



HAL
open science

A parsimonious model for mass-univariate vertex-wise analysis

Baptiste Couvy-Duchesne, Futao Zhang, Kathryn E Kemper, Julia Sidorenko, Naomi R Wray, Peter M Visscher, Olivier Colliot, Jian Yang

► **To cite this version:**

Baptiste Couvy-Duchesne, Futao Zhang, Kathryn E Kemper, Julia Sidorenko, Naomi R Wray, et al..
A parsimonious model for mass-univariate vertex-wise analysis. 2022. hal-03118366v2

HAL Id: hal-03118366

<https://hal.science/hal-03118366v2>

Preprint submitted on 12 May 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A parsimonious model for mass-univariate vertex-wise analysis

Baptiste Couvy-Duchesne,^{a,b,*} Futao Zhang,^a Kathryn E. Kemper,^a Julia Sidorenko,^a
Naomi R. Wray,^{a,†} Peter M. Visscher,^{a,†} Olivier Colliot,^{b,†} Jian Yang^{a,c,d,†}

^aInstitute for Molecular Bioscience, the University of Queensland, St Lucia, QLD

^bSorbonne University, Paris Brain Institute (ICM), CNRS, Inria, Inserm, AP-HP, Hôpital de la Pitié Salpêtrière, F-75013, Paris, France

^cSchool of Life Sciences, Westlake University, Hangzhou, Zhejiang, China

^dWestlake Laboratory of Life Sciences and Biomedicine, Hangzhou, Zhejiang, China

[†]These authors contributed equally

Abstract

Purpose: Covariance between grey-matter measurements can reflect structural or functional brain networks though it has also been shown to be influenced by confounding factors (e.g. age, head size, scanner), which could lead to lower mapping precision (increased size of associated clusters) and create distal false positives associations in mass-univariate vertex-wise analyses.

Approach: We evaluated this concern by performing state-of-the-art mass-univariate analyses (general linear model, GLM) on traits simulated from real vertex-wise grey matter data (including cortical and subcortical thickness and surface area). We contrasted the results with those from linear mixed models (LMMs), which have been shown to overcome similar issues in omics association studies.

Results: We showed that when performed on a large sample (N=8,662, UK Biobank), GLMs yielded greatly inflated false positive rate (cluster false discovery rate>0.6). We showed that LMMs resulted in more parsimonious results: smaller clusters and reduced false positive rate but at a cost of increased computation. Next, we performed mass-univariate association analyses on five real UKB traits (age, sex, BMI, fluid intelligence and smoking status) and LMM yielded fewer and more localised associations. We identified 19 significant clusters displaying small associations with age, sex and BMI, which suggest a complex architecture of at least dozens of associated areas with those phenotypes.

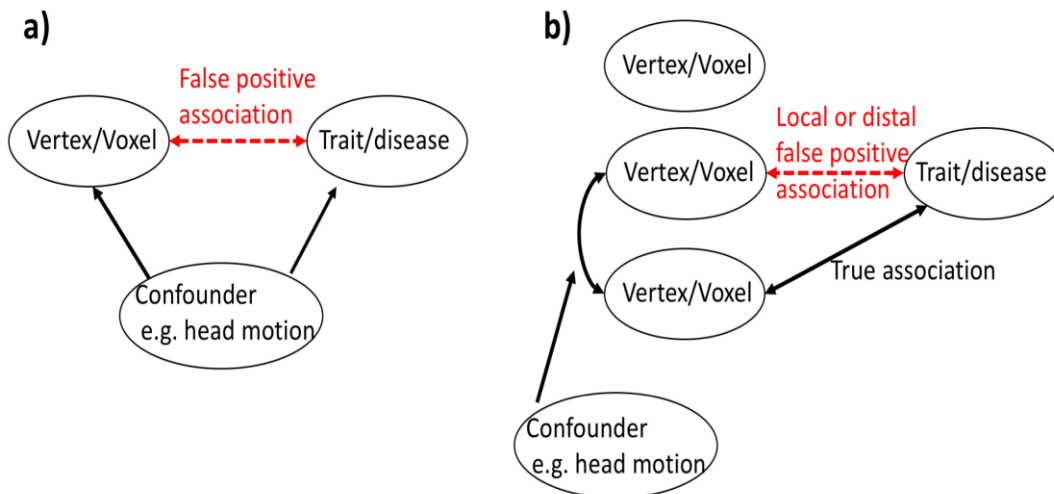
Conclusions: The published literature could contain a large proportion of redundant (possibly confounded) associations, that are largely prevented using LMMs. The parsimony of LMMs results from controlling for the joint effect of all vertices, which prevents local and distal redundant associations from reaching significance.

Keywords: structural brain MRI, vertex-wise processing, linear mixed model, association, brain mapping.

* Baptiste Couvy-Duchesne (b.couvyduchesne@uq.edu.au)

1 Introduction

43 Brain MRI scans can generate hundreds of thousands of vertex/voxel-wise measurements per
 44 individual, which can be linked to other measured traits/diseases using mass univariate
 45 vertex/voxel-wise association analyses. Results of association analyses (and subsequent follow-
 46 up analyses) can shed light on the brain networks or cell composition relevant for the
 47 trait/disease and may be leveraged for brain-feature based phenotype prediction. However, brain
 48 measurements may exhibit a pattern of correlation, owing to factors (e.g. head size, MRI
 49 scanner/artefact (1) or demographics (2)) which can generate confounded brain-trait associations.
 50 Induced local correlations with a true brain-biomarker can generate a smear of association (i.e. a
 51 cluster of associated vertices) which may limit the precise localisation of the directly associated
 52 regions. On the other hand, long-range vertex correlations caused or inflated by factors irrelevant
 53 to the trait of interest, may be more prejudicial, as they can yield distal false positives (**Figure 1**).



54

55 **Figure 1:** Illustration of the traditional confounding paradigm a) and of the confounding that may arise in association studies performed across
 56 correlated brain features b).

57 One sided arrows represent a causal effect, and two-sided arrows a correlation.

58

59 Two approaches can be used to limit the inflation of false positives described above. One is to
60 control for the confounders in the association testing, although it requires knowledge and
61 measurement of the factors influencing (or more generally associated with) the covariance
62 between brain measurements. Note that these factors can overlap with traditional confounders of
63 neuroimaging studies (e.g. head size, age, sex, head motion), and additional confounders are
64 being identified as sample sizes increase (3). Another correction strategy is to control for the
65 other vertices in the association testing, in order to remove the signal that could be attributed to
66 another brain vertex or region. The difficulty of such approach is that typically, the number of
67 vertex/voxel-wise measurements (p) far exceeds the number of participants (N) in the study. The
68 $p \gg N$ paradigm implies that the marginal joint associations with all p vertices cannot be
69 estimated in a single general linear model (GLM).

70 Statistically, the challenge of mass univariate vertex-wise analyses resembles that of genome-
71 wide association studies (GWAS) or methylation-wide associations studies (MWAS), which aim
72 to identify genomic regions associated with a phenotype in the presence of correlated features
73 (i.e., genetic variants or DNA methylation probes). Several studies have demonstrated that
74 feature correlation (i.e., Linkage Disequilibrium (LD) or population structure in genetics) can
75 result in inflated false positive rate (4-6), even more so when the sample size increases (5). This
76 led GLMs to be replaced by linear mixed models (LMMs) (6-8) which co-varies out all features
77 by fitting them as random effects. LMMs have been shown to better control the inflation of false
78 positive associations arising from LD or correlation between probes and to minimise the
79 occurrence of false positives in both GWAS and MWAS (6, 7, 9).

80 LMMs are commonly used in neuroimaging to model longitudinal data(10). Instead, we rely
81 here on a novel formulation that allows fitting the high-dimensional brain image as a single

82 random-effect. Such LMMs allow estimation of the overall degree of association between a trait
83 and a high-dimensional brain image, coined “morphometricity” in the context of structural brain
84 measurements(11, 12). Recently, we have shown that a single LMM framework was suited to
85 estimate morphometricity in large datasets, to draw links between traits through their
86 associations with similar brain structure (grey-matter correlation) and to build brain-based
87 predictors(12). The LMMs we propose here complement our previous work by identifying the
88 vertices/voxels that contribute to the morphometricity and phenotype prediction.

89 Here, we sought to evaluate whether the inflation of false positives observed in omics data is
90 also present in neuroimaging data. In the first part of the analysis, we performed extensive
91 simulations of continuous phenotypes from real grey-matter data to quantify false positive rate as
92 well as statistical power, mapping precision and prediction accuracy achieved from mass-
93 univariate analyses. We compared the performances of the current state-of-the-art GLMs to that
94 of LMMs inspired by omics association studies. In the second part, we sought to characterise the
95 brain regions associated with real phenotypes (i.e., age, sex, BMI, fluid IQ, and smoking status)
96 that previously exhibited significant morphometricity(12), in order to confirm the results
97 obtained on simulated traits. Our analyses relied on 14,451 MRI images collected by the UK
98 Biobank (UKB), one of the largest brain imaging initiative (13).

99 *1.1 Novelties and contribution*

100 The novelties and contributions of our paper are as follows:

- 101 • We propose novel linear mixed models for brain mapping, inspired from those using in
102 genetics, which aims at overcoming false positive issues found in standard analyses.

- 103 • By controlling for all brain measurements (fitted as a random effect) the LMMs remove
104 redundant associations leading to more parsimonious results.
- 105 • We demonstrate that, compared to the current state-of-the-art, the LMMs minimise false
106 positive rate while also maximising power, mapping precision and prediction accuracy.

107

108 **2. Material and methods**

109 *2.1. Models of mass-univariate vertex wise analyses*

110 First, we considered five GLMs that differ in term of covariates used when estimating the
111 association (b_i) between the trait and the i th (standardised) vertex-wise measurement (\mathbf{X}_i). They
112 can be written under the form:

$$113 \quad \mathbf{y} = \mathbf{Z}\mathbf{c} + \mathbf{X}_i b_i + \boldsymbol{\varepsilon} \quad (1)$$

114 with \mathbf{y} the vector of phenotype for the N individuals, \mathbf{Z} a matrix of size $N \times q$ of q covariates
115 and \mathbf{c} a vector of the q fixed effects.

116 The five GLMs are differentiated as follows: 1) GLM with no covariates (“no covariates”), 2)
117 GLM including the most commonly used covariates in similar analyses: age, sex and intra-
118 cranial volume (ICV) (“age, sex, ICV corrected”), 3) & 4) GLMs including 5 and 10 principal
119 components (PCs) of grey-matter variation, respectively (“5 global PCs”, “10 global PCs”), 5)
120 GLM including 10 PCs specific to the measurement type (cortical thickness, cortical surface,
121 subcortical thickness or subcortical surface area), referred to as “10 modality specific PCs”.
122 Grey-matter PCs capture the main axes of covariations between vertices, and we expect that by
123 controlling for them we may be able to remove unmeasured or unknown factors contributing to
124 long-range correlation between vertices (which might include demographics, MRI machine, head

125 motion, software update, processing option *etc.*). Note that PCs from genetic data are commonly
 126 used in GWAS in order to limit the false positive rate of GLMs analyses (14) but are rarely used
 127 in neuroimaging analyses. The difficulties of PC correction are to determine the optimal number
 128 of PCs, which controls for confounding effects without removing signals of interest. In practice,
 129 this may prove extremely difficult considering that the optimal number of PCs could depend on
 130 the trait/variable of interest, and that PCs are notoriously hard to interpret and have not been
 131 comprehensively investigated on these data. Thus, we arbitrarily chose two scenarios with the
 132 first 5 or 10 PCs. In addition, GLMs without covariates are also very rare, but worth considering
 133 in order to appreciate the effect of including covariates.

134 Finally, we considered three LMMs that can be seen as extensions of the previous approaches
 135 in that they further control for all vertex-wise measurements. The first LMM model (“LMM
 136 global BRM”), analogous to the MOA (MLM-based Omic Association) model (6), can be
 137 written as:

$$138 \quad \mathbf{y} = \mathbf{X}_i \mathbf{b}_i + \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (2)$$

139 Here, \mathbf{X} is the $N \times p$ matrix of all standardised vertex-wise measurements, $\boldsymbol{\beta}$ is the $p \times 1$ vector
 140 of joint vertex-trait associations. $\boldsymbol{\beta}$ is a vector of random effects, allowing for $p > N$, with
 141 $\boldsymbol{\beta} \sim \mathcal{N}(0, \mathbf{I} \sigma_{\boldsymbol{\beta}}^2)$, and $\boldsymbol{\varepsilon}$ is the error term assumed to follow $\boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I} \sigma_{\boldsymbol{\varepsilon}}^2)$. $\sigma_{\boldsymbol{\beta}}^2$ and $\sigma_{\boldsymbol{\varepsilon}}^2$ are the
 142 variances of the random effects $\boldsymbol{\beta}$ and $\boldsymbol{\varepsilon}$. The variance-covariance matrix for \mathbf{y} is $\text{var}(\mathbf{y}) = \mathbf{V} =$
 143 $\mathbf{X} \mathbf{X}' \sigma_{\boldsymbol{\beta}}^2 + \mathbf{I} \sigma_{\boldsymbol{\varepsilon}}^2 = \mathbf{B} p \sigma_{\boldsymbol{\beta}}^2 + \mathbf{I} \sigma_{\boldsymbol{\varepsilon}}^2$. Here, we regard $\mathbf{B} = \mathbf{X} \mathbf{X}' / p$ as the brain relatedness matrix and
 144 $p \sigma_{\boldsymbol{\beta}}^2$ the morphometricity (proportion of phenotypic variance captured by all vertices) (15).

145 We considered a second LMM (“LMM with covariates”) that includes known covariates (age,
146 sex and ICV) fitted as fixed effects. Thus, we can separate the effect of the random effects from
147 that of the known covariates on the results. The model becomes:

$$148 \quad \mathbf{y} = \mathbf{Zc} + \mathbf{X}_i \mathbf{b}_i + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (3)$$

149 Our third LMM (“LMM multi. BRM”) includes 4 random effects ($\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \boldsymbol{\beta}_3, \boldsymbol{\beta}_4$), each
150 corresponding to a type of vertices (cortical thickness, cortical surface area, subcortical thickness
151 and subcortical surface area).

$$152 \quad \mathbf{y} = \mathbf{X}_i \mathbf{b}_i + \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{X}_2 \boldsymbol{\beta}_2 + \mathbf{X}_3 \boldsymbol{\beta}_3 + \mathbf{X}_4 \boldsymbol{\beta}_4 + \boldsymbol{\varepsilon} \quad (4)$$

153 This more general LMM allows the distribution of effect sizes to differ based on vertex type,
154 rather than enforcing a single distribution over all types of measurements (15). Note that each
155 random effect takes up a single degree of freedom meaning that LMMs and GLMs have a
156 comparable (large) numbers of degrees of freedom given the same sample size.

157

158 *2.2. Statistical testing and multiple comparison*

159 We performed a χ^2 test of the association between a vertex (\mathbf{X}_i) and the phenotype using that, for
160 large sample size N , $\left(\frac{b_i}{SE(b_i)}\right)^2 \sim \chi_1^2$ under the null hypothesis of no association. In each model
161 (GLM or LMM), we accounted for multiple testing over the vertices using Bonferroni correction,
162 thus setting a brain-wide significance threshold of $0.05/652,283=7.6e-8$. We chose the
163 straightforward Bonferroni correction over random field theory (RFT)(16) as RFT requires
164 stationarity and a smooth mesh of vertex-wise residuals, which is unlikely to be the case here
165 (we did not apply kernel smoothing on the data as it reduced the estimated morphometricity of
166 the UKB phenotypes (15)). In addition, RFT is not currently implemented to be performed using

167 residuals of LMMs or across several surfaces and type of measurements. Bonferroni correction is
168 expected to be conservative under the null hypothesis (no association) because the correlations
169 between vertices means that the effective number is tests lower than the number of tests
170 conducted and used for the Bonferroni correction.

171 *2.3. MRI Image processing*

172 MRI images were mostly collected in Cheadle (for 96% of the sample) and Newcastle using a
173 3T Siemens Skyra machine (software platform VD13) and a 32-channel head coil (13) (see
174 **Supp. 1**, for MRI sequence details).

175 We processed the T1w and T2 FLAIR images together to enhance the tissue segmentation in
176 FreeSurfer 6.0 (17), which should result in a more precise skull stripping and pial surfaces
177 definition. When the T2 FLAIR was not acquired or not usable, we processed the T1w image
178 alone, though a recent report showed this results in systematic differences in cortical thickness
179 (18). This may represent a source of noise in the data, albeit it was limited in term of number of
180 individuals (see quality control, **Supp. 1**). We extracted vertex-wise data mapping cortical
181 surface area and thickness (“recon-all” processing in FreeSurfer) and used the maximal
182 resolution allowed by the software (fsaverage atlas - unsmoothed). In short, FreeSurfer segments
183 the grey/white and grey/cerebrospinal fluid borders, which delimitate the grey-matter. Surfaces
184 are mapped onto a spherical atlas to align the cortical folding patterns of the individuals, and a
185 tessellation is applied. Cortical thickness is calculated as the closest distance from the two grey-
186 matter boundaries, for each vertex on the tessellated surface (19). Surface area is measured as the
187 mean area of all faces that meet at a particular vertex, on the grey/white matter surface(20). We
188 previously showed that this cortical processing maximised the morphometricity for a wide range

189 of phenotypes (15). In other words, this cortical processing maximised the information retained
190 by the processed MRI images. In addition, we applied the ENIGMA-shape processing (21, 22),
191 where subcortical structures segmented in FreeSurfer are projected onto spherical atlases to
192 quantify vertex-wise radial thickness and log Jacobian determinant (21, 22), which is analogous
193 to a surface area (23). This yielded a vertex-wise characterization of the hippocampus, putamen,
194 amygdala, thalamus, caudate, pallidum and accumbens. Overall, the imaging data used in the
195 analyses comprised 652,283 vertex measurements per individual: 299,009 for cortical thickness,
196 another 299,034 for cortical surface area, 27,120 for subcortical thickness and 27,120 for
197 subcortical surface area.

198 In a post-hoc analysis, we also utilised smoothed cortical data (surface based kernel with
199 FWHM=20mm), in order to evaluate the robustness of our results to variation in the MRI
200 processing.

201 *2.4. Main sample for simulation and discovery*

202 Our final sample comprised 9,890 adults with complete cortical and subcortical data, aged 62.5
203 on average (SD=7.5, range 44.6–79.6) with slightly more (52.4%) female participants (see Supp.
204 1 for participant inclusion and exclusion). Of note, 341 participants did not have an exploitable
205 T2 image.

206 We performed a stringent quality control (QC) to exclude one of each pair of individuals whose
207 brains were too similar or dissimilar relative to most other individuals, resulting in 1,228
208 exclusions (12.4% of the sample). The main reason for this exclusion was to prevent bias in the
209 LMM estimates, although it should also remove individuals flagged as outliers by other QC
210 criteria (e.g. 80.6% of the participants processed using T1w only, spike-like cortical parcellation

211 in FreeSurfer)(12) (see **Supp. 1** for more details on QC). Importantly, all analyses were
212 performed on the same list of individuals (post QC) to ensure that performance of the models
213 would be comparable.

214 *2.5. Independent samples for prediction and replication*

215 Our first independent sample included an additional 4,942 participants of the UKB with a T1w
216 image (downloaded in May 2018, most participants also had an exploitable T2w). The final
217 sample (N=4,160 after processing and QC) was on average 63.1 years old (SD=7.46, range 46.1-
218 80.3) with 52.1% of females.

219 In addition, we used the OASIS3 (Open Access Series of Imaging Studies) sample (24) to
220 evaluate the generalizability of the prediction. The OASIS3 dataset gathers several longitudinal
221 MRI studies conducted in the Washington University Knight Alzheimer Disease Research
222 Center over the past 15 years. Our final sample included 1,006 unique participants after
223 processing based on T1w images and QC. When several visits were available for a participant,
224 we selected the one with the most phenotypic information. Participants were 71.1 years old on
225 average (SD=9.18, range 42.6-95.7) and mostly female (55.5%). Almost a quarter of the
226 participants (23.6%) had a diagnosis of Alzheimer's disease at the time of imaging.

227 *2.6. Mass-univariate analyses on simulated phenotypes*

228 *2.6.1. Simulation of phenotypic traits from real grey-matter data*

229 We simulated phenotypic traits from the UKB processed (standardised) grey-matter data, instead
230 of relying on synthetic/simulated images. This novel approach ensures the vertex-wise data
231 retains a realistic correlation structure. In addition, our framework includes simulation of the

232 phenotype under not only the null hypothesis (“H0”) that no vertex is associated with the
233 phenotype but also the alternative hypothesis (“H1”) that a set of vertices are truly associated
234 with the phenotype.

235 First, we randomly selected a set of associated vertices and drew their relative effects from a
236 normal distribution. We then calculated the simulated phenotypes as a linear combination of the
237 individuals’ vertex values and noise (6). We considered three scenarios that differ in term of
238 number of associated vertices and total association with the phenotype. This global association
239 between grey-matter measurements and a trait has been coined morphometricity (11, 15) and
240 may be expressed as the proportion of the trait variance (R^2) captured by the vertex-wise
241 measurement. Our scenarios were: i) 10 associated vertices accounting for a phenotype
242 morphometricity of $R^2=0.20$ (i.e. 20% of the trait variance); ii) 100 associated vertices with
243 $R^2=0.50$; iii) 1000 vertices with $R^2=40\%$. For each scenario, we simulated 100 phenotypes.

244 In follow-up analyses, we simulated phenotypes using the same parameters, this time
245 restricting the associated vertices to a single type of measurement. This allowed evaluation of the
246 specificity of each type of measurement, which possess a unique correlation pattern. In addition,
247 this ensures our phenotypes were not associated with cortical vertices only, which represent 90%
248 of the vertex-wise measurements.

249 To evaluate the effect of smoothing on our results, we simulated phenotypes from smoothed
250 brain maps. For the ease of computation, we restricted the analysis of smoothed data to the case
251 of 10 associated vertices ($R^2=0.2$). We kept the same associated vertices (and weights) as in the
252 previous simulation from unsmoothed data. Finally, we randomly simulated 100 “null” traits, in
253 order to evaluate the calibration of the models under the null hypothesis of no association. All
254 simulations were generated using the OSCA software (6).

255

256 *2.6.2. Inflation of test statistics*

257 First, we compared the empirical distribution of χ^2 statistics to the expected distribution,
258 which is assumed to follow a $\chi^2(1)$ for non-associated (null) vertices. We considered the ratio of
259 empirical over expected median χ^2 , known as the inflation factor (λ), which is expected to be
260 equal to one across non-associated vertices. We also used the nominal false positive rate (FPR)
261 defined as the proportion of null vertices with p-values < 0.05 (expected to be 0.05). Correlation
262 between associated and null vertices (e.g. due to confounding factors) typically result in an
263 inflation of test statistics, which may cause null vertices to reach significance in mass-univariate
264 analyses.

265 *2.6.3. Discoverability and mapping precision*

266 First, we quantified the model discoverability using the true positive rate (TPR) defined as the
267 proportion of truly associated vertices reaching significance (after Bonferroni correction).
268 Importantly, the TPR is dependent on the false positive rate, which can limit comparison across
269 models (see statistical power below). In addition, we quantified the mapping precision of mass-
270 univariate analyses by reporting the median size of the true positive (TP) clusters. We defined TP
271 clusters as sets of significant contiguous vertices of the mesh that contain a true positive vertex.

272 *2.6.4. False positives and statistical power*

273 We reported the Family-Wise Error Rate (FWER) defined as the proportion of replicates with
274 at least one false positive vertex (null vertex significant after Bonferroni correction). In the
275 presence of strong correlation between neighboring vertices, it is statistically difficult to separate

276 a true positive vertex from the flanking ones, thus we can expect a FWER greater than 5%.
277 Hence, we also reported the cluster FWER defined as the proportion of replicates with at least
278 one false positive cluster. FWER is more stringent than False Discovery Rate (FDR), implying
279 that any false positives that remain after FWER correction would also be observed using FDR.

280 To account for the models' differences in FWER, we further reported the statistical power,
281 defined as the TPR for a set risk alpha. We chose cluster FWER<0.2, which was easier to
282 achieve than the traditional FWER<0.05, as we enforced comparable FWER by iteratively
283 lowering the significance threshold, for each of the models (**Appendix 2**). The choice of risk
284 alpha does not impact the relative performance of the models, and we can expect models best
285 powered for FWER<0.2 to also be best powered at other FWER levels.

286 Finally, we reported the proportion of false positive clusters out of all significant clusters
287 (cluster FDR). We labelled false positive clusters, the groups of significantly associated,
288 contiguous vertices that did not contain a true positive association.

289 In follow up analyses, we simulated associations on a single type of vertex-wise
290 measurements, in order to evaluate the probability of false positive (FWER) arising on the same
291 type of measurements, other types of measurements as well as contra-lateral regions.

292 *2.6.5. Prediction from significant vertices*

293 We evaluated the prediction accuracy achieved from the brain regions reaching significance,
294 in the different mass-univariate models. We used prediction as a meta-criterion to compare the
295 model performances, as it is dependent on power, true and false positives, and association effect
296 sizes. We selected the most significant vertex in each cluster and constructed a linear predictor
297 using association weights (\hat{b}_i , see (1) and (2)) estimated from the different mass-univariate

298 analyses. Because some significant clusters might contain several independent signals, we also
299 built predictors that included all significant vertices. We evaluated the prediction of in the
300 independent UKB and OASIS3 samples.

301 *2.6.6. Mass-univariate analyses of UK Biobank phenotypes*

302 Next, we performed mass-univariate vertex-wise analyses on five UKB phenotypes that
303 showed significant replicated morphometricity (15): age, sex, BMI, smoking status and fluid
304 intelligence. We used the raw fluid intelligence score provided by the UKB,
305 a non-standard test which has demonstrated some reliability in a test-retest analysis (25).

306 For each UKB phenotype and model, we reported the number of significant vertices, number
307 of significant clusters as well as their sizes. We defined significance using a Bonferroni
308 significance threshold of $0.05/(652283*5)=1.5e-8$, which accounts for the total number of tests
309 performed. For those phenotypes, the true pattern of association is unknown which prevents
310 evaluation of the false positive rate (or power) of the different approaches. However, false
311 positives or redundant associations should not improve prediction accuracy. In this regard, we
312 evaluated each GLM or LMM model in both the UKB replication and OASIS3 datasets. As
313 above, we used linear predictors, and reported the prediction accuracy (correlation) controlling
314 for age, sex, ICV and site. In OASIS3, we also corrected for clinical status (Alzheimer's disease
315 and mild cognitive impairment).

316 **3. Results**

317 *3.1. Phenotypes simulated under H0*

318 We found that all GLM and LMM models behaved well under the null hypothesis, as
319 indicated by no inflation of test statistic, FPR, or of false positive rate (FWER). As expected
320 under a stringent Bonferroni correction, all approaches were conservative as indicated by
321 $FWER < 3\%$ (**SFig. 1**).

322 *3.2. Phenotypes simulated under H1*

323 *3.2.1. Inflation of test statistics*

324 First, we quantified whether we could observe an inflation of test statistics on the vertices not
325 associated with the simulated phenotypes. As expected in presence of correlation between truly
326 associated and null vertices, we observed a global inflation of (median) test statistics when using
327 GLMs (**Figure 2, STable 1**). This was confirmed by an FPR greater than 5% for all GLM
328 models even though controlling for covariates or PCs, reduced the inflation of test-statistics
329 compared to the “no covariates” GLM. In comparison, LMMs appropriately controlled the
330 inflation of test statistics on null vertices ($\lambda < 1$ and $FDR < 5\%$; **Figure 2, STable 1**).

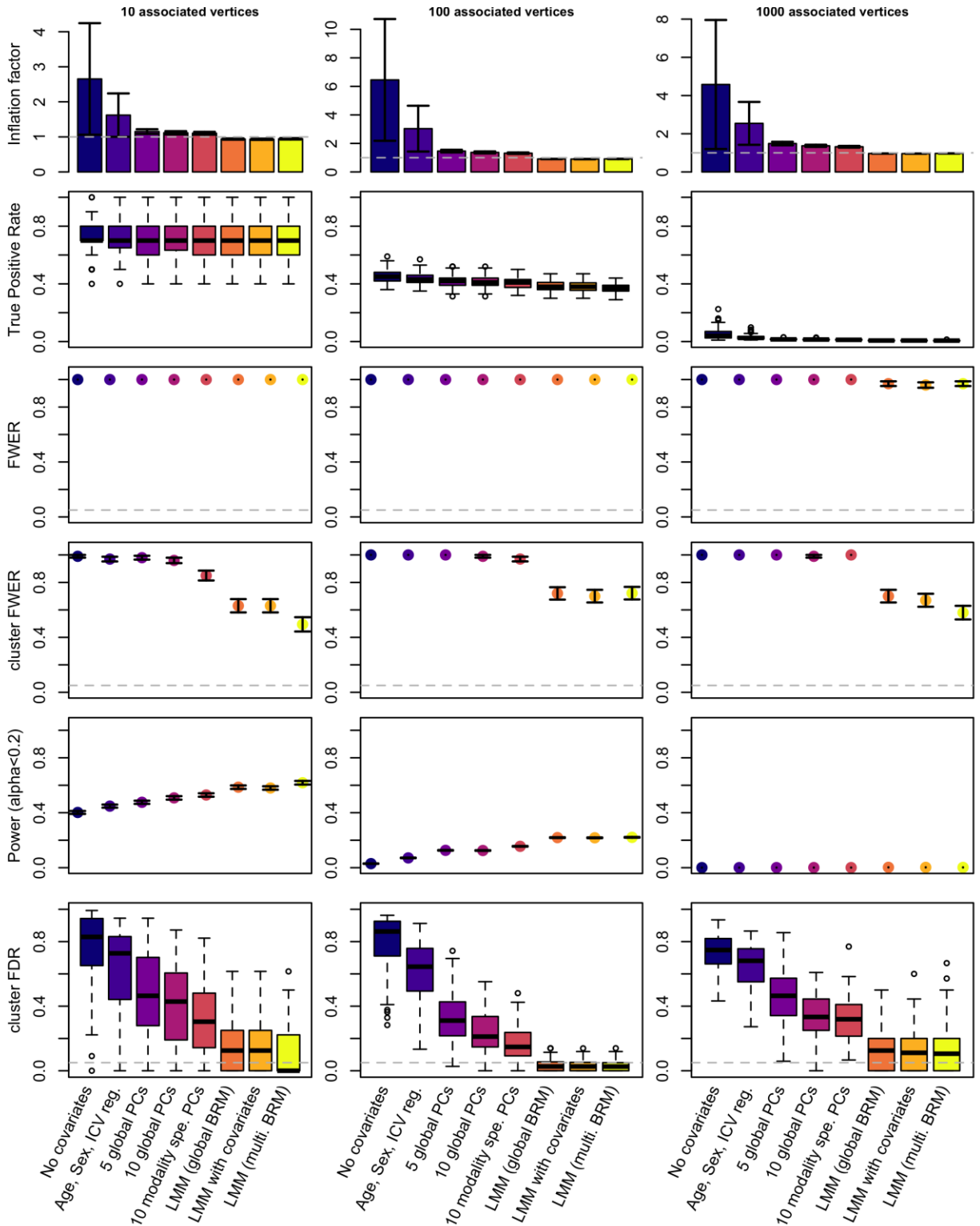
331 *3.2.2. True Positive Rate*

332 First, we confirmed that the TPR (after Bonferroni correction) was dependent on the scenarios
333 which corresponded to different effect sizes for the vertices. For example, about 70% of the truly
334 associated brain regions were detected in the case of a simple trait (10 associated vertices each
335 accounting for 2% of the phenotypic variance on average). On the other hand, less than 5% of

336 the associated brain regions were identified for the most complex phenotypes (scenario 3, 1000
337 vertices each accounting for 0.04% of variance, **Figure 2, STable 1**).

338 Across all scenarios, LMMs exhibited a slightly reduced TPR compared to the GLMs (**Figure**
339 **2, STable1**). We investigated this result using phenotypes simulated from a single type of
340 measurement. We found TPR of LMMs to be especially reduced on subcortical thickness and
341 surface area (**SFig. 2**).

342



343
 344 **Figure 2:** Performance of GLMs and LMMs for mass-univariate vertex-wise analyses: test inflation, statistical power and false positive rate.
 345 The columns correspond to the different scenarios considered when simulating traits. We simulated 100 phenotypic traits for each scenario. Bars
 346 represent +/- SE across the 100 replicates. Clusters are composed of groups of contiguous vertices each significantly associated with the
 347 phenotype (after Bonferroni correction). We labelled them as false positives if they did not include a true positive association.
 348

349

350 *3.2.3 False positives*

351 Here, we evaluated the occurrence of false positive vertices or clusters from our simulations. We
352 found that every single simulation yielded at least 1 false positive vertex after Bonferroni
353 correction (FWER=1, **Figure 2**). We noted that the FWER of 0.97 (SE=0.02) found for LMMs
354 in the scenario of “1000 associated vertices”, came from three simulations returning no
355 significant associations.

356 When evaluating the results at a cluster level, we found that using GLMs almost always
357 resulted in one or more false positive cluster (**Figure 2, STable 1**), leading to cluster
358 FWER>85%. Cluster FWER was reduced to 49-72% by using LMMs (**Figure 2, STable 1**).
359 Despite this improvement, no model ensured a cluster-FWER below 5%. LMMs also minimised
360 the proportion of false positive clusters (cluster FDR), compared to the GLM approaches. At the
361 extreme, more than 70% of the significant clusters were false positives using GLMs without
362 covariate. This reduced to about 60% when controlling for age, sex and ICV and further reduced
363 to less than 17% using LMMs (**Figure 2, STable1**).

364 Next, we simulated phenotypes associated with a single type of measurement and reported the
365 FWER for each type of measurement in **SFig. 3-6**. This allowed evaluation of whether false
366 positives could appear as a result of associations with vertices from other types of measurements.
367 We found that using GLMs resulted in contamination of signal between all the different types of
368 measurements, as indicated by FWER>5% (**SFig. 3-6**). In comparison, LMMs always minimised
369 the probability of false positives appearing on non-associated types of measurement. In

370 particular, LMMs ensured that associations on the cortex did not inflate the false positive rate on
371 subcortical structures, and vice versa (FWER<5%).

372

373 *3.2.4. Statistical power*

374 We found that the models differ in terms of false positive rate, which limits the direct
375 comparison of TPR. Instead, we reported the statistical power, which consists in the TPR for a
376 fixed level of FWER (cluster FWER<0.2). We found the LMMs to be more powerful than the
377 GLMs (**Figure 2, Supp. 2**).

378 *3.2.5. Mapping precision*

379 We defined mapping precision as the median size of the true positive clusters. LMMs led to a
380 more precise localisation of the associations by minimising the size of true positive clusters
381 (whether we looked at clusters median or maximal size, **Figure 3, STable1**).

382 The median size of true positive clusters was reduced by a factor greater than ten on subcortical
383 measurements, and by a factor greater than two on cortical thickness when using LMMs (**STable**
384 **1**). Of note, positive clusters on cortical surface area were particularly small (most clusters were
385 composed of a single vertex), independent of the model used, **Figure 3, STable 1**). However,
386 LMMs still offered a greater precision than the GLMs when considering the maximal cluster size
387 (**STable 1, SFig. 3-6**).

388 *3.2.6. Prediction accuracy from significant vertices*

389 As a way of aggregating the previous metrics of performance, we compared prediction
390 accuracy achieved from significant vertices, using the UKB replication sample. Across all
391 models and scenarios, selecting the top vertex per significant cluster maximised prediction

392 accuracy, compared to including all significant vertices. This was expected, as significant
393 vertices from the same cluster tag likely redundant information, leading to overweight the
394 prediction signal coming from large clusters.

395 In simulation scenarios 1 and 2, we found that including more covariates in the GLMs
396 resulted in greater prediction accuracy despite that predictors included fewer vertices (**Figure 3,**
397 **STable 1**). In addition, LMMs yielded marginally better prediction accuracy than the best GLM
398 using even fewer vertices (**Figure 3, STable 1**), consistent with observation from previous
399 studies (6, 9). For the third simulation scenario, the prediction accuracy was comparable and
400 limited for all models (**Figure 3, STable 1**).

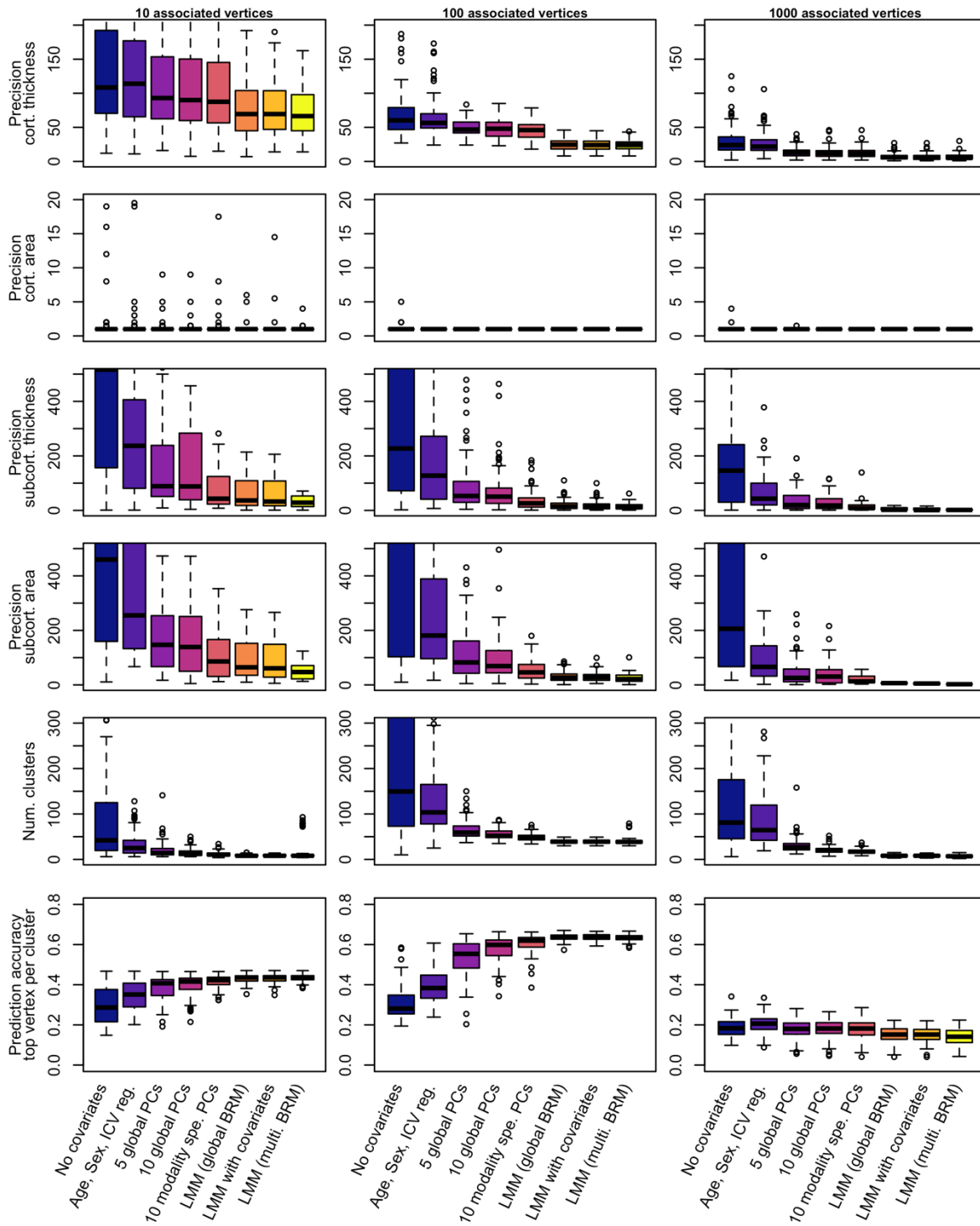
401 *3.2.7. Analyses using smoothed cortical surfaces*

402 We repeated the analysis using smoothed cortical meshes of surface and thickness
403 (FWHM=20mm), which is more commonly used in the literature than unsmoothed meshes
404 (STable 4-8). We sought to investigate how robust our results were to such variation of MRI
405 processing.

406 Overall, smoothing did not change the results of the model comparison. LMMs resulted again
407 in a reduced false positive rate (lower cluster FWER and cluster FDR) as well as reduced power
408 (seemingly more important than in the unsmoothed case). LMMs maximised mapping precision
409 and prediction accuracy, despite relying on fewer significant clusters (**SFig. 7**). Of note,
410 performing analyses on smoothed data decreased the mapping precision, leading to true positive
411 clusters roughly ten times larger on cortical meshes (**Figure 2, SFig. 7**).

412 Data smoothing resulted in a large inflation of test statistic and FPR for GLMs (**Figure 2,**
413 **SFig. 7**), which is to be expected as smoothing increases the amount of correlation between

414 vertices. We noticed that smoothing led to an increase of cluster FWER for the GLM with 10
415 PCs, while it decreased cluster FWER for the LMMs (despite the associated vertices and effect
416 sizes remaining the same). This result warrants a more fined-grained evaluation of the
417 associations. We can only hypothesise that the 20mm (FWHM) smoothing can induce medium-
418 range correlations (hence medium range false positives in GLMs) while it also increases local
419 correlation which might aggregate false positive clusters in LMMs.



420

421
422
423
424
425

Figure 3: Mapping precision and prediction accuracy from significant vertices between the different models of mass-univariate analyses. The columns correspond to the different simulation scenarios. We simulated 100 phenotypic traits for each scenario. Bars represent \pm SE across the 100 replicates. Clusters are composed of groups of contiguous vertices each significantly associated with the phenotype (after Bonferroni correction). We labelled them as true positives if they included a true positive association. (Mapping) precision refers to the median size of the true positive clusters.

426 *3.3 Morphometricity of the phenotypes*

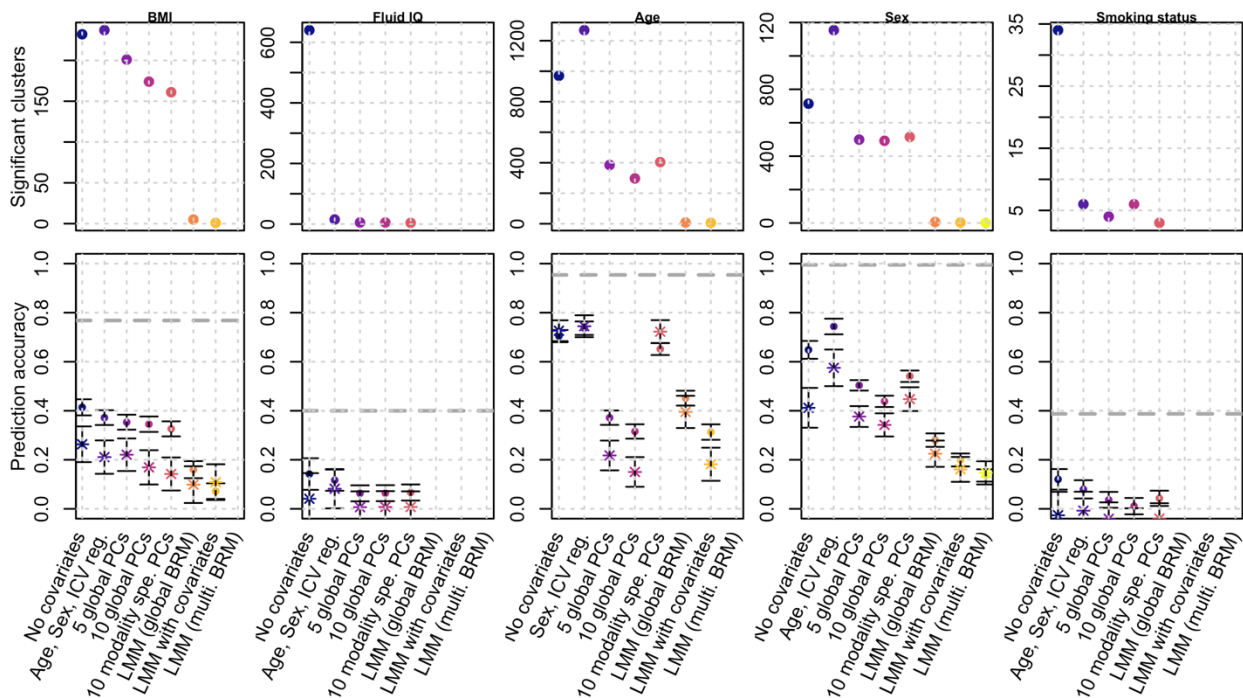
427 First, we confirmed that the morphometricity estimates of our simulated traits matched the values
428 chosen in simulations (**SFig. 8**). For the five UKB phenotypes, we also found consistent
429 morphometricity using the three LMM models (**STable 2**), suggesting associations across all
430 types of vertex measurements.

431 BMI and fluid intelligence exhibited large and moderate morphometricity ($R^2=0.51$ ($SE=0.031$)
432 and $R^2=0.17$ ($SE=0.034$)) but only a limited association with age, sex or the first 10 principal
433 components from vertex-wise data (adjusted R^2 with ten PCs: $R^2=0.032$ for fluid intelligence,
434 $R^2=0.033$ for BMI), which resembles the case of our simulations. Age and sex displayed high
435 morphometricity ($R^2=0.83$ ($SE=0.026$) and $R^2=0.99$ ($SE=0.024$)) and large associations with the
436 first ten PCs (adjusted $R^2=0.41$ for age, $R^2=0.43$ for sex). Smoking status is a discrete variable
437 (non-smoker, former smoker, still smoking) with a morphometricity of $R^2=0.12$ ($SE=0.029$), and
438 adjusted $R^2=9.2e-3$ with first 10 PCs (**STable 2**). Note that the morphometricity estimates are
439 slightly larger than the ones reported previously (15), which had mean cortical thickness and area
440 regressed out.

441 *3.4. Analysis of UK Biobank phenotypes*

442 We sought to confirm the differences in model performance by applying them to real phenotypic
443 traits. Using GLM without covariates resulted in many vertices and clusters reaching
444 significance (**Figure 4, STable 2**). Unsurprisingly, correcting for covariates which account for a
445 large fraction of the phenotypic variance (see adjusted R^2 with covariates and PCs, **STable 2**),
446 drastically reduced the number of associations in the GLMs. For example, correcting for ten PCs
447 in mass-univariate analyses of age and sex reduced the number of associated vertices by a factor

448 8-13, compared to the GLM without covariates (**Figure 4, STable 2**). For smoking status, the
 449 number of significant vertices and clusters also dropped despite a negligible association with PCs
 450 (**Figure 4, STable 2**). Similarly, for fluid intelligence, correcting for the top 10 PCs did not
 451 remove much of the trait variance over controlling for age sex and ICV (adjusted $R^2=0.030$ with
 452 age, sex, ICV, adjusted $R^2=0.034$ when further controlling for PCs) though it greatly reduced the
 453 number of associations. In addition, the more covariates we corrected for, the smaller the size of
 454 the associated clusters, suggesting they do remove confounding effects.



455 **Figure 4:** Number of significant clusters and prediction accuracy for the real UKB phenotypes
 456 Bars represent the 95% confidence intervals of the prediction accuracy (correlations). Dots indicate prediction accuracy in the UKB replication
 457 sample, while stars correspond to the prediction achieved in the OASIS3 sample. Prediction accuracy is reported controlling for age, sex (when
 458 pertinent), ICM, site/machine. In the OASIS3 dataset, we further controlled for clinical status. The dashed lines correspond to the estimated
 459 morphometricity, which corresponds to the theoretical maximum prediction accuracy achievable from a linear predictor.
 460
 461
 462

463 We found that across all phenotypes, LMMs resulted in a more parsimonious pattern of
 464 associations (**Figure 4, STable 2**). Thus, using the LMM with a single random-effect
 465 component, we identified 5 clusters associated with BMI, 8 with age and 6 with sex (**STable 2**).

466 LMM with covariates yielded fewer associations, while LMM with multiple random-effects was
467 the most conservative (**STable2**).

468 Next, we compared the prediction accuracy achieved from the vertices reaching significance
469 using each model (**Figure 4, STable 2**). Predicting our traits of interest allows evaluation of how
470 power and false positive rate of the different models may counterbalance each other. In addition,
471 prediction into independent samples quantifies the generalizability of findings obtained in the
472 different mass-univariate approaches. For BMI, we found that prediction accuracy from GLMs in
473 the UKB replication sample was greater than that in the OASIS3 sample, which suggests that
474 GLMs based predictors capture information that is sample specific (e.g., the same confounders
475 are more likely to be shared in the same cohort than across different cohorts). In contrast, the
476 prediction accuracy from LMMs was comparable between the UKB and OASIS3 samples,
477 pointing towards a better generalizability of the prediction. This suggests that the higher
478 prediction accuracy in the UKB replication sample for GLM is likely to be driven by
479 confounding factors shared between UKB data sets. The comparable performance of GLM and
480 LMM seen on OASIS3 for BMI aligns with our simulations.

481 For age and sex prediction, prediction accuracy of LMMs was sometimes inferior to that
482 achieved from GLMs, in particular those from the simplest models (“no Covariates” and “age,
483 sex, ICV”). Overall, prediction based on LMMs generalised well (comparable accuracy in the
484 UKB and OASIS3), while the GLMs often displayed heterogeneous performances across the test
485 samples (in particular for the GLMs with PCs, which may suffer from PCs being different
486 between samples).

487 Regarding fluid IQ and smoking status, no LMM predictor was available, and the different
488 GLMs resulted in comparable, albeit limited prediction accuracy.

489 3.5. *Description of associated regions*

490 We listed the significant associations identified using LMM (global BRM) in **STable 3**
491 (**SFig.9-11** for Manhattan plots, **SFig. 12-16** for brain plots). The significant associations were in
492 the range of $R^2=0.5-1\%$. Most associations were observed with subcortical volumes though the
493 top cluster for sex was spatially located at the border of the lateral-orbitofrontal and medial
494 orbitofrontal gyri (based on the Desikan atlas(26)). Out of the 85 vertices associated with age,
495 sex and BMI, 68 replicated in an independent UKB sample ($p<0.05/85$, **Table 2**). In particular,
496 4/11 associations replicated for BMI, 43/47 for age, and 21/27 for sex. The replication rate was
497 slightly lower in the OASIS3 dataset, where none of the vertices reached significance for BMI,
498 15/47 associations were replicated for age, and 12/27 for sex. Overall, the sign of the
499 associations was consistent across the 3 datasets (**STable 3**).

500 **4. Discussion**

501 Using extensive and realistic simulations, we evaluated the statistical power, false positive
502 rate and precision of GLMs and LMMs for vertex-wise grey-matter association studies. In
503 particular, we evaluated the different models in the context of big-data neuroimaging (large
504 sample size but even greater number of correlated brain vertices) (27). We consistently found
505 that using state-of-the-art GLMs resulted in a large number of false positive associations and
506 clusters, whether we used smoothed or not-smoothed grey-matter surfaces. Thus, across all
507 scenarios tested, more than 60% of the significant clusters were false positives using a standard
508 GLM that controlled for age, sex and ICV. In comparison, false discovery rate was below 17%
509 using LMMs, though still greater than the 5% expectation (**STable 1, Figure 2, SFig. 7**). In
510 addition, we showed that unlike GLMs, LMMs could appropriately separate cortical from

511 subcortical associations, even though signal contamination between thickness and surface still
512 occurred (**SFig. 2-5**).

513 Our results suggest that previously reported results from mass univariate vertex-wise analyses
514 obtained using standard GLM approaches could contain many redundant associations, some of
515 which are likely to be false positives induced by confounding factors that cause correlation
516 between vertices (e.g. (28-31), see also Figure 1b). Note that albeit redundant in term of
517 association and prediction, some of the brain regions identified using GLM may correspond to
518 indirect manifestations of the trait/disease of interest, which may be relevant to understand the
519 dynamics of grey-matter structure. Importantly, the type 1 error (greater than 5%) we observed in
520 simulations also warns against taking for granted results from LMMs.

521 The increased false positive rate for GLMs has been well documented in omics association
522 analyses studies (e.g. GWAS (8, 14) or MWAS (6, 9)) and has been attributed to proximal and
523 distal correlations between features, caused by factors independent of the trait of interest (e.g.
524 genetic ancestry in genetics, (14), cell composition of the biological sample and smoking status
525 in DNA methylation (6, 32)). On the other hand, LMMs can reduce the probability of generating
526 false positives, by fitting all other vertices as random effects which accounts for the complex
527 correlation structure between vertices within and between individuals. In brain imaging, more
528 work is needed to identify the factors that contribute to local and distal correlations between
529 vertices, hence inducing a correlation between true associations and “null” vertices, beyond the
530 usual covariates or confounders used in neuroimaging (e.g. MRI scanner/artefact (1) or
531 demographics (2)).

532 LMMs yielded fewer true positive associations, using the Bonferroni adjusted significance
533 threshold (in particular for the simulated associations on the subcortical nuclei (**SFig. 2**)).

534 However, this result must be interpreted with caution as it may be partly due to a more stringent
535 control of false positives, resulting in overall fewer vertices reach significance (**Figure 2, STable**
536 **2**). To better compare the models performances, we estimated statistical power (i.e. TPR for a set
537 false positive rate) and noted that the LMMs were more powerful than the GLMs (**Figure 2,**
538 **STable 2, Supp. 2**).

539 Despite this, LMMs are known to suffer from a power reduction, which arises from the double
540 fitting of the vertex of interest, once as fixed effect and again as a random effect (eq. 2)(7, 33).
541 For subcortical structures, the effect of double-fitting could be exacerbated by the high level of
542 correlation between vertices. A workaround (7) is to exclude the candidate vertex (and vertices
543 strongly correlated) from the BRM calculation (33), though this requires computation of the
544 BRM p times (complexity is $O(pN^3)$, with N the sample size and p the number of vertices),
545 which becomes impractical for large sample sizes (7, 33). In comparison, the current LMM
546 implementation makes our analysis scalable to samples sizes of tens of thousands (computational
547 complexity of $O(pN^2 + N^3 + pN)$) (6). It should be noted that Restricted Maximum Likelihood
548 (REML) estimation approach used in LMMs requires substantially more computational resources
549 than the GLMs and thus requires the use of high performance clusters.

550 Beyond power and false positive rate, we observed from simulations that LMMs could
551 pinpoint the grey-matter association with greater precision (smaller clusters of true positives,
552 **Figure 3**). Lastly, we found that prediction achieved from clusters reaching significance in
553 LMMs was on par with that from the best GLMs (**Figure 3**), despite fewer vertices included in
554 the predictor. This suggests a higher specificity of the LMMs. Overall, our simulations indicate
555 that LMM with a single random effect currently offers a good trade-off between power and false

556 positive rate. However, it still fails to ensure a cluster FWER below 5% (also reported on
557 MWAS (6)), despite a stringent Bonferroni correction to account for multiple testing.

558 Next, we applied the mass-univariate vertex-wise models to five real phenotypes of the UKB:
559 age, sex, BMI, smoking status and fluid IQ. As in the simulations, the LMMs identified fewer
560 vertices and clusters than the GLMs (**Figure 4, STable 2**). The LMM with multiple random-
561 effect components was the most stringent (a single cluster of association), consistent with
562 simulations which showed it had the lowest FWER and statistical power. In contrast, the LMM
563 with a single random-effect component identified several cortical and subcortical associations
564 with BMI, age and sex (**STable 3**). Most (12/19) of the top vertices in the associated cluster
565 replicated in the UKB left out sample, and 6 replicated in the OASIS3 sample (**STable 3**). The
566 lower replication rate in OASIS3 may be due to a lower power even though we cannot rule out
567 that the same confounders might act similarly on the two UKB data sets. Overall, replication
568 may be warranted to conclude about an association in future studies, considering the inflation of
569 false positives (even when using LMMs, **Figure 2**). The top associated vertices with age, sex and
570 BMI each captured less than 1% of the phenotypic variance, suggesting that many more small
571 associations are likely to account for the full morphometricity of the phenotypes (**STable 2**). Our
572 results echo the warning against the risk of small associations being confounded (e.g. by
573 artefacts) in big-data neuroimaging (27), which was confirmed by a recent exploratory study of
574 putative MRI confounders in the UKB (3). Note that LMMs can reduce false positive
575 associations caused by correlations across and within the different types of measurements
576 (**Figure 2, 3**). Finally, unlike in our simulations (**Figure 3**), LMMs often resulted in lower
577 prediction accuracy than GLMs in the UKB left out sample (**Figure 4**). Nonetheless, prediction
578 from LMMs generalised better in the OASIS3 dataset (**Figure 4**) (24). This suggests that LMMs

579 result in a more robust and parsimonious predictor, less sensitive to sample specific vertex-wise
580 patterns and confounders.

581 In the past years, many studies have been published on the association between grey-matter
582 structure and our phenotypes of interest (see **STable 4-8** for a selective review of publications).
583 Our simulation and empirical results suggest that some of these studies could report a substantial
584 number of false positive or redundant associations. Nevertheless, due to the limitations outlined
585 below, it is unclear which of these studies suffer from this issue and to which extent.

586 Firstly, it has been shown in the omics literature, that power of LMM may be reduced for
587 phenotypes strongly associated with the covariation between features (7, 34). This is likely the
588 case for age and sex as indicated by their strong association with the PCs calculated from vertex-
589 wise data (**STable 2**). This may be an important limitation for phenotypes associated with a
590 cascade of changes in grey-matter, for which LMM would be over conservative.

591 In addition, LMM assumes a normal distribution of random effects, which may not be
592 realistic for all phenotypes studied. It is equivalent to assuming highly regionalised and
593 specialised brain regions, each displaying a small association with the phenotype. Thus, LMM
594 may be sub-optimal under some architectures of association, such as if only a specific but sizable
595 brain region is associated with the trait. Several models have been proposed to relax the LMM
596 hypothesis, for example, to include large/outlying associations as fixed effects (stepwise
597 LMM(35)), break down the feature list into sets of small and large associations (data driven
598 approach: MOMENT(6)), or consider more complex distributions using Bayesian LMMs
599 (Bayesian alphabet (34, 36)). They remain to be evaluated in the context of vertex-wise analyses.
600 More simulations are warranted, to study other trait architectures, different trait distributions
601 (e.g. skewed, discrete) or to evaluate more sophisticated models. Of note, we limited our trait

602 complexity to 1,000 associated brain regions, even if the true pattern of association might be
603 more complex. Our simulations suggest that LMM outperform GLM independently of the trait
604 complexity, but also that larger samples are required to study traits with more complex
605 architecture (**Figure 2**). Our framework of simulation may be easily adapted for such
606 investigations, and offers the advantage of estimation of statistical power as well as false positive
607 rate, which are not often reported at the same time (37, 38).

608 The nature of the grey-matter regions identified in our GLM analyses of real phenotypes (for
609 which the truth is unknown) can be a matter of debate, which depends on the (also unknown)
610 nature of the correlation between vertices. Two key scenarios can explain the correlation but the
611 data currently available to us does not allow to differentiate between them. First, the correlation
612 could be solely due to confounders (e.g. **Figure 1b**), in which case the distal associations are
613 false positives. Second, the correlation between vertices could reflect dynamic brain pathways
614 relevant to the trait of interest. In this case, one could describe the GLM associations not found
615 using LMM as redundant rather than false positives. Since we cannot differentiate between these
616 two important causes of between-vertex correlation, we chose to label LMM models as
617 parsimonious, until we understand better the effect of confounders on the vertices correlation
618 structure as well as the longitudinal changes in grey-matter and their relationship with the
619 phenotypes.

620 Finally, some additional limitations are worthy of note as they may limit the interpretation of
621 mass-univariate vertex-wise analyses (compared to GWAS results). First, grey-matter
622 associations may be both causes or consequences of the phenotype studied, unlike GWAS
623 findings, which can impact how to consider redundant associations. At one end of the spectrum
624 are phenotypes such as age for which the direction of the causality is obvious (nothing causes

625 chronological age). When describing which parts of the brain are affected by aging, one may be
626 interested in reporting all associations, including all indirect and redundant. Though, there is no
627 guarantee that those brain regions correctly map the brain pathway of ageing as they might also
628 reach significance due to confusion factors. On the other hand, for many other phenotypes, the
629 direction of causality is unclear (e.g. smoking, BMI) and one may prefer a more parsimonious
630 and robust brain mapping. Second, grey-matter vertices are semi-arbitrary features which may be
631 defined and measured in different ways (e.g. different cortical meshes in FreeSurfer). For
632 instance, the resolution of the cortical tessellation is arbitrary and thus so is the number of local
633 vertices which are found to be significant. Hence, the results presented might differ if one were
634 to use a different MRI processing or vertex definition (e.g., volume processing from SPM,
635 coarser surface mesh). In addition, we used Bonferroni to control for multiple testing, although
636 approaches based on RFT are more commonly used (**STable 4-8**)(16). RFT based correction is
637 reportedly less stringent than Bonferroni, (at least for smoothed data, on which the RFT
638 hypotheses are more likely to be met)(39), which suggests that RFT would also suffer from the
639 inflation of test statistics that we reported.

640 Furthermore, we did not consider all possible covariates in GLM analysis, focussing on the
641 more commonly used in previous analyses (age, sex, ICV, **STable 4-8**). More work is needed to
642 evaluate the extended set(s) of covariates which have been recently proposed, from a large-scale
643 study of the UKB data(3). Finally, mass-univariate results may depend on the study sample used,
644 which raises the question of generalisability into to samples from different age or ethnic groups
645 or with different MRI qualities for instance.

646 In summary, we found that results obtained using the current state-of-the-art models (GLMs)
647 used in MRI-trait association analyses likely suffer from a large inflation of false positive or

648 redundant associations due to the unaccounted correlation between vertices. In contrast, LMMs
649 allow to control for all vertices fitted as a random effect, which result in a more parsimonious,
650 robust and conservative characterisation of the localised associations between a phenotype and
651 grey-matter structure. However, LMM results should still be interpreted with caution as our
652 simulations show that the false positive rate remains higher than the standard type 1 error of 5%,
653 even after Bonferroni correction.

654 **Disclosures**

655 The authors declare no conflict of interest.
656

657 **Acknowledgments**

658 Informed consent was obtained from all UK Biobank participants. Procedures are controlled by a
659 dedicated Ethics and Guidance Council (ukbiobank.ac.uk/ethics), with the Ethics and
660 Governance Framework available at [ukbiobank.ac.uk/wp-](http://ukbiobank.ac.uk/wp-content/uploads/2011/05/EGF20082.pdf)
661 [content/uploads/2011/05/EGF20082.pdf](http://ukbiobank.ac.uk/wp-content/uploads/2011/05/EGF20082.pdf). IRB approval was also obtained from the North West
662 Multi-centre Research Ethics Committee. This research has been conducted using the UK
663 Biobank Resource under Application Number 12505.

664 All necessary patient/participant consent has been obtained and the appropriate institutional
665 forms have been archived by the OASIS team.

666 This research was supported by the Australian National Health and Medical Research Council
667 (1078037, 1078901, 1113400, 1161356 and 1107258), the Australian Research Council
668 (FT180100186 and FL180100072), the Sylvia & Charles Viertel Charitable Foundation, the
669 program “Investissements d’avenir” ANR-10-IAIHU-06 (Agence Nationale de la Recherche-10-

670 IA Institut Hospitalo-Universitaire-6) and reference ANR-19-P3IA-0001 (PRAIRIE 3IA
671 Institute), the European Union H2020 program (project EuroPOND, grant number 666992, the
672 joint NSF/NIH/ANR program “Collaborative Research in Computational Neuroscience” (project
673 HIPLAY7, grant number ANR-16-NEUC-0001-01), the ICM Big Brain Theory Program
674 (project DYNAMO, project PredictICD), and the Abeona Foundation (project Brain@Scale).
675 The OASIS3 data were provided by Principal Investigators: T. Benzinger, D. Marcus, J. Morris
676 supported by NIH grants: P50AG00561, P30NS09857781, P01AG026276, P01AG003991,
677 R01AG043434, UL1TR000448, R01EB009352.

678 We used R (40) (v3.6.2) for analyses not performed using OSCA (6) and for plots. We used the
679 colour-blind friendly R palette *viridis* (41), *ukbtools* (42) to facilitate UKB phenotype
680 manipulation, *Morpho* and *Rvcg* (43) to identify clusters, *rgl*(44) to generate brain plots of
681 associations. Other packages used include, *dplyr* (45), *readr* (46), *rmarkdown* (47), *matrixStats*
682 (48), *RcolorBrewer* (49), *gridExtra* (50), *ggplot2*(51), *png* (52), *epuRate* (53).

683 We would like to thank the Research Computing Centre (RCC) at the University of Queensland
684 for their support with high performance computing, data handling, storage and processing.

685 **Code, Data, and Materials Availability**

686 Data used in this manuscript are held and distributed by the OASIS and UKB teams. We have
687 released the code used in image processing and analyses to facilitate replication and
688 dissemination of the results (<https://baptistecd.github.io/Brain-Mapping-LMM/>). Upon
689 publication, we will also release, the summary statistics of mass-univariate analyses performed
690 on the UKB phenotypes (<https://cnsgenomics.com/content/data>). Supplementary figures (**SFig.**

691 **1-16)**, tables (**STable 1-7**) and sections (**Supp. 1-2**) may be found on the GitHub repository:
692 https://github.com/baptisteCD/Brain-Mapping-LMM/blob/main/Supp_Article_JMI.pdf .

693

694 **5. References**

695

- 696 1. A. A. Chen et al., "Removal of Scanner Effects in Covariance Improves Multivariate
697 Pattern Analysis in Neuroimaging Data," *bioRxiv* 858415 (2020).
- 698 2. M. Montembeault et al., "The impact of aging on gray matter structural covariance
699 networks," *NeuroImage* **63**(2), 754-759 (2012).
- 700 3. F. Alfaro-Almagro et al., "Confound modelling in UK Biobank brain imaging,"
701 *NeuroImage* 117002 (2020).
- 702 4. L. R. Cardon, and L. J. Palmer, "Population stratification and spurious allelic
703 association," *Lancet* **361**(9357), 598-604 (2003).
- 704 5. J. Marchini et al., "The effects of human population structure on large genetic association
705 studies," *Nat Genet* **36**(5), 512-517 (2004).
- 706 6. F. Zhang et al., "OSCA: a tool for omic-data-based complex trait analysis," *Genome Biol*
707 **20**(1), 107 (2019).
- 708 7. J. Yang et al., "Advantages and pitfalls in the application of mixed-model association
709 methods," *Nat Genet* **46**(2), 100-106 (2014).
- 710 8. A. L. Price et al., "New approaches to population stratification in genome-wide
711 association studies," *Nature reviews. Genetics* **11**(7), 459-463 (2010).
- 712 9. M. F. Nabais et al., "Significant out-of-sample classification from methylation profile
713 scoring for amyotrophic lateral sclerosis," *NPJ Genom Med* **5**(10) (2020).
- 714 10. J. L. Bernal-Rusiel et al., "Statistical analysis of longitudinal neuroimage data with
715 Linear Mixed Effects models," *NeuroImage* **66**(249-260) (2013).
- 716 11. M. R. Sabuncu et al., "Morphometricity as a measure of the neuroanatomical signature of
717 a trait," *Proceedings of the National Academy of Sciences of the United States of America*
718 **113**(39), E5749-E5756 (2016).
- 719 12. B. Couvy-Duchesne et al., "A unified framework for association and prediction from
720 vertex-wise grey-matter structure," *Human Brain Mapping* **n/a**(n/a), (2020).
- 721 13. K. L. Miller et al., "Multimodal population brain imaging in the UK Biobank prospective
722 epidemiological study," *Nat Neurosci* **19**(11), 1523-1536 (2016).
- 723 14. A. L. Price et al., "Principal components analysis corrects for stratification in genome-
724 wide association studies," *Nature Genetics* **38**(8), 904-909 (2006).
- 725 15. B. Couvy-Duchesne et al., "Widespread associations between grey matter structure and
726 the human phenome," *bioRxiv* 696864 (2019).
- 727 16. T. Nichols, and S. Hayasaka, "Controlling the familywise error rate in functional
728 neuroimaging: a comparative review," *Statistical methods in medical research* **12**(5),
729 419-446 (2003).
- 730 17. B. Fischl, "FreeSurfer," *NeuroImage* **62**(2), 774-781 (2012).

- 731 18. H. Lindroth et al., "Examining the identification of age-related atrophy between T1 and
732 T1 + T2-FLAIR cortical thickness measurements," *Sci Rep* **9**(1), 11288 (2019).
- 733 19. B. Fischl, and A. M. Dale, "Measuring the thickness of the human cerebral cortex from
734 magnetic resonance images," *Proceedings of the National Academy of Sciences* **97**(20),
735 11050-11055 (2000).
- 736 20. A. M. Winkler et al., "Measuring and comparing brain cortical surface area and other
737 areal quantities," *NeuroImage* **61**(4), 1428-1443 (2012).
- 738 21. B. A. Gutman et al., "A Family of Fast Spherical Registration Algorithms for Cortical
739 Shapes," in *Multimodal Brain Image Analysis: Third International Workshop, MBIA
740 2013, Held in Conjunction with MICCAI 2013, Nagoya, Japan, September 22, 2013,
741 Proceedings* L. Shen, T. Liu, P.-T. Yap, H. Huang, D. Shen, and C.-F. Westin, Eds., pp.
742 246-257, Springer International Publishing, Cham (2013).
- 743 22. B. A. Gutman et al., "Shape Matching with Medial Curves and 1-D Group-Wise
744 Registration," *2012 9th Ieee International Symposium on Biomedical Imaging (Isbi)* 716-
745 719 (2012).
- 746 23. G. V. Roshchupkin et al., "Heritability of the shape of subcortical brain structures in the
747 general population," *Nat Commun* **7**(13738) (2016).
- 748 24. P. J. LaMontagne et al., "OASIS-3: Longitudinal Neuroimaging, Clinical, and Cognitive
749 Dataset for Normal Aging and Alzheimer Disease," *medRxiv* 2019.2012.2013.19014902
750 (2019).
- 751 25. C. Fawns-Ritchie, and I. J. Deary, "Reliability and validity of the UK Biobank cognitive
752 tests," *PloS one* **15**(4), e0231627-e0231627 (2020).
- 753 26. R. S. Desikan et al., "An automated labeling system for subdividing the human cerebral
754 cortex on MRI scans into gyral based regions of interest," *NeuroImage* **31**(3), 968-980
755 (2006).
- 756 27. S. M. Smith, and T. E. Nichols, "Statistical Challenges in "Big Data" Human
757 Neuroimaging," *Neuron* **97**(2), 263-268 (2018).
- 758 28. C. K. Tamnes et al., "Development of the Cerebral Cortex across Adolescence: A
759 Multisample Study of Inter-Related Longitudinal Changes in Cortical Volume, Surface
760 Area, and Thickness," *J Neurosci* **37**(12), 3402-3412 (2017).
- 761 29. S. J. Ritchie et al., "Sex Differences in the Adult Human Brain: Evidence from 5216 UK
762 Biobank Participants," *Cereb Cortex* **28**(8), 2959-2975 (2018).
- 763 30. S. R. Cox et al., "Associations between vascular risk factors and brain MRI indices in UK
764 Biobank," *European Heart Journal* **40**(28), 2290-2300 (2019).
- 765 31. F. J. Navas-Sanchez et al., "Cortical morphometry in frontoparietal and default mode
766 networks in math-gifted adolescents," *Hum Brain Mapp* **37**(5), 1893-1902 (2016).
- 767 32. A. E. Jaffe, and R. A. Irizarry, "Accounting for cellular heterogeneity is critical in
768 epigenome-wide association studies," *Genome Biol* **15**(2), R31 (2014).
- 769 33. J. Listgarten et al., "Improved linear mixed models for genome-wide association studies,"
770 *Nat Methods* **9**(6), 525-526 (2012).
- 771 34. L. R. Lloyd-Jones et al., "Inference on the Genetic Basis of Eye and Skin Color in an
772 Admixed Population via Bayesian Linear Mixed Models," *Genetics* **206**(2), 1113-1126
773 (2017).
- 774 35. V. Segura et al., "An efficient multi-locus mixed-model approach for genome-wide
775 association studies in structured populations," *Nat Genet* **44**(7), 825-830 (2012).

- 776 36. G. Moser et al., "Simultaneous Discovery, Estimation and Prediction Analysis of
777 Complex Traits Using a Bayesian Mixture Model," *Plos Genetics* **11**(4), (2015).
- 778 37. A. Eklund, T. E. Nichols, and H. Knutsson, "Cluster failure: Why fMRI inferences for
779 spatial extent have inflated false-positive rates," *P Natl Acad Sci USA* **113**(28), 7900-
780 7905 (2016).
- 781 38. S. Noble, D. Scheinost, and R. T. Constable, "Cluster failure or power failure? Evaluating
782 sensitivity in cluster-level inference," *NeuroImage* **209**(116468 (2020)).
- 783 39. K. J. Worsley, "An improved theoretical P value for SPMs based on discrete local
784 maxima," *NeuroImage* **28**(4), 1056-1062 (2005).
- 785 40. R Development Core Team, "R: A Language and Environment for Statistical
786 Computing," R Foundation for Statistical Computing, Vienna, Austria (2012).
- 787 41. S. Garnier, "viridis: Default Color Maps from 'matplotlib'," (2018).
- 788 42. K. Hanscombe, "ukbtools: Manipulate and Explore UK Biobank Data," (2017).
- 789 43. S. Schlager, "Chapter 9 - Morpho and Rvcg – Shape Analysis in R: R-Packages for
790 Geometric Morphometrics, Shape Analysis and Surface Manipulations," in *Statistical
791 Shape and Deformation Analysis* G. Zheng, S. Li, and G. Székely, Eds., pp. 217-256,
792 Academic Press (2017).
- 793 44. D. Adler, and D. Murdoch, "rgl: 3D Visualization Using OpenGL," (2020).
- 794 45. H. Wickham, and R. Francois, "dplyr: A Grammar of Data Manipulation," (2015).
- 795 46. H. H. Wickham, J.; Francois, R., "readr: Read Rectangular Text Data," (2017).
- 796 47. J. X. Allaire, Yihui.; McPherson, Jonathan.; Luraschi, Javier.; Ushey, Kevin.; Atkins,
797 Aron.; Wickham, Hadley.; Cheng, Joe.; Chang, Winston. , "rmarkdown: Dynamic
798 Documents for R," (2018).
- 799 48. H. Bengtsson, "matrixStats: Functions that Apply to Rows and Columns of Matrices (and
800 to Vectors)," (2019).
- 801 49. E. Neuwirth, "RColorBrewer: ColorBrewer Palettes," (2014).
- 802 50. B. Auguie, "gridExtra: Miscellaneous Functions for "Grid" Graphics," (2017).
- 803 51. H. Wickham, *ggplot2: Elegant Graphics for Data Analysis*, Springer Publishing
804 Company, Incorporated (2009).
- 805 52. S. Urbanek, "png: Read and write PNG images," (2013).
- 806 53. Y. Holtz, "epuRate: A clean template for R Markdown documents," (2020).
- 807

808 **Baptiste Couvy-Duchesne** is a CJ Martin fellow (National Health Medical Research Council,
809 Australia) and INRIA researcher (National Institute for Research in Digital Science and
810 Technology, France). He obtained his PhD in 2017 from the University of Queensland, working
811 on complex-trait genetics and MRI brain imaging. Since, BCD has been working at adapting and
812 applying some of the methods used in big-data genetics to the analysis of high-dimensional brain

813 MRI. He is a member of the ENIGMA (Enhancing Neuro-Imaging Genetics using Meta-
814 Analyses) and PGC (Psychiatry Genetics) consortia.

815
816 **Caption List**

817
818 **Figure 1:** Illustration of the traditional confounding paradigm a) and of the confounding that
819 may arise in association studies performed across correlated brain features b).
820 One sided arrows represent a causal effect, and two-sided arrows a correlation.

821
822 **Figure 2:** Performance of GLMs and LMMs for mass-univariates vertex-wise analyses: test
823 inflation, statistical power and false positive rate.
824 The columns correspond to the different scenarios considered when simulating traits. We
825 simulated 100 phenotypic traits for each scenario. Bars represent +/- SE across the 100
826 replicates. Clusters are composed of groups of contiguous vertices each significantly associated
827 with the phenotype (after Bonferroni correction). We labelled them as false positives if they did
828 not include a true positive association.

829
830 **Figure 3:** Mapping precision and prediction accuracy from significant vertices between the
831 different models of mass-univariate analyses
832 The columns correspond to the different simulation scenarios. We simulated 100 phenotypic
833 traits for each scenario. Bars represent +/- SE across the 100 replicates. Clusters are composed of
834 groups of contiguous vertices each significantly associated with the phenotype (after Bonferroni
835 correction). We labelled them as true positives if they included a true positive association.
836 (Mapping) precision refers to the median size of the true positive clusters.

837
838 **Figure 4:** Number of significant clusters and prediction accuracy for the real UKB phenotypes
839 Bars represent the 95% confidence intervals of the prediction accuracy (correlations). Dots
840 indicate prediction accuracy in the UKB replication sample, while stars correspond to the
841 prediction achieved in the OASIS3 sample. Prediction accuracy is reported controlling for age,
842 sex (when pertinent), ICM, site/machine. In the OASIS3 dataset, we further controlled for
843 clinical status. The dashed lines correspond to the estimated morphometricity, which corresponds
844 to the theoretical maximum prediction accuracy achievable from a linear predictor.

845
846 **Author Bios:**

847 Baptiste Couvy-Duchesne is an INRIA researcher and CJ Martin Fellow, working for the Paris
848 Brain Institute and the University of Queensland. His research focuses on developing and
849 applying novel statistical methods to analyse large scale brain MRI data.

850
851 Futao Zhang is a post-doctoral researcher with the Institute for Molecular Biosciences (IMB) at
852 the University of Queensland. He develops performant methods and software for the analysis of
853 genomic and other large scale datasets.

854
855 Kathryn Kemper is a post-doctoral researcher with the Institute for Molecular Biosciences (IMB)
856 at the University of Queensland. Her research in statistical genetics focuses on examining the
857 causes of variation within human populations for traits such as height and body mass index.
858

859 Julia Sidorenko is a post-doctoral research assistant with the Institute for Molecular Biosciences
860 (IMB) at the University of Queensland. Her research in statistical genetics focuses on examining
861 the causes of variation within human populations. She also processes and manages genetic and
862 genomic datasets, including the UKBiobank.
863

864 Naomi Wray holds joint appointments at the Institute for Molecular Bioscience (IMB) and the
865 Queensland Brain Institute (QBI) within the University of Queensland. She is a National Health
866 and Medical Research Council (NHMRC) Leadership Fellow, a Fellow of the Australian
867 Academy of Science and a Fellow of the Australian Academy of Health and Medical Science.
868 Her research focusses on development of quantitative genetics and genomics methodology with
869 application to psychiatric and neurological disorders.
870

871 Peter Visscher joined the University of Queensland in 2011, where he is Professor of
872 Quantitative Genetics. He is a Laureate Fellow of the Australian Research Council. Visscher was
873 elected a Fellow of the Australian Academy of Science in 2010, a Fellow of the Royal Society
874 (London) in 2018 and a Foreign Member of the Royal Netherlands Academy of Arts and
875 Sciences in 2018. Visscher's research is about genetic variation for complex traits (including
876 quantitative traits and disease) in populations, with the broad aim to understand and quantify the
877 causes and consequences of human trait variation.
878

879
880 Olivier Colliot is a Research Director at CNRS and the co-head of the ARAMIS Lab
881 (www.aramislab.fr), a joint laboratory between CNRS, Inria, Inserm, Sorbonne University and
882 the Paris Brain Institute. He also holds a chair at the PRAIRIE Institute for Artificial
883 Intelligence. He is Conference Chair of SPIE Medical Imaging Image Processing conference and
884 an Associate Editor of Medical Image Analysis and Frontiers in Brain Imaging Methods. His
885 research interests include machine learning, medical image analysis, and their applications to
886 neurological disorders.
887

888 Jian Yang is a Professor of Statistical Genetics at the School of Life Sciences, Westlake
889 University, China. He received his PhD in 2008 from Zhejiang University, China, before
890 undertaking postdoctoral research at the QIMR Berghofer Medical Research Institute in
891 Australia (2008-2011). He moved to The University of Queensland (UQ), Australia, as a
892 Research Fellow in 2012 and was reappointed as a Senior Research Fellow and Group Leader in
893 January 2014. He was promoted to be an Associate Professor in December 2014, and then a
894 Professor in 2017 at UQ. He joined Westlake University in 2020. His primary research interests
895 are focused on understanding the genomic variations among individuals within and between
896 populations and the links of genomic variations with health outcomes.
897

898

899

900

901