

Sample size requirement for achieving multisite harmonization using structural brain MRI features

Authors:

Parekh, Pravesh^{a, c, d, *}

Bhalerao, Gaurav Vivek^{b, c, d, *}

the ADBS consortium[^]

John, John P^{a, c, d, #}

Venkatasubramanian, G^{b, c, d, #}

* equal contribution; # Corresponding authors

Keywords: neuroimaging; harmonization; sample size; multisite; Mahalanobis distance; cross-validation

Affiliations:

^a Multimodal Brain Image Analysis Laboratory, National Institute of Mental Health and Neurosciences (NIMHANS), Bangalore, India

^b Translational Psychiatry Lab, National Institute of Mental Health and Neurosciences (NIMHANS), Bangalore, India

^c ADBS Neuroimaging Centre, National Institute of Mental Health and Neurosciences (NIMHANS), Bangalore, India

^d Department of Psychiatry, National Institute of Mental Health and Neurosciences (NIMHANS), Bangalore, India

[^] Accelerator Program for Discovery in Brain disorders using Stem cells (ADBS) Consortium:

Biju Viswanath¹, Naren P. Rao¹, Janardhanan C. Narayanaswamy¹, Palanimuthu T. Sivakumar¹, Arun Kandasamy¹, Muralidharan Kesavan¹, Urvakhsh Meherwan Mehta¹, Odity Mukherjee², Meera Purushottam¹, Ramakrishnan Kannan¹, Bhupesh Mehta¹, Thennarasu Kandavel¹, B. Binukumar¹, Jitender Saini¹, Deepak Jayarajan¹, A. Shyamsundar¹, Sydney Moirangthem¹, K. G. Vijay Kumar¹, Jayant Mahadevan¹, Bharath Holla¹, Jagadisha Thirthalli¹, Prabha S. Chandra¹, Bangalore N. Gangadhar¹, Pratima Murthy¹, Mitradas M. Panicker³, Upinder S. Bhalla³, Sumantra Chattarji³, Vivek Benegal¹, Mathew Varghese¹, Janardhan Y. C. Reddy¹, Padinjat Raghu³, Mahendra Rao³, and Sanjeev Jain

¹National Institute of Mental Health and Neurosciences (NIMHANS); ²Institute for Stem Cell Biology and Regenerative Medicine (InStem); ³National Center for Biological Sciences (NCBS)

Address for correspondence:

Dr. John P. John
Multimodal Brain Image Analysis Laboratory
Department of Psychiatry,
National Institute of Mental Health and
Neurosciences (NIMHANS),
Bangalore - 560029, India
jpjnimhans@gmail.com; jpj@nimhans.ac.in
(+91) 080 2699 5329

Dr. Ganesan Venkatasubramanian
Translational Psychiatry Lab
Department of Psychiatry,
National Institute of Mental Health and
Neurosciences (NIMHANS),
Bangalore - 560029, India
venkat.nimhans@gmail.com; gvs@nimhans.ac.in
(+91) 080 2699 5256

1 Abstract

2 When data is pooled across multiple sites, the extracted features are confounded by site effects.
3 Harmonization methods attempt to correct these site effects while preserving the biological
4 variability within the features. However, little is known about the sample size requirement for
5 effectively learning the harmonization parameters and their relationship with the increasing
6 number of sites. In this study, we performed experiments to find the minimum sample size
7 required to achieve multisite harmonization (using neuroHarmonize) using volumetric and
8 surface features by leveraging the concept of learning curves. Our first two experiments show
9 that site-effects are effectively removed in a univariate and multivariate manner; however, it is
10 essential to regress the effect of covariates from the harmonized data additionally. Our
11 following two experiments with actual and simulated data showed that the minimum sample
12 size required for achieving harmonization grows with the increasing average Mahalanobis
13 distances between the sites and their reference distribution. We conclude by positing a general
14 framework to understand the site effects using the Mahalanobis distance. Further, we provide
15 insights on the various factors in a cross-validation design to achieve optimal inter-site
16 harmonization.

1 Introduction

2 With the advent of standardized data sharing structures such as the Brain Imaging Data
3 Structure (Gorgolewski et al., 2016) and the availability of open-source platforms for data
4 sharing, such as OpenNeuro (Markiewicz et al., 2021), it is increasingly common to share
5 different kinds of neuroimaging data. Performing analyses by pooling samples across these
6 datasets allows for increased sample size, better representation of geographical diversity, and
7 the potential to develop robust, generalizable models. However, using multi-site data is
8 challenging owing to factors like different scanners/hardware, differences in acquisition
9 protocols, conditions of data acquisition (such as subject-positioning, eyes-open vs. eyes-
10 closed during resting-state functional magnetic resonance imaging (fMRI)), variations in terms
11 of image quality, etc. When pooling samples across scanners, it is essential to correct for site-
12 related variations, which can otherwise influence outcome measurements like cortical
13 thickness and brain volumes (for example, see (Fortin et al., 2018; Lee et al., 2019; Liu et al.,
14 2020; Medawar et al., 2021; Takao et al., 2014; Wittens et al., 2021)). These “batch effects”
15 are well-known in fields like microarray technology, where methods have been developed to
16 correct for these effects (see (Johnson et al., 2007) and (Leek et al., 2010) for a review).

17 Multiple harmonization methods have been proposed for correcting scanner-related differences
18 in neuroimaging. A regression-based procedure to correct for site-effects is to add dummy-
19 coded scanner/site variables (for example, (Bruin et al., 2019; Fennema-Notestine et al., 2007;
20 Pardoe et al., 2008; Rozycki et al., 2018; Segall et al., 2009; Stonnington et al., 2008)). Another
21 popular method for correcting site-specific effects is to use ComBat (Fortin et al., 2018;
22 Johnson et al., 2007). The ComBat harmonization method models the location and scale of the
23 variables to be harmonized, accounting for the additive and multiplicative effects of the site on
24 the variables and preserving the site-specific biological variability (Fortin et al., 2018, 2017).
25 Various extensions to ComBat have been proposed, such as ComBat-GAM, which models the

1 non-linear effect of the age (Pomponio et al., 2020), CovBat, which additionally models the
2 covariance in the data (Chen et al., n.d.), ComBat for longitudinal data (Beer et al., 2020), etc.
3 A newly developed method, NeuroHarmony, attempts to generalize to unseen scanners/sites,
4 taking into account the image quality metrics (Garcia-Dias et al., 2020).

5 Previous work has shown that harmonization methods like ComBat can remove site-related
6 effects. For example, (Zavaliangos-Petropulu et al., 2019) examined differences in ROI-level
7 diffusion measures and found only one remaining ROI showing significant protocol-related
8 differences post harmonization. In (Fortin et al., 2017), the authors showed that ComBat
9 effectively removed site-related differences from voxel-level diffusion scalar maps and ROI-
10 level diffusion measures. Similarly, in (Fortin et al., 2018), the authors showed that the site-
11 related effects on cortical thickness were removed using Combat. In addition to univariate
12 results, the authors also demonstrated the removal of site-related effects in a multivariate
13 manner: a support vector machine (SVM) classifier was unable to predict the site after
14 harmonization.

15 The sample used for “learning” harmonization parameters must adequately capture site-related
16 effects to achieve inter-site harmonization. This becomes critical in situations where
17 harmonization needs to be carried out in a cross-validation manner – learning of the
18 harmonization parameters happens from the training set, and these parameters are then applied
19 to the test set (for example, in machine learning). Therefore, a central question in such
20 paradigms is finding the minimum sample size required to eliminate the site effects.
21 Additionally, it is essential to assess the sample size requirement in the context of a potential
22 multivariate relationship between the variables and the site. In this paper, we attempt to address
23 this lacuna by leveraging the concept of learning curves to find the minimum sample size
24 required to remove site-related effects. By iteratively increasing the sample size per site and

1 training a machine learning classifier to predict the site, we attempt to find the sample size at
2 which the classifier prediction reduces to chance.

3 2. Methodology

4 2.1 Datasets

5 For this study, we selected T1-weighted MRI scans of healthy subjects from four publicly
6 available datasets and supplemented them with scans from our labs. We restricted our selection
7 to datasets that had scans of at least 300 healthy subjects acquired on the same scanner. The
8 first dataset was the Southwest University Adult Lifespan Dataset (SALD) (Wei et al., 2018),
9 a cross-sectional sample of 494 subjects. The second dataset consisted of the scans acquired at
10 the Guy's Hospital and available as part of the IXI dataset (available at [https://brain-](https://brain-development.org/ixi-dataset/)
11 [development.org/ixi-dataset/](https://brain-development.org/ixi-dataset/)) and composed of 322 subjects (henceforth referred to as “Guys”
12 dataset). The third dataset was the Amsterdam Open MRI Collection (AOMIC) (Snoek et al.,
13 2021a, 2021b) and consisted of 928 subjects. The fourth dataset was a pooled¹ version of four
14 different datasets: the Beijing Normal University (BNU) dataset 1 (Lin et al., 2015) ($n = 57$;
15 available at: https://fcon_1000.projects.nitrc.org/indi/CoRR/html/bnu_1.html), BNU dataset 2
16 (Huang et al., 2016) ($n = 61$; available at:
17 https://fcon_1000.projects.nitrc.org/indi/CoRR/html/bnu_2.html), BNU dataset 3 ($n = 48$;
18 available at: https://fcon_1000.projects.nitrc.org/indi/CoRR/html/bnu_3.html), from the
19 Consortium for Reliability and Reproducibility (CoRR) dataset (Zuo et al., 2014), and the
20 “Beijing_Zang” dataset ($n = 198$) from the 1000 Functional Connectomes Project (Biswal et
21 al., 2010). The CoRR dataset is test-retest reliability scans, and for the present work, we only
22 considered the baseline scans for all subjects (we henceforth refer to this combined dataset as

¹ We pooled these datasets as they were acquired on the same scanner; (Huang et al., 2016) mentions that the BNU series was acquired on the same scanner and Beijing_Zang was confirmed to have been acquired on the same scanner (Y.F. Zang, personal communication, December 04, 2021)

1 “BNUBeijing”). The final dataset (“NIMHANS” dataset) consisted of 372 subjects collected
 2 at the National Institute of Mental Health and Neurosciences (NIMHANS) as part of two
 3 different research projects.

4 2.2 Image Acquisition and Processing

5 We have summarized the critical acquisition parameters for the datasets in **Table 1**. For each
 6 image, we set the origin (i.e., (0,0,0) coordinate) to correspond to the anterior commissure (AC)
 7 using *acpcdetect* v2 (Ardekani, 2018; Ardekani et al., 1997; Ardekani and Bachman, 2009)
 8 (available at: <https://www.nitrc.org/projects/art/>). We then, visually examined the images and
 9 manually set the origin to the AC using the display utility in SPM12 v7771
 10 (<https://www.fil.ion.ucl.ac.uk/spm/software/spm12/>) in case of *acpcdetect* failure. Then, we
 11 ran the segmentation and surface pipeline for all the images using the Computational Anatomy
 12 Toolbox (CAT) 12.7 v 1727 (<http://www.neuro.uni-jena.de/cat/>) with SPM12 v7771 in the
 13 background, running on MATLAB R2020a (The MathWorks, Natick, USA;
 14 <https://www.mathworks.com>).

15 **Table 1:** Summary of crucial acquisition parameters across the five datasets

Resolution (mm)	Slices	TR (ms)	TE (ms)	TI (ms)	Flip Angle (degrees)
SALD: Siemens Trio – 3T ($n = 494$)					
1.0000 × 1.0000 × 1.0000	176	1900	2.52	900	90
Guys: Philips Intera – 1.5T ($n = 322$)					
1.2000 × 0.9375 × 0.9375	130/140/150	9.813	4.603	-	8
AOMIC: Philips Intera – 3T ($n = 928$)					
1.0000 × 1.0000 × 1.0000	160	81	37	-	8
BNU-1: Siemens Trio – 3T ($n = 57$)					
1.3300 × 1.0000 × 1.0000	144	2530	3.39	1100	7
BNU-2: Siemens Trio – 3T ($n = 61$)					
1.3300 × 1.0000 × 1.0000	128	2530	3.39	1100	7

BNU-3: Siemens Trio – 3T ($n = 48$)

1.3300 × 1.0000 × 1.0000 127/128 2530 3.39 1100 7

Beijing_Zang: Siemens Trio – 3T ($n = 198$)

1.3300 × 1.0000 × 1.0000 128/176 - - - -

NIMHANS: Siemens Skyra – 3T ($n = 372$)

1.0000 × 0.9727 × 0.9727 192 1900 2.4 900 9

1.0000 × 1.0039 × 1.0039

1.0000 × 0.9609 × 0.9609

1.0000 × 0.9766 × 0.9766

1.0000 × 1.0000 × 1.0000

1

2 2.3 Quality Check

3 First, we rejected the data of any subject with an age of less than 18 years. Next, we flagged
4 the data of any subject whose CAT12 report had a noise rating of “D” or below. In the next
5 step, we eliminated the scans of any subject where the CAT12 quantified white matter
6 hyperintensity exceeded a volume of 10 cm³ (while this criterion is arbitrary, it helped eliminate
7 any scans with a potential of underlying white matter abnormalities). Additionally, to identify
8 scans with improper segmentation, we calculated the Dice coefficient of the (binarized)
9 modulated normalized gray matter image with the (binarized) template gray matter image
10 (from CAT12) and flagged the images with a Dice coefficient less than 0.9. In addition, the
11 NIMHANS dataset has undergone a thorough visual quality check as part of ongoing research
12 work at our labs.

13 2.4 Selected Sample

14 After the quality check, we had 489 subjects in the SALD dataset, 302 in the Guys dataset, 928
15 in the AOMIC dataset, 341 in the BNUBeijing dataset, and 318 in the NIMHANS dataset.
16 Since the lowest number was 302 subjects in the Guys dataset, we decided to randomly subset
17 300 subjects from each dataset for further experiments. Further, we rounded the age to the

1 nearest integer. The socio-demographic details of these 1500 subjects are summarized in **Table**
 2 **2**.

3 **Table 2:** Summary of socio-demographics details for each dataset after quality check

Dataset	# Males	# Females	Age (Males):	Age (Females):	Age (Overall):
			Mean \pm SD,	Mean \pm SD,	Mean \pm SD,
			Min – Max	Min – Max	Min – Max
SALD	117	183	45.62 \pm 17.82, 20 – 80	45.28 \pm 17.46, 19 – 78	45.41 \pm 17.57, 19 – 80
Guys	132	168	48.33 \pm 16.44, 20 – 86	51.41 \pm 14.96, 21 – 80	50.05 \pm 15.68, 20 – 86
AOMIC	144	156	23.07 \pm 1.74, 20 – 26	23.03 \pm 1.79, 20 – 26	23.05 \pm 1.77, 20 – 26
BNUBeijing	131	169	21.79 \pm 1.86, 18 – 27	21.53 \pm 1.89, 18 – 29	21.64 \pm 1.88, 18 – 29
NIMHANS	188	112	27.20 \pm 5.10, 18 – 49	27.28 \pm 6.80, 18 – 50	27.23 \pm 5.78, 18 – 50

4
 5 **2.5 Features**

6 We extracted regional gray matter volumes using the Hammers atlas (CAT12 version;
 7 (Faillenot et al., 2017; Gousias et al., 2008; Hammers et al., 2003) as the primary features. This
 8 version of the Hammers atlas consists of 68 regions of interest (ROIs); from these, we excluded
 9 the bilateral parcels of corpus callosum, brainstem, and the ventricles, resulting in a total of 60
 10 gray matter volumes for each subject. In addition to volumetric features, we also performed
 11 experiments using regional cortical thickness, fractal complexity, sulcal depth, and gyrification
 12 index. For these surface features, we used the Desikan-Killiany (DK40) atlas (Desikan et al.,
 13 2006) consisting of 72 parcellations; from these, we excluded the bilateral corpus callosum and
 14 the unknown parcels, resulting in a total of 68 surface estimates for each subject.

1 2.6 Harmonization

2 For all the experiments, we used the neuroHarmonize (Pomponio et al., 2020) (available at:
3 <https://github.com/rpomponio/neuroHarmonize/>) toolbox for harmonizing the features across
4 scanners. For each experiment, each model (see below), we preserved the effects of age, total
5 intracranial volume (TIV), and sex (dummy coded as 1 for females).

6 2.7 Experiments

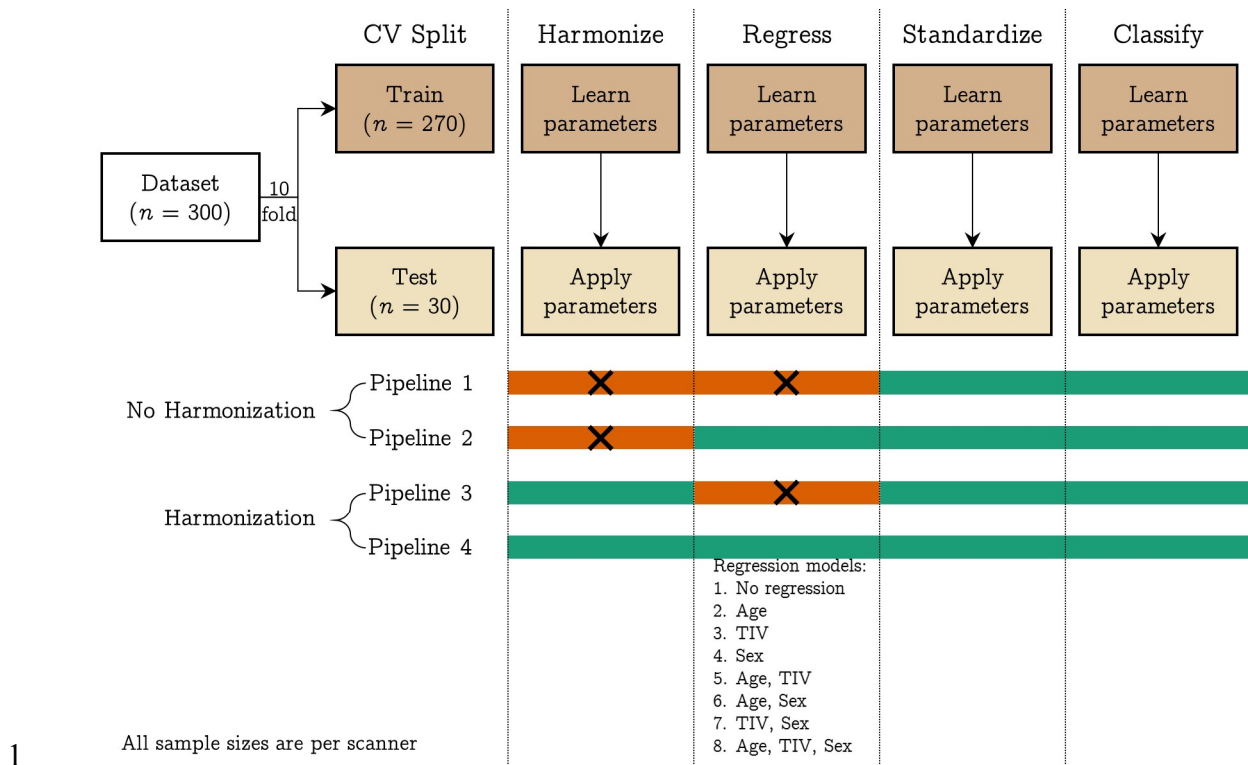
7 The primary motivation behind this study is to estimate the minimum sample size required to
8 achieve inter-site harmonization, such that the multivariate mapping between features and
9 scanners is removed. To achieve this, we first show that univariate (experiment 1) and
10 multivariate (experiment 2) mapping exists between features and scanners; we also show that
11 this mapping is removed after performing harmonization. Then, in experiment 3, we use the
12 concept of learning curves to estimate the minimum sample size required to remove the
13 (multivariate) site effects. Finally, in experiment 4, we extended the learning curve experiment
14 using simulated data. This was done to examine the sample size requirement under a wide-
15 range of Mahalanobis distances (see below) and number of sites. Critically, we note that for
16 experiments 2-4, we used the same seeds to ensure the comparability of the results.

17 **2.7.1 Experiment 1: univariate differences**

18 Following (Garcia-Dias et al., 2020), we performed two-sample Kolmogorov-Smirnov (KS)
19 tests between pairs of scanners for each ROI's brain volume and surface features to test if the
20 distribution of the features were significantly different before and after harmonization.
21 However, it should be noted that age, TIV, and sex will confound ROI features. Therefore, for
22 each pair of scanners, for each ROI, we performed a linear regression to correct for these
23 confounding variables. Then, the residuals from the linear regression were used for the pairwise
24 two-sample KS test. Similar to the earlier work (Garcia-Dias et al., 2020), we did not perform
25 any correction for multiple comparisons and examined our results at $\alpha = 0.05$.

1 **2.7.2 Experiment 2: multivariate differences**

2 We fit a linear SVM model to predict the scanner using volumetric features. We repeated this
3 before and after harmonization using 10-fold cross-validation considering a pair of scanners at
4 a time (10 combinations), three scanners at a time (10 combinations), four scanners at a time
5 (five combinations), and all five scanners at the same time. We used the one vs. one coding
6 method in MATLAB for multiclass classification. We performed the following operations
7 within a cross-validation framework: regression of age, TIV, and sex (independently and
8 combinations thereof; eight combinations), standardization of features, and training of linear
9 SVM (using the default hyperparameter $C = 1$). In each of these steps, the parameters were
10 learned from the training data of each fold and then applied to the test data. In the case of
11 harmonization, the training data was harmonized, and harmonization parameters were applied
12 to the test data before the regression step for each cross-validation fold. This was done to ensure
13 that the regression coefficients were not confounded by site effects; a similar approach has
14 been followed in (Pomponio et al., 2020). We repeated the 10-fold cross-validation 50 times
15 and performed an additional 50 repeats of permutation testing (i.e., 100 repetitions for
16 permutation testing). For permutation testing, we permuted the class labels of the entire data
17 and calculated the p values as described in (Ojala and Garriga, 2010). This procedure is
18 illustrated in **Figure 1**.



2 **Figure 1:** Pipelines implemented in experiment 2: we trained a linear SVM classifier to predict the scanner from
 3 raw and harmonized structural features; additionally, we explored eight different regression models where we
 4 regressed the effect of different confounding variables from the structural features; the four pipelines have four
 5 different modules: harmonization, regression, standardization, and classification; the steps indicated with orange
 6 color were not performed in that pipeline. The 10-fold cross-validation was repeated 50 times and an additional
 7 50 repeats of permutation testing (i.e., 100 repeats of permutation) were performed to assess whether the
 8 classification performance was above chance level. [color version of this figure is available online]

9 2.7.3 Quantifying the multivariate site-effect

10 In order to quantify the multivariate site-effect prevalent in our dataset, we used the
 11 Mahalanobis distance (MD) (Mahalanobis, 1936), which is a multivariate extension of Cohen's
 12 d (Cohen, 1988) and can be used for calculating the standardized mean differences between
 13 groups (Del Giudice, 2009). We first calculated a reference distribution using the overall mean
 14 and pooled covariance matrix from the individual distributions. Then, we calculated the MD
 15 from each individual distribution to the reference distribution as:

$$16 \quad MD_i = \sqrt{(\mu_i - \mu_R)^T C^{-1} (\mu_i - \mu_R)}$$

17 where μ_i indicates the means of all variables in that individual distribution, μ_R indicates the
 18 means of all variables in the reference distribution, C indicates the overall pooled variance-

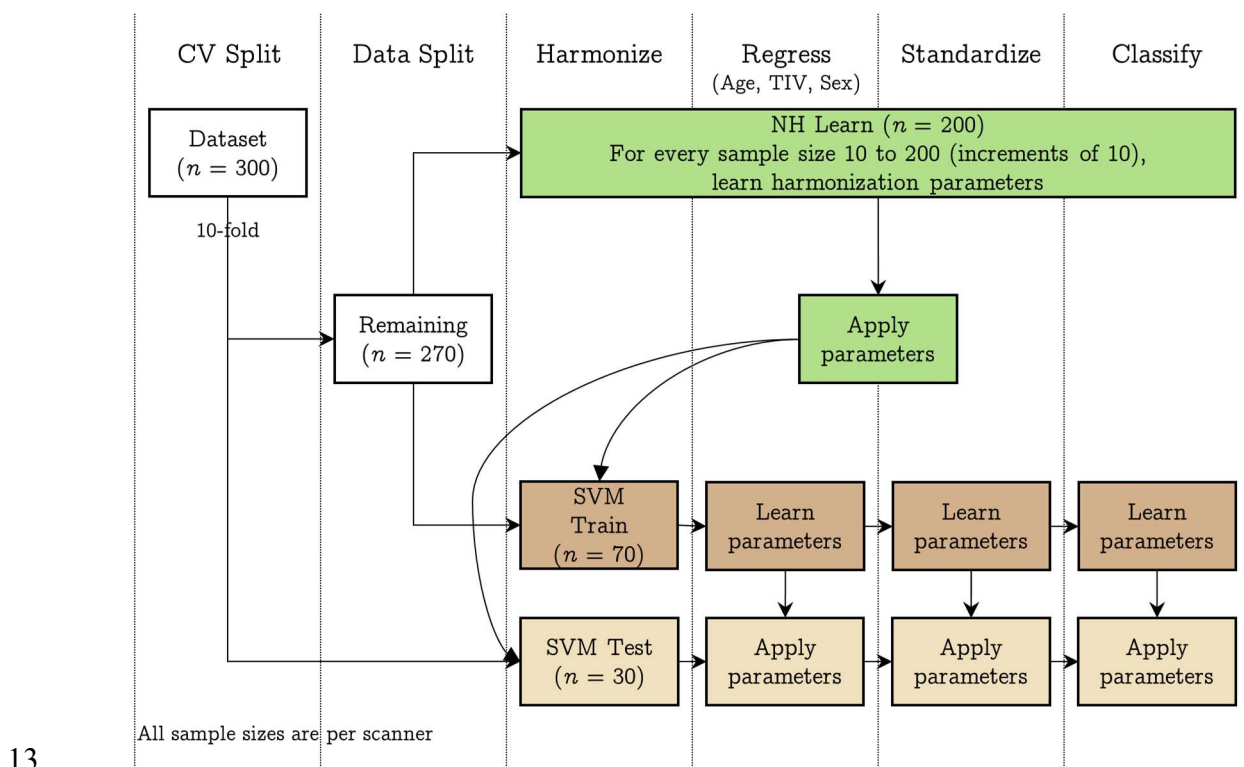
1 covariance matrix, and MD_i indicates the calculated Mahalanobis distance of that individual
2 distribution from the reference. We then average the individual D_i s to get an overall measure
3 of the site-effect. We calculated two versions of MD – one without regressing any covariates
4 and one after regressing the effect of TIV, age, and sex. When performing the regression, we
5 regressed the effect of the covariates on the entire data (i.e., linearly regressing the effect of
6 covariates across all sites in a single model) before calculating the MD.

7 **2.7.4 Experiment 3: learning curves**

8 In order to estimate the minimum sample size required to remove the site-effects, we leveraged
9 the concept of learning curves. Stated simply, a learning curve can be used to assess the
10 performance of a learner (SVM, in our case) as a function of increasing sample size. Within a
11 cross-validation framework, when small sample sizes are used for learning harmonization
12 parameters, we do not expect multivariate site-effects to be effectively removed from the test
13 set. This can be tested by dividing the datasets into three parts: the first part used for learning
14 harmonization parameters (and applying to the other parts), the second part used for learning
15 the multivariate mapping between features and scanners, and the third part used for evaluating
16 the classifier. If the site-effects are removed by harmonization, then a classifier will not be able
17 to learn the mapping between features and scanners, and therefore the test accuracy of this
18 classifier will be as good as chance level (which can be assessed by permutation testing).

19 The experimental design is illustrated in **Figure 2**. For this experiment, we first performed a
20 10-fold split on the dataset ($n = 1500$, 300 samples per scanner), resulting in 270 samples and
21 30 samples (“SVM test”) per scanner; the 270 samples per scanner were then split into 200
22 (“NH learn”) and 70 samples (“SVM train”) per scanner. The *NH learn* sample was used for
23 learning harmonization parameters, the *SVM train* sample was used for training SVM, and the
24 *SVM test* sample was used for testing the SVM performance. For the *NH learn* sample, we
25 iteratively increased the sample size from 10 samples per scanner to 200 samples per scanner

1 in increments of 10 samples (i.e., within each fold, 20 different sample sizes were used to learn
 2 harmonization parameters). For each sample size, we learnt the harmonization parameters and
 3 then applied to *SVM train* and *SVM test* samples. The harmonized *SVM train* sample was then
 4 used for training a linear SVM classifier to predict the scanner and the classifier performance
 5 assessed on harmonized *SVM test* sample. The entire process was repeated 50 times and an
 6 additional 50 repeats were performed for permutation testing (i.e., a total of 100 repeats for
 7 permutation testing). For this experiment, we learnt the regression coefficients (for the effect
 8 of age, TIV, and sex) from the *SVM train* samples and applied the coefficients to *SVM test*
 9 samples. Similarly, the standardization parameters were learnt from the *SVM train* samples and
 10 applied to the *SVM test* samples. Similar to experiment 2, we performed experiment 3 by taking
 11 a pair of scanners at a time (10 combinations), three scanners at a time (10 combinations), four
 12 scanners at a time (five combinations), and all five scanners at the same time.



14 **Figure 2:** Pipeline implemented in experiment 3: we trained a linear SVM classifier to predict the scanner after
 15 using different samples sizes to achieve harmonization of structural features; first, we performed a 10-fold split
 16 on the data resulting in 30 samples per scanner (SVM Test) and 270 samples per scanner. The 270 samples were
 17 next split into 200 samples per scanner (NH learn) and 70 samples per scanner (SVM Train). For every sample

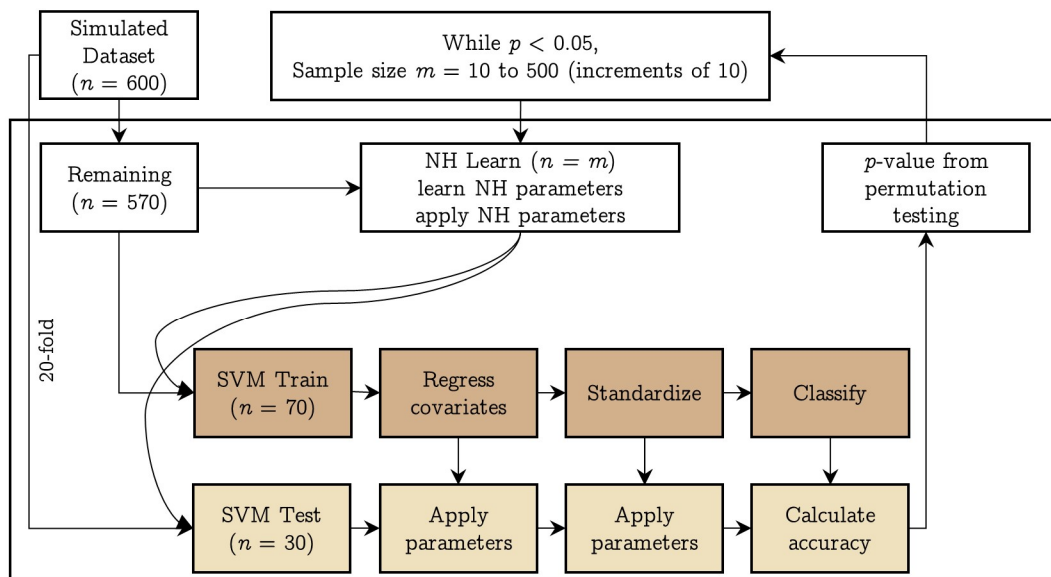
1 size 10 to 200, at increments of 10, we learnt the harmonization parameters using NH learn and applied it to SVM
2 Train and SVM Test samples. Then, after regressing the effect of age, TIV, and sex, we standardized the SVM
3 Train data (and applied the regression and standardization parameters to SVM Test) and trained a linear SVM
4 classifier to predict the scanner. Model performance was assessed on SVM Test dataset. The 10-fold cross-
5 validation was repeated 50 times and an additional 50 repeats of permutation testing (i.e., 100 repeats of
6 permutation) were performed to assess whether the classification performance was above chance level. [color
7 version of this figure is available online]

8 **2.7.5 Experiment 4: simulation-based approach**

9 We created 600 simulated samples per site by sampling from a multivariate normal distribution
10 with means and covariances from the fractal dimension features from the actual data.
11 Specifically, since we had five sites, we took their means and covariances and combined them
12 to have 25 different combinations of means and covariances; using these, we obtained 25
13 simulated datasets having 600 samples in each site. This strategy allowed us to get datasets
14 covering a wide range of MDs between pairs of sites. We then repeated the learning curve
15 experiment (experiment 3) using these simulated datasets with the difference that the k -fold
16 split was 20 fold, resulting in 30 samples per site for *SVM Test*, a holdout of 70 samples per
17 site for *SVM Train*, and 500 samples for *NH Learn* (which we iterated over from 10 samples
18 per site to 500 samples per site in increments of 10, stopping when the *SVM Test* performance
19 was not above chance). The design of this experiment is shown in **Figure 3**.

20 Having obtained the simulated data for 25 sites, we took combinations of two, three, and four
21 sites taken at a time and ran experiment 3 with the data splits as mentioned above. In this
22 simulation-based experiment, we did not attempt to simulate the covariates, in view of the
23 complex multivariate relationship within the covariates and between the features and the
24 covariates. Given the large number of possible combinations and the computational costs
25 involved, we only ran the experiment by randomly² selecting 100 two-, three-, and four-site
26 combinations.

² We randomly selected 15 MDs ≤ 0.5 , 70 MDs between 0.5 and 1.0, and 15 MDs > 1.0



All sample sizes are per scanner

1

2 **Figure 3:** Pipeline implemented in experiment 4: using simulated data, we trained a linear SVM classifier to
3 predict the scanner after using different samples sizes to achieve harmonization of structural features; we
4 performed a 20-fold split on the data resulting in 30 samples per scanner (SVM Test) and 570 remaining samples
5 per scanner. The 570 samples were next split into 500 samples per scanner (NH learn) and 70 samples per scanner
6 (SVM Train). For every sample size 10 to 500, at increments of 10, we learnt the harmonization parameters using
7 NH learn and applied it to SVM Train and SVM Test samples. Then, after regressing the effect of age, TIV, and
8 sex, we standardized the SVM Train data (and applied the regression and standardization parameters to SVM
9 Test) and trained a linear SVM classifier to predict the scanner. Model performance was assessed on SVM Test
10 dataset. The 20-fold cross-validation was repeated 50 times and an additional 50 repeats of permutation testing
11 (i.e., 100 repeats of permutation) were performed to assess whether the classification performance was above
12 chance level. The whole process was repeated for every sample size in NH learn till the classification performance
13 was above chance level. [color version of this figure is available online]

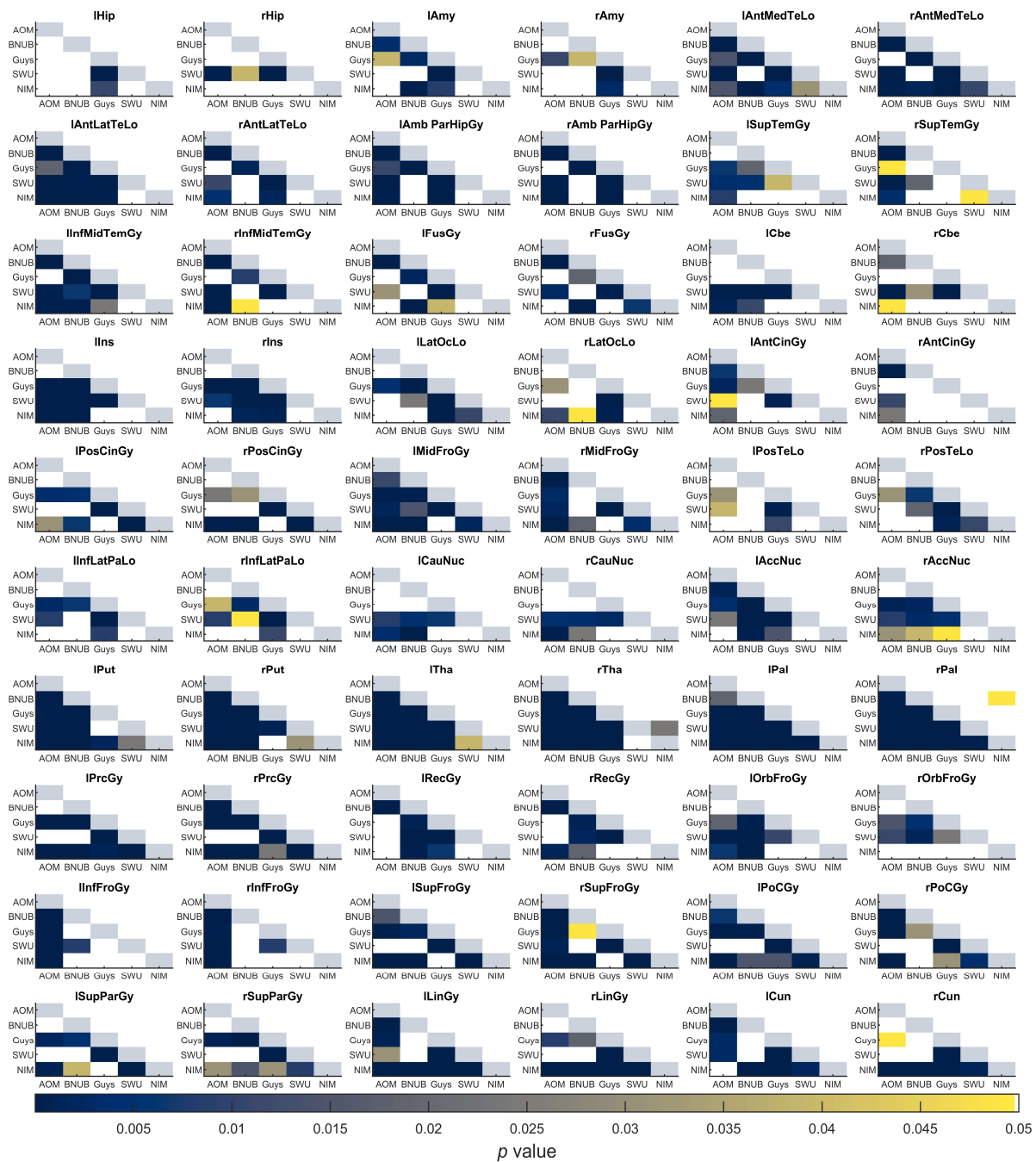
14 2.8 Data/code availability statement

15 The neuroimaging data used for all sites (except NIMHANS) is publicly available. The code
16 for the experiments is available at <https://github.com/parekhpravesh/HarmonizationPaper>. The
17 demographics and feature information for the NIMHANS dataset can be requested from the
18 corresponding authors on presentation of a reasonable data analysis request. The other datasets
19 used in this study are publicly available.

1 3. Results

2 3.1 Experiment 1: univariate difference

3 For every ROI, we compared the distribution of volumes and surface features between pairs of
4 scanners before and after harmonization using a two-sample KS test at $\alpha = 0.05$ (after
5 regressing the linear effect of age, TIV, and sex). Before harmonization, several of the gray
6 matter volumes across ROIs were statistically different between pairs of scanners; after
7 harmonization, most of these were not significantly different, except for the right pallidum
8 (BNUBeijing vs. NIMHANS, $p = 0.0491$) and right thalamus (SWU vs. NIMHANS, $p =$
9 0.0243) (see **Figure 4**). For cortical thickness, most ROIs were statistically different before
10 harmonization; after harmonization, only left pericalcarine region was statistically significant
11 (Guys vs. NIMHANS, $p = 0.0113$) (see Figure S1). For fractal dimension, several ROIs were
12 significantly different before harmonization; after harmonization, the fractal dimension of only
13 the right frontal pole was statistically significant (AOMIC vs. NIMHANS, $p = 0.0391$) (see
14 Figure S2). For sulcal depth, several ROIs were significantly different before harmonization
15 but after harmonization, none of the ROIs showed a statistically significant difference (see
16 Figure S3). For gyrification index, several ROIs were significantly different before
17 harmonization; after harmonization, the left pars orbitalis showed a statistically significant
18 difference (AOMIC vs. BNUBeijing, $p = 0.0309$ and AOMIC vs. SWU, $p = 0.0189$) (see Figure
19 S4). Additionally, we note that AOMIC vs. Guys and Guys vs. NIMHANS did not have many
20 significantly different ROIs (less than 10) before harmonization for fractal dimension, sulcal
21 depth, and gyrification index. Similarly, BNUBeijing vs. SWU only had four significantly
22 different ROIs (fractal dimension). The overall number of significantly different ROIs before
23 and after harmonization across feature categories is summarized in **Table 3**.



1

2 **Figure 4:** Summary of p -values from two sample Kolmogorov-Smirnov (KS) test between pairs of scanners for
 3 gray matter volumes. Each sub-plot indicates the p -values before (lower triangle) and after harmonization (upper
 4 triangle) between all pairs of scanners; the diagonal elements are shaded in a constant color to help distinguish
 5 lower and upper triangles. Each cell is color coded based on their p -value and only values smaller than 0.05 are
 6 shown. See Table S1 for the full names of the ROIs. Note that AOMIC dataset has been abbreviated to “AOM”,
 7 BNUBeijing dataset has been abbreviated to “BNUB”, and “NIMHANS” dataset has been abbreviated to “NIM”.
 8 [color version of this figure is available online]

9

10 **Table 3:** Number of statistically significant ROIs before and after harmonization using a two sample Kolmogorov-
 11 Smirnov (KS) test between pairs of scanners for different feature categories at $\alpha = 0.05$

Scanner	Grey matter volumes	Cortical thickness	Fractal dimension	Sulcal depth	Gyrification index
---------	---------------------	--------------------	-------------------	--------------	--------------------

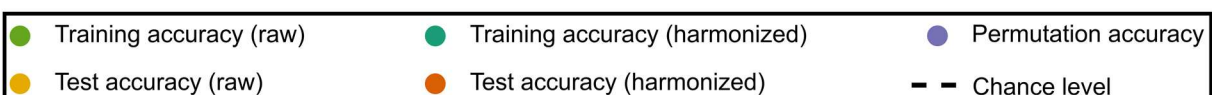
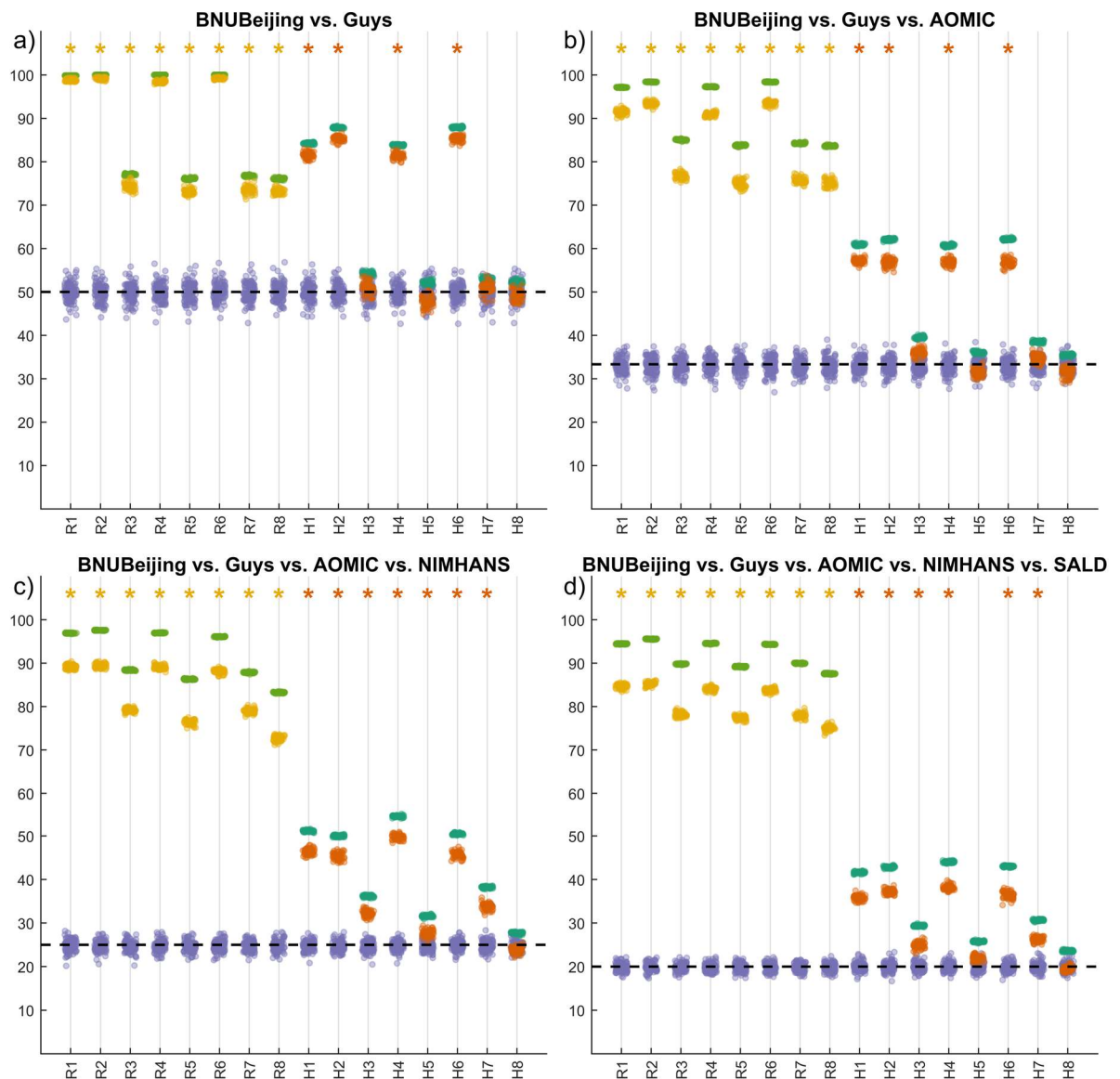
	# BH	# AH	# BH	# AH	# BH	# AH	# BH	# AH	# BH	# AH
AOMIC – BNUBeijing	35	0	53	0	44	0	59	0	36	1
AOMIC – Guys	44	0	51	0	2	0	9	0	6	0
AOMIC – SWU	41	0	32	0	33	0	45	0	30	1
AOMIC – NIMHANS	44	0	29	0	23	1	31	0	13	0
BNUBeijing – Guys	44	0	38	0	19	0	40	0	15	0
BNUBeijing – SWU	28	0	35	0	4	0	18	0	13	0
BNUBeijing – NIMHANS	40	1	51	0	21	0	43	0	31	0
Guys – SWU	56	0	55	0	35	0	40	0	36	0
Guys – NIMHANS	35	0	47	1	8	0	8	0	7	0
SWU – NIMHANS	27	1	41	0	24	0	27	0	27	0

1 **BH:** before harmonization; **AH:** after harmonization

2 3.2 Experiment 2: multivariate differences

3 For volumetric features, for all site combinations, before harmonization, the SVM models were
4 always able to predict the sites above chance level, irrespective of which covariates were
5 regressed. When harmonization was performed (within cross-validation), for all site
6 combinations where no regression of covariates was performed, or when TIV alone was
7 regressed, or when TIV and sex were regressed, the SVM model predictions were above chance
8 level. When only sex was regressed, only AOMIC vs. BNUBeijing SVM model predictions
9 were not statistically above chance level. When age alone, or TIV and age, or age and sex were
10 regressed, SVM prediction for certain site combinations remained statistically significant. Only
11 when TIV, age, and sex were regressed, the SVM predictions for all site combinations were
12 statistically not significant. An example plot showing these accuracies before and after
13 harmonization for all categories of covariate regression (for volumetric features) is shown in
14 **Figure 5**. Overall, all harmonized SVM accuracies were lower than before harmonization,
15 irrespective of which covariates were regressed. This indicates that harmonization does remove

1 site effects; however, it is important to additionally regress the confounding variables of TIV,
2 age, and sex after harmonization to eliminate site-effects completely. For surface features, we
3 saw a similar trend, albeit some differences (see supplementary material); the overall trend of
4 SVM accuracies being non-significant after regression of TIV, age, and sex was consistent
5 across the four feature categories – cortical thickness, fractal dimensions, sulcal depth, and
6 gyrification index.



Regression models – raw (R)/harmonized (H):

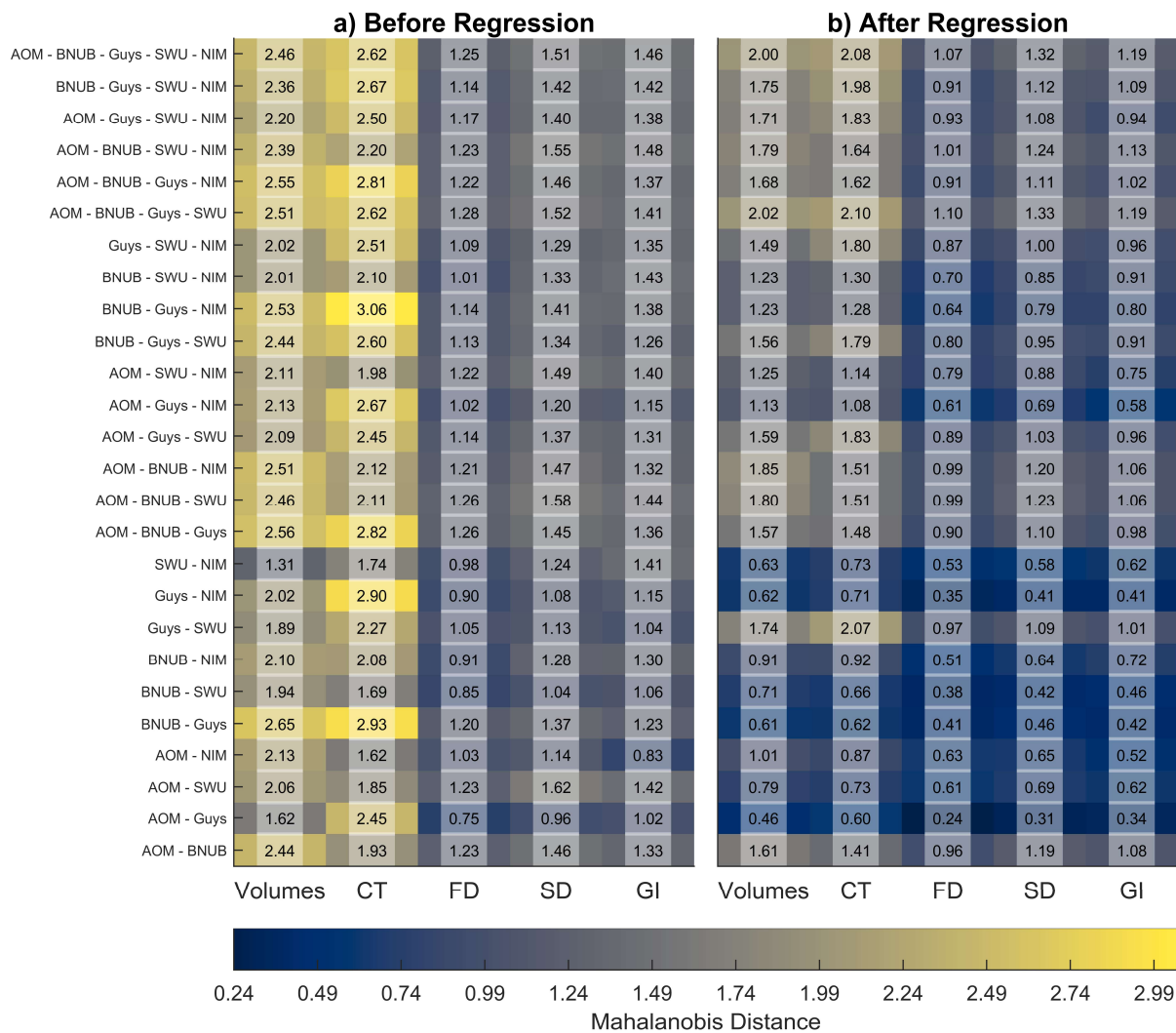
- | | | | |
|------------------|----------------|---------------------|--------------------------|
| 1. No regression | 3. Regress TIV | 5. Regress age, TIV | 7. Regress TIV, sex |
| 2. Regress age | 4. Regress sex | 6. Regress age, sex | 8. Regress age, TIV, sex |

1

2 **Figure 5:** Summary of experiment 2 for representative cases for a) two sites, b) three sites, c) four sites, and d)
 3 five sites taken at a time; The x-axis indicates the model type – raw (R) data and harmonized (H) data with different
 4 combinations of covariates being regressed while the y-axis indicates the 10-fold cross-validated percentage
 5 accuracy of SVM classifier; the training and test accuracy points are across 50 repeats of 10-fold cross-validation
 6 while the permutation accuracy points are across 100 repeats of 10-fold cross-validation; the asterisk mark
 7 indicates models where the permutation testing p -value was less than 0.05; the theoretical chance level
 8 accuracy is indicated with a dashed black line [color version of this figure is available online]

1 3.3 Quantifying the multivariate site-effect

2 We used the MD to quantify the site-effect before and after regression. These results are
3 summarized in **Figure 6**. Overall, we observed a reduction in MD after the regression of
4 confounding variables of age, TIV, and sex. Further, for fractal dimension, sulcal depth, and
5 gyrification index, we observed that the MD were smaller than volumetric and cortical
6 thickness features. We observed the smallest effect sizes in SWU vs. NIMHANS (volumes),
7 AOMIC vs. NIMHANS (cortical thickness and gyrification index), AOMIC vs. Guys (fractal
8 dimension and sulcal depth) before regression; after regressing of the covariates, the smallest
9 effect sizes were in AOMIC vs. Guys (for all feature categories). The largest effect sizes were
10 seen in BNUBeijing vs. Guys (volumes), BNUBeijing vs. Guys vs. NIMHANS (cortical
11 thickness), AOMIC vs. BNUBeijing vs. Guys vs. SWU (fractal dimension), AOMIC vs. SWU
12 (sulcal depth), and AOMIC vs. BNUBeijing vs. Guys vs. SWU vs. NIMHANS (gyrification
13 index) before regression; after regression, the largest effect sizes were seen in AOMIC vs.
14 BNUBeijing vs. Guys vs. SWU (volumes, cortical thickness, fractal dimension, and sulcal
15 depth), and in AOMIC vs. BNUBeijing vs. Guys vs. SWU vs. NIMHANS (gyrification index).



1

2 **Figure 6:** Average Mahalanobis distances between combinations of sites before and after regression of age, TIV,
 3 and sex for raw data; for any site combination, we first created a reference distribution using the overall mean and
 4 the pooled covariance; then, we calculated the distances of each site from this reference distribution and
 5 summarized it as the overall average; the x-axes indicate the different feature categories – grey matter volumes,
 6 cortical thickness (CT), fractal dimension (FD), sulcal depth (SD), and gyrification index (GI). Note that the
 7 AOMIC dataset has been abbreviated to “AOM,” BNUBeijing dataset has been abbreviated to “BNUB”, and the
 8 NIMHANS dataset has been abbreviated to “NIM”. [color version of this figure is available online]

9 3.4 Experiment 3: learning curves

10 When examining learning curves, we looked for the sample size at which the SVM classifier
 11 performance was no different than chance level (“convergence”) – we considered the sample
 12 size at which we first observed $p \geq 0.05$ as the required minimum sample size to remove site-
 13 effects completely. For volumetric features, for two-site combinations, all combinations except
 14 AOMIC vs. BNUBeijing and Guys vs. SWU converged and the sample size required ranged
 15 from 120 per site (Guys vs. NIMHANS) to 180 per site (AOMIC vs. NIMHANS). When

1 considering more than two-site combinations, we did not see convergence up to 200 samples
2 per site (see Figures S5 and S6). The two-site results for volumetric data is summarized in
3 **Figure 7**.

4 For cortical thickness, for two-site combinations, all combinations except BNUBeijing vs.
5 NIMHANS and Guys vs. SWU converged (see Figure S7). The sample size required ranged
6 from 110 per site (Guys vs. NIMHANS) to 200 per site (AOMIC vs. BNUBeijing). We did not
7 see convergence up to 200 samples per site when considering more than two-site combinations
8 (see Figures S8 and S9).

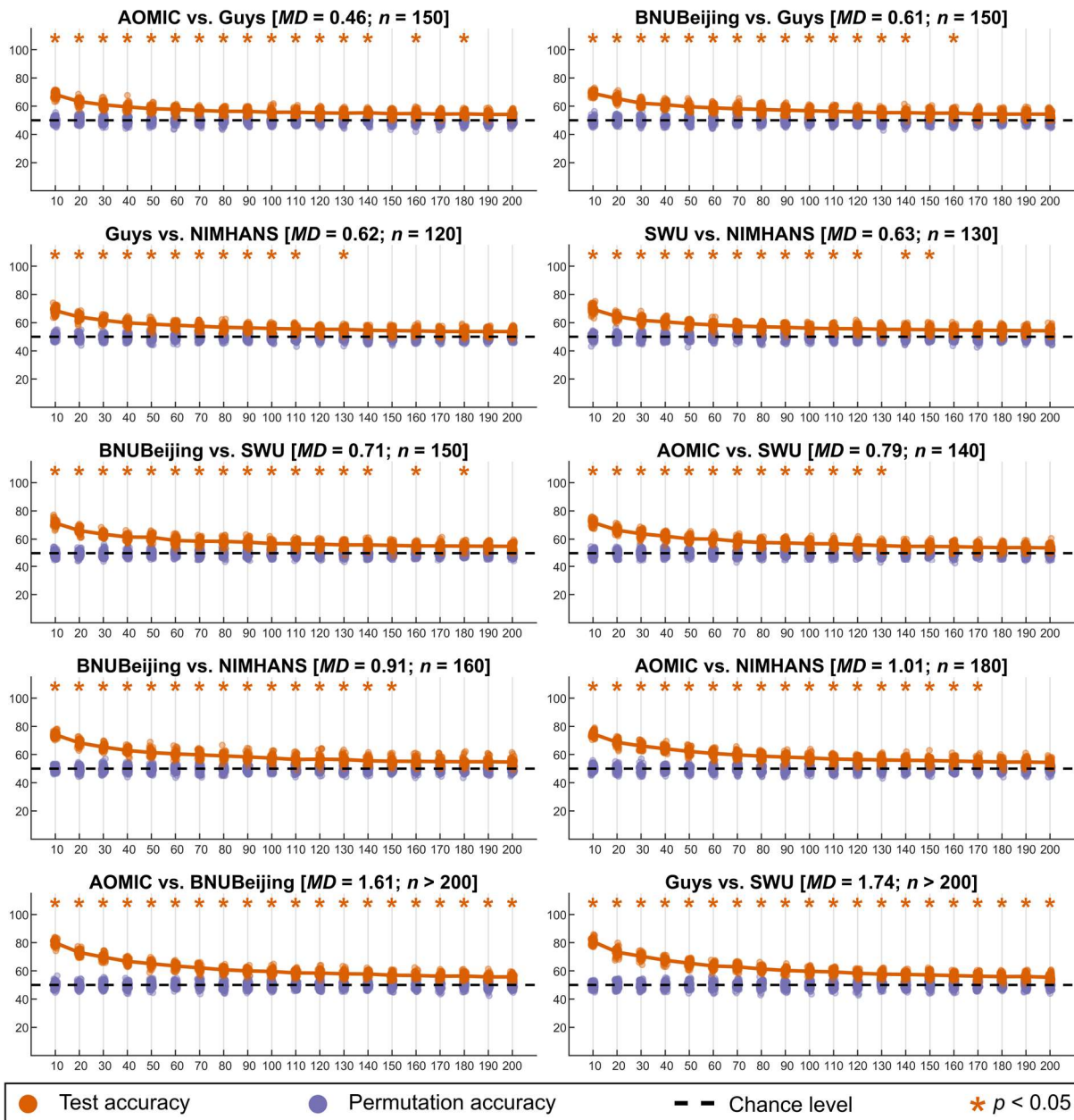
9 For fractal dimension, for two-site combinations, all combinations converged, and the sample
10 size required ranged from 50 samples per site (AOMIC vs. Guys) to 200 samples per site (Guys
11 vs. SWU) (see Figure S10). For three-site combinations, we observed convergence for AOMIC
12 vs. Guys vs. NIMHANS (160 samples per site), BNUBeijing vs. SWU vs. NIMHANS (190
13 samples per site), and BNUBeijing vs. Guys vs. NIMHANS (200 samples per site) (see Figure
14 S11). The other multi-site combinations did not show convergence for up to 200 samples per
15 site (see Figure S12).

16 For sulcal depth, for two-site combinations, all combinations except Guys vs. SWU converged,
17 and the required sample sizes ranged between 70 samples per site (AOMIC vs. Guys) and 160
18 samples per site (AOMIC vs. SWU and AOMIC vs. BNUBeijing) (see Figure S13). For three-
19 site combinations, only AOMIC vs. Guys vs. NIMHANS converged (180 samples per site)
20 (see Figure S14); all other multi-site combinations did not show convergence for up to 200
21 samples per site (see Figure S15).

22 For the gyrification index, for two-site combinations, all combinations except Guys vs. SWU
23 converged, and the required sample sizes ranged between 70 samples per site (AOMIC vs.
24 Guys) and 200 (AOMIC vs. BNUBeijing) (see Figure S16). For three-site combinations,

1 AOMIC vs. BNUBeijing vs. Guys showed convergence at 200 samples per site (see Figure
2 S17); all other multi-site combinations did not show convergence for up to 200 samples per
3 site (see Figure S18).

4 In general, for every site-combination, we observed that the minimum sample size required for
5 convergence increased with increasing average MDs. We observed that low MDs generally led
6 to lower required sample size (for example, across all feature categories, the lowest MDs was
7 ~0.24 for fractal dimension between AOMIC and Guys, which showed a convergence at a mere
8 50 samples per site). We calculated the correlation between the average MDs (across all feature
9 categories) and the minimum sample sizes required per site (excluding the site combinations
10 which did not converge) and found a strong association between these ($r = 0.82, p < 0.00$).



1

2

3

4

5

6

7

8

9

10

11

12

13

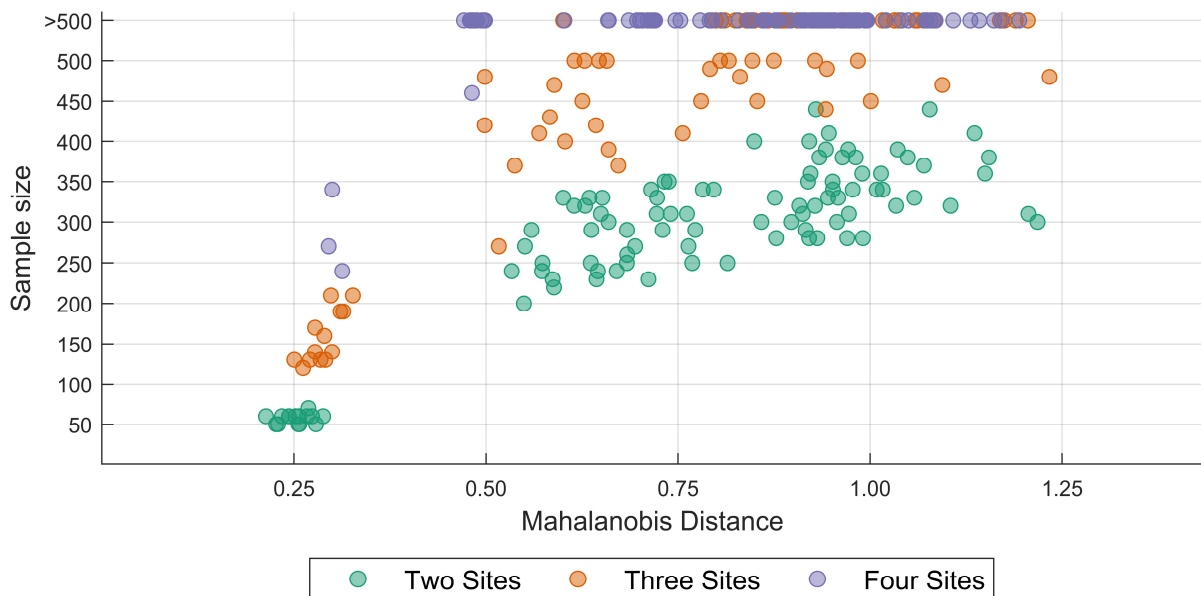
Figure 7: Summary of learning curves for volumetric features for two-site combinations; the orange points indicate the test accuracy of the SVM classifier (50 repeats of 10-fold cross-validation), the purple points indicate the permutation test accuracy of the SVM classifier (100 repeats of 10-fold cross-validation), while the dashed black line indicates the theoretical chance accuracy level; the x -axis indicates the sample size used for learning harmonization parameters (“*NHLearn*”) while the y -axis indicates the test accuracy in percentage. The title of each figure indicates the site-combinations, the average Mahalanobis distance (MD) of the two sites from the reference, and the sample size required for learning harmonization parameter (n) such that the SVM classifier performance was no different than chance level; the accuracies that were above chance are marked with an orange asterisk mark [color version of this figure is available online]

3.5 Experiment 4: simulation-based approach

Similar to our learning curve experiment, when examining the sample size requirement for a

variety of MD s, we observed that the sample size requirement increased with increasing MD s.

1 When comparing the sample size requirement across a number of sites for comparable MDs,
2 we observed that the sample size requirement typically increased with an increasing number of
3 sites. A summary of the results from the simulation-based approach is shown in **Figure 8**.



4
5 **Figure 8:** Plot of the sample size required for achieving inter-site harmonization for a range of Mahalanobis
6 distances for two-, three-, and four-site scenarios; the features were simulated using means and covariances from
7 fractal dimension features from real data (see text for details) [color version of this figure is available online]

8 4. Discussion

9 This study aimed to estimate the minimum sample size required for achieving inter-site
10 harmonization of volumetric and surface measures. The first step in this quest was to establish
11 whether site effects are effectively removed after harmonization. The first experiment showed
12 that most measures across ROIs showed site-effects before harmonization; almost all the site-
13 effects were removed (in a univariate manner) after harmonization. The few ROIs which were
14 statistically different after harmonization are likely false positives as we did not correct for
15 multiple comparisons. However, it is not enough to show that site-effects are corrected in a
16 univariate manner. Therefore, we employed a machine learning approach to assess whether the
17 site effects are eliminated in a multivariate manner. The results of this experiment showed that
18 (in most cases) site effects are effectively removed only when performing harmonization and

1 regressing the confounding effect of TIV, age, and sex. This reveals an interesting aspect
2 related to the site-effect. Given that the TIV is calculated from the images, the TIV itself is
3 confounded by site-effects. If the TIV is “preserved” during harmonization, then some residual
4 site-effects will remain in the harmonized features. Thus, when TIV is regressed after
5 harmonization, there is better correction of site-effects. Similarly, if other covariates are
6 strikingly different between sites, the harmonized feature values will also show these effects.
7 Therefore, it becomes essential to regress the effect of additional confounding factors from the
8 harmonized data.

9 Next, when we examined the learning curves to find the minimum sample size needed for
10 effective removal of site-effects (after regressing the confounding effects of TIV, age, and sex
11 from harmonized data), we observed that the required sample size exceeded the available
12 sample size in several cases, especially when harmonizing features from more than two sites.
13 This is evident when examining the average MDs – as the value of MDs increases, the sample
14 size needed to remove site-effects increases (correlation between MDs and sample size = 0.82).
15 For example, ROI volumes and cortical thicknesses had, in general, larger values for MDs and
16 did not show convergence for more than two sites taken at a time. In contrast, fractal dimension,
17 on average, had smaller values for MDs and therefore showed convergence for a few of the
18 three site combinations.

19 We observed a similar trend with the simulation approach, where we observed convergence at
20 a sample size greater than 200 for larger MDs. This provides insights into why we did not see
21 convergence (at a maximum sample size of 200) for several cases in real data. This experiment
22 also reveals an interesting property about the sample size requirement – in general, for the same
23 MDs, the sample size requirement increases with an increasing number of sites. However, we
24 note that it is possible that our simulation approach overestimated the sample size requirement
25 given that we did not simulate (and control for) the effect of covariates. In experiment 2, we

1 showed that the site-effect reduces when covariates are controlled for – since we did not
2 simulate and control the covariate in the simulated data, the uncorrected residual site-effects
3 might have led to an overestimate of the sample size requirement.

4 4.1 Recommendations and future directions

5 **Mahalanobis distance:** In this study, we have provided a framework to examine the site-
6 effects as the MD between each site and the reference distribution. The reference distribution
7 is created using the overall mean and pooled covariance matrix from the individual sites. When
8 considering more than two sites, we have computed the average of the MDs between each site
9 and the reference distribution. However, this average MD value is just a mid-point summary
10 statistic which will not entirely reflect the distances between the individual sites and the
11 reference. For example, for three sites, the MD between one site and the reference may be far
12 greater than the other two sites and the reference. Therefore, when considering the
13 harmonization of data from more than two sites, it may additionally be useful to examine the
14 range of the MD (i.e., the difference between the maximum MD and the minimum MD) as the
15 sample size requirement may not only be driven by the average MD. We note that (Zhang et
16 al., 2018) provide a data harmonization approach that allows one to specify a reference site.
17 When using this approach, the MD can, then, be calculated with respect to the reference site
18 rather than creating a reference distribution. In addition, future work should also examine the
19 relationship between the minimum sample size and other distance measures. For example,
20 distance measures which are bounded within a specific range (for example, the Jensen-Shannon
21 divergence) could be useful in generalizing the distances across any number of sites as this
22 might enable direct comparison of the distance measure independent of number of sites and
23 number of features.

24 **Number of sites and number of features:** In our experiments, we examined the minimum
25 sample size required for removing site effects for up to five site combinations. Further, we

1 restricted our analyses to 60 volumetric features and 68 surface features. However, it is possible
2 that as the number of sites and the number of features increase, the MD increases, and thereby
3 the minimum sample size required for achieving inter-site harmonization increases. This
4 would, of course, be dependent on the distribution of the features and the correlation among
5 them. Future work should explore the sample size requirement by densely sampling a 3-
6 dimensional grid of a number of sites \times number of features \times MDs.

7 **Effect of covariates:** In our second experiment, we examined the effect of preserving age,
8 TIV, and sex as covariates and then regressed the linear effect of these covariates after
9 harmonizing the data. These results showed that this regression step is essential to achieving
10 complete inter-site harmonization in a multivariate sense. We have previously remarked about
11 the relation between TIV and site effects. However, when simulating data, we have not
12 accounted for the covariates given the multivariate relationship within the covariates and
13 between the covariates and the features. The minimum sample size required to achieve inter-
14 site harmonization may change depending on the number of covariates preserved during
15 harmonization and may additionally depend on the type of the covariate
16 (categorical/continuous), along with the MDs, the number of sites, and the number of features.
17 Further, when harmonizing data across groups of subjects (for example, healthy subjects and
18 patients), it might be better to have representative samples of all the groups in the
19 harmonization training set and add the group effect as a covariate to be preserved.

20 **Cross-validation:** In all experiments (experiment 2-4), all the steps (including regression of
21 covariates and standardization of data) were performed within a cross-validation framework
22 (see **Figure 1** and **Figure 2** for experimental designs). This is an often overlooked step where
23 confound correction is performed on the entire data before any machine learning, leading to
24 information leakage (see, for example, (Snoek et al., 2019)). Similarly, other steps, including
25 data harmonization and standardization, should be performed within a cross-validation

1 framework to prevent any information leakage. Therefore, the design of the experiment
2 (including the number of folds and the type of cross-validation) is a critical factor. It is
3 important to ensure that the training sample is representative of the actual dataset; in addition,
4 our results show that the number of folds is an important consideration as it will directly control
5 the number of samples available for learning harmonization parameters.

6 **Alternate methods for harmonization:** As mentioned before, in addition to ComBat, other
7 extensions to ComBat like CovBat (Chen et al., 2021) have been proposed. Similarly, recently
8 proposed methods like NeuroHarmony (Garcia-Dias et al., 2020) need to be tested and
9 evaluated in terms of their sample size requirement for achieving inter-site harmonization.
10 Additionally, it would be interesting to explore the use of site-specific regression of covariates
11 within cross-validation and mixed-effects modeling (where covariates like TIV, age, and sex
12 are fixed-effects and site is a random-effect) as alternate methods for data harmonization.

13 **Conclusion:**

14 For multi-site studies that involve any form of cross-validation, it is important to carefully
15 design the experiment such that there is enough number of samples (per site) available in the
16 training dataset to learn the harmonization parameters adequately. Examples of such situations
17 include machine learning classification and prediction studies and model-building studies to
18 apply the model to new data. In this study, we have provided a framework utilizing
19 Mahalanobis distance to quantify the site effects. Through real data and simulations, we have
20 shown the estimated minimum sample size required to remove site effects completely. We have
21 attempted to provide some rules of thumb for this sample size requirement under different
22 circumstances (see **Figure 8**). However, our work indicates that further research needs to be
23 carried out in this area while accounting for various previously enlisted factors.

1 5. Acknowledgments

2 This work is funded by the Department of Biotechnology, Government of India by grant
3 number: BT/PR17316/MED/31/326/2015. The images in the “NIMHANS” dataset were
4 acquired with funding support from Department of Science and Technology, Government of
5 India by grant number: DST/JPJ/CSI/043 (to JPJ) and the Department of Biotechnology –
6 Wellcome Trust India Alliance Senior Fellowship Grant (500236/Z/11/Z) (to GV). G.V.
7 acknowledges the support of Department of Biotechnology (DBT)-Wellcome Trust India
8 Alliance (IA/CRC/19/1/610005) and Department of Biotechnology, Government of India
9 (BT/HRD-NBA-NWB/38/2019-20(6)). Several of the figures use the `tight_subplot` function
10 (Pekka Kumpulainen (2020). `tight_subplot(Nh, Nw, gap, marg_h, marg_w)`
11 ([https://www.mathworks.com/matlabcentral/fileexchange/27991-tight_subplot-nh-nw-gap-](https://www.mathworks.com/matlabcentral/fileexchange/27991-tight_subplot-nh-nw-gap-marg_h-marg_w)
12 `marg_h-marg_w`), MATLAB Central File Exchange. Retrieved January 27, 2020). The *cividis*
13 color scheme used in **Figure 6** is from (Nuñez et al., 2018). The color schemes of several of
14 the figures are from ColorBrewer 2.0 (<http://colorbrewer2.org/>) by Cynthia A. Brewer,
15 Geography, Pennsylvania State University (accessed 25-Oct-2019).

16 6. Author contributions

17 **Pravesh Parekh** and **Gaurav Vivek Bhalerao**: Conceptualization, Methodology, Software,
18 Formal analysis, Data Curation, Writing - Original Draft, Writing - Review & Editing,
19 Visualization

20 **The ADBS Consortium**: Computing resources

21 **John P. John** and **Ganesan Venkatasubramanian**: Resources, Writing - Review & Editing,
22 Supervision, Project administration, Funding acquisition

1 7. Ethics statement

2 The NIMHANS dataset was acquired as part of two research projects which received ethical
3 clearance from the Institute Ethics Committee, prior to data collection. No additional ethical
4 clearance was requested as the other datasets are already publicly available.

5 8. Declaration of interest

6 None

7 9. Role of funding agency

8 The funding agency had no role in study design, collection, analysis, and interpretation of data,
9 writing the report, and deciding to submit the article for publication.

10 10. References

- 11 1. Ardekani, B.A., 2018. A New Approach to Symmetric Registration of Longitudinal Structural MRI of the
12 Human Brain. *bioRxiv*. <https://doi.org/10.1101/306811>
- 13 2. Ardekani, B.A., Bachman, A.H., 2009. Model-based automatic detection of the anterior and posterior
14 commissures on MRI scans. *NeuroImage* 46, 677–682. <https://doi.org/10.1016/j.neuroimage.2009.02.030>
- 15 3. Ardekani, B.A., Kershaw, J., Braun, M., Kanuo, I., 1997. Automatic detection of the mid-sagittal plane in 3-
16 D brain images. *IEEE Transactions on Medical Imaging* 16, 947–952. <https://doi.org/10.1109/42.650892>
- 17 4. Beer, J.C., Tustison, N.J., Cook, P.A., Davatzikos, C., Sheline, Y.I., Shinohara, R.T., Linn, K.A., 2020.
18 Longitudinal ComBat: A method for harmonizing longitudinal multi-scanner imaging data. *NeuroImage* 220,
19 117129. <https://doi.org/10.1016/j.neuroimage.2020.117129>
- 20 5. Biswal, B.B., Mennes, M., Zuo, X.-N., Gohel, S., Kelly, C., Smith, S.M., Beckmann, C.F., Adelstein, J.S.,
21 Buckner, R.L., Colcombe, S., Dogonowski, A.-M., Ernst, M., Fair, D., Hampson, M., Hoptman, M.J., Hyde,
22 J.S., Kiviniemi, V.J., Kötter, R., Li, S.-J., Lin, C.-P., Lowe, M.J., Mackay, C., Madden, D.J., Madsen, K.H.,
23 Margulies, D.S., Mayberg, H.S., McMahon, K., Monk, C.S., Mostofsky, S.H., Nagel, B.J., Pekar, J.J., Peltier,
24 S.J., Petersen, S.E., Riedl, V., Rombouts, S.A.R.B., Rypma, B., Schlaggar, B.L., Schmidt, S., Seidler, R.D.,
25 Siegle, G.J., Sorg, C., Teng, G.-J., Veijola, J., Villringer, A., Walter, M., Wang, L., Weng, X.-C., Whitfield-
26 Gabrieli, S., Williamson, P., Windischberger, C., Zang, Y.-F., Zhang, H.-Y., Castellanos, F.X., Milham, M.P.,
27 2010. Toward discovery science of human brain function. *PNAS* 107, 4734–4739.
28 <https://doi.org/10.1073/pnas.0911855107>
- 29 6. Bruin, W., Denys, D., van Wingen, G., 2019. Diagnostic neuroimaging markers of obsessive-compulsive
30 disorder: Initial evidence from structural and functional MRI studies. *Progress in Neuro-
31 Psychopharmacology and Biological Psychiatry, Promising neural biomarkers and predictors of treatment
32 outcomes for psychiatric disorders: Novel neuroimaging approaches* 91, 49–59.
33 <https://doi.org/10.1016/j.pnpbp.2018.08.005>
- 34 7. Chen, A.A., Beer, J.C., Tustison, N.J., Cook, P.A., Shinohara, R.T., Shou, H., Initiative, T.A.D.N., 2021.
35 Mitigating site effects in covariance for machine learning in neuroimaging data. *Human Brain Mapping* n/a.
36 <https://doi.org/10.1002/hbm.25688>

- 1 8. Cohen, J., 1988. *Statistical Power Analysis for the Behavioral Sciences*, Second. ed. Routledge Academic,
2 New York, NY. <https://doi.org/10.1016/C2013-0-10517-X>
- 3 9. Del Giudice, M., 2009. On the Real Magnitude of Psychological Sex Differences. *Evol Psychol* 7,
4 147470490900700220. <https://doi.org/10.1177/147470490900700209>
- 5 10. Desikan, R.S., Ségonne, F., Fischl, B., Quinn, B.T., Dickerson, B.C., Blacker, D., Buckner, R.L., Dale, A.M.,
6 Maguire, R.P., Hyman, B.T., Albert, M.S., Killiany, R.J., 2006. An automated labeling system for subdividing
7 the human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage* 31, 968–980.
8 <https://doi.org/10.1016/j.neuroimage.2006.01.021>
- 9 11. Faillenot, I., Heckemann, R.A., Frot, M., Hammers, A., 2017. Macroanatomy and 3D probabilistic atlas of
10 the human insula. *NeuroImage* 150, 88–98. <https://doi.org/10.1016/j.neuroimage.2017.01.073>
- 11 12. Fennema-Notestine, C., Gamst, A.C., Quinn, B.T., Pacheco, J., Jernigan, T.L., Thal, L., Buckner, R., Killiany,
12 R., Blacker, D., Dale, A.M., Fischl, B., Dickerson, B., Gollub, R.L., 2007. Feasibility of Multi-site Clinical
13 Structural Neuroimaging Studies of Aging Using Legacy Data. *Neuroinform* 5, 235–245.
14 <https://doi.org/10.1007/s12021-007-9003-9>
- 15 13. Fortin, J.-P., Cullen, N., Sheline, Y.I., Taylor, W.D., Aselcioglu, I., Cook, P.A., Adams, P., Cooper, C., Fava,
16 M., McGrath, P.J., McInnis, M., Phillips, M.L., Trivedi, M.H., Weissman, M.M., Shinohara, R.T., 2018.
17 Harmonization of cortical thickness measurements across scanners and sites. *NeuroImage* 167, 104–120.
18 <https://doi.org/10.1016/j.neuroimage.2017.11.024>
- 19 14. Fortin, J.-P., Parker, D., Tunç, B., Watanabe, T., Elliott, M.A., Ruparel, K., Roalf, D.R., Satterthwaite, T.D.,
20 Gur, R.C., Gur, R.E., Schultz, R.T., Verma, R., Shinohara, R.T., 2017. Harmonization of multi-site diffusion
21 tensor imaging data. *NeuroImage* 161, 149–170. <https://doi.org/10.1016/j.neuroimage.2017.08.047>
- 22 15. Garcia-Dias, R., Scarpazza, C., Baecker, L., Vieira, S., Pinaya, W.H.L., Corvin, A., Redolfi, A., Nelson, B.,
23 Crespo-Facorro, B., McDonald, C., Tordesillas-Gutiérrez, D., Cannon, D., Mothersill, D., Hernaus, D.,
24 Morris, D., Setien-Suero, E., Donohoe, G., Frisoni, G., Tronchin, G., Sato, J., Marcelis, M., Kempton, M.,
25 van Haren, N.E.M., Gruber, O., McGorry, P., Amminger, P., McGuire, P., Gong, Q., Kahn, R.S., Ayesa-
26 Arriola, R., van Amelsvoort, T., Ortiz-García de la Foz, V., Calhoun, V., Cahn, W., Mechelli, A., 2020.
27 Neuroharmony: A new tool for harmonizing volumetric MRI data from unseen scanners. *NeuroImage* 220,
28 117127. <https://doi.org/10.1016/j.neuroimage.2020.117127>
- 29 16. Gorgolewski, K.J., Auer, T., Calhoun, V.D., Craddock, R.C., Das, S., Duff, E.P., Flandin, G., Ghosh, S.S.,
30 Glatard, T., Halchenko, Y.O., Handwerker, D.A., Hanke, M., Keator, D., Li, X., Michael, Z., Maumet, C.,
31 Nichols, B.N., Nichols, T.E., Pellman, J., Poline, J.-B., Rokem, A., Schaefer, G., Sochat, V., Triplett, W.,
32 Turner, J.A., Varoquaux, G., Poldrack, R.A., 2016. The brain imaging data structure, a format for organizing
33 and describing outputs of neuroimaging experiments. *Scientific Data* 3, 160044.
34 <https://doi.org/10.1038/sdata.2016.44>
- 35 17. Gousias, I.S., Rueckert, D., Heckemann, R.A., Dyet, L.E., Boardman, J.P., Edwards, A.D., Hammers, A.,
36 2008. Automatic segmentation of brain MRIs of 2-year-olds into 83 regions of interest. *NeuroImage* 40, 672–
37 684. <https://doi.org/10.1016/j.neuroimage.2007.11.034>
- 38 18. Hammers, A., Allom, R., Koepp, M.J., Free, S.L., Myers, R., Lemieux, L., Mitchell, T.N., Brooks, D.J.,
39 Duncan, J.S., 2003. Three-dimensional maximum probability atlas of the human brain, with particular
40 reference to the temporal lobe. *Hum. Brain Mapp.* 19, 224–247. <https://doi.org/10.1002/hbm.10123>
- 41 19. Huang, L., Huang, T., Zhen, Z., Liu, J., 2016. A test-retest dataset for assessing long-term reliability of brain
42 morphology and resting-state brain activity. *Sci Data* 3, 160016. <https://doi.org/10.1038/sdata.2016.16>
- 43 20. Johnson, W.E., Li, C., Rabinovic, A., 2007. Adjusting batch effects in microarray expression data using
44 empirical Bayes methods. *Biostatistics* 8, 118–127. <https://doi.org/10.1093/biostatistics/kxj037>
- 45 21. Lee, H., Nakamura, K., Narayanan, S., Brown, R.A., Arnold, D.L., 2019. Estimating and accounting for the
46 effect of MRI scanner changes on longitudinal whole-brain volume change measurements. *NeuroImage* 184,
47 555–565. <https://doi.org/10.1016/j.neuroimage.2018.09.062>
- 48 22. Leek, J.T., Scharpf, R.B., Bravo, H.C., Simcha, D., Langmead, B., Johnson, W.E., Geman, D., Baggerly, K.,
49 Irizarry, R.A., 2010. Tackling the widespread and critical impact of batch effects in high-throughput data.
50 *Nat Rev Genet* 11, 733–739. <https://doi.org/10.1038/nrg2825>

- 1 23. Lin, Q., Dai, Z., Xia, M., Han, Z., Huang, R., Gong, G., Liu, C., Bi, Y., He, Y., 2015. A connectivity-based
2 test-retest dataset of multi-modal magnetic resonance imaging in young healthy adults. *Sci Data* 2, 150056.
3 <https://doi.org/10.1038/sdata.2015.56>
- 4 24. Liu, S., Hou, B., Zhang, Y., Lin, T., Fan, X., You, H., Feng, F., 2020. Inter-scanner reproducibility of brain
5 volumetry: influence of automated brain segmentation software. *BMC Neuroscience* 21, 35.
6 <https://doi.org/10.1186/s12868-020-00585-1>
- 7 25. Mahalanobis, P.C., 1936. On the generalized distance in statistics. *Proceedings of the National Institute of*
8 *Sciences (Calcutta)* 2, 49–55.
- 9 26. Markiewicz, C.J., Gorgolewski, K.J., Feingold, F., Blair, R., Halchenko, Y.O., Miller, E., Hardcastle, N.,
10 Wexler, J., Esteban, O., Goncavles, M., Jwa, A., Poldrack, R., 2021. The OpenNeuro resource for sharing of
11 neuroscience data. *eLife* 10, e71774. <https://doi.org/10.7554/eLife.71774>
- 12 27. Medawar, E., Thieleking, R., Manuilova, I., Paerisch, M., Villringer, A., Witte, A.V., Beyer, F., 2021.
13 Estimating the effect of a scanner upgrade on measures of grey matter structure for longitudinal designs.
14 *PLOS ONE* 16, e0239021. <https://doi.org/10.1371/journal.pone.0239021>
- 15 28. Nuñez, J.R., Anderton, C.R., Renslow, R.S., 2018. Optimizing colormaps with consideration for color vision
16 deficiency to enable accurate interpretation of scientific data. *PLOS ONE* 13, e0199239.
17 <https://doi.org/10.1371/journal.pone.0199239>
- 18 29. Ojala, M., Garriga, G.C., 2010. Permutation Tests for Studying Classifier Performance. *Journal of Machine*
19 *Learning Research* 11, 1833–1863.
- 20 30. Pardoe, H., Pell, G.S., Abbott, D.F., Berg, A.T., Jackson, G.D., 2008. Multi-site voxel-based morphometry:
21 Methods and a feasibility demonstration with childhood absence epilepsy. *NeuroImage* 42, 611–616.
22 <https://doi.org/10.1016/j.neuroimage.2008.05.007>
- 23 31. Pomponio, R., Erus, G., Habes, M., Doshi, J., Srinivasan, D., Mamourian, E., Bashyam, V., Nasrallah, I.M.,
24 Satterthwaite, T.D., Fan, Y., Launer, L.J., Masters, C.L., Maruff, P., Zhuo, C., Völzke, H., Johnson, S.C.,
25 Fripp, J., Koutsouleris, N., Wolf, D.H., Gur, Raquel, Gur, Ruben, Morris, J., Albert, M.S., Grabe, H.J.,
26 Resnick, S.M., Bryan, R.N., Wolk, D.A., Shinohara, R.T., Shou, H., Davatzikos, C., 2020. Harmonization of
27 large MRI datasets for the analysis of brain imaging patterns throughout the lifespan. *NeuroImage* 208,
28 116450. <https://doi.org/10.1016/j.neuroimage.2019.116450>
- 29 32. Rozycki, M., Satterthwaite, T.D., Koutsouleris, N., Erus, G., Doshi, J., Wolf, D.H., Fan, Y., Gur, R.E., Gur,
30 R.C., Meisenzahl, E.M., Zhuo, C., Yin, H., Yan, H., Yue, W., Zhang, D., Davatzikos, C., 2018. Multisite
31 Machine Learning Analysis Provides a Robust Structural Imaging Signature of Schizophrenia Detectable
32 Across Diverse Patient Populations and Within Individuals. *Schizophr Bull* 44, 1035–1044.
33 <https://doi.org/10.1093/schbul/sbx137>
- 34 33. Segall, J.M., Turner, J.A., van Erp, T.G.M., White, T., Bockholt, H.J., Gollub, R.L., Ho, B.C., Magnotta, V.,
35 Jung, R.E., McCarley, R.W., Schulz, S.C., Lauriello, J., Clark, V.P., Voyvodic, J.T., Diaz, M.T., Calhoun,
36 V.D., 2009. Voxel-based Morphometric Multisite Collaborative Study on Schizophrenia. *Schizophr Bull* 35,
37 82–95. <https://doi.org/10.1093/schbul/sbn150>
- 38 34. Snoek, L., Miletić, S., Scholte, H.S., 2019. How to control for confounds in decoding analyses of
39 neuroimaging data. *NeuroImage* 184, 741–760. <https://doi.org/10.1016/j.neuroimage.2018.09.074>
- 40 35. Snoek, L., van der Miesen, M.M., Beemsterboer, T., van der Leij, A., Eigenhuis, A., Scholte, S.H., 2021a.
41 The Amsterdam Open MRI Collection, a set of multimodal MRI datasets for individual difference analyses.
42 *Sci Data* 8, 85. <https://doi.org/10.1038/s41597-021-00870-6>
- 43 36. Snoek, L., van der Miesen, M.M., van der Leij, A., Beemsterboer, T., Eigenhuis, A., Scholte, S.H., 2021b.
44 AOMIC-ID1000.
- 45 37. Stonnington, C.M., Tan, G., Klöppel, S., Chu, C., Draganski, B., Jack, C.R., Chen, K., Ashburner, J.,
46 Frackowiak, R.S.J., 2008. Interpreting scan data acquired from multiple scanners: A study with Alzheimer's
47 disease. *NeuroImage* 39, 1180–1185. <https://doi.org/10.1016/j.neuroimage.2007.09.066>
- 48 38. Takao, H., Hayashi, N., Ohtomo, K., 2014. Effects of study design in multi-scanner voxel-based
49 morphometry studies. *NeuroImage* 84, 133–140. <https://doi.org/10.1016/j.neuroimage.2013.08.046>
- 50 39. Wei, D., Zhuang, K., Ai, L., Chen, Q., Yang, W., Liu, W., Wang, K., Sun, J., Qiu, J., 2018. Structural and
51 functional brain scans from the cross-sectional Southwest University adult lifespan dataset. *Sci Data* 5,
52 180134. <https://doi.org/10.1038/sdata.2018.134>

- 1 40. Wittens, M.M.J., Allemeersch, G.-J., Sima, D.M., Naeyaert, M., Vanderhasselt, T., Vanbinst, A.-M., Buls,
2 N., De Brucker, Y., Raeymaekers, H., Franssen, E., Smeets, D., van Hecke, W., Nagels, G., Bjerke, M., de
3 Mey, J., Engelborghs, S., 2021. Inter- and Intra-Scanner Variability of Automated Brain Volumetry on Three
4 Magnetic Resonance Imaging Systems in Alzheimer's Disease and Controls. *Frontiers in Aging Neuroscience*
5 13, 641. <https://doi.org/10.3389/fnagi.2021.746982>
- 6 41. Zavaliangos-Petropulu, A., Nir, T.M., Thomopoulos, S.I., Reid, R.I., Bernstein, M.A., Borowski, B., Jack Jr.,
7 C.R., Weiner, M.W., Jahanshad, N., Thompson, P.M., 2019. Diffusion MRI Indices and Their Relation to
8 Cognitive Impairment in Brain Aging: The Updated Multi-protocol Approach in ADNI3. *Frontiers in*
9 *Neuroinformatics* 13, 2. <https://doi.org/10.3389/fninf.2019.00002>
- 10 42. Zhang, Y., Jenkins, D.F., Manimaran, S., Johnson, W.E., 2018. Alternative empirical Bayes models for
11 adjusting for batch effects in genomic studies. *BMC Bioinformatics* 19, 262. [https://doi.org/10.1186/s12859-](https://doi.org/10.1186/s12859-018-2263-6)
12 [018-2263-6](https://doi.org/10.1186/s12859-018-2263-6)
- 13 43. Zuo, X.-N., Anderson, J.S., Bellec, P., Birn, R.M., Biswal, B.B., Blautzik, J., Breitner, J.C.S., Buckner, R.L.,
14 Calhoun, V.D., Castellanos, F.X., Chen, A., Chen, B., Chen, J., Chen, X., Colcombe, S.J., Courtney, W.,
15 Craddock, R.C., Di Martino, A., Dong, H.-M., Fu, X., Gong, Q., Gorgolewski, K.J., Han, Y., He, Ye, He,
16 Yong, Ho, E., Holmes, A., Hou, X.-H., Huckins, J., Jiang, T., Jiang, Y., Kelley, W., Kelly, C., King, M.,
17 LaConte, S.M., Lainhart, J.E., Lei, X., Li, H.-J., Li, Kaiming, Li, Kuncheng, Lin, Q., Liu, D., Liu, J., Liu, X.,
18 Liu, Y., Lu, G., Lu, J., Luna, B., Luo, J., Lurie, D., Mao, Y., Margulies, D.S., Mayer, A.R., Meindl, T.,
19 Meyerand, M.E., Nan, W., Nielsen, J.A., O'Connor, D., Paulsen, D., Prabhakaran, V., Qi, Z., Qiu, J., Shao,
20 C., Shehzad, Z., Tang, W., Villringer, A., Wang, H., Wang, K., Wei, D., Wei, G.-X., Weng, X.-C., Wu, X.,
21 Xu, T., Yang, N., Yang, Z., Zang, Y.-F., Zhang, L., Zhang, Q., Zhang, Zhe, Zhang, Zhiqiang, Zhao, K., Zhen,
22 Z., Zhou, Y., Zhu, X.-T., Milham, M.P., 2014. An open science resource for establishing reliability and
23 reproducibility in functional connectomics. *Sci Data* 1, 140049. <https://doi.org/10.1038/sdata.2014.49>
- 24