

1 Evaluation of ComBat harmonization for reducing across- 2 tracer differences in regional amyloid PET analyses

3 Running title: Evaluating ComBat for amyloid PET analyses

4 Author information

5 Braden Yang¹ (0000-0002-2558-4132), Tom Earnest¹ (0000-0001-8671-8424), Sayantan
6 Kumar¹ (0000-0001-7213-0734), Deydeep Kothapalli¹, Tammie Benzinger¹ (0000-0002-8114-
7 0552), Brian Gordon¹ (0000-0003-2109-2955), Aristeidis Sotiras^{1,2} (0000-0003-0795-8820)

8 ¹ Mallinckrodt Institute of Radiology, Washington University School of Medicine in St. Louis, St. Louis,
9 MO, USA 63110

10 ² Institute for Informatics, Data Science and Biostatistics, Washington University School of Medicine in St.
11 Louis, St. Louis, MO, USA 63110

12 Please address all correspondence to: Braden Yang (b.y.yang@wustl.edu), Aristeidis Sotiras
13 (aristeidis.sotiras@wustl.edu)

14 Abstract

15 Introduction

16 Differences in amyloid positron emission tomography (PET) radiotracer pharmacokinetics and
17 binding properties lead to discrepancies in amyloid- β uptake estimates. Harmonization of tracer-
18 specific biases is crucial for optimal performance of downstream tasks. Here, we investigated
19 the efficacy of ComBat, a data-driven harmonization model, for reducing tracer-specific biases
20 in regional amyloid PET measurements from [¹⁸F]-florbetapir (FBP) and [¹¹C]-Pittsburgh
21 Compound-B (PiB).

22 **Methods**

23 One-hundred-thirteen head-to-head FBP-PiB scan pairs, scanned from the same subject within
24 ninety days, were selected from the Open Access Series of Imaging Studies 3 (OASIS-3)
25 dataset. The Centiloid scale, ComBat with no covariates, ComBat with biological covariates, and
26 GAM-ComBat with biological covariates were used to harmonize both global and regional
27 amyloid standardized uptake value ratios (SUVR). Variants of ComBat, including longitudinal
28 ComBat and PEACE, were also tested. Intraclass correlation coefficient (ICC) and mean
29 absolute error (MAE) were computed to measure the absolute agreement between tracers.
30 Additionally, longitudinal amyloid SUVRs from an anti-amyloid drug trial were simulated using
31 linear mixed effects modeling. Differences in rates-of-change between simulated treatment and
32 placebo groups were tested, and change in statistical power/Type-I error after harmonization
33 was quantified.

34 **Results**

35 In the head-to-head tracer comparison, ComBat with no covariates was the best at increasing
36 ICC and decreasing MAE of both global summary and regional amyloid PET SUVRs between
37 scan pairs of the same group of subjects. In the clinical trial simulation, harmonization with both
38 Centiloid and ComBat increased statistical power of detecting true rate-of-change differences
39 between groups and decreased false discovery rate in the absence of a treatment effect. The
40 greatest benefit of harmonization was observed when groups exhibited differing FBP-to-PiB
41 proportions.

42 **Conclusion**

43 ComBat outperformed the Centiloid scale in harmonizing both global and regional amyloid
44 estimates. Additionally, ComBat improved the detection of rate-of-change differences between

45 clinical trial groups. Our findings suggest that ComBat is a viable alternative to Centiloid for
46 harmonizing regional amyloid PET analyses.

47 **Keywords**

48 Positron emission tomography, amyloid- β , harmonization, Centiloid, ComBat

49 **Introduction**

50 Positron emission tomography (PET) is widely used in clinical and research settings for
51 measuring and monitoring amyloid- β deposition *in vivo* in the brain for patients who are at risk of
52 developing or who already present with Alzheimer's disease (AD). In clinical trials for anti-
53 amyloid drugs, PET is an important tool for screening appropriate candidates who have
54 undergone significant amyloidosis in the brain [1]. Moreover, PET has also been used for
55 monitoring the progression of global amyloid burden longitudinally within these trials, which
56 along with measures of cognitive function serves as a crucial secondary endpoint [2,3]. In
57 research settings, PET is able to resolve the spatial distribution of amyloid within specific
58 regions of the brain, enabling the design of multivariable statistical analyses and predictive
59 models of AD using voxel-wise [4,5] or region-of-interest (ROI) based [6–8] PET biomarkers as
60 multidimensional features.

61 Several PET radiotracers for imaging brain amyloid pathology have been developed. The first
62 amyloid PET tracer developed for human imaging studies was [^{11}C]-Pittsburgh compound B
63 (PiB) [9], but due to its short half-life requires an on-site cyclotron to produce. Consequently, PiB
64 is not accessible by many sites and not appropriate for use in clinical trials. Alternatively,
65 amyloid measurements obtained from ^{18}F -based tracers such as [^{18}F]-florbetapir (FBP) [10–12],
66 [^{18}F]-florbetaben [13] and [^{18}F]-flutemetamol [14,15] have been shown to correlate well with PiB.

67 Coupled with a much longer half-life than PiB, these tracers are a much more suitable option for
68 clinical trials due to their accessibility and ability to be distributed off-site.

69 Nonetheless, previous studies that performed a head-to-head comparison of amyloid PET
70 tracers have demonstrated significant disparities in dynamic range and non-specific binding
71 properties between tracers [10,13,16]. Subsequently, this makes it difficult to compare
72 quantitative amyloid measurements between images acquired using different tracers. This may
73 also negatively impact the performance of downstream tasks such as detecting significant
74 treatment effects in anti-amyloid drug trials [17].

75 To address this, Klunk *et al.* introduced the Centiloid scale [18], which linearly transforms the
76 dynamic range of a global estimate of amyloid burden to a common scale and converts it to
77 Centiloid (CL) units. This involves calibrating the scale to a preselected cohort of amyloid-
78 negative healthy controls and amyloid-positive typical AD patients, where the average global
79 burden of the two groups are set to 0 CL and 100 CL, respectively. However, the calibration
80 process requires at least two PET scans from the same subject within a short time period in
81 order to calibrate conversion equations. Additionally, a single equation is usually derived to
82 operate on the global amyloid estimate, but this cannot address local disparities in amyloid PET
83 signal between tracers. Other methods for tracer harmonization that are based on data-driven
84 and/or machine learning techniques such as principal component analysis [19], non-negative
85 matrix factorization [20], and deep learning [21,22] have been proposed, but like Centiloid they
86 focus on the global amyloid burden.

87 Alternatively, ComBat [23] is a data-driven harmonization model which has been widely applied
88 in magnetic resonance imaging (MRI) analyses to adjust for differences in scanners and
89 acquisition protocols. It has been used to correct regional volume and cortical thickness
90 measurements from MRI [24–27], and has more recently been applied to [¹⁸F]-
91 fluorodeoxyglucose-PET [28] and amyloid PET [29] biomarkers. Much of the current literature

92 on applying ComBat has focused on reducing scanner-level and institutional-level biases.
93 However, it remains unclear whether ComBat is applicable for mitigating across-tracer variance,
94 specifically in regional amyloid PET measurements.

95 Here, we aimed to evaluate the efficacy of ComBat for harmonizing standardized uptake value
96 ratios (SUVR) from amyloid PET across two tracers - PiB and FBP. Specifically, we addressed
97 two primary inquiries. Firstly, we investigated whether ComBat harmonization may increase the
98 agreement between regional SUVRs obtained from the two tracers. This was accomplished
99 through a head-to-head comparison of PiB and FBP. We selected a set of PiB-FBP scan pairs
100 acquired from the same subject in a short time period and compared measures of the absolute
101 agreement between regional SUVRs before and after ComBat harmonization. Secondly, we
102 explored the utility of ComBat harmonization in the context of clinical tasks. This was examined
103 by simulating a multi-tracer anti-amyloid drug trial where two different amyloid tracers were used
104 to measure brain amyloid deposition, under the assumption that different sites have access to
105 different tracers. We generated longitudinal amyloid PET data of hypothetical treatment and
106 placebo groups with a known underlying treatment effect, and assigned each group a specific
107 proportion of PiB-to-FBP scans. We then gauged whether ComBat harmonization improves the
108 statistical power of detecting the underlying treatment effect when using two different tracers.

109 **Materials & Methods**

110 **Participants and data**

111 Data for this study were acquired from the Open Access Series of Imaging Studies 3 (OASIS-3)
112 dataset [30], which consisted of 1098 total participants and their longitudinal imaging data. Of
113 these, we selected 997 who underwent PiB and/or FBP amyloid PET imaging. All available PET
114 scans, including the initial baseline scan and any follow-up scans, were utilized in this study, for

115 a total of 678 FBP scans and 1157 PiB scans. Additionally, each subject's age at scan, sex and
116 apolipoprotein-ε4 (APOE) allele carriership were extracted. Subjects who were missing any of
117 these variables were excluded from further analyses.

118 **Image acquisition and processing**

119 All amyloid PET imaging from OASIS-3 were acquired at Washington University in St. Louis
120 using one of four Siemens scanner models: Biograph mMR PET/MR 3T, Biograph 40 PET/CT,
121 Biograph 128 Vision Edge PET/CT, and ECAT HR+ 962 PET. For PiB PET, participants
122 received a bolus injection of 6-20 mCi of PiB, and a 60-minute dynamic scan was acquired. For
123 FBP PET, participants received a bolus injection of 10 mCi of FBP, and either a 70-minute
124 dynamic scan was acquired, or a 20-minute dynamic scan was acquired at 50-minutes post-
125 injection. Additionally, T1-weighted MRI scans were acquired and utilized for PET processing.
126 All MRI imaging from the OASIS-3 dataset were acquired at Washington University in St. Louis
127 using one of three Siemens scanner models: Vision 1.5T, TIM Trio 3T, and Biograph mMR
128 PET/MR 3T.

129 PET images were processed using the PET Unified Pipeline (<https://github.com/ysu001/PUP>),
130 described in [31]. Briefly, raw PET images were smoothed to 8mm spatial resolution, corrected
131 for inter-frame motion, and coregistered to the T1 MRI scan acquired closest in time using a
132 vector-gradient algorithm. T1 images were segmented and parcellated into cortical and
133 subcortical ROIs using FreeSurfer 5.0 or 5.1 for 1.5T scans or FreeSurfer 5.3 for 3T scans. For
134 each ROI, regional SUVRs were computed from the peak time windows of each tracer (30-to-60
135 minutes post-injection for PiB, 50-to-70 minutes post-injection for FBP). The average of the left
136 and right cerebellar cortex was used as the reference region. Additionally, a summary estimate
137 of global amyloid burden was derived by computing the SUVR of a meta-ROI comprised of
138 lateral and medial orbitofrontal, middle and superior temporal, superior frontal, rostral middle

139 frontal, and precuneus ROIs from both hemispheres. For subsequent analyses, we chose to
140 focus on 68 cortical, 16 subcortical regions, and the global summary region. The full list of
141 regions is given in Supp. Table S1.

142 **Data harmonization**

143 Five harmonization methods were investigated in the current study: Centiloid [18], ComBat
144 [23,27], GAM-ComBat [26], longitudinal ComBat [32], and Probabilistic Estimation for Across-
145 batch Compatibility Enhancement (PEACE) [29]. These methods are briefly described below.

146 **Centiloid**

147 The Centiloid scale [18] is a method of linearly transforming global amyloid burden estimates
148 from SUVRs to a scale that is standardized across tracers. Centiloid ranges from 0 to 100, with
149 0 corresponding to the average amyloid burden of a group of healthy controls, and 100
150 corresponding to the average amyloid burden of typical AD patients. Note that Centiloids are
151 allowed to fall above 100 CL or below 0 CL.

152 Although Centiloid is calibrated against and primarily used to harmonize the global summary
153 SUVR, it can also be applied to regional or voxel-wise SUVRs [17,18]. To convert regional
154 SUVRs to Centiloid, we utilized the conversion equations that were previously validated for the
155 OASIS-3 cohort [10,33]:

$$\begin{aligned} 156 \quad CL_{PiB} &= 111.8 SUVR_{PiB} - 119.3 \\ CL_{FBP} &= 163.6 SUVR_{FBP} - 181.0 \end{aligned} \quad (1)$$

157 **ComBat**

158 ComBat is a data-driven method for adjusting data with batch-specific effects [23], where
159 batches refer to any nominal variable(s) which may contribute confounding biases in the target
160 measurement. It utilizes a multivariable linear regression to model measurements in terms of

161 batch-specific shift and scale parameters, as well as other covariates which model variance due
162 to biologically relevant effects. For batch effect i , subject j and feature k , ComBat models the
163 measurement y_{ijk} as:

$$164 \quad y_{ijk} = \alpha_k + X_j \beta_k + \gamma_{ik} + \delta_{ik} \epsilon_{ijk} \quad (2)$$

165 where α_k is the mean measurement across all subjects and all batches, X_j is the vector of
166 biological covariates associated with subject j , and β_k is the vector of coefficients for X_j . The
167 batch-specific shift (additive) and scale (multiplicative) parameters are represented by γ_{ik} and
168 δ_{ik} respectively. These modify the measurement from the group average to account for batch-
169 specific biases. ϵ_{ijk} is the error term, which is assumed to be normally distributed with zero
170 mean and unit variance. γ_{ik} and δ_{ik} are estimated using an empirical Bayesian approach, and
171 once estimated, the measurement without batch effects can be recovered by the following:

$$172 \quad y_{ijk}^* = \frac{y_{ijk} - \alpha_k - X_j \beta_k - \gamma_{ik}}{\delta_{ik}} + \alpha_k + X_j \beta_k \quad (3)$$

173 This adjustment ensures that only variance due to the batch effects is corrected for, while
174 variance due to the covariates is preserved, which is a unique advantage of ComBat over other
175 batch-adjusting techniques. In subsequent experiments, we selected age, sex, and APOE
176 carriership as the covariates of interest to preserve.

177 **GAM-ComBat**

178 A limitation of the ComBat model is that it is only able to model covariates as linearly related to
179 the target variable. To address this, Pomponio *et al.* [26] developed GAM-ComBat, a variant of
180 ComBat which can model continuous covariates non-linearly using generalized additive models
181 (GAM):

$$182 \quad y_{ijk} = f(X_j) + \gamma_{ik} + \delta_{ik} \epsilon_{ijk} \quad (4)$$

183 where f is the GAM. In subsequent experiments, we explored modeling the age covariate non-
184 linearly using GAMs.

185 **Other ComBat variants**

186 We also tested two variants of ComBat: longitudinal ComBat [32] and PEACE [29]. Longitudinal
187 ComBat incorporates a subject-specific random intercept term to the original ComBat model:

$$188 \quad y_{ijk} = \alpha_k + X_j(t)\beta_k + \gamma_{ik} + \eta_{jk} + \delta_{ik}\epsilon_{ijk}(t) \quad (5)$$

189 where η_{jk} is the subject-specific random intercept, which is preserved after harmonization.
190 Additionally, the covariate design matrix X_j and error term ϵ_{ijk} are parameterized by time t ,
191 allowing for these terms to be varied across time. This model is appropriate for harmonizing
192 data consisting of multiple repeated measurements of the same subjects at different time points.

193 PEACE differs from ComBat in two aspects: (1) it models target measurements using a bimodal
194 Gaussian mixture model to estimate two clusters of the data, then estimates the batch-specific
195 parameters independently of the cluster assignments; and (2) rather than estimating model
196 parameters and hyperparameters in an empirical Bayesian manner, it employs a fully Bayesian
197 approach where these parameters are assumed to be distributed by fixed priors. PEACE
198 addresses ComBat's limitation of assuming that the target features, after residualizing covariate
199 terms, are distributed normally. However, this is often not the case for amyloid imaging data,
200 where the distribution of amyloid burden across subjects often exhibits a bimodal pattern of low
201 amyloid (amyloid-negative) and high amyloid (amyloid-positive) clusters [29,34].

202 The PEACE model is given by the following:

$$203 \quad y_{ijk} = \alpha_{z_{jk}} + X_j\beta_{z_{jk}} + \gamma_{ik} + \delta_{ik}\epsilon_{ijk} \quad (6)$$

204 where z_j indicates the cluster assignment of subject j . PEACE is fit using Hamiltonian Monte
205 Carlo Markov chain sampling. In our experiments, we ran 4 separate Markov chains where 50
206 warmup iterations were sampled followed by 50 samples drawn from the posterior, for a total of
207 200 samples after warmup. We then averaged over these 200 samples to obtain the final
208 harmonized data.

209 The original implementation of PEACE only considered a single covariate. We modified the
210 model to allow for either no covariates or multiple covariates. In subsequent experiments, we
211 selected age, sex, and APOE carriership as the covariates of interest to preserve. We further
212 configured PEACE to train on one dataset and be applied on held-out data.

213 **Statistical analysis**

214 To evaluate Centiloid and ComBat for harmonizing inter-tracer differences for both global and
215 regional amyloid PET features, we performed two experiments: a head-to-head comparison of
216 FBP and PiB to evaluate absolute agreement, and a clinical trial simulation to evaluate the
217 clinical utility of harmonization. A summary of the pipelines for the two experiments and the data
218 used for each is illustrated in Fig. 1.

219 **Tracer head-to-head comparison**

220 We performed a head-to-head comparison of FBP and PiB measurements and evaluated their
221 absolute agreement after harmonization. We identified 113 FBP-PiB scan pairs across 99
222 subjects which were acquired within 90 days. All remaining scans were used to train ComBat
223 models, which included 565 FBP and 1044 PiB scans.

224 Centiloid, three different configurations of ComBat, and PEACE were applied to the global
225 summary SUVR and 84 regional SUVRs from the head-to-head dataset to harmonize tracer
226 differences. We tested ComBat without any covariates, ComBat with age, sex and APOE- ϵ 4

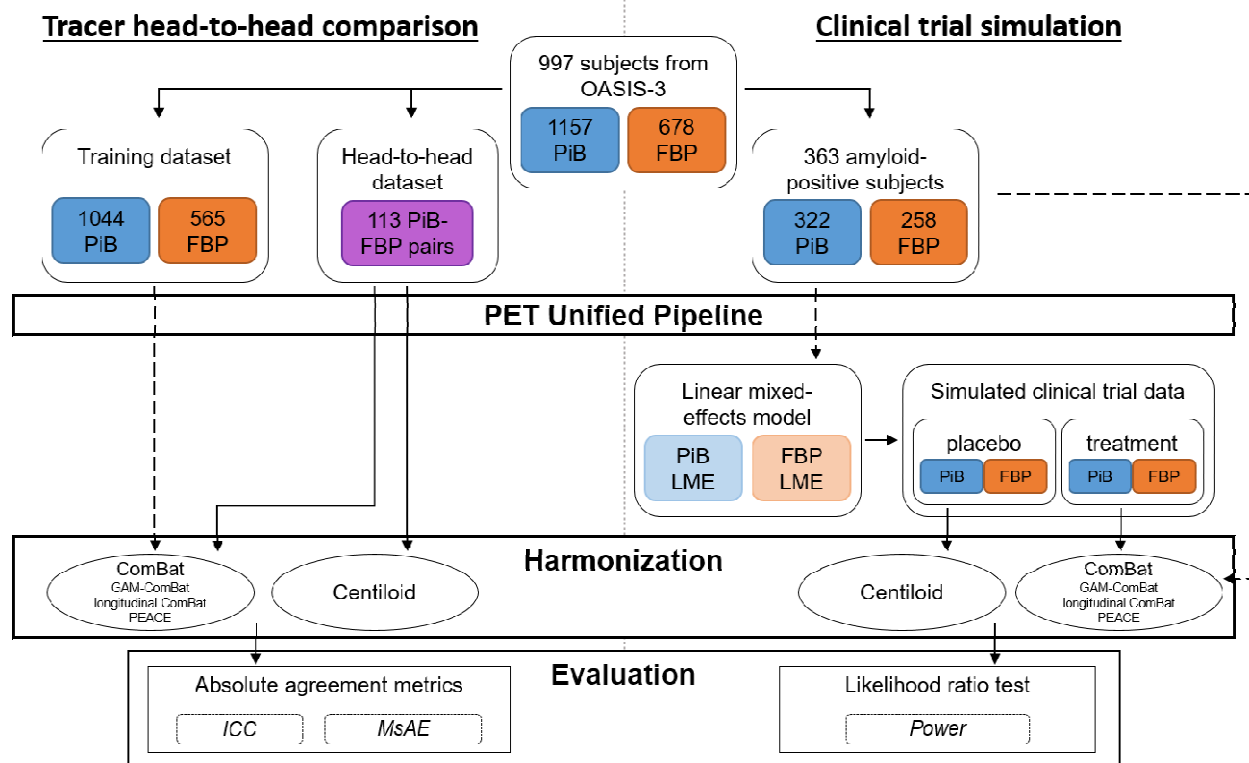


Fig. 1 Flowchart of data for the tracer head-to-head comparison (left) and clinical trial simulation (right). Dotted arrows indicate where data was used to train ComBat or linear mixed effects models.

227 carriership as linear covariates, and GAM-ComBat with sex and APOE- 4 carriership as linear
 228 covariates and age as a non-linear covariate. Additionally, we tested PEACE without any
 229 covariates and PEACE with age, sex and APOE- 4 carriership as linear covariates. Note that
 230 we omitted longitudinal ComBat from this analysis since the head-to-head data was treated as
 231 cross-sectional.

232 To evaluate the absolute agreement between FBP and PiB measurements, two metrics were
 233 computed. Firstly, intraclass correlation coefficient (ICC) using a fixed rater, single
 234 measurement model (i.e., ICC3) was estimated. ICC is roughly the ratio of intraclass variance to
 235 total variance, and values closer to 1 indicate better agreement between the two tracers.
 236 Secondly, the absolute error (AE) between FBP and PiB measurements was computed:

237 (7)

238 where PiB_i and FBP_i are the measurements made with PiB and FBP from scan pair i . We also
239 computed the mean absolute error (MAE) across all scan pairs:

$$240 \quad MAE = \frac{1}{N} \sum_{i=1}^N AE_i \quad (8)$$

241 where N is the number of scan pairs. To facilitate comparisons of absolute errors between
242 (un)harmonized SUVRs and Centiloids, we scaled Centiloids to a similar dynamic range as
243 SUVRs. Utilizing the Centiloid conversion equations in Equation 1, we computed the $SUVR_{FBP}$
244 and $SUVR_{PiB}$ that result in 0 CL and 100 CL (denoted as $SUVR_{\{tracer\},0CL}$ and $SUVR_{\{tracer\},100CL}$,
245 respectively). Then, using the average of FBP and PiB SUVRs at these anchor points, we
246 linearly mapped Centiloids back to SUVR using the following equation:

$$247 \quad SUVR_{CL} = CL \left(\frac{\frac{SUVR_{FBP,100} + SUVR_{PiB,100}}{2} - \frac{SUVR_{FBP,0} + SUVR_{PiB,0}}{2}}{100} \right) + \frac{SUVR_{FBP,0} + SUVR_{PiB,0}}{2} \quad (9)$$

248 Substituting the respective SUVRs ($SUVR_{FBP,100} = 1.718$, $SUVR_{PiB,100} = 1.961$, $SUVR_{FBP,0} =$
249 1.106 , $SUVR_{PiB,0} = 1.067$) into the above equation yields the following:

$$250 \quad SUVR_{CL} = 0.007528 * CL + 1.0867 \quad (10)$$

251 We will refer to this as the scaled Centiloid in the remainder of the text.

252 Paired t-tests were performed to test for significant differences in the distributions of ICC and
253 MAE between unharmonized SUVRs and each of the four harmonization methods. Additionally,
254 we further subdivided each FreeSurfer region into three groups – regions belonging to the
255 global summary meta-ROI, other cortical regions not part of the summary meta-ROI, and
256 subcortical regions. We then performed paired t-tests for each group separately to compare
257 across harmonization methods. In all statistical tests, Bonferroni correction was applied to
258 correct for multiple comparisons.

259 **Clinical trial simulation**

260 We evaluated Centiloid and ComBat in the context of improving detection of treatment effects in
261 an anti-amyloid drug trial setting, with the assumption that multiple amyloid PET tracers were
262 used due to pooling of data from multiple institutions. To accomplish this, we modeled a
263 simulation experiment after those described in Chen *et al.* [17] to generate data of placebo and
264 treatment groups. We varied the proportion of FBP-to-PiB scans of each group, then tested for
265 group differences of amyloid rate-of-change.

266 We selected subjects who presented as PET amyloid-positive at least once during their
267 participation in OASIS-3. To mark scans as amyloid-positive, we used a global summary SUVR
268 threshold of 1.31 for PiB and 1.24 for FBP. These thresholds were previously validated for the
269 OASIS-3 cohort [10]. From these criteria, we identified 363 amyloid-positive subjects, from
270 which 258 FBP and 322 PiB scans were selected.

271 For each tracer and for each region-of-interest (including the global summary region), a linear
272 mixed effects (LME) model was fit on the selected scans to predict longitudinal SUVR. Sex,
273 APOE carriership, baseline age, and time-from-baseline were specified as fixed effects. A
274 random intercept grouped by subject was specified as the only random effect. This resulted in
275 two LME models – one fitted to FBP data and one to PiB data – for each region-of-interest.
276 Fitted LME models were then used to generate new longitudinal data of placebo and treatment
277 groups. For the placebo group, the models were applied as is to generate SUVRs that follow the
278 natural longitudinal trajectory among amyloid-positive subjects in OASIS-3. For the treatment
279 group, we added a negative rate-of-change term to the LME equation to mimic a treatment
280 effect. We tested multiple values of the treatment effect from 0 to -0.03 SUVR, varying in
281 increments of -0.01 SUVR. These values were chosen based off of previously reported clinical
282 trial effect sizes [2,17].

283 To simulate a single subject's data, we randomly sampled the empirical distributions of the
284 number of longitudinal scans, age at baseline scan, and interval between scans among the
285 OASIS-3 amyloid-positive cohort to generate longitudinal time points. Each simulated subject
286 was randomly assigned sex and APOE carriership based on the empirical distributions of these
287 covariates. We then allocated a tracer (either PiB or FBP) to each time point, with the following
288 constraints: (1) the proportion of tracers across all scans from all subjects approximated a
289 prespecified proportion; (2) a subject could only switch tracers once during their clinical trial
290 participation, reflecting a realistic scenario where multiple tracers are utilized in a single study.
291 Time points assigned to PiB or FBP were then input into the corresponding trained LME model
292 to obtain simulated SUVR measurements. For our experiments, we varied the percentage of
293 FBP scans from 0.1 to 0.9 in increments of 0.2 for both clinical trial groups independently. We
294 fixed the total number of subjects to 50 per group, with the number of scans per subject ranging
295 from 2 to 6, the mean (\pm standard deviation) time in between scans being 3.35 ± 1.46 years,
296 and the mean (\pm standard deviation) baseline age being 69.4 ± 9.2 years.

297 The simulated data was harmonized using one of four methods: Centiloid, ComBat, PEACE, or
298 longitudinal ComBat. We tested ComBat, PEACE and longitudinal ComBat both without
299 covariates and with age, sex, and APOE- ϵ 4 carriership as linear covariates. Note that we
300 omitted GAM-ComBat from this analysis, since the simulated data was generated using a linear
301 age term in the LME. ComBat and PEACE were trained using all available amyloid-positive
302 scans. For longitudinal ComBat, since the random intercept terms are indexed by subjects in the
303 training dataset, it is not possible to use this model to harmonize data for new subjects.
304 Therefore, we trained longitudinal ComBat on the simulated subjects' data, and trained models
305 were applied to harmonize this data.

306 To test for group differences in the rate-of-change in amyloid SUVR between placebo and
307 treatment groups, we first fitted the same LME described previously, but with three additional

308 terms – clinical trial group, interaction of time-from-baseline with trial group, and tracer. We then
309 tested for statistical significance of the time-from-baseline and clinical trial group interaction
310 term, which would indicate whether the two groups exhibit different rates-of-change. Likelihood
311 ratio tests were used to compare the fit of the full model with a nested model that excludes this
312 term, and significance was determined using $\alpha = 0.05$. A one-way test was used, meaning that
313 we considered a rate-of-change difference to be statistically significant only if the treatment
314 group had a lower rate-of-change than the placebo. Simulations were repeated 1000 times for
315 each permutation of tracer mixing proportions and treatment effect. Statistical power was
316 computed as the proportion of simulation iterations which resulted in a significant finding. Note
317 that for a treatment effect of zero, i.e. the absence of a ground truth treatment effect, this
318 corresponds to the Type-I error rate.

319 **Results**

320 **Demographics**

321 Descriptive statistics of each cohort are listed in Table 1. A two-tailed t-test was used to test for
322 differences in age at scan, and Fisher's exact test was used to test for differences in sex,
323 APOE- ϵ 4, and Clinical Dementia Rating® (CDR). Significant differences in age and CDR were
324 observed between the head-to-head cohort and the single-tracer FBP cohort ($p < 0.005$). Age
325 was also significantly different between the head-to-head and the mixed-tracer PiB cohorts ($p <$
326 0.05), and CDR was significantly different between the head-to-head and single-tracer PiB
327 cohorts ($p < 0.005$). For the tracer-vs-tracer comparison, age was significantly different in both
328 single- and mixed-tracer cohorts in the training dataset ($p < 0.005$), and for only the mixed-
329 tracer cohort in the simulation dataset ($p < 0.005$). Lastly, sex and APOE- ϵ 4 carriership were
330 significantly different between FBP and PiB in the single-tracer simulation dataset ($p < 0.01$).

Table 1 Demographics of each cohort. The training dataset of the head-to-head comparison was split into subjects who were scanned with only one tracer (single-tracer) and those who were scanned with multiple tracers (mixed-tracer). Statistically significant differences are denoted with asterisks and crosses. Asterisks indicate comparisons between head-to-head and each of the 4 training sub-cohorts, whereas crosses indicate FBP vs. PiB comparisons within a single/mixed tracer cohort. The number of symbols indicates the significance level (1 = $p < 0.05$, 2 = $p < 0.01$, 3 = $p < 0.005$, 4 = $p < 1e-4$). Note that the reported CDRs are counted by scans, where the closest CDR score in time was assigned to every scan.

	Tracer head-to-head comparison					Simulation			
	head-to-head	training				single-tracer	PiB	mixed-tracer	
		single-tracer	mixed-tracer		FBP			PiB	
		FBP	PiB	FBP	PiB	FBP	PiB	FBP	PiB
Number of subjects	99	301	452	219	219	153	150	60	60
Number of scans	113	325	656	240	388	173	200	85	122
Mean age at scan (\pm sd) ¹	69 \pm 8.48	72.3 \pm 7.63 ^{***}	70.2 \pm 9.57 ^{***}	70.6 \pm 8.4	66.9 \pm 8.66 ^{*†***}	73.1 \pm 7.24	74.5 \pm 7.27	74.4 \pm 7.02	70.8 \pm 6.9 ^{***}
Number of males/females	45/54	145/156	205/247	78/141	78/141	55/98	84/66 ^{***}	22/38	22/38
APOE noncarriers/carriers ²	67/32	179/122	283/169	158/61	158/61	78/75	52/98 ^{**}	26/34	26/34
Number of CDR = 0/CDR = 0.5/CDR > 1 ³	108/4/1	250/58/17 [*]	527/104/25 ^{***}	230/9/1	373/14/1	114/45/14	116/63/21	75/8/2	116/5/1
Number of scans per subject (\pm sd)	1.1 \pm 0.38	1.1 \pm 0.27	1.5 \pm 0.76	1.1 \pm 0.31	1.8 \pm 0.9	1.1 \pm 0.39	1.3 \pm 0.67	1.4 \pm 0.65	2 \pm 1.12
Mean years between scans (\pm sd)	1.6 \pm 0.53	3 \pm 0.82	3.5 \pm 1.81	3 \pm 0.92	4.2 \pm 2.17	2.5 \pm 1.25	3.2 \pm 1.49	1.8 \pm 0.78	3.4 \pm 1.54

* = significance of head-to-head vs. training comparison

† = significance of tracer vs. tracer comparison

¹ sd = standard deviation

² APOE = apolipoprotein-E4

³ CDR = Clinical Dementia Rating[®]

331 Tracer head-to-head comparison

332 We evaluated the ability of Centiloid, ComBat and PEACE to improve the absolute agreement
 333 between FBP and PiB using the head-to-head dataset. For the global summary region, absolute
 334 agreement, as measured by ICC, increased after harmonization with either Centiloid ($ICC =$
 335 0.912) or ComBat with no covariates ($ICC = 0.916$), compared to the unharmonized SUVR

336 ($ICC = 0.882$) (Table 2, Fig. 2). ICC also increased slightly after harmonization with ComBat with

Fig. 2 Global summary measures computed from PiB and FBP scans in the tracer head-to-head dataset. The red line indicates the best fit line from ordinary least squares linear regression, the gray area indicates the confidence interval of the slope, and the black line represents the identity line. Intraclass correlation coefficient (ICC) is reported on the bottom right of each scatterplot.

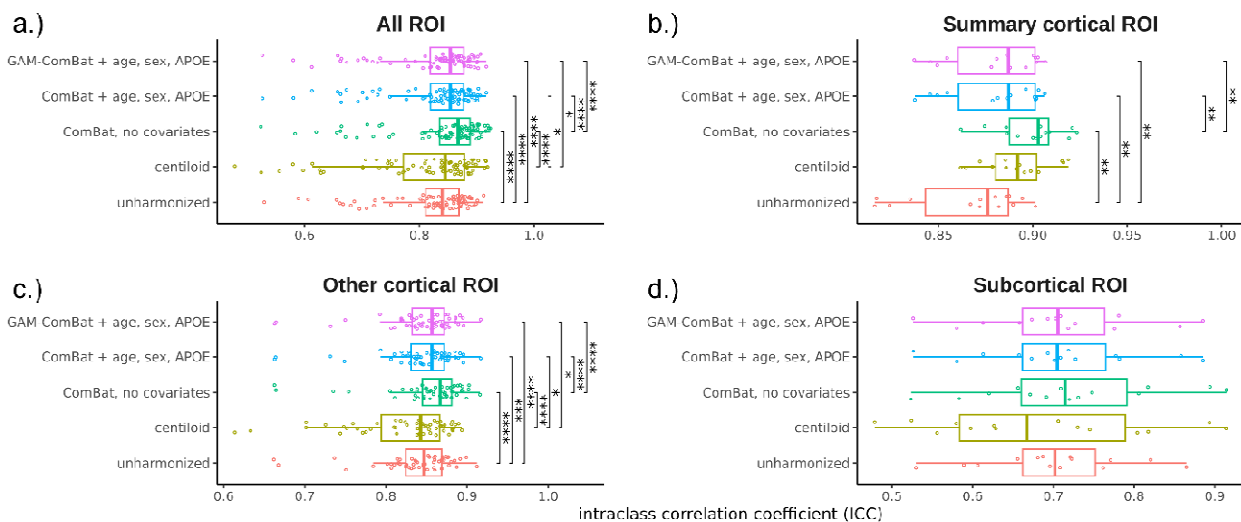
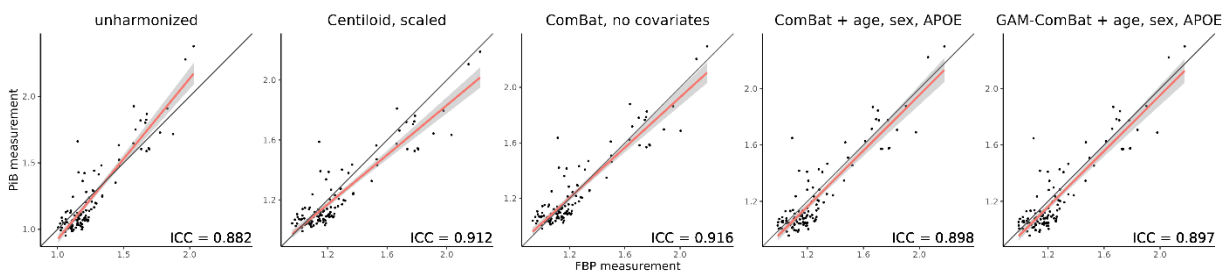


Fig. 3 Distribution of regional ICC across all ROIs, grouped by harmonization method and ROI subgroup. Each point represents a single ROI. Significance levels from paired t-tests with Bonferroni correction are indicated for each pair of harmonization methods (* = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.005$, **** = $p < 1e-4$).

337 covariates () and GAM-ComBat (), albeit not to the same degree as

338 ComBat with no covariates. Similarly, PEACE, either with ($ICC = 0.884$) or without ($ICC = 0.897$)

339 covariates, did not perform as well as ComBat with no covariates or Centiloid in increasing ICC

340 (Table 2, Supp. Fig. S1). For ROI measurements, all three ComBat harmonization methods led

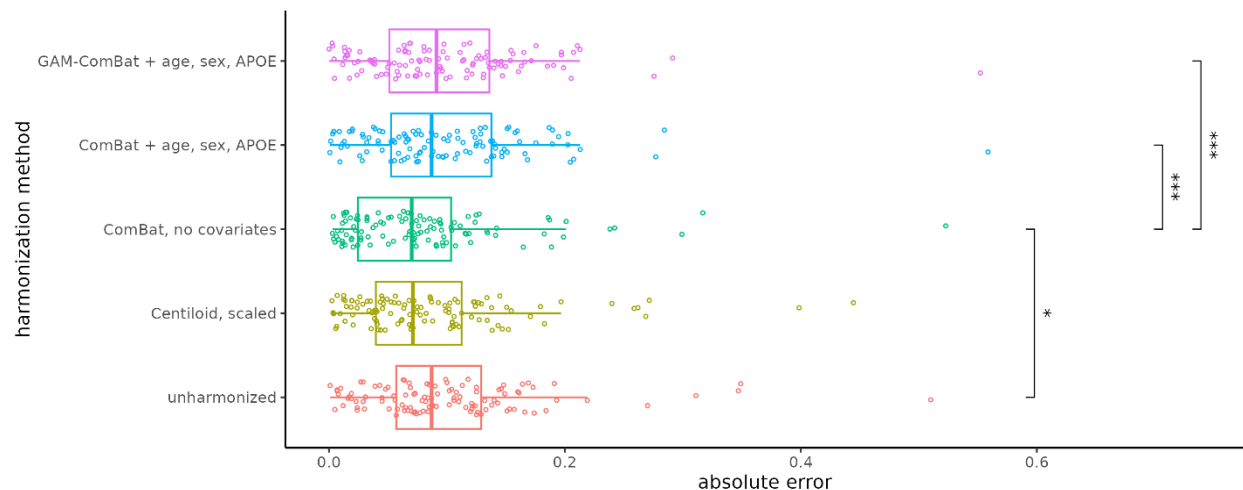


Fig. 4 Absolute error of the global summary measure from each PiB and FBP scan pair in the tracer head-to-head dataset. Each point represents a single scan pair. Significance levels from paired t-tests with Bonferroni correction are indicated for each pair of harmonization methods (* = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.005$, **** = $p < 1e-4$).

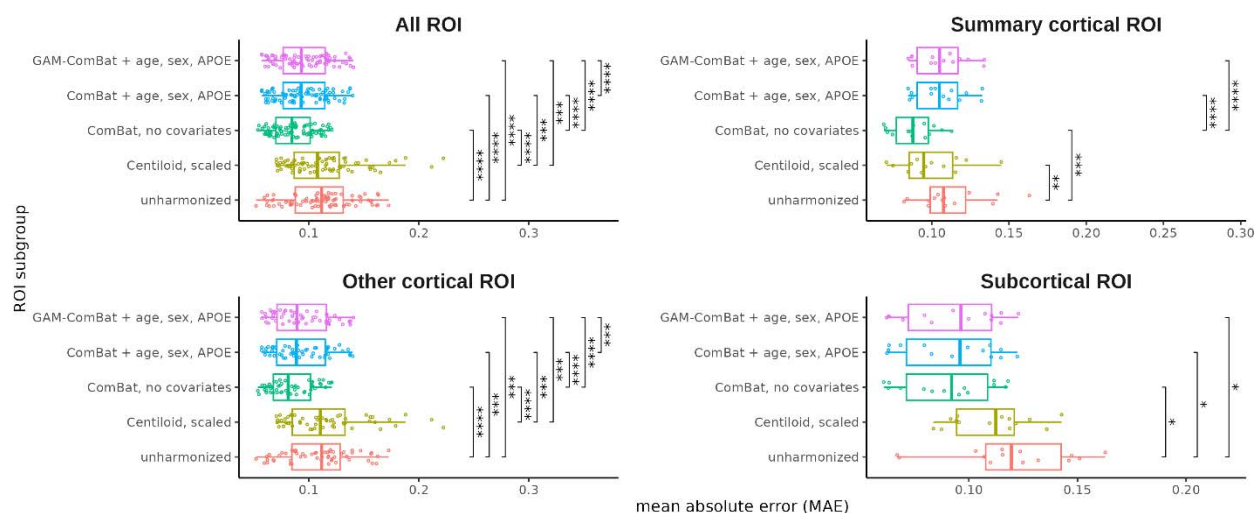


Fig. 5 Distribution of regional MAE across all ROIs, grouped by harmonization method and ROI subgroup. Each point represents a single ROI. Significance levels from paired t-tests with Bonferroni correction are indicated for each pair of harmonization methods (* = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.005$, **** = $p < 1e-4$).

341 to a statistically significant increase in average ICC among all ROIs compared to unharmonized
 342 SUVR ($p < 1e-4$), with ComBat with no covariates again performing the best ($\overline{ICC} = 0.838$)
 343 (Table 2, Fig. 3a). Additionally, ComBat with no covariates performed the best within the
 344 summary cortical ROIs ($\overline{ICC} = 0.899$) and other cortical ROIs ($\overline{ICC} = 0.853$) (Table 2, Fig. 3b-c).
 345 No method was effective at improving cross-tracer agreement for the subcortical ROIs, with

Table 2 Mean \pm standard deviation of ICC and MAE from the tracer head-to-head comparison. The mean across scan pairs was computed for the global summary region, whereas the mean across ROIs was computed for the FreeSurfer ROI. Bold values indicate the best performing harmonization method. Asterisks indicate statistical significance of paired t-tests comparing each harmonization method with unharmonized (* = $p < 0.05$, ** = $p < 0.01$, *** = $p < 0.005$, **** = $p < 1e-4$).

346 mean ICC of less than 0.75 even after harmonization (Table 2, Fig. 3d). PEACE did not improve
 347 the mean ICC with statistical significance in any of the groups except for other cortical ROIs
 348 (Supp. Fig. S2), and PEACE both with and without covariates achieved a lower mean ICC than
 349 ComBat with no covariates (Table 2). When plotting region-wise ICC on the inflated brain
 350 surface, we observed that Centiloid resulted in a decrease in ICC in the bilateral occipital and
 351 sensorimotor regions, and

Harmonization method	Global summary		FreeSurfer ROI							
	ICC	MAE	All ROI		Summary cortical ROI		Other cortical ROI		Subcortical ROI	
			ICC	MAE	ICC	MAE	ICC	MAE	ICC	MAE
unharmonized	0.882	0.101 \pm 0.076	0.819 \pm 0.077	0.111 \pm 0.029	0.866 \pm 0.03	0.113 \pm 0.023	0.838 \pm 0.048	0.109 \pm 0.031	0.701 \pm 0.089	0.12 \pm 0.028
Centiloid, scaled	0.912	0.088 \pm 0.075	0.811 \pm 0.098	0.112 \pm 0.032	0.891 \pm 0.018	0.099 \pm 0.02**	0.822 \pm 0.063	0.116 \pm 0.037	0.686 \pm 0.137	0.11 \pm 0.019
ComBat, no covariates	0.916	0.08 \pm 0.075	0.838 \pm 0.083 ****	0.086 \pm 0.019 ****	0.899 \pm 0.018 **	0.089 \pm 0.014 **	0.853 \pm 0.05 ****	0.084 \pm 0.019****	0.72 \pm 0.111	0.09 \pm 0.021
ComBat + age, sex, APOE	0.898	0.099 \pm 0.075	0.827 \pm 0.08****	0.095 \pm 0.023****	0.881 \pm 0.024**	0.106 \pm 0.017	0.844 \pm 0.049****	0.093 \pm 0.024**	0.707 \pm 0.099	0.092 \pm 0.022
GAM-ComBat + age, sex, APOE	0.897	0.1 \pm 0.075	0.827 \pm 0.08****	0.095 \pm 0.023****	0.881 \pm 0.024**	0.107 \pm 0.017	0.844 \pm 0.049****	0.093 \pm 0.024**	0.707 \pm 0.099	0.092 \pm 0.022
PEACE, no covariates	0.897	0.092 \pm 0.077**	0.813 \pm 0.087	0.097 \pm 0.022****	0.875 \pm 0.027	0.103 \pm 0.019	0.831 \pm 0.058	0.093 \pm 0.021**	0.684 \pm 0.094	0.104 \pm 0.027
PEACE + age, sex, APOE	0.884	0.083 \pm 0.095	0.822 \pm 0.079	0.086 \pm 0.021 ****	0.867 \pm 0.028	0.094 \pm 0.018	0.842 \pm 0.05	0.083 \pm 0.021 ****	0.698 \pm 0.087	0.092 \pm 0.022

352 in the left temporal and parietal cortices (Supp. Fig. S3b). PEACE also exhibited a decrease in
 353 ICC in similar regions (Supp. Fig. S3f-g). In contrast, none of the ComBat variants led to such
 354 decrease in these regions (Supp. Fig. S3c-e), with ComBat with no covariates having the
 355 highest magnitude of change in ICC across multiple regions.

356 When assessing the absolute error between FBP and PiB measurements of the global summary
 357 region (Table 2, Fig. 4, Supp. Fig. S4), ComBat with no covariates ($MAE = 0.080$) and PEACE,
 358 both with ($MAE = 0.083$) and without ($MAE = 0.092$) covariates, were the only methods that

359 reduced the mean absolute error with statistical significance ($p < 0.05$) compared to
360 unharmonized SUVRs ($MAE = 0.101$). ComBat with no covariates lowered the absolute error
361 the greatest between the three methods. For ROI measurements (Table 2, Fig. 5, Supp. Fig.
362 S5), all methods except Centiloid significantly reduced the average MAE among all ROIs ($p <$
363 $1e-4$), but ComBat with no covariates and PEACE with covariates resulted in the greatest
364 reduction ($\overline{MAE} = 0.086$). ComBat with no covariates also performed the best within the
365 summary cortical ROIs ($\overline{MAE} = 0.089$) and subcortical ROIs ($\overline{MAE} = 0.09$), while PEACE with
366 covariates performed the best within other cortical ROIs ($\overline{MAE} = 0.083$), although ComBat with
367 no covariates performed comparably. When plotting region-wise MAE on the inflated brain
368 surface, we observed that both ComBat with no covariates and PEACE with covariates led to
369 the greatest reduction in MAEs in ROIs across the bilateral frontal, parietal and occipital regions
370 (Supp. Fig. S6).

371 **Clinical trial simulation**

372 We performed simulations to test for group differences in amyloid rate-of-change between
373 treatment and placebo groups in a hypothetical clinical trial, and evaluated whether
374 harmonization improved the ability of detecting these differences. For the global summary
375 SUVR, both Centiloid and ComBat without covariates resulted in overall increases in statistical
376 power after harmonization in the presence of a treatment effect (i.e., for rate-of-change \in
377 $\{-0.01, -0.02, -0.03\}$), primarily when the placebo group had high FBP composition and the
378 treatment group had low FBP composition (Fig. 6). A similar increase in power was observed
379 when using ComBat with covariates, albeit at a lesser magnitude. However, PEACE led to
380 widespread decreases in power, and although longitudinal ComBat led to increases in power in
381 certain tracer composition configurations, the effect was not as consistent as either Centiloid or
382 ComBat with no covariates (Supp. Fig. S7). In the absence of a treatment effect (i.e., rate-of-
383 change = 0), Centiloid, ComBat and longitudinal ComBat with no covariates achieved large

384 decreases in Type-I error, primarily when the placebo group had low FBP composition and the

385 treatment group had high FBP composition. Slight decreases in power were also observed in

386 the case of a low treatment effect (i.e., rate-of-change = -0.01). Supplementary Table S2 shows

387 the mean power computed across all 25 configurations of treatment and placebo FBP

388 compositions. Out of all methods, Centiloid achieved the largest mean power for detecting the

389 treatment effect in every rate-of-change except for -0.01 (longitudinal ComBat performed the

390 best in this case), and Centiloid also achieved the lowest mean Type-I error. However, ComBat

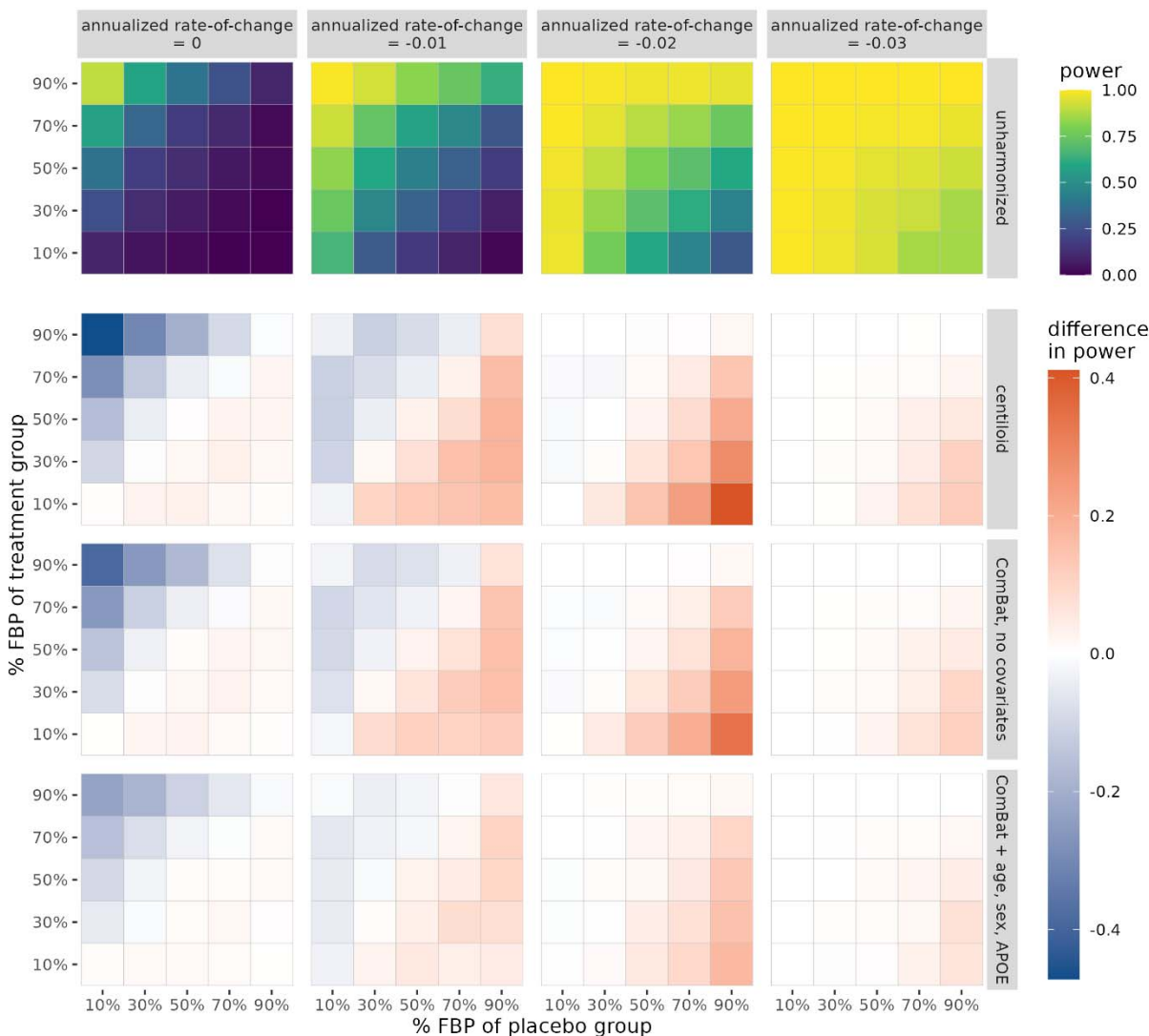


Fig. 6 Statistical power of detecting group differences in rate-of-change of the global summary SUVR between treatment and placebo groups, computed as the proportion of significant findings over 1000 iterations. Power is plotted for unharmonized SUVR, while the difference in power relative to unharmonized is plotted for all harmonization methods. The true underlying rate-of-change is varied across columns. The proportion of FBP scans in the placebo and treatment groups is varied across the horizontal and vertical axes of each heatmap, respectively. Note that for annualized rate-of-change equal to zero, the proportion of significant findings corresponds to Type-I error rate.

391 with no covariates was often the second-best performing method.

392 For regional SUVRs, similar patterns of change in power after harmonization were observed.
393 Across all ROIs, mean power was increased/Type-I error was decreased after harmonization
394 with either Centiloid or ComBat, with both methods producing comparable changes in power
395 (Supp. Table S3). Additionally, these same changes in regional power after harmonization were
396 consistent across all three ROI subgroups. Surface plots revealed that regions which exhibited
397 relatively low power of detecting treatment effects (such as the frontal and medial parietal
398 cortices and putamen in the case of rate-of-change = -0.02) experienced a high increase in
399 power after harmonization with either Centiloid or ComBat with no covariates (Supp. Fig. S8-9).
400 Again, harmonization with PEACE led to widespread decreases in statistical power across
401 multiple regions, and longitudinal ComBat showed worse performance in increasing power
402 compared to Centiloid or ComBat with no covariates, with the exception of when rate-of-change
403 = -0.01 (Supp. Fig. S8-9, Supp. Table S3).

404 **Discussion**

405 We demonstrated that ComBat may effectively harmonize amyloid PET measurements across
406 FBP and PiB. Notably, ComBat with no covariates outperformed Centiloid in increasing absolute
407 agreement between tracers in both the global summary and regional measurements, and
408 resulted in a comparable improvement in detecting group differences in the simulated clinical
409 trial. As more studies shift focus from using a global summary metric of amyloid burden to using
410 the spatial distribution of regional amyloid as features [6,7], harmonization techniques like
411 ComBat that can be applied to multiple regions become appealing for pooling PET data across
412 multiple tracers.

413 ComBat poses several methodological advantages over Centiloid. Firstly, calibration of Centiloid
414 requires *a priori* selection of representative individuals from dichotomous groups, namely
415 healthy control and typical AD cohorts. This *a priori* cohort selection may introduce bias into the

416 calibration process. Especially if the selected sample is small and/or captures only a subset of
417 the overall population (e.g., biased towards a single ethnicity group), Centiloid may not
418 generalize well to heterogeneous or out-of-sample datasets. ComBat circumvents this
419 requirement, which allows it to learn a robust harmonization on a wider spectrum of data which
420 consists of controls, AD patients and “in-between” subjects. Furthermore, much like Centiloid, a
421 trained ComBat model may then be used to harmonize out-of-sample data. However, one
422 should still carefully consider the subject selection process and ensure that the proportions of
423 cohorts are not skewed towards a particular cohort (e.g., many more healthy controls than AD
424 patients), since this may lead to a suboptimal harmonization model [29]. Secondly, Centiloid
425 requires at least two PET scans of different tracers for each subject in the calibration cohort,
426 one of which should be acquired using PiB. In contrast, ComBat can train using just one scan
427 per subject and does not require PiB to be used. Indeed, we expect ComBat to generalize well
428 to harmonizing two (or more) [^{18}F]-based tracers without the need for PiB, although further
429 investigation using head-to-head data of these tracers should be conducted to verify this.
430 Thirdly, a region-specific harmonization is important for addressing sources of tracer bias which
431 variably affect different regions, such as non-specific binding [35]. However, as suggested by
432 Klunk et al. [18], a region-specific Centiloid calibration is not ideal, since it would fix different
433 SUVRs of different regions to the same Centiloid value. In contrast, ComBat independently
434 removes the tracer-specific variance from the target measurements without scaling the dynamic
435 range of each region to fixed points, making it a more suitable technique for regional
436 harmonization. A caveat to this is that ComBat effectively centers the harmonized
437 measurements on the global mean and variance of the data. As such, these measurements can
438 no longer be interpreted as belonging to the scale of a particular batch, but rather on an
439 aggregated scale. Alternative approaches such as modified ComBat [36] allow users to choose
440 a “gold standard” reference batch to which all other batches are adjusted to, which may aid in
441 improving the interpretability of harmonized measurements.

442 Lastly, ComBat has the advantage of being able to preserve covariate relationships in the target
443 measurement, which may be useful in downstream analyses such as in predictive models which
444 take into account biologically-related variance to make accurate predictions. However, in the
445 context of purely evaluating the absolute agreement between tracers, we observed that
446 including covariates into the ComBat model led to worse ICC and absolute errors. This may be
447 partially due to differences in covariate distributions between the training and head-to-head
448 cohorts, of which age and CDR differed with statistical significance. Although CDR was not
449 explicitly included as a ComBat covariate, it may have indirectly contributed to a biased ComBat
450 model which does not generalize well to testing data with different covariate characteristics.

451 It was noted that no harmonization method investigated in this study performed well for the
452 subcortical regions. Notably, these regions lie close to white matter regions, and thus may be
453 affected by non-specific binding more so than cortical regions. This may contribute to more
454 noise in the subcortical regions, which batch harmonization methods such as ComBat are not
455 able to mitigate. One potential area of investigation is to evaluate whether partial volume
456 correction [35] would have an effect on regional harmonization of PET SUVRs, especially for
457 regions which experience high amounts of signal spill-over from neighboring white matter
458 regions.

459 Our simulation experiments revealed the importance of harmonization in settings where multiple
460 tracers are utilized to track brain amyloid deposition in clinical trial participants. Particularly,
461 harmonization was the most beneficial when trial groups exhibited differing proportions of tracer
462 data. In these “off-diagonal” cases, tracer biases contributed to a substantial confounding effect
463 across clinical trial groups, resulting in either a reduction of power in detecting the true
464 underlying treatment effect, or an increase in Type-I error in the case when no treatment effect
465 exists. Harmonization effectively served to mitigate these confounding effects due to tracer
466 differences. This was consistent with previous reports that found significant differences in

467 amyloid rates-of-change across different tracers within real clinical trial groups, and that these
468 differences were subsequently removed after harmonization [17]. Interestingly, we observed an
469 asymmetric effect where harmonization led to changes in power in one off-diagonal, but not in
470 the other. This was most likely because we utilized a one-sided statistical test to test for rate-of-
471 change differences. In the case of low FBP% in the placebo group and high FBP% in the
472 treatment group, and when there was no ground-truth treatment effect introduced, a high
473 amount of Type-I error suggested that FBP contributed to a greater rate-of-change compared to
474 PiB due to tracer biases alone. However, in the opposite off-diagonal, these biases did not
475 contribute to any Type-I error. While tracer biases were still present in the overall data, this
476 indicated that they did not interfere with the detection of one-way group differences.

477 It is important to note, however, that scenarios of high imbalance of tracer proportions between
478 clinical trial arms are very unlikely to occur in a real-world setting, assuming proper
479 randomization. In the more realistic case where trial groups exhibit an equal proportion of tracer
480 data, a much lower change in power was observed compared to the off-diagonal cases. This is
481 likely because the same tracer bias would affect both groups equally, which statistically would
482 not influence the detection of group differences.

483 We investigated PEACE and longitudinal ComBat, which were previously validated for scanner-
484 wise harmonization, for specifically tracer harmonization in the current study. Both methods led
485 to mixed results in the head-to-head comparison as well as in the clinical trial simulation. In the
486 head-to-head comparison, PEACE with covariates performed similarly to ComBat with no
487 covariates in terms of MAE, but not in terms of ICC. Additionally, in the clinical trial simulation,
488 PEACE failed to improve statistical power of detecting the treatment effect. This can partly be
489 explained by our decision to simulate SUVRs following a unimodal Gaussian distribution, which
490 was motivated by the fact that a true clinical trial will only enroll amyloid-positive participants and
491 exclude amyloid-negative individuals. Longitudinal ComBat resulted in mixed improvements in

492 power, but these improvements were not as consistent compared to Centiloid or ComBat with
493 no covariates. We speculate that increased model complexity may have contributed to models
494 which were less robust to the data at hand. Ultimately, we found that ComBat with no
495 covariates, which is the simplest model with the fewest number of parameters to estimate and
496 the fewest assumptions made, consistently performed either comparably to or better than
497 PEACE or longitudinal ComBat.

498 There are several limitations to the current work. Firstly, on the basis of purely increasing tracer
499 agreement, there are no clear recommendations on the choice of including covariates in
500 ComBat. One caveat to using ComBat is that, unless explicitly accounted for in the covariates, it
501 will assume that any biases due to real biological differences between tracer cohorts are batch
502 differences, which are subsequently removed. Therefore, one should carefully examine the
503 composition of the data at hand and consider whether it is necessary to model known biological
504 factors via the covariate terms. Secondly, data from the simulation experiment were generated
505 from models trained on a cohort of amyloid-positive subjects from OASIS-3 instead of data from
506 an actual anti-amyloid drug trial. Although simulations were set up to mimic data that would be
507 collected in a successful trial, it remains to be seen whether our hypotheses would hold on real-
508 world clinical trial data. Thirdly, our conclusions on the performance of ComBat for tracer
509 harmonization are limited to PiB and FBP. Although we expect ComBat to be robust to other
510 ^{18}F -based amyloid tracers such as ^{18}F -florbetaben and ^{18}F -flutemetamol, future work is
511 required to validate this using head-to-head data from these tracers. Finally, our simulation
512 analysis only focused on early amyloid-positive individuals, which we assumed to exhibit
513 temporal amyloid accumulation in a roughly linear fashion. However, to draw conclusions on a
514 cohort of both amyloid-negative and positive individuals (and even late-stage individuals with
515 plateaued amyloidosis), a sigmoidal or piecewise linear model should instead be used in order
516 to model the non-linearities of amyloid accumulation across the broader AD spectrum [37,38].

517 **Conclusion**

518 Harmonization of amyloid PET radiotracers is imperative for removing tracer-specific biases in
519 amyloid burden measurements for optimal performance of downstream tasks, such as
520 enhancing statistical power and reducing false discoveries in clinical trials. In the current study,
521 we demonstrated that ComBat is effective for harmonizing both global and regional amyloid
522 measurements in an entirely data-driven way. Our experimental results suggest that ComBat
523 not only increases the absolute agreement of measurements made within scan pairs of the
524 same group of subjects by different tracers, but also provides a significant benefit to the
525 performance of detecting true treatment effects in anti-amyloid drug trials. ComBat thus
526 presents as a viable technique for harmonizing regional-based analyses of amyloid PET.

527 **List of Abbreviations**

528 PET: positron emission tomography; FBP: [¹⁸F]-florbetapir; PiB: [¹¹C]-Pittsburgh Compound-B;
529 OASIS-3: Open Access Series of Imaging Studies 3; SUVR: standardized uptake value ratio;
530 ICC: intraclass correlation coefficient; MAE: mean absolute error; AD: Alzheimer's disease; ROI:
531 region-of-interest; CL: Centiloid; MRI: magnetic resonance imaging; APOE: apolipoprotein-ε4;
532 GAM: generalized additive model; LME: linear mixed effects; CDR: Clinical Dementia Rating;
533 PEACE: Probabilistic Estimation for Across-batch Compatibility Enhancement.

534 **Supplementary Information**

535 **Additional file 1.docx** - Supplementary figures and tables. **Supp. Table S1**: List of FreeSurfer
536 regions-of-interest used in the study, along with their subgroupings. **Supp. Fig. S1**: Scatterplots
537 of the global summary SUVR from the tracer head-to-head comparison, with results of PEACE.
538 **Supp. Fig. S2**: Boxplots of regional SUVRs from the tracer head-to-head comparison, with

539 results of PEACE. **Supp. Fig. S3:** Regional ICCs plotted on the surface from the tracer head-to-
540 head comparison. **Supp. Fig. S4:** Boxplots of absolute errors of the global summary SUVR from
541 the tracer head-to-head comparison, with results of PEACE. **Supp. Fig. S5:** Boxplots of mean
542 absolute errors of the regional SUVR from the tracer head-to-head comparison, with results of
543 PEACE. **Supp. Fig. S6:** Regional MAEs plotted on the surface from the tracer head-to-head
544 comparison. **Supp. Fig. S7:** Heatmaps of statistical power of the clinical trial simulation when
545 run on the global summary SUVR, with results of PEACE and longitudinal ComBat. **Supp. Fig.**
546 **S8:** Brain surface plots of mean statistical power from the simulation experiment. **Supp. Fig. S9:**
547 Subcortical plots of mean statistical power from the simulation experiment. **Supp. Table S2:**
548 Mean statistical power of detecting significant rate-of-change differences between treatment and
549 placebo groups in the simulation experiment for the global summary amyloid estimate. **Supp.**
550 **Table S3:** Mean statistical power of detecting significant rate-of-change differences between
551 treatment and placebo groups in the simulation experiment for the ROI amyloid measurements.

552 **Key Points**

- 553 • ComBat is a data driven harmonization method which, unlike Centiloid, does not require
554 *a priori* selection and stratification of training cohorts and is able to harmonize regional
555 amyloid PET estimates.
- 556 • ComBat with no covariates performed the best in increasing the absolute agreement of
557 regional amyloid PET measurements made within scan pairs of the same group of
558 subjects using two different radiotracers.
- 559 • ComBat increased the statistical power of detecting treatment effects and decreased
560 Type-I error of falsely detecting effects in a simulated anti-amyloid drug trial.

561 **Declarations**

562 **Ethics approval and consent to participate**

563 Ethics approvals were obtained by the OASIS-3 dataset. All participants were consented into
564 Knight ADRC-related projects in accordance with the Declaration of Helsinki and following
565 procedures approved by the Institutional Review Board of Washington University School of
566 Medicine in St. Louis. For more details, we refer the reader to the OASIS-3 reference.

567 **Consent for publication**

568 Not applicable

569 **Availability of data and materials**

570 Data utilized in this study were obtained from the OASIS-3 open access dataset. Data can be
571 requested at <https://sites.wustl.edu/oasisbrains>.

572 Code for this study will be made publicly available at <https://github.com/sotiraslab>. All statistical
573 analyses and simulation experiments were implemented using R version 4.4.0 and python
574 version 3.10.10. The *neuroharmonize* python package
575 (<https://github.com/rpomponio/neuroHarmonize>) was used to train and apply ComBat and GAM-
576 ComBat models, and the *longCombat* R package (<https://github.com/jcbeer/longCombat>) was
577 used to train and apply longitudinal ComBat.

578 **Competing Interests**

579 AS reported receiving personal fees from BrightFocus for serving as a grant reviewer and stock
580 from TheraPanacea outside the submitted work. All remaining authors have no conflicting
581 interests to report.

582 **Funding**

583 BY was supported by the Imaging Science Pathways NIH T32 EB014855 and BrightFocus
584 Foundation grant ADR A2021042S. AS was supported by NIH award R01 AG067103 and
585 BrightFocus Foundation grant ADR A2021042S.

586 Computations were performed using the facilities of the Washington University Research
587 Computing and Informatics Facility, which were partially funded by NIH grants S10OD025200,
588 1S10RR022984-01A1 and 1S10OD018091-01. Additional support is provided by The
589 McDonnell Center for Systems Neuroscience.

590 **Authors' contributions**

591 All authors contributed to the conceptualization and design of the study. BY implemented all
592 data analyses and experiments and wrote the first draft of the manuscript. AS, BG and TB
593 contributed to the interpretation of data. TE, SK and DK provided technical support. All authors
594 were involved with manuscript revision, and all approved of the final draft.

595 **Acknowledgements**

596 Acknowledgement is made to the donors of the ADR A2021042S, a program of the BrightFocus
597 Foundation, for support of this research. Data were provided by OASIS-3: Longitudinal
598 Multimodal Neuroimaging (Principal Investigators: T. Benzinger, D. Marcus, J. Morris). OASIS-3
599 was supported by the following funding sources: NIH P50 AG00561, P30 NS09857781, P01
600 AG026276, P01 AG003991, R01 AG043434, UL1 TR000448, R01 EB009352. AV-45 doses
601 were provided by Avid Radiopharmaceuticals, a wholly owned subsidiary of Eli Lilly.

602 **References**

603 1. Chapleau M, Iaccarino L, Soleimani-Meigooni D, Rabinovici GD. The Role of Amyloid PET in
604 Imaging Neurodegenerative Disorders: A Review. J Nucl Med. 2022;63:13S-19S.

- 605 2. Swanson CJ, Zhang Y, Dhadda S, Wang J, Kaplow J, Lai RYK, et al. A randomized, double-
606 blind, phase 2b proof-of-concept clinical trial in early Alzheimer's disease with lecanemab, an
607 anti-A β protofibril antibody. *Alz Res Therapy*. 2021;13:80.
- 608 3. Shcherbinin S, Evans CD, Lu M, Andersen SW, Pontecorvo MJ, Willis BA, et al. Association
609 of Amyloid Reduction After Donanemab Treatment With Tau Pathology and Clinical Outcomes:
610 The TRAILBLAZER-ALZ Randomized Clinical Trial. *JAMA Neurology*. 2022;79:1015–24.
- 611 4. Mathotaarachchi S, Pascoal TA, Shin M, Benedet AL, Kang MS, Beaudry T, et al. Identifying
612 incipient dementia individuals using machine learning and amyloid imaging. *Neurobiology of*
613 *Aging*. 2017;59:80–90.
- 614 5. Choi H, Jin KH. Predicting cognitive decline with deep learning of brain metabolism and
615 amyloid imaging. *Behavioural Brain Research*. 2018;344:103–9.
- 616 6. Pfeil J, Hoenig MC, Doering E, van Eimeren T, Drzezga A, Bischof GN. Unique regional
617 patterns of amyloid burden predict progression to prodromal and clinical stages of Alzheimer's
618 disease. *Neurobiology of Aging*. 2021;106:119–29.
- 619 7. Pascoal TA, Therriault J, Mathotaarachchi S, Kang MS, Shin M, Benedet AL, et al.
620 Topographical distribution of A β predicts progression to dementia in A β positive mild cognitive
621 impairment. *Alzheimer's & Dementia: Diagnosis, Assessment & Disease Monitoring*.
622 2020;12:e12037.
- 623 8. Ezzati A, Abdulkadir A, Jack CR, Thompson PM, Harvey DJ, Truelove-Hill M, et al. Predictive
624 value of ATN biomarker profiles in estimating disease progression in Alzheimer's disease
625 dementia. *Alzheimer's & Dementia*. 2021;17:1855–67.

- 626 9. Klunk WE, Engler H, Nordberg A, Wang Y, Blomqvist G, Holt DP, et al. Imaging brain amyloid
627 in Alzheimer's disease with Pittsburgh Compound-B. *Annals of Neurology*. 2004;55:306–19.
- 628 10. Su Y, Flores S, Wang G, Hornbeck RC, Speidel B, Joseph-Mathurin N, et al. Comparison of
629 Pittsburgh compound B and florbetapir in cross-sectional and longitudinal studies. *Alzheimer's &*
630 *Dementia: Diagnosis, Assessment & Disease Monitoring*. 2019;11:180–90.
- 631 11. Wolk DA, Zhang Z, Boudhar S, Clark CM, Pontecorvo MJ, Arnold SE. Amyloid imaging in
632 Alzheimer's disease: comparison of florbetapir and Pittsburgh compound-B positron emission
633 tomography. *J Neurol Neurosurg Psychiatry*. 2012;83:923–6.
- 634 12. Landau SM, Breault C, Joshi AD, Pontecorvo M, Mathis CA, Jagust WJ, et al. Amyloid- β
635 Imaging with Pittsburgh Compound B and Florbetapir: Comparing Radiotracers and
636 Quantification Methods. *Journal of Nuclear Medicine*. 2013;54:70–7.
- 637 13. Villemagne VL, Mulligan RS, Pejoska S, Ong K, Jones G, O'Keefe G, et al. Comparison of
638 ^{11}C -PiB and ^{18}F -florbetaben for A β imaging in ageing and Alzheimer's disease. *Eur J Nucl*
639 *Med Mol Imaging*. 2012;39:983–9.
- 640 14. Adamczuk K, Schaefferbeke J, Nelissen N, Neyens V, Vandebulcke M, Goffin K, et al.
641 Amyloid imaging in cognitively normal older adults: comparison between ^{18}F -flutemetamol and
642 ^{11}C -Pittsburgh compound B. *Eur J Nucl Med Mol Imaging*. 2016;43:142–51.
- 643 15. Mountz JM, Laymon CM, Cohen AD, Zhang Z, Price JC, Boudhar S, et al. Comparison of
644 qualitative and quantitative imaging characteristics of ^{11}C PiB and ^{18}F flutemetamol in normal
645 control and Alzheimer's subjects. *NeuroImage: Clinical*. 2015;9:592–8.

- 646 16. Landau SM, Thomas BA, Thurfjell L, Schmidt M, Margolin R, Mintun M, et al. Amyloid PET
647 imaging in Alzheimer's disease: a comparison of three radiotracers. *Eur J Nucl Med Mol*
648 *Imaging*. 2014;41:1398–407.
- 649 17. Chen CD, McCullough A, Gordon B, Joseph-Mathurin N, Flores S, McKay NS, et al.
650 Longitudinal head-to-head comparison of 11C-PiB and 18F-florbetapir PET in a Phase 2/3
651 clinical trial of anti-amyloid- β monoclonal antibodies in dominantly inherited Alzheimer's disease.
652 *Eur J Nucl Med Mol Imaging*. 2023;50:2669–82.
- 653 18. Klunk WE, Koeppe RA, Price JC, Benzinger TL, Devous Sr. MD, Jagust WJ, et al. The
654 Centiloid Project: Standardizing quantitative amyloid plaque estimation by PET. *Alzheimer's &*
655 *Dementia*. 2015;11:1-15.e4.
- 656 19. Pegueroles J, Montal V, Bejanin A, Vilaplana E, Aranha M, Santos-Santos MA, et al. AMYQ:
657 An index to standardize quantitative amyloid load across PET tracers. *Alzheimer's & Dementia*.
658 2021;17:1499–508.
- 659 20. Bourgeat P, Doré V, Doecke J, Ames D, Masters CL, Rowe CC, et al. Non-negative matrix
660 factorisation improves Centiloid robustness in longitudinal studies. *NeuroImage*.
661 2021;226:117593.
- 662 21. Chen K, Ghisays V, Luo J, Chen Y, Lee W, Wu T, et al. Harmonizing florbetapir and PiB
663 PET measurements of cortical A β plaque burden using multiple regions-of-interest and machine
664 learning techniques: An alternative to the Centiloid approach. *Alzheimer's & Dementia*.
665 2024;20:2165–72.
- 666 22. Liu H, Nai Y-H, Saridin F, Tanaka T, O' Doherty J, Hilal S, et al. Improved amyloid burden
667 quantification with nonspecific estimates using deep learning. *Eur J Nucl Med Mol Imaging*.
668 2021;48:1842–53.

- 669 23. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using
670 empirical Bayes methods. *Biostatistics*. 2007;8:118–27.
- 671 24. Richter S, Winzeck S, Correia MM, Kornaropoulos EN, Manktelow A, Outtrim J, et al.
672 Validation of cross-sectional and longitudinal ComBat harmonization methods for magnetic
673 resonance imaging data on a travelling subject cohort. *Neuroimage: Reports*. 2022;2:100136.
- 674 25. Sun D, Rakesh G, Haswell CC, Logue M, Baird CL, O’Leary EN, et al. A comparison of
675 methods to harmonize cortical thickness measurements across scanners and sites.
676 *NeuroImage*. 2022;261:119509.
- 677 26. Pomponio R, Erus G, Habes M, Doshi J, Srinivasan D, Mamourian E, et al. Harmonization
678 of large MRI datasets for the analysis of brain imaging patterns throughout the lifespan.
679 *NeuroImage*. 2020;208:116450.
- 680 27. Fortin J-P, Cullen N, Sheline YI, Taylor WD, Aselcioglu I, Cook PA, et al. Harmonization of
681 cortical thickness measurements across scanners and sites. *NeuroImage*. 2018;167:104–20.
- 682 28. Leithner D, Schöder H, Haug A, Vargas HA, Gibbs P, Häggström I, et al. Impact of ComBat
683 Harmonization on PET Radiomics-Based Tissue Classification: A Dual-Center PET/MRI and
684 PET/CT Study. *Journal of Nuclear Medicine*. 2022;63:1611–6.
- 685 29. Bilgel M. Probabilistic estimation for across-batch compatibility enhancement for amyloid
686 PET. *Alzheimer’s & Dementia: Diagnosis, Assessment & Disease Monitoring*. 2023;15:e12436.
- 687 30. LaMontagne PJ, Benzinger TLS, Morris JC, Keefe S, Hornbeck R, Xiong C, et al. OASIS-3:
688 Longitudinal Neuroimaging, Clinical, and Cognitive Dataset for Normal Aging and Alzheimer
689 Disease. *medRxiv*. 2019;2019.12.13.19014902.

- 690 31. Su Y, D'Angelo GM, Vlassenko AG, Zhou G, Snyder AZ, Marcus DS, et al. Quantitative
691 Analysis of PiB-PET with FreeSurfer ROIs. PLOS ONE. 2013;8:e73377.
- 692 32. Beer JC, Tustison NJ, Cook PA, Davatzikos C, Sheline YI, Shinohara RT, et al. Longitudinal
693 ComBat: A method for harmonizing longitudinal multi-scanner imaging data. NeuroImage.
694 2020;220:117129.
- 695 33. Su Y, Flores S, Hornbeck RC, Speidel B, Vlassenko AG, Gordon BA, et al. Utilizing the
696 Centiloid scale in cross-sectional and longitudinal PiB PET studies. NeuroImage: Clinical.
697 2018;19:406–16.
- 698 34. Properzi MJ, Buckley RF, Chhatwal JP, Donohue MC, Lois C, Mormino EC, et al. Nonlinear
699 Distributional Mapping (NoDiM) for harmonization across amyloid-PET radiotracers.
700 NeuroImage. 2019;186:446–54.
- 701 35. Su Y, Blazey TM, Snyder AZ, Raichle ME, Marcus DS, Ances BM, et al. Partial volume
702 correction in quantitative amyloid imaging. NeuroImage. 2015;107:55–64.
- 703 36. Stein CK, Qu P, Epstein J, Buros A, Rosenthal A, Crowley J, et al. Removing batch effects
704 from purified plasma cell gene expression microarrays with modified ComBat. BMC
705 Bioinformatics. 2015;16:63.
- 706 37. Jack CR, Holtzman DM. Biomarker Modeling of Alzheimer's Disease. Neuron.
707 2013;80:1347–58.
- 708 38. Bilgel M, An Y, Zhou Y, Wong DF, Prince JL, Ferrucci L, et al. Individual estimates of age at
709 detectable amyloid onset for risk factor assessment. Alzheimer's & Dementia. 2016;12:373–9.
- 710