pe-
804 tersen, S. E., Nov. 2011. Functional Network Organization of the Human
805 Brain. *Neuron* 72 (4), 665–678.
806 URL <http://dx.doi.org/10.1016/j.neuron.2011.09.006>
- 807 Shehzad, Z., Kelly, C. M., Reiss, P. T., Gee, D. G., Gotimer, K., Uddin, L. Q.,
808 Lee, S. H. H., Margulies, D. S., Roy, A. K. K., Biswal, B. B., Petkova, E.,
809 Castellanos, F. X., Milham, M. P., Oct. 2009. The resting brain: uncon-
810 strained yet reliable. *Cerebral cortex* (New York, N.Y. : 1991) 19 (10), 2209–
811 2229.
812 URL <http://dx.doi.org/10.1093/cercor/bhn256>
- 813 Sheline, Y. I., Price, J. L., Yan, Z., Mintun, M. A., 2010. Resting-state func-
814 tional mri in depression unmasks increased connectivity between networks via
815 the dorsal nexus. *Proceedings of the National Academy of Sciences* 107 (24),
816 11020–11025.
- 817 Sutton, B. P., Goh, J., Hebrank, A., Welsh, R. C., Chee, M. W., Park, D. C.,
818 2008. Investigation and validation of intersite fmri studies using the same
819 imaging hardware. *Journal of Magnetic Resonance Imaging* 28 (1), 21–28.
- 820 Tam, A., Dansereau, C., Badhwar, A., Orban, P., Belleville, S., Chertkow, H.,
821 Dagher, A., Hanganu, A., Monchi, O., Rosa-Neto, P., Shmuel, A., Wang,
822 S., Breitner, J., Bellec, P., 2015. Common effects of amnestic mild cognitive
823 impairment on resting-state connectivity across four independent studies.
824 *Frontiers in Aging Neuroscience* 7 (242).
825 URL [http://www.frontiersin.org/aging_neuroscience/10.3389/
826 fnagi.2015.00242/abstract](http://www.frontiersin.org/aging_neuroscience/10.3389/fnagi.2015.00242/abstract)
- 827 van den Heuvel, M., Mandl, R., Hulshoff Pol, H., Apr. 2008. Normalized Cut
828 Group Clustering of Resting-State fMRI Data. *PLoS ONE* 3 (4), e2001+.
829 URL <http://dx.doi.org/10.1371/journal.pone.0002001>
- 830 Van Dijk, K. R., Hedden, T., Venkataraman, A., Evans, K. C., Lazar, S. W.,
831 Buckner, R. L., Jan. 2010. Intrinsic functional connectivity as a tool for human
832 connectomics: theory, properties, and optimization. *Journal of neurophysiol-
833 ogy* 103 (1), 297–321.
834 URL <http://dx.doi.org/10.1152/jn.00783.2009>

- 835 Van Dijk, K. R., Sabuncu, M. R., Buckner, R. L., 2012. The influence of head
836 motion on intrinsic functional connectivity mri. Neuroimage 59 (1), 431–438.
- 837 Vanhoutte, G., Verhoye, M., der Linden, A. V., May 2006. Changing body
838 temperature affects the t2* signal in the rat brain and reveals hypothalamic
839 activity. Magn Reson Med 55 (5), 1006–1012.
840 URL <http://dx.doi.org/10.1002/mrm.20861>
- 841 Wang, D., Liu, B., Qin, W., Wang, J., Zhang, Y., Jiang, T., Yu, C., Dec. 2012.
842 KIBRA gene variants are associated with synchronization within the default-
843 mode and executive control networks. NeuroImage.
844 URL <http://dx.doi.org/10.1016/j.neuroimage.2012.12.022>
- 845 Worsley, K. J., Friston, K. J., Sep. 1995. Analysis of fMRI Time-Series
846 Revisited—Again. NeuroImage 2 (3), 173–181.
847 URL <http://dx.doi.org/10.1006/nimg.1995.1023>
- 848 Yan, C., Liu, D., He, Y., Zou, Q., Zhu, C., Zuo, X., Long, X., Zang, Y., May
849 2009. Spontaneous Brain Activity in the Default Mode Network Is Sensitive to
850 Different Resting-State Conditions with Limited Cognitive Load. PLoS ONE
851 4 (5), e5743+.
852 URL <http://dx.doi.org/10.1371/journal.pone.0005743>
- 853 Yan, C.-G. G., Craddock, C. C., Zuo, X.-N. N., Zang, Y.-F. F., Milham, M. P.,
854 Oct. 2013. Standardizing the intrinsic brain: towards robust measurement
855 of inter-individual variation in 1000 functional connectomes. NeuroImage 80,
856 246–262.
857 URL <http://view.ncbi.nlm.nih.gov/pubmed/23631983>
- 858 Yeo, B. T., Krienen, F. M., Sepulcre, J., Sabuncu, M. R., Lashkari, D.,
859 Hollinshead, M., Roffman, J. L., Smoller, J. W., Zöllei, L., Polimeni, J. R.,
860 Fischl, B., Liu, H., Buckner, R. L., Sep. 2011. The organization of the human
861 cerebral cortex estimated by intrinsic functional connectivity. Journal of neu-
862 rophysiology 106 (3), 1125–1165.
863 URL <http://dx.doi.org/10.1152/jn.00338.2011>
- 864 Zuo, X.-N., Anderson, J. S., Bellec, P., Birn, R. M., Biswal, B. B., Blautzik, J.,
865 Breitner, J. C. S., Buckner, R. L., Calhoun, V. D., Castellanos, F. X., Chen,
866 A., Chen, B., Chen, J., Chen, X., Colcombe, S. J., Courtney, W., Craddock,
867 R. C., Di Martino, A., Dong, H.-M., Fu, X., Gong, Q., Gorgolewski, K. J.,
868 Han, Y., He, Y., He, Y., Ho, E., Holmes, A., Hou, X.-H., Huckins, J., Jiang,
869 T., Jiang, Y., Kelley, W., Kelly, C., King, M., LaConte, S. M., Lainhart, J. E.,
870 Lei, X., Li, H.-J., Li, K., Li, K., Lin, Q., Liu, D., Liu, J., Liu, X., Liu, Y., Lu,
871 G., Lu, J., Luna, B., Luo, J., Lurie, D., Mao, Y., Margulies, D. S., Mayer,
872 A. R., Meindl, T., Meyerand, M. E., Nan, W., Nielsen, J. A., O'Connor,
873 D., Paulsen, D., Prabhakaran, V., Qi, Z., Qiu, J., Shao, C., Shehzad, Z.,
874 Tang, W., Villringer, A., Wang, H., Wang, K., Wei, D., Wei, G.-X., Weng,
875 X.-C., Wu, X., Xu, T., Yang, N., Yang, Z., Zang, Y.-F., Zhang, L., Zhang,

⁸⁷⁶ Q., Zhang, Z., Zhang, Z., Zhao, K., Zhen, Z., Zhou, Y., Zhu, X.-T., Milham,
⁸⁷⁷ M. P., 9 Dec. 2014. An open science resource for establishing reliability and
⁸⁷⁸ reproducibility in functional connectomics. *Sci Data* 1, 140049.

879 *Supplementary Material – Statistical power and prediction accuracy in*
880 *multisite resting-state fMRI connectivity*

881 Submitted to Neuroimage.

882 C. Dansereau^{1,2}, Y. Benhajali^{1,3}, C. Risterucci⁴, E. Merlo Pich⁴, P. Orban¹, D.
883 Arnold⁵, P. Bellec^{1,2}

884 ¹Centre de Recherche de l’Institut Universitaire de Gériatrie de Montréal,
885 Montréal, CA

886 ²Department of Computer Science and Operations Research, University of Mon-
887 treal, Montreal, CA

888 ³Département d’anthropologie, Université Montréal, Montréal, CA

889 ⁴Clinical Imaging, pRED, F.Hoffman-La Roche, Basel, CH

890 ⁵NeuroRx inc., Montréal, CA

891
892 For all questions regarding the paper, please address correspondence to
893 Pierre Bellec, CRIUGM, 4545 Queen Mary, Montreal, QC, H3W 1W5, Canada.
894 Email: pierre.bellec (at) criugm.qc.ca.

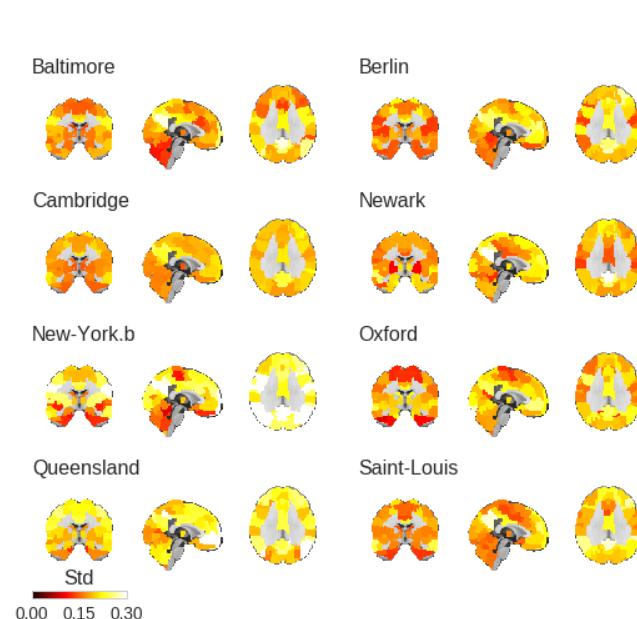


Figure S1: Standard deviation of resting-state connectivity across subjects, in the DMN, for each site, superimposed on the MNI152 template.

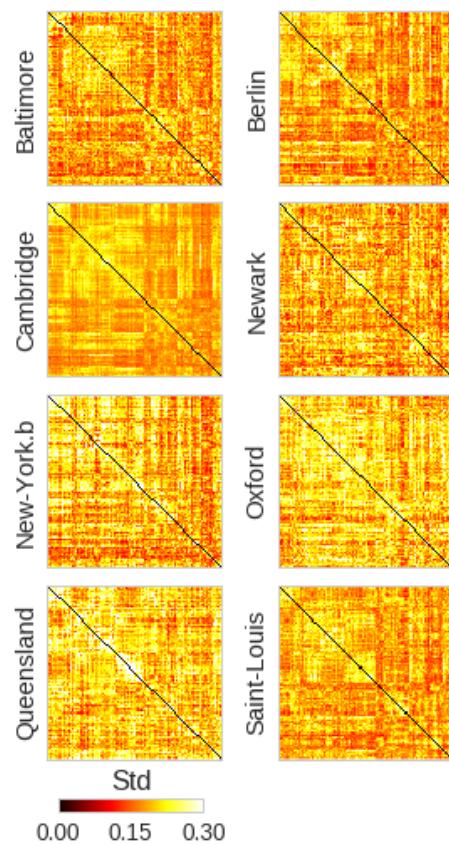


Figure S2: Standard deviation of resting-state connectivity across subjects, for the full connectome and each site.

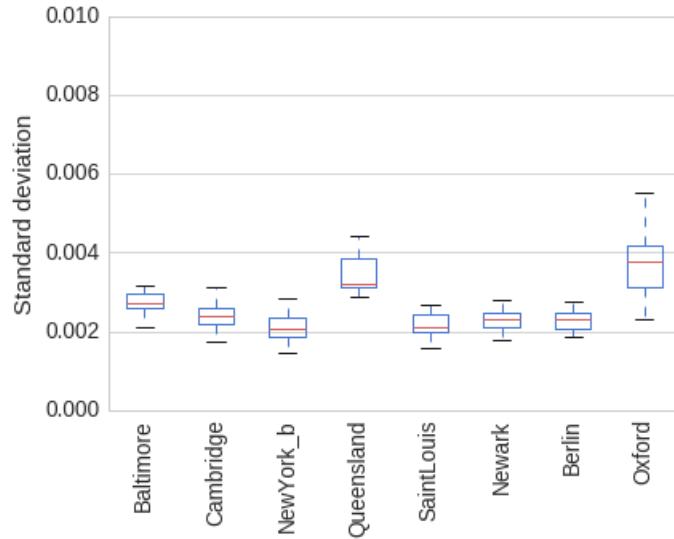


Figure S3: Standard deviation of resting-state time-series across subjects, averaged across all connections, at each site.

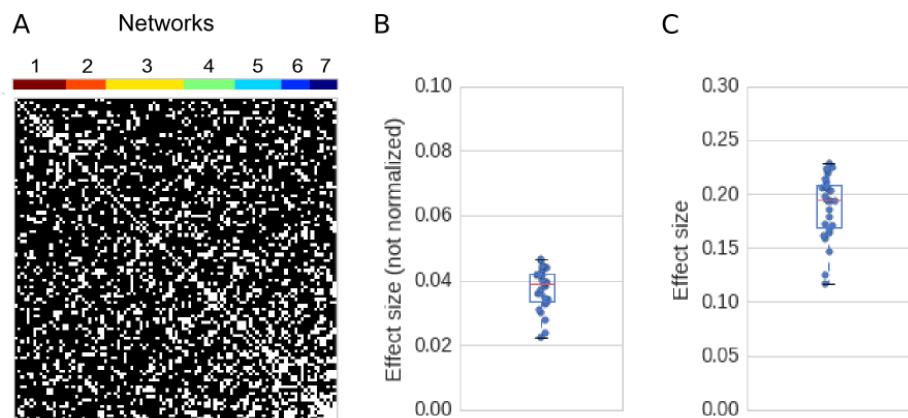


Figure S4: Panel A shows the results of a White test for homoscedasticity, across sites. Panel B show the average absolute difference in standard deviation between any pair of sites, and Panel C show the same difference, relative to the average of the standard deviation at the two sites.

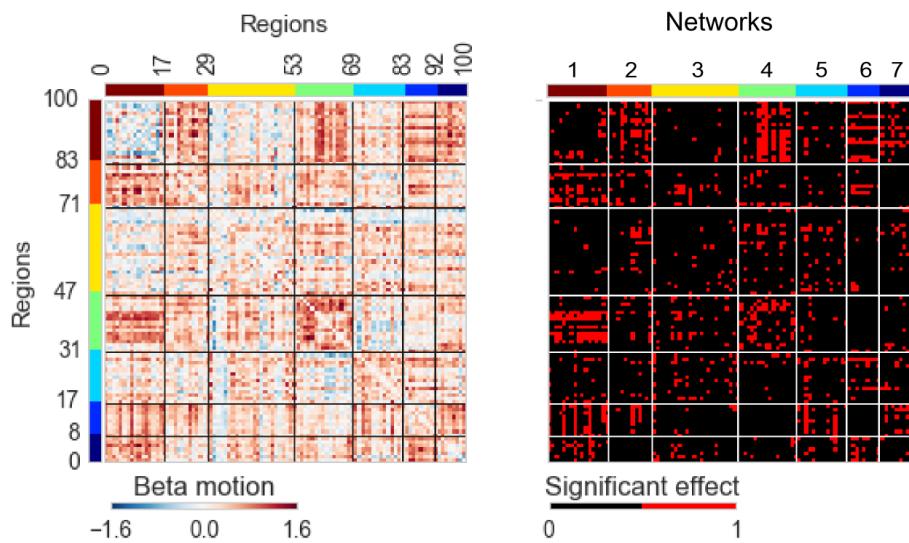


Figure S5: The figure shows average connectomes across all sites, as well as connections with a significant motion effect.

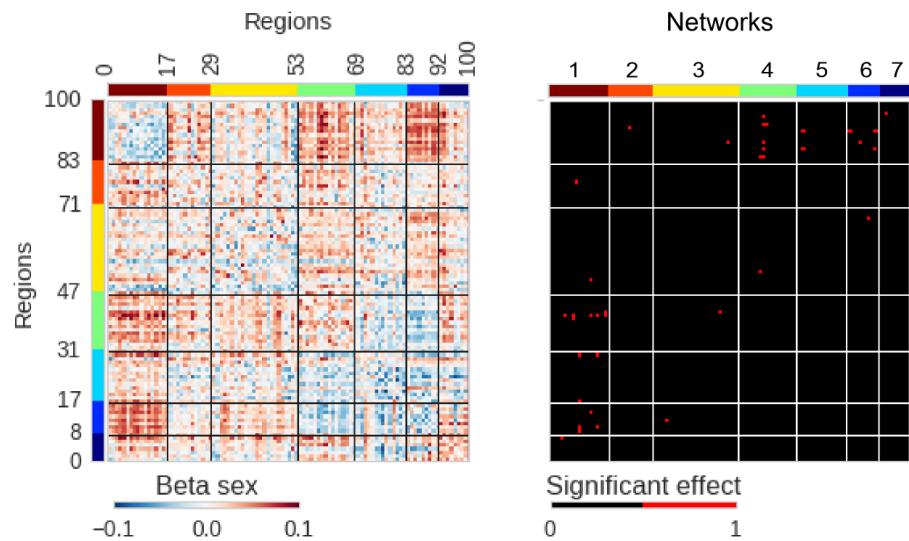


Figure S6: The figure shows average connectomes across all sites, as well as connections with a significant sex effect.

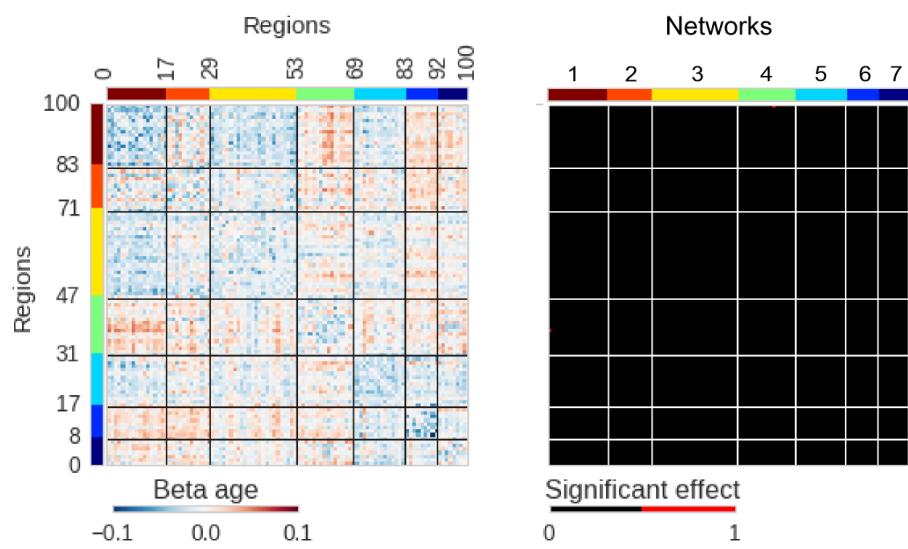


Figure S7: The figure shows average connectomes across all sites, as well as connections with a significant age effect.

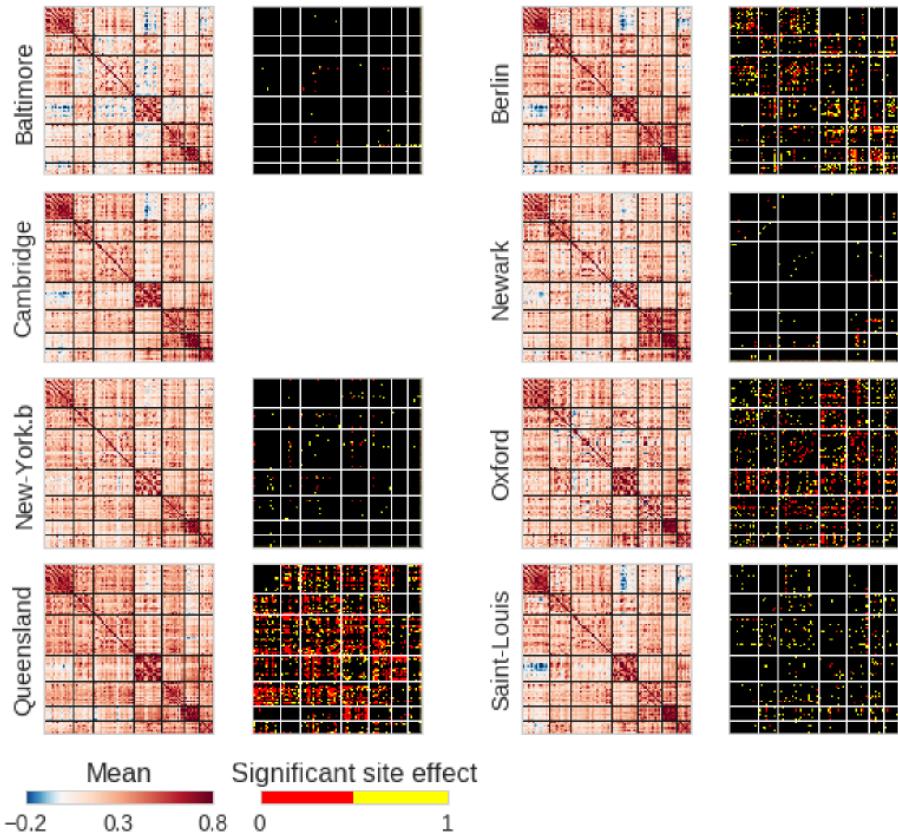


Figure S8: Average connectomes for individual sites, as well as connections with a significant site effect. This Figure is identical to Figure 2 in the paper, with the difference that the Cambridge site was excluded from the analysis. The intersection (\cap) of the significant site effects are shown in red and the symmetric difference (Δ) of the significant site effects are shown in yellow. Baltimore $\cap : 9, \Delta : 16$, Berlin $\cap : 318, \Delta : 333$, Newark $\cap : 23, \Delta : 36$, New-York.b $\cap : 25, \Delta : 45$, Oxford $\cap : 377, \Delta : 251$, Queensland $\cap : 946, \Delta : 389$, Saint-Louis $\cap : 49, \Delta : 162$