

Controlling for Human Population Stratification in Rare Variant Association Studies

Matthieu Bouaziz^{1,2}, Jimmy Mullaert^{1,2,3,4}, Benedetta Bigio⁵, Yoann Seeleuthner^{1,2}, Jean-Laurent Casanova^{1,2,5,6,7}, Alexandre Alcais^{1,2}, Laurent Abel^{1,2,5,¶}, Aurélie Cobat^{1,2,¶*}.

1. Laboratory of Human Genetics of Infectious Diseases, Necker Branch, INSERM U1163, Paris, France, EU.

2. Paris Descartes University, Imagine Institute, 75015 Paris, France, EU.

3. Université de Paris, IAME, INSERM, F-75018 Paris, France

4. AP-HP, Hôpital Bichat, DEBRC, F-75018 Paris, France

5. St. Giles Laboratory of Human Genetics of Infectious Diseases, Rockefeller Branch, The Rockefeller University, New York, USA.

6. Howard Hughes Medical Institute, New-York, NY, USA

7. Pediatric Hematology and Immunology Unit, Necker Hospital for Sick Children, 75015 Paris, France, EU.

¶ These authors contributed equally to this work

* Corresponding author : aurelie.cobat@inserm.fr (AC)

21 **Abstract**

22 Population stratification is a strong confounding factor in human genetic association studies. In
23 analyses of rare variants, the main correction strategies based on principal components (PC) and linear
24 mixed models (LMM), may yield conflicting conclusions, due to both the specific type of structure
25 induced by rare variants and the particular statistical features of association tests. Studies evaluating
26 these approaches generally focused on specific situations with limited types of simulated structure
27 and large sample sizes. We investigated the properties of several correction methods in the context
28 of a large simulation study using real exome data, and several within- and between- continent
29 stratification scenarios. We also considered different sample sizes, with situations including as few
30 as 50 cases, to account for the analysis of rare disorders. In this context, we focused on a genetic
31 model with a phenotype driven by rare deleterious variants well suited for a burden test. For analyses
32 of large samples, we found that accounting for stratification was more difficult with a continental
33 structure than with a worldwide structure. LMM failed to maintain a correct type I error in many
34 scenarios, whereas PCs based on common variants failed only in the presence of extreme continental
35 stratification. When a sample of 50 cases was considered, an inflation of type I errors was observed
36 with PC for small numbers of controls (≤ 100), and with LMM for large numbers of controls (≥ 1000).
37 We also tested a promising novel adapted local permutation method (LocPerm), which maintained a
38 correct type I error in all situations. All approaches capable of correcting for stratification properly
39 had similar powers for detecting actual associations pointing out that the key issue is to properly
40 control type I errors. Finally, we found that adding a large panel of external controls (*e.g.* extracted
41 from publicly available databases) was an efficient way to increase the power of analyses including
42 small numbers of cases, provided an appropriate stratification correction was used.

43

44

45 **Author Summary**

46

47 Genetic association studies focusing on rare variants using next generation sequencing (NGS) data
48 have become a common strategy to overcome the shortcomings of classical genome-wide association
49 studies for the analysis of rare and common diseases. The issue of population stratification remains
50 however a substantial question that has not been fully resolved when analyzing NGS data. In this
51 work, we propose a comprehensive evaluation of the main strategies to account for stratification, that
52 are principal components and linear mixed model, along with a novel approach based on local
53 permutations (LocPerm). We compared these correction methods in many different settings,
54 considering several types of population structures, sample sizes or types of variants. Our results
55 highlighted important limitations of some classical methods as those using principal components (in
56 particular in small samples) and linear mixed models (in several situations). In contrast, LocPerm
57 maintained a correct type I error in all situations. Also, we showed that adding a large panel of external
58 controls, *e.g* coming from publicly available databases, is an efficient strategy to increase the power
59 of an analysis including a low number of cases, as long as an appropriate stratification correction is
60 used. Our findings provide helpful guidelines for many researchers working on rare variant
61 association studies.

62

63 **Introduction**

64

65 Genetic association studies focusing on rare variants have become a popular approach to
66 analyzing rare and common diseases. The advent of next-generation sequencing (NGS) and the
67 development of new statistical approaches have rendered possible the comprehensive investigation
68 of rare genetic variants, overcoming the shortcomings of classical genome-wide association studies
69 (GWAS) [1, 2]. The main methods for testing rare variants for association do not test single variants
70 against a phenotype, as in GWAS, but generally use an aggregation strategy within a genetic unit,
71 usually a gene. These gene-based tests can be divided into two main categories: burden and variance-
72 component tests [1-4]. Population stratification occurs when study subjects, usually cases and
73 controls, are recruited from genetically heterogeneous populations. This problem is well known in
74 association studies with common variants, causing an inflation of the type I error rate and reducing
75 power. Several statistical approaches can be used to account for population stratification in GWAS.
76 The most widely used are based on Principal Components (PC) analysis [5, 6] and Linear Mixed
77 Models (LMM) [7-10].

78 Population stratification also affects association studies including rare variants [11-13].
79 However, it remains unclear whether the same correction methods can be applied to rare variant
80 association studies [12, 14], particularly as rare and common variants may induce different types of
81 population structure [12, 15]. Many studies have investigated the bias introduced by population
82 stratification in the analysis of rare variants and have highlighted the need for corrective approaches
83 to obtain meaningful results [12, 16, 17]. The performance of the correction method depends on the
84 study setting and the method used to analyze the variants [11, 12, 18-21]. PC has been widely
85 investigated [5, 6, 22-25] and shown to yield satisfactory correction at large geographic scales, but
86 not at finer scales [20]. LMM have also been studied [19, 26] and shown to account for stratification
87 well if variance-component approaches are used to test for association [19]. Most of these studies
88 used simulated genetic data that did not completely reproduce the complexity of real exome
89 sequences, and limited types of population structures. In addition, they used large numbers of cases
90 (*e.g.* generally more than 500), which may not always be possible in practice, particularly in studies
91 focusing on rare diseases.

92 We aimed at addressing such limitations of classical comparative studies with the
93 comprehensive evaluation study proposed in this article. We investigated the main correction methods
94 for rare variant association studies in the context of limited sample sizes, as in studies of rare disorders.
95 For an accurate assessment of the different approaches, we used real NGS data from two sources:
96 1000 Genomes data [27] and our in-house cohort, with data for > 5,000 exomes [28]. We focused on

97 two population structure scenarios: within-continent stratification (recent separation) and between-
98 continent stratification (ancient separation). We also considered different sample sizes, including
99 situations with as few as 50 cases, which have, to our knowledge, never been extensively investigated
100 in this manner. We focused on a classic genetic model for a rare disease with a phenotype driven by
101 rare deleterious variants well suited for a burden test, such as the cohort allelic sums test (CAST) [3].
102 We tested two classical correction methods, PC and LMM, a promising novel correction method
103 called adapted local permutations (LocPerm) [29] and considered an uncorrected CAST-like test as a
104 reference. Our global objective here is to provide useful practical insight into how best to account for
105 population stratification in rare variant association studies.

106

107 **Materials and methods**

108 **Simulation study**

109

110 **Exome data.** For a realistic comparison of the correction approaches, we used two real exome
111 datasets rather than program-based simulated exomes. Simulated data tend to provide erroneous site
112 frequency spectra or LD structures [30]. The first dataset used was our HGID (Human Genetic of
113 Infectious Diseases) database, containing 3,104 samples of in-house WES data generated with the
114 SureSelect Human All Exon V4+UTRs exome capture kit (<https://agilent.com>). All study participants
115 provided written informed consent for the use of their DNA in studies aiming to identify genetic risk
116 variants for disease. IRB approval was obtained from The Rockefeller University and Necker
117 Hospital for Sick Children, along with a number of collaborating institutions. The second dataset used
118 was the 2,504 whole genomes from 1000 Genomes phase 3 (<http://www.internationalgenome.org/>)
119 reduced with the same capture kit. We merged all the exomes from these two databases into a single
120 large dataset before selecting samples. We performed quality control, retaining only coding variants
121 with a depth of coverage (DP) > 8, a genotype quality (GQ) > 20, a minor read ratio (MRR) > 0.2 and
122 call-rate > 95% [31]. We then excluded all related individuals based on the kinship coefficient (King's
123 kinship $2K > 0.1875$) [32, 33], resulting in a final set of 4,887 unrelated samples. From these samples,
124 we created two types of samples, as comparable as possible to those used in practice in association
125 studies. The first sample, the “European” sample, consisted of samples from patients of European
126 ancestry, and was used to assess stratification at the continental level. The second, the “Worldwide”
127 sample, consisted of samples from European individuals together with North-African, Middle-
128 Eastern, and South-Asian samples, for the assessment of intercontinental stratification.

129

130 **European sample.** We selected samples from individuals of European ancestry based on a reference

131 sample and genetic distance. We first picked a European sample (sample HG00146 from the GBR
132 population of 1000 Genomes, Figure 1A) and calculated its genetic distance to all other samples in
133 the combined dataset. We used a Euclidean distance based on the first 10 PCs: the distance between
134 individuals i and j is calculated as $d_{ij}^2 = \sum_{k=1}^{10} \lambda_k |PC_{ki}^{CV} - PC_{kj}^{CV}|^2$, where PC^{CV} is the matrix of
135 principal components calculated on common variants and λ_k is the eigenvalue corresponding to the
136 k -th principal component PC_k^{CV} . We considered that a sample could be “European” if its distance to
137 the reference sample was below a certain threshold. This threshold was empirically chosen to ensure
138 that all individuals of known European ancestry from the 1000 Genomes and our in-house HGID
139 cohorts were included. The final sample consisted of 1,523 individuals, and included all the European
140 samples from 1000 Genomes. We empirically separated the samples into three groups on the basis of
141 ancestry (Figure 1B): Northern ancestry (including principally the FIN samples from 1000 Genomes),
142 Middle-Europe ancestry (including the CEU and GBR samples from 1000 Genomes) and Southern
143 ancestry (including the TSI and IBS samples from 1000 Genomes). The sample size for each
144 subpopulation is shown in Table S1. After removal of the 102,219 private variants, the final sample
145 contained 328,989 biallelic SNPs (Table S2).

146

147 **Worldwide sample.** The Worldwide sample was created in a similar manner. We selected four
148 different reference samples of European (sample HG00146 from the GBR population of 1000
149 Genomes), South-Asian (sample NA20847 from the GIH population of 1000 Genomes), Middle-
150 Eastern and North-African (samples from our in-house sample with a reported and verified Middle-
151 Eastern or North-African ancestry) ancestry (Figure 2A). The genetic distances between each sample
152 and the four reference samples were calculated as previously described. Thresholds were applied such
153 that each sample with a reported ancestry of interest was assigned to the correct population and there
154 was no overlap between the subpopulations (Figure 2B). The final Worldwide sample included 1,967
155 individuals separated into four subpopulations (Table S1). Note that all the European samples of this
156 sample were also present in the European sample. This sample contained 483,762 biallelic SNPs after
157 removal of the 132,565 private variants (Table S2).

158

159 **Stratification scenarios.** We first assessed the various correction approaches on case/control samples
160 with large sample sizes (*i.e.* with the whole European or Worldwide sample). We used the same three
161 stratification scenarios for both samples. In each scenario, we considered a fixed proportion of 15%
162 cases and 85% controls. Thus, in all our scenarios, the case/control ratio was unbalanced, as is often
163 the case in practice. Comparison studies generally consider balanced scenarios with large numbers of

164 cases and controls, corresponding to the ideal situation for most correction approaches, and their
 165 performance in more realistic conditions may therefore be overestimated. We considered a first
 166 scenario without stratification (No PS), in which we randomly selected 15% of the samples in each
 167 subpopulation as cases, the rest being used as controls. The second scenario corresponded to moderate
 168 stratification (Moderate PS), with the cases selected mostly from certain subpopulations. The third
 169 scenario was an extreme situation (High PS), in which all the cases were selected from a single
 170 subpopulation. The distribution of cases for the European and the Worldwide samples is shown, for
 171 each scenario, in Table 1.

172

173 **Table 1: Distribution of the cases in the sub-populations of the European and the Worldwide**
 174 **samples for the different population stratification (PS) scenarios.**

European sample				
Scenario	Northern-Europe (n=127)	Middle-Europe (n=651)	Southern-Europe (n=745)	
No PS	19 (15 %) ^a	98 (15 %)	112 (15 %)	
Moderate PS	6 (5 %)	45 (7 %)	177 (24 %)	
High PS	0 (0 %)	0 (0 %)	228 (30 %)	
Worldwide sample				
Scenario	Europe (n=700)	South-Asia (n=543)	North-Africa (n=359)	Middle-East (n=365)
No PS	105 (15 %) ^a	81 (15 %)	53 (15 %)	54 (15 %)
Moderate PS	177 (25 %)	60 (11 %)	29 (8 %)	29 (7 %)
High PS	294 (42 %)	0 (0 %)	0 (0 %)	0 (0 %)

175 ^a#cases (% of the sub-population)

176

177 In practice, the samples used in rare variant association studies are frequently not very large. This is
 178 particularly true for rare diseases, for which only small numbers of cases are available. Case numbers
 179 may also be small as a consequence of the WES cost. The usual analysis strategy involves matching
 180 the controls to the cases. One key question is whether the addition of unmatched controls could
 181 increase the power of the analysis when population stratification is taken into account properly. Such
 182 controls are now available in large cohorts, such as the 1000 Genomes (Genomes Project, Auton (27)),
 183 UK10K [34], and UK Biobank [35] cohorts. We decided to investigate such strategies, by considering
 184 several scenarios with 50 cases and various numbers of controls of similar or different ancestries
 185 (Table 2). We considered three possible types of cases: 50 cases from the rather homogeneous

186 Southern-Europe subpopulation (50SE), 50 cases from the more heterogeneous whole European
 187 population (50E) and 50 cases selected Worldwide (50W). Four types of controls were considered:
 188 100 controls from the same population as the cases (100SE, 100E, 100W), 1000 controls from the
 189 total European sample (1000E), 1000 controls randomly chosen from the total Worldwide sample
 190 (1000W) and 2000 controls randomly chosen from the total Worldwide sample (2000W).

191

192 **Table 2: Stratification scenarios for the small size study.** The first 4 scenarios correspond to cases
 193 from the Southern-Europe sub-population (SE), the following 4 scenarios to cases from whole
 194 European sample (E) and the final 4 to cases from the Worldwide population (W). Controls are
 195 randomly drawn among the Southern-European, European or Worldwide populations.

Scenario	Cases	Controls
50SE-100SE	50 from Southern-Europe	100 from Southern-Europe
50SE-1000E	50 from Southern-Europe	1000 from all Europe
50SE-1000W	50 from Southern-Europe	1000 Worldwide
50SE-2000W	50 from Southern-Europe	2000 Worldwide
50E-100E	50 from all Europe	100 from all Europe
50E-1000E	50 from all Europe	1000 from all Europe
50E-1000W	50 from all Europe	1000 Worldwide
50E-2000W	50 from all Europe	2000 Worldwide
50W-100W	50 Worldwide	100 Worldwide
50W-1000E	50 Worldwide	1000 from all Europe
50W-1000W	50 Worldwide	1000 Worldwide
50W-2000W	50 Worldwide	2000 Worldwide

196

197

198 **Type I error rate evaluation.** For each type of sample and stratification scenario, the type I error
 199 rate was estimated under the null hypothesis of no association between a gene and the phenotype
 200 (H_0). We therefore simulated phenotypes, for the large sample, by randomly assigning the case and
 201 control states according to the stratification proportions provided in Table 1, respecting a fixed
 202 proportion of cases of 15%. Each protein-coding gene was then tested for association with the
 203 phenotype by the various statistical approaches described in the Statistical methods section. The rare
 204 variants included in these tests were biallelic variants with a $MAF \leq 5\%$ in the sample analyzed. We
 205 included only genes with at least 10 rare variant carriers, resulting in 17,619 genes being studied in

206 the European sample, and 17,854 genes in the Worldwide sample. A similar simulation process was
207 applied to the small samples, according to the proportions of cases and controls described in Table 2.
208 In these scenarios, the number of genes with at least 10 mutation carriers retained depended on sample
209 size (Table S3). This procedure was repeated 10 times for each sample, to account for sampling
210 variation. The type I error rate at the nominal level α was evaluated by assessing the quantity $fp =$
211 $\frac{\#\{p\text{-value}_i \leq \alpha, i = 1, \dots, G\}}{G}$ where G is the total number of genes tested. We decided to provide an adjusted
212 prediction interval (PI), accounting for the large number of methods investigated, with the type I error
213 rate as suggested in previous studies [19]. The bounds of this interval are $fp \pm Z_{0.975/\#(methods)}$
214 $\sqrt{fp(1-fp)/G}$ where $Z_{0.975/\#(methods)}$ replaces the usual 97.5 percentile of the normal distribution
215 $Z_{0.975}$ after adjustment for the number of methods investigated. An approach was considered to
216 provide a good correction if its type I error rate was found within this interval.

217
218 **Power studies.** Power was estimated under the alternative hypothesis of an association between a
219 gene and the phenotype (H_1). We selected a subset of 10 genes for the power analysis. All these genes
220 had a cumulative frequency of rare variants (*i.e.* with $MAF \leq 5\%$) of $\sim 10\%$ (*i.e.* $\sim 20\%$ of carriers) and
221 at least 10 mutation carriers. In addition, we considered $\sim 50\%$ of the rare variants of each gene to be
222 causal, with the same direction of effect, and used the presence of at least one of these variants to
223 define the binary genetic score described in the Statistical method section. This implies that there was
224 no cumulative effect of carrying several causal variants, and that the relative risk is defined at the
225 gene level. Table S4 provides details of the 10 genes selected and their causal variants for the
226 European and Worldwide samples. For each gene tested, a phenotype was simulated, using a binomial
227 distribution and penetrance as parameters. For each stratification scenario, penetrance was calculated
228 from the proportion of cases and controls, the frequency of carriers, and the relative risk ($RR=1, 2, 3, 4$).
229 An example is presented in Table S5 for the first gene tested. Tests of association between the genes
230 and the simulated phenotypes were performed 500 times per gene, and power was estimated by
231 evaluating the same quantity as for the type I error rate averaged over the 10 genes and the 500
232 replicates.

233 234 **Statistical methods**

235
236 **Association test.** Let us now consider an association study including n individuals. The binary
237 phenotype is denoted $\mathbf{Y} = (y_1, \dots, y_n)$, where y_i is the status of individual i coded 0 (healthy) or 1
238 (affected). We call $\mathbf{X} = (x_{ij})_{i=1\dots n, j=1\dots p}$ the $n \times p$ genotype matrix for n individuals and p markers. Each

239 term x_{ij} corresponds to the genotype of sample i at marker j and is coded 0, 1 or 2 according to the
240 number of minor alleles. We also introduce the normalized genotype matrix $\tilde{\mathbf{X}} = (\tilde{x}_{ij})_{i=1\dots n, j=1\dots p}$, where
241 each term is $\tilde{x}_{ij} = \frac{x_{ij} - \mu_j}{\sqrt{f_j(1-f_j)}}$ with μ_j the column mean and f_j the observed allele frequency of each
242 marker.

243
244 Several routine statistical tests are available for assessing the association between rare variants and a
245 phenotype. Considering our focus on a small number of cases with phenotypes driven by the presence
246 of at least one causal variant, the most appropriate approach is that based on the CAST method [3].
247 This approach collapses variants into a single genetic score that takes a value of 0 if there are no rare
248 variants in the region or 1 if there is at least one variant. Considering a given genetic region g , in our
249 case a gene, the score for this region is denoted $\mathbf{Z}_g = (z_{g1}, \dots, z_{gn})$, where $z_{gi} = I$ (at least one rare
250 variant in the region g for individual i), $I()$ being the indicator function.

251
252 The corresponding association test can be expressed in a logistic regression framework.

$$\text{logit}(P(\mathbf{Y} = 1)) = \alpha + \beta_g \mathbf{Z}_g$$

253
254
255
256 Where α and β_g are the model parameters for the intercept and the genetic score. Under the null
257 hypothesis of no association $\{\beta_g = 0\}$ the likelihood ratio test (LRT) statistics follow a χ_{1df}^2
258 distribution.

259
260 **Genetic similarity.** Certain methods, including PC and LMM, account for population stratification
261 by using a large number of single-nucleotide polymorphisms (SNPs) to derive genetic similarity
262 matrices (also called relatedness matrices). Considering a set H of p_H SNPs, a normalized similarity
263 matrix $\mathbf{S}^H = \tilde{\mathbf{X}}^H \tilde{\mathbf{X}}^{H'}$ can be derived, where $\tilde{\mathbf{X}}^H$ is the normalized genotype matrix reduced to the
264 markers of set H . Each term $s_{ik, i=1\dots n, k=1\dots n}$ represents the genetic similarity between samples i and k
265 based on the SNPs of set H .

266
267 With whole-exome sequencing (WES) data, a broad range of SNPs are now available, and it is usual
268 to separate them into categories based on their minor allele frequencies (MAFs) [18, 19, 24]. We will
269 consider four categories of variants, based on the MAFs calculated for the total sample: rare variants
270 (RVs; $0\% < MAF < 1\%$), low-frequency variants (LFVs; $1\% \leq MAF < 5\%$), common variants (CVs;
271 $MAF \geq 5\%$) and all variants (ALLVs; the union of RVs, LFVs and CVs). We excluded private variants

272 from these sets of variants, because their sparse distribution tends to have a strong influence on the
273 calculation of similarity matrices. We also pruned all these sets to remove variants with a pairwise r^2
274 < 0.2 , to reduce the effect of linkage disequilibrium. We investigated the effect of using these different
275 sets of SNPs $H \in \{RVs, LFVs, CVs, ALLVs\}$ to derive PC-based or LMM corrections.

276
277 **Principal component (PC) approach.** PC analysis creates new variables from SNP data, the
278 principal components, corresponding to axes of genetic variation. These variables can be included,
279 as covariates, in a statistical model, such as the one described above to adjust for population
280 stratification. Principal components $PC^H = (PC_1^H, \dots, PC_{n-1}^H)$ are based on a given set of SNPs H and
281 are derived from the singular vector decomposition of the normalized similarity matrix S^H . After
282 adjustment for the first m principal components, the corresponding logistic model becomes:

$$283 \logit(P(Y = 1)) = \alpha + \beta_g Z_g + \gamma_1 PC_1^H + \dots + \gamma_m PC_m^H$$

284
285 where $\gamma_1, \dots, \gamma_m$ are new model parameters for the PCs.

286
287 Under the null hypothesis of no association $\{\beta_g = 0\}$, the LRT statistics follow a χ_{1df}^2 distribution.
288 We investigated correction based on the first 3, 5, 10 or 50 PCs, calculated on the four possible sets
289 of variants, RVs, CVs, LFVs and ALLVs. In the following, we use a notation such that PC3_{CV}, for
290 example, indicates that the first three PCs based on common variants were used.

291
292 **Linear mixed models (LMM).** Linear mixed models were initially developed to alleviate the effect
293 of familial relatedness in association analyses, and have also been used to correct for population
294 stratification in GWAS. This regressive approach considers both fixed and random effects and uses
295 a genetic similarity matrix to improve estimation of the parameters of interest. Using the previous
296 CAST regression framework, the LMM model becomes:

$$297 Y = \alpha + \beta_g Z_g + \mathbf{u} + \epsilon$$

298
299 where $\mathbf{u} \sim MVN(0, \tau S^H)$ is a vector of random effects based on the similarity matrix S^H and an
300 additional variance parameter τ . Under the null hypothesis of no association $\{\beta_g = 0\}$, the LRT
301 statistics follow a χ_{1df}^2 distribution. We focus here on LMM based on the relatedness matrices
302 constructed with the four sets of variants previously described, and with for instance the notation
303
304

305 LMM_{CV} indicating that common variants were used.

306

307 **Adapted local permutations (LocPerm).** Permutation strategies have been designed to derive p -
308 values when the 'true' null distribution of the test statistic T_0 is unknown [36]. This is the case for
309 population stratification, which creates a bias that cannot be numerically derived. The rationale
310 behind permutation procedures is to simulate several test statistics (T_1, \dots, T_B) under the null
311 hypothesis, to derive an approximated distribution as close as possible to the unknown true null
312 distribution, and to use these statistics to estimate a p -value. With the classical permutation approach,
313 the simulation of test statistics under H_0 is achieved by randomly resampling phenotypes (*i.e.*
314 exchanging them between individuals). Adapted local permutations are based on the observation that,
315 in the presence of population structure, not all phenotypes are exchangeable [29]. A given sample has
316 a higher chance of sharing its phenotype with another sample of the same ancestry. The principle is,
317 therefore, to establish, for each sample, a neighborhood, *i.e.* a set of samples between which it is
318 reasonable to exchange phenotypes. These neighborhoods are established according to a genetic
319 distance derived from the first 10 PCs:

320

$$321 \quad d_{ij}^2 = \sum_{k=1}^{10} \lambda_k |PC_{ki}^{CV} - PC_{kj}^{CV}|^2$$

322

323 where PC^{CV} is the matrix of principal components calculated on the set of common variants and λ_k
324 is the eigenvalue corresponding to the k -th principal component PC_k^{CV} . This distance is used to create
325 a neighborhood of 30 individuals around each sample [29]. Permutations can then be performed for
326 each sample, within its neighborhood.

327

328 A straightforward empirical way to derive a p -value for the permutation test is to assess the quantity
329 $pv = \#\{T_i \geq T_0\}/B$ where $\#$ is the cardinal function and B is the number of permutations. This method
330 is dependent on the number of permutations computed, and a large number of permutations is required
331 for the accurate estimation of small p -values. Mullaert et al. proposed an alternative semi-parametric
332 approach, in which a limited number of resampled statistics are used to estimate the mean (m) and
333 standard deviation (σ) of the test statistic under H_0 . The previously described CAST-like LRT
334 statistics are used, through their square roots with a sign attributed according to the direction of the
335 effect, $T_i = \text{sign}(\text{effect})\sqrt{|LRT|}$, to estimate the $N(m, \sigma^2)$ distribution parameters and then calculate
336 the p -value. We evaluated both the semi-parametric approach using 500 local permutations and the
337 full empiric approach using 5000 local permutations. These two approaches yielded very similar

338 results. We therefore present here only the results for the semi-empiric approach.

339

340 **Implementation of the simulations and methods.** We used R software (<https://www.R-project.org/>)
341 to code the comparison pipeline and implement the logistic and permutation models. Principal
342 components and similarity matrices were obtained with Plink2 software ([https://www.cog-
344 genomics.org/plink/2.0/](https://www.cog-
343 genomics.org/plink/2.0/)), and GEMMA was used for the LMM method [10, 37].

345 **Results**

346

347 **Large study size**

348

349 The results of the simulation study under the null hypothesis for the European sample of 1,523
350 individuals are presented in Table 3 (for $\alpha=0.001$) and Table S6 (for $\alpha=0.01$). In the absence of
351 stratification, the four methods had correct type I error rates, within the 95% PI bounds (Table 3A,
352 Table S6A). This was the case for PC3 and LMM, regardless of the type of variant considered. In the
353 presence of moderate stratification (Table 3B, Table S6B), the unadjusted CAST approach displayed
354 the expected inflation of type I error rate (0.00163 at $\alpha=0.001$). The PC3 method corrected properly
355 regardless of the type of variant at $\alpha=0.001$, but a slight inflation of type I error was observed for RVs
356 and LFVs at $\alpha=0.01$. The use of LMM led to an inflation of type I error rates at $\alpha=0.001$, unless all
357 variants were considered, which gave rates within the 95% PI at $\alpha=0.01$. LocPerm had a correct type
358 I error rate at both α levels. In the presence of strong stratification (Table 3C, Table S6C), the
359 unadjusted CAST method gave a strong inflation of type I error rate, to 0.00359 at $\alpha=0.001$. The PC
360 and LMM approaches also led to inflated type I errors (between 0.00133 and 0.00175 at $\alpha=0.001$),
361 the lowest level of inflation being observed when CVs or all variants were considered. For the PC
362 approach, increasing the number of PCs did not improve the correction, consistent with previous
363 findings reported by Persyn et al. (2018). The use of 50 PCs resulted in an inflation of type I error
364 whatever the level of stratification, probably due to an overadjustment of the regression model (Table
365 S7). Thus, in the presence of strong population structure, classical methods were unable to handle the
366 stratification properly. The adapted local permutations approach was the only method able to correct
367 for stratification in this scenario, with a slightly conservative result of 0.00863 at $\alpha=0.01$ (Table S6C).

368

369 **Table 3: Type I error rates of the different approaches for the large size European sample.** The
370 nominal level alpha considered is $\alpha = 0.001$ and the corresponding 95%PI adjusted for the 10

371 methods is [0,00079-0,00121]. Type I error rates under the lower bound of the 95%PI are displayed
 372 in italic and above the upper bound of the 95%PI in bold.

373

A - No Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.00106	0.00108	0.00118	0.00082
LFVs		0.0011	0.00119	
CVs		0.00104	0.00118	
ALLVs		0.00108	0.00116	
B – Moderate Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.00163	0.00117	0.00141	0.00095
LFVs		0.00101	0.00125	
CVs		0.001	0.00124	
ALLVs		0.00102	0.00117	
C – High Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.00359	0.00157	0.00175	0.00087
LFVs		0.00137	0.00176	
CVs		0.00136	0.00161	
ALLVs		0.00133	0.00145	

374

375 The results of the simulation study under H_0 for the Worldwide sample of 1,967 individuals are
 376 presented in Table 4 (for $\alpha=0.001$) and Table S8 (for $\alpha=0.01$). In the absence of stratification, none
 377 of the main approaches had a significantly inflated type I error rate (Table 4A and Table S8A). At α
 378 $=0.01$, LMM corrections were slightly conservative. The presence of moderate or strong stratification
 379 led to extremely inflated type I errors at $\alpha=0.001$ for the unadjusted CAST approach, with values of
 380 0.00681 and 0.137, respectively. For PC3 and LMM, a satisfactory correction was obtained at
 381 $\alpha=0.001$ with CVs, whereas, at $\alpha=0.01$, PC gave a slight inflation of type I error and LMM results
 382 were slightly conservative. The three other types of variants could not properly account for
 383 stratification for PC3 and LMM. Increasing the number of PCs did not improve the results obtained
 384 with PC3 (Table S9) for the Worldwide sample. LocPerm maintained a correct type I error rate in
 385 both scenarios, with values of 0.00096 and 0.00113 at $\alpha=0.001$ for moderate and strong stratification,

386 respectively. Overall, the analyses under the null hypothesis within the European and Worldwide
 387 samples showed that accounting for stratification was generally more difficult with a continental
 388 structure than with a worldwide structure. PC3 and LMM based on CVs were capable of maintaining
 389 a correct type I error rate in most of the situations considered, with the exception of high levels of
 390 stratification in Europe, and LocPerm correctly accounted for stratification in all the situations
 391 considered.

392
 393 **Table 4: Type I error rates of the different approaches for the large size Worldwide sample.**
 394 The nominal level alpha considered is $\alpha = 0.001$ and the corresponding 95%PI adjusted for the 10
 395 methods is [0,00079-0,00121]. Type I error rates under the lower bound of the 95%PI are displayed
 396 in italic and above the upper bound of the 95%PI in bold.

A - No Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.00085	0.00099	0.00093	0.00087
LFVs		0.00099	0.00094	
CVs		0.00099	0.00093	
ALLVs		0.00099	0.00093	
B – Moderate Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.00681	0.00259	0.00456	0.00096
LFVs		0.00109	0.00123	
CVs		0.00105	0.00117	
ALLVs		0.00128	0.00162	
C – High Stratification				
	CAST	PC3	LMM	LocPerm
RVs	0.13698	0.00662	0.01834	0.00113
LFVs		0.0012	0.00163	
CVs		0.00119	0.00115	
ALLVs		0.00127	0.00266	

397
 398 With respect to the results of the simulation under H_0 , we focused the power studies on the methods
 399 providing satisfactory correction (*i.e.* PC3_{CV}, LMM_{CV} and LocPerm), in addition to the unadjusted
 400 CAST. Only powers derived from a correct type I error rate under H_0 are presented in the main text.

401 Adjusted powers accounting for inflated type I error rates are provided in the Supplementary figures
402 for information. The results of the power study for the European sample are presented in Figure 3 and
403 Figure S1. In situations with no stratification or moderate stratification, all approaches had similar
404 powers, of about 50% at $\alpha=0.001$ for a relative risk of 3, for example (Figure 3). In the presence of
405 strong stratification, only LocPerm was able to correct for confounding and to maintain power levels
406 (Figure 3C). The adjusted powers (Figure S1) indicate that all three correction methods provide very
407 similar powers when type I error is controlled. The results of the power study for the Worldwide
408 sample are presented in Figure 4 and Figure S2. As for the European sample, all methods had similar
409 powers in the absence of stratification or the presence of moderate stratification. In the presence of
410 strong stratification, LocPerm was slightly less powerful than the other methods (Figure 4C) with for
411 a RR of 3 at $\alpha=0.001$, a power of 64% as opposed to the powers of 69 and 72% obtained for PC3_{CV},
412 and LMM_{CV}, respectively. It is also interesting to compare the power of each method, separately,
413 between the different stratification scenarios (Figures S3 for the European sample and S4 for the
414 Worldwide sample). Power was very similar for any given technique in the different stratification
415 scenarios, indicating that the correction methods maintained the level of power observed in the
416 absence of stratification.

417

418 **Small study size**

419

420 The results of the simulation study under the null hypothesis for a small sample size, based on 50
421 cases, are presented in Table 5 (for $\alpha=0.001$) and Table S10 (for $\alpha=0.01$). Only PC3_{CV}, LMM_{CV} and
422 LocPerm, which provided a satisfactory correction for stratification in the large sample study, were
423 investigated for small sample sizes. In scenarios without stratification (*i.e.* controls and cases of the
424 same origin), an inflation of type I errors was observed: 1) with PC3 (about 0.0015 at $\alpha=0.001$) when
425 the number of controls was low (100), and, to a lesser extent, with CAST (about 0.0012 at $\alpha=0.001$),
426 and 2) with LMM (about 0.002 at $\alpha=0.001$) when the number of controls was high (1000 or 2000).
427 In the presence of stratification (*i.e.* a large number of controls with an origin different from that of
428 the cases), a strong inflation of type I error rates was observed for CAST. This was also the case for
429 LMM_{CV}, albeit to a lesser extent, particularly for stratification within Europe or when the cases came
430 from the Worldwide sample and the controls from Europe only. Both PC3_{CV} and LocPerm provided
431 correct type I error rates in all the scenarios considered with small numbers of cases and a large
432 number of controls.

433

434 **Table 5: Type I error rates of the different approaches for the small size sample scenarios.** The

435 nominal level alpha considered is $\alpha = 0.001$. Type I error rates under the lower bound of the 95%PI
 436 are displayed in italic and above the upper bound of the 95%PI in bold.
 437 STable 3 provides the adjusted 95%PI for the different number of genes tested in each scenario.
 438

Scenario	CAST	PC3 _{CV}	LMM _{CV}	LocPerm
50SE-100SE	0.0012	0.0015	0.0012	0.0009
50SE-1000E	0.0016	0.0012	0.0028	0.0008
50SE-1000W	0.0046	0.0011	0.0015	0.0010
50SE-2000W	0.0046	0.0010	0.0016	0.0011
50E-100E	0.0014	0.0015	0.0012	0.0010
50E-1000E	0.0010	0.0010	0.0021	0.0009
50E-1000W	0.0051	0.001	0.0014	0.0010
50E-2000W	0.0050	0.0009	0.0015	0.0011
50W-100W	0.0013	0.0015	0.0012	0.0010
50W-1000E	0.0077	<i>0.0007</i>	0.0053	0.0010
50W-1000W	0.0009	0.0010	0.0021	0.0009
50W-2000W	0.0009	0.0009	0.002	0.0010

439
 440 A power study was performed for PC3_{CV} and LocPerm with small numbers of cases (Figure 5). Both
 441 approaches gave a correct type I error rate and similar results, but power was slightly higher for
 442 LocPerm than for PC3 when the 50 cases came from Europe as a whole or from the Worldwide
 443 sample. For cases were from Southern Europe, considering 1000 controls from the whole of Europe
 444 gave a power twice that obtained when considering 100 controls of the same origin as the cases
 445 (Figure 5A). For example, for a RR of 4 and at $\alpha=0.001$, the power increased from 15% to 34% under
 446 these conditions with LocPerm. A smaller increase was observed if 1000 controls from the Worldwide
 447 sample were used, increasing to a similar level with the use of 2000 Worldwide controls. When the
 448 cases were from anywhere in Europe, a similar increase in power was observed with 1000 European
 449 and with 1000 Worldwide controls, whereas the use of 2000 Worldwide resulted in no greater a power
 450 than the use of 1000 Worldwide controls. Finally, when the cases were selected from the Worldwide
 451 sample, the use of 1000 Worldwide controls gave a power almost double that achieved with 100
 452 Worldwide controls, whereas the use of 1000 controls from Europe did not substantially increase the
 453 power. These results indicate that using a large panel of worldwide controls to increase sample size
 454 is a good strategy for increasing the power of a study while correcting for stratification with

455 approaches such as PC3_{CV} or LocPerm.

456

457 **Computational considerations**

458

459 We also assessed the computing time required for the different approaches. While the unadjusted
460 CAST method does not imply the computation of any particular matrix, the same covariance matrix
461 is necessary for PC3_{CV}, LMM and LocPerm and additional specific permutation matrices are required
462 for LocPerm only. We ran each method separately, CAST, PC3_{CV} and LocPerm with R, and LMM_{CV}
463 with GEMMA, on the 1,523 individuals and the 17,619 genes of the European sample, under a
464 hypothesis of no association. We broke down the runtime of each method into a pretreatment phase
465 (covariance and permutation matrices) and a gene-testing phase (see Table S11). The pretreatment
466 runtime was dependent only on the number of individuals (and the set of SNPs used for the
467 calculations) and this part of the analysis was performed only once. The runtime of the gene-testing
468 phase depended on the number of individuals and the number of genes tested, and could be repeated
469 for different analyses (*e.g.* for different MAF thresholds). PC3_{CV} and LMM_{CV} had similar
470 pretreatment times, markedly shorter than that for LocPerm, which also requires the calculation of
471 permutation matrices. However, the need to calculate these matrices only once decreases the relative
472 disadvantage of the LocPerm method. In terms of gene-testing time, LMM_{CV} was the fastest approach
473 when used with GEMMA, but this may not be the case for other programs that have not been
474 optimized. A comparison of the methods implemented with R showed that the adjustment on PCs and
475 LocPerm took 1.4x and 2.5x longer, respectively, than the unadjusted test. These comparisons were
476 run on a 64-bit Intel Xeon Linux machine with a CPU of 3.70 GHz and 64 GB of RAM.

477

478 **Discussion**

479

480 We performed a large simulation study based on real exomes data to investigate the ability of several
481 approaches (*i.e.* PCs, LMM and LocPerm) to account for population stratification in rare-variant
482 association studies of a binary trait. In our simulation study, the efficiency of PCs and LMM to
483 correct for population stratification was dependent on the type of variant used to derive the similarity
484 matrices, the best performance being obtained with CVs. It was generally not possible to correct the
485 stratification bias with RVs, even with the exclusion of private variants for the calculation of the
486 matrices. Private variants have very sparse distributions, which may lead to difficulties in calculation,
487 and their inclusion resulted in an even lower efficiency of correction for population structure (data
488 not shown). Other studies evaluating different types of variants reached the same conclusions [24,

489 25] although one reported better performances for PC based on RVs [14]. However, this study was
490 based on simulated NGS data, which may have led to an unrealistic rare variant distribution. Our
491 results also indicate that CVs or ALLVs were the best sets of variants for the LMM approach applied
492 to CAST, confirming the results of Luo et al. based on the SKAT test [19]. Variant selection remains
493 an area in which there are perspectives for improving the corrections provided by strategies such as
494 PC or LMM [12, 26], although the use of CVs appeared to be a good choice in most situations.

495

496

497 With the optimal set of variants, PC generally corrected for population stratification more efficiently
498 than LMM. This is consistent with benefits of the PC approach over LMM observed in the presence
499 of spatially confined confounders [38], which is often the case with rare variants. For large sample
500 sizes, both PC and LMM controlled for stratification better at larger geographic scales than at finer
501 scales. In small samples (50 cases and 100 controls), PC approaches gave inflated type I errors even
502 in the absence of population stratification, as previously reported [18, 29, 39]. This inflation
503 disappeared when the sample included additional controls, whatever their ethnic origin, even with a
504 highly unbalanced case-control ratio. By contrast, the type I error of LMM was inflated in samples
505 with highly unbalanced case-control ratios, whatever the level of population stratification, as
506 previously noted in the context of GWAS [40]. Finally, the adapted local permutations procedure
507 recently proposed by Mullaert et al. [29] gave very promising results, as it fully corrected for
508 population stratification, regardless of the scale over which the stratification occurred, sample size
509 and case-control ratio. When valid under H_0 , the three correction methods had similar powers. For a
510 given setting, power was similar in the different stratification situations, indicating that the correction
511 method could maintain the power it would have in the absence of stratification. These results are in
512 partial agreement with several studies reporting a small loss of power for PC-adjusted logistic
513 regression in the presence of stratification relative to an absence of stratification [13, 20].

514

515 We also investigated the specific situation in which only a very small number of cases are available,
516 which is particularly relevant in the context of rare disorders. In this setting, we show that PC and
517 LocPerm provide correct type I errors when the number of controls is large, regardless of the ethnic
518 origin of the controls. In addition, the strategy of adding controls, even of worldwide origin, provided
519 a substantial gain of power for PC and LocPerm when few cases were available. This is an important
520 finding, highlighting the potential interest of using publicly available controls, such as those of the
521 1000G project, to increase the power of a study with a small sample size. We also investigated an
522 additional scenario in which all cases were strictly from our in-house HGID cohort and the controls

523 were obtained from both the HGID and 1000 Genomes cohorts (data not shown). This scenario gave
524 identical results to those presented here, indicating that, even in the presence of heterogeneity in the
525 types of exome data considered for cases and controls (*e.g.* in terms of kit or technology used), the
526 conclusions drawn here still apply. Overall, these results validate a strategy of using additional
527 external controls to increase the power of a study, provided that an efficient stratification correction
528 approach is used.

529
530 We focused on the investigation of diseases caused by a few deleterious variants, for which the
531 CAST-like approach is particularly appropriate. Additional studies are required to investigate more
532 complex genetic models, such as the presence of both risk and protective variants of a given gene, for
533 which other association tests, such as variant-component approaches, may be more appropriate.
534 Different results can be expected, as the effect of population stratification differs between testing
535 strategies [17, 20]. In addition, the novel LocPerm strategy has not been evaluated in combination
536 with other association tests. In the situations we considered, our study highlighted several useful
537 conclusions for rare variant association studies in the presence of stratification: 1) the key issue is to
538 properly control type I errors as powers are comparable, 2) population stratification can be corrected
539 by PC_{3CV} in most instances, unless there is a high degree of intracontinental stratification and a small
540 sample size, 3) LocPerm proposes a satisfying correction in all instances, and 4) strategies based on
541 the inclusion of a large number of additional controls (*e.g.* from publicly available databases) provide
542 a substantial gain of power provided that stratification is controlled for correctly.

543 **Acknowledgments**

544

545 We thank both branches of the Laboratory of Human Genetics of Infectious Diseases for helpful
546 discussions and support. The Laboratory of Human Genetics of Infectious Diseases was supported in
547 part by grants from the French National Agency for Research (ANR) under the “Investissement
548 d’avenir” program (grant number ANR-10-IAHU-01), the TBPATHEGEN project (ANR-14-CE14-
549 0007-01), the MYCOPARADOX project (ANR-16-CE12-0023-01), the Integrative Biology of
550 Emerging Infectious Diseases Laboratory of Excellence (grant number ANR-10-LABX-62-IBEID),
551 the St. Giles Foundation, the National Center for Research Resources and the National Center for
552 Advancing Sciences (NCATS), and the Rockefeller University.

553

554

555 Bibliography

- 556 1. Auer PL, Lettre G. Rare variant association studies: considerations, challenges and
557 opportunities. *Genome Med.* 2015;7(1):16. doi: 10.1186/s13073-015-0138-2. PubMed PMID:
558 25709717; PubMed Central PMCID: PMC4337325.
- 559 2. Lee S, Abecasis GR, Boehnke M, Lin X. Rare-variant association analysis: study designs and
560 statistical tests. *Am J Hum Genet.* 2014;95(1):5-23. doi: 10.1016/j.ajhg.2014.06.009. PubMed
561 PMID: 24995866; PubMed Central PMCID: PMC4085641.
- 562 3. Morgenthaler S, Thilly WG. A strategy to discover genes that carry multi-allelic or mono-
563 allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutat Res.* 2007;615(1-2):28-
564 56. doi: 10.1016/j.mrfmmm.2006.09.003. PubMed PMID: 17101154.
- 565 4. Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for
566 sequencing data with the sequence kernel association test. *Am J Hum Genet.* 2011;89(1):82-93. doi:
567 10.1016/j.ajhg.2011.05.029. PubMed PMID: 21737059; PubMed Central PMCID:
568 PMC3135811.
- 569 5. Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet.*
570 2006;2(12):e190. doi: 10.1371/journal.pgen.0020190. PubMed PMID: 17194218; PubMed Central
571 PMCID: PMC1713260.
- 572 6. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal
573 components analysis corrects for stratification in genome-wide association studies. *Nat Genet.*
574 2006;38(8):904-9. doi: 10.1038/ng1847. PubMed PMID: 16862161.
- 575 7. Kang HM, Sul JH, Service SK, Zaitlen NA, Kong SY, Freimer NB, et al. Variance component
576 model to account for sample structure in genome-wide association studies. *Nat Genet.*
577 2010;42(4):348-54. doi: 10.1038/ng.548. PubMed PMID: 20208533; PubMed Central PMCID:
578 PMC3092069.
- 579 8. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D. FaST linear mixed
580 models for genome-wide association studies. *Nat Methods.* 2011;8(10):833-5. doi:
581 10.1038/nmeth.1681. PubMed PMID: 21892150.
- 582 9. Listgarten J, Lippert C, Kadie CM, Davidson RI, Eskin E, Heckerman D. Improved linear
583 mixed models for genome-wide association studies. *Nat Methods.* 2012;9(6):525-6. doi:
584 10.1038/nmeth.2037. PubMed PMID: 22669648; PubMed Central PMCID: PMC3597090.
- 585 10. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association studies.
586 *Nat Genet.* 2012;44(7):821-4. doi: 10.1038/ng.2310. PubMed PMID: 22706312; PubMed Central
587 PMCID: PMC3386377.
- 588 11. Jiang Y, Epstein MP, Conneely KN. Assessing the impact of population stratification on
589 association studies of rare variation. *Hum Hered.* 2013;76(1):28-35. doi: 10.1159/000353270.
590 PubMed PMID: 23921847; PubMed Central PMCID: PMC4406348.
- 591 12. Mathieson I, McVean G. Differential confounding of rare and common variants in spatially
592 structured populations. *Nat Genet.* 2012;44(3):243-6. doi: 10.1038/ng.1074. PubMed PMID:
593 22306651; PubMed Central PMCID: PMC3303124.
- 594 13. O'Connor TD, Kiezun A, Bamshad M, Rich SS, Smith JD, Turner E, et al. Fine-scale patterns
595 of population stratification confound rare variant association tests. *PLoS One.* 2013;8(7):e65834.
596 doi: 10.1371/journal.pone.0065834. PubMed PMID: 23861739; PubMed Central PMCID:
597 PMC3701690.
- 598 14. Liu Q, Nicolae DL, Chen LS. Marbled inflation from population structure in gene-based
599 association studies with rare variants. *Genet Epidemiol.* 2013;37(3):286-92. doi:
600 10.1002/gepi.21714. PubMed PMID: 23468125; PubMed Central PMCID: PMC3716585.
- 601 15. De la Cruz O, Raska P. Population structure at different minor allele frequency levels. *BMC*
602 *Proc.* 2014;8(Suppl 1 Genetic Analysis Workshop 18Vanessa Olmo):S55. doi: 10.1186/1753-6561-
603 8-S1-S55. PubMed PMID: 25519390; PubMed Central PMCID: PMC4143691.
- 604 16. Moore CB, Wallace JR, Wolfe DJ, Frase AT, Pendergrass SA, Weiss KM, et al. Low

- 605 frequency variants, collapsed based on biological knowledge, uncover complexity of population
606 stratification in 1000 genomes project data. *PLoS Genet.* 2013;9(12):e1003959. doi:
607 10.1371/journal.pgen.1003959. PubMed PMID: 24385916; PubMed Central PMCID:
608 PMCPMC3873241.
- 609 17. Zawistowski M, Reppell M, Wegmann D, St Jean PL, Ehm MG, Nelson MR, et al. Analysis of
610 rare variant population structure in Europeans explains differential stratification of gene-based tests.
611 *Eur J Hum Genet.* 2014;22(9):1137-44. doi: 10.1038/ejhg.2013.297. PubMed PMID: 24398795;
612 PubMed Central PMCID: PMCPMC4135410.
- 613 18. Babron MC, de Tayrac M, Rutledge DN, Zeggini E, Genin E. Rare and low frequency variant
614 stratification in the UK population: description and impact on association tests. *PLoS One.*
615 2012;7(10):e46519. doi: 10.1371/journal.pone.0046519. PubMed PMID: 23071581; PubMed
616 Central PMCID: PMCPMC3465327.
- 617 19. Luo Y, Maity A, Wu MC, Smith C, Duan Q, Li Y, et al. On the substructure controls in rare
618 variant analysis: Principal components or variance components? *Genet Epidemiol.* 2018;42(3):276-
619 87. doi: 10.1002/gepi.22102. PubMed PMID: 29280188; PubMed Central PMCID:
620 PMCPMC5851819.
- 621 20. Persyn E, Redon R, Bellanger L, Dina C. The impact of a fine-scale population stratification
622 on rare variant association test results. *PLoS One.* 2018;13(12):e0207677. doi:
623 10.1371/journal.pone.0207677. PubMed PMID: 30521541; PubMed Central PMCID:
624 PMCPMC6283567.
- 625 21. Wang C, Zhan X, Bragg-Gresham J, Kang HM, Stambolian D, Chew EY, et al. Ancestry
626 estimation and control of population stratification for sequence-based association studies. *Nat*
627 *Genet.* 2014;46(4):409-15. doi: 10.1038/ng.2924. PubMed PMID: 24633160; PubMed Central
628 PMCID: PMCPMC4084909.
- 629 22. Baye TM, He H, Ding L, Kurowski BG, Zhang X, Martin LJ. Population structure analysis
630 using rare and common functional variants. *BMC Proc.* 2011;5 Suppl 9:S8. doi: 10.1186/1753-
631 6561-5-S9-S8. PubMed PMID: 22373300; PubMed Central PMCID: PMCPMC3287920.
- 632 23. Sha Q, Zhang K, Zhang S. A Nonparametric Regression Approach to Control for Population
633 Stratification in Rare Variant Association Studies. *Sci Rep.* 2016;6:37444. doi: 10.1038/srep37444.
634 PubMed PMID: 27857226; PubMed Central PMCID: PMCPMC5114546.
- 635 24. Zhang Y, Guan W, Pan W. Adjustment for population stratification via principal components
636 in association analysis of rare variants. *Genet Epidemiol.* 2013;37(1):99-109. doi:
637 10.1002/gepi.21691. PubMed PMID: 23065775; PubMed Central PMCID: PMCPMC4066816.
- 638 25. Zhang Y, Shen X, Pan W. Adjusting for population stratification in a fine scale with principal
639 components and sequencing data. *Genet Epidemiol.* 2013;37(8):787-801. doi: 10.1002/gepi.21764.
640 PubMed PMID: 24123217; PubMed Central PMCID: PMCPMC3864649.
- 641 26. Listgarten J, Lippert C, Heckerman D. FaST-LMM-Select for addressing confounding from
642 spatial structure and rare variants. *Nat Genet.* 2013;45(5):470-1. doi: 10.1038/ng.2620. PubMed
643 PMID: 23619783.
- 644 27. Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, et al. A global
645 reference for human genetic variation. *Nature.* 2015;526(7571):68-74. doi: 10.1038/nature15393.
646 PubMed PMID: 26432245; PubMed Central PMCID: PMCPMC4750478.
- 647 28. Boisson-Dupuis S, Ramirez-Alejo N, Li Z, Patin E, Rao G, Kerner G, et al. Tuberculosis and
648 impaired IL-23-dependent IFN-gamma immunity in humans homozygous for a common TYK2
649 missense variant. *Sci Immunol.* 2018;3(30). doi: 10.1126/sciimmunol.aau8714. PubMed PMID:
650 30578352; PubMed Central PMCID: PMCPMC6341984.
- 651 29. Mullaert J, Bouaziz M, Seeleuthner Y, Bigio B, Casanova J-L, Alcais A, et al. Taking
652 population stratification into account by local permutations in rare-variant association studies on
653 small samples. *bioRxiv.* 2020:2020.01.29.924977. doi: 10.1101/2020.01.29.924977.
- 654 30. Moutsianas L, Agarwala V, Fuchsberger C, Flannick J, Rivas MA, Gaulton KJ, et al. The

- 655 power of gene-based rare variant methods to detect disease-associated variation and test hypotheses
656 about complex disease. *PLoS Genet.* 2015;11(4):e1005165. doi: 10.1371/journal.pgen.1005165.
657 PubMed PMID: 25906071; PubMed Central PMCID: PMC4407972.
- 658 31. Belkadi A, Bolze A, Itan Y, Cobat A, Vincent QB, Antipenko A, et al. Whole-genome
659 sequencing is more powerful than whole-exome sequencing for detecting exome variants. *Proc Natl*
660 *Acad Sci U S A.* 2015;112(17):5473-8. doi: 10.1073/pnas.1418631112. PubMed PMID: 25827230;
661 PubMed Central PMCID: PMC4418901.
- 662 32. Anderson CA, Pettersson FH, Clarke GM, Cardon LR, Morris AP, Zondervan KT. Data
663 quality control in genetic case-control association studies. *Nat Protoc.* 2010;5(9):1564-73. doi:
664 10.1038/nprot.2010.116. PubMed PMID: 21085122; PubMed Central PMCID: PMC3025522.
- 665 33. Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. Robust relationship
666 inference in genome-wide association studies. *Bioinformatics.* 2010;26(22):2867-73. doi:
667 10.1093/bioinformatics/btq559. PubMed PMID: 20926424; PubMed Central PMCID:
668 PMC3025716.
- 669 34. Consortium UK, Walter K, Min JL, Huang J, Crooks L, Memari Y, et al. The UK10K project
670 identifies rare variants in health and disease. *Nature.* 2015;526(7571):82-90. doi:
671 10.1038/nature14962. PubMed PMID: 26367797; PubMed Central PMCID: PMC4773891.
- 672 35. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank: an open
673 access resource for identifying the causes of a wide range of complex diseases of middle and old
674 age. *PLoS Med.* 2015;12(3):e1001779. doi: 10.1371/journal.pmed.1001779. PubMed PMID:
675 25826379; PubMed Central PMCID: PMC4380465.
- 676 36. Rudolph PE. Good, Ph.: Permutation Tests. A Practical Guide to Resampling Methods for
677 Testing Hypotheses. Springer Series in Statistics, Springer-Verlag, Berlin — Heidelberg — New
678 York: 1994, x, 228 pp., DM 74,00; öS 577.20; sFr 74.—. ISBN 3-540-94097-9. *Biometrical Journal.*
679 1995;37(2):150-. doi: 10.1002/bimj.4710370203.
- 680 37. Zhou X, Stephens M. Efficient multivariate linear mixed model algorithms for genome-wide
681 association studies. *Nat Methods.* 2014;11(4):407-9. doi: 10.1038/nmeth.2848. PubMed PMID:
682 24531419; PubMed Central PMCID: PMC4211878.
- 683 38. Zhang Y, Pan W. Principal component regression and linear mixed model in association
684 analysis of structured samples: competitors or complements? *Genet Epidemiol.* 2015;39(3):149-55.
685 doi: 10.1002/gepi.21879. PubMed PMID: 25536929; PubMed Central PMCID: PMC4366301.
- 686 39. Zhang X, Basile AO, Pendergrass SA, Ritchie MD. Real world scenarios in rare variant
687 association analysis: the impact of imbalance and sample size on the power in silico. *BMC*
688 *Bioinformatics.* 2019;20(1):46. doi: 10.1186/s12859-018-2591-6. PubMed PMID: 30669967;
689 PubMed Central PMCID: PMC6343276.
- 690 40. Zhou W, Nielsen JB, Fritsche LG, Dey R, Gabrielsen ME, Wolford BN, et al. Efficiently
691 controlling for case-control imbalance and sample relatedness in large-scale genetic association
692 studies. *Nat Genet.* 2018;50(9):1335-41. doi: 10.1038/s41588-018-0184-y. PubMed PMID:
693 30104761; PubMed Central PMCID: PMC6119127.

695

696 **Supporting information**

697

698 **Table S1. Distribution of the samples in the European and Worldwide sub-populations**

699

700 **Table S2. Distribution of the variants in the European and the Worldwide samples according**
701 **to their MAFs as described in the Material and Methods section.**

702

703 **Table S3. Number of genes tested and 95%PI in each scenario of the small size sample study.**

704 Prediction intervals are adjusted on the 4 methods tested.

705

706 **Table S4. Details of the genes selected for the power analysis in the European and the**
707 **Worldwide samples.** Freq() indicates the cumulative frequency of the causal variants.

708

709 **Table S5. Example of penetrances used for the power estimation of the gene ADAMTS4 in the**
710 **European sample.** F_0 and F_1 represent the penetrances for non-carriers and carriers considering a
711 relative risk RR of 3, a total proportion of 15% of cases and proportions of carriers of 9% in Northern-
712 Europe (NE) and Middle-Europe (ME) and 8% in Southern-Europe (SE). Within a given sample these
713 penetrances are calculate by $F_0 = n_{cases}/(n_{non-carriers}+RR.n_{carriers})$ and $F_1=RR.F_0$

714

715 **Table S6. Type I error rates of the different approaches for the large size European sample.**

716 The nominal level alpha considered is $\alpha = 0.01$ and the corresponding 95%PI adjusted for the 10
717 methods is [0.00933-0.01067]. Type I error rates under the lower bound of the 95%PI are displayed
718 in italic and above the upper bound of the 95%PI in bold.

719

720 **Table S7. Type I error rates of the PC approach with 3, 5, 10 or 50 PCs for the large size**

721 **European sample.** The nominal level alpha considered is $\alpha = 0.001$ and the corresponding 95%PI
722 adjusted for the 16 methods is [0.00078-0.00122]. Type I error rates under the lower bound of the
723 95%PI are displayed in italic and above the upper bound of the 95%PI in bold.

724

725 **Table S8. Type I error rates of the different approaches for the large size Worldwide sample.**

726 The nominal level alpha considered is $\alpha = 0.01$ and the corresponding 95%PI adjusted for the 10
727 methods is [0.00933-0.01067]. Type I error rates under the lower bound of the 95%PI are displayed
728 in italic and above the upper bound of the 95%PI in bold.

729

730 **Table S9. Type I error rates of the PC approach with 3, 5, 10 or 50 PCs for the large size**
731 **Worldwide sample.** The nominal level alpha considered is $\alpha = 0.001$ and the corresponding 95%PI
732 adjusted for the 16 methods is [0.00078-0.00122]. Type I error rates under the lower bound of the
733 95%PI are displayed in italic and above the upper bound of the 95%PI in bold.

734

735 **Table S10. Type I error rates of the different approaches for the small size sample scenarios.**
736 The nominal level alpha considered is $\alpha = 0.01$. Type I error rates under the lower bound of the
737 95%PI are displayed in italic and above the upper bound of the 95%PI in bold.

738 Table S3 provides the adjusted 95%PI for the different number of genes tested in each scenario.

739

740 **Table S11. Runtime of each method calculated on 1,523 individuals and 17,619 genes of the**
741 **large size European sample under the null hypothesis.** Note that if the analyses are conducted
742 several times, with for instance different MAF thresholds or modes of inheritance, the pre-treatment
743 part does not have to be performed again.

744

745 **Figure S1. Histogram of adjusted powers of the correction methods for the large size European**
746 **sample (n=1,523) at the level $\alpha = 0.001$.** (A) Without stratification. (B) With moderate
747 stratification. (C) With high stratification. Relative risks considered vary from 2 to 4 on the x-axis.

748

749 **Figure S2. Histogram of adjusted powers for the correction methods for the large size**
750 **Worldwide sample at the level $\alpha = 0.001$.** (A) Without stratification. (B) With moderate
751 stratification. (C) With high stratification. Relative risks considered vary from 2 to 4 on the x-axis.

752

753 **Figure S3. Histogram of powers for methods with a correct type I error rate for the large size**
754 **European sample (n=1,523) at the level $\alpha = 0.001$.** (A) Principal components. (B) Linear Mixed
755 Models. (C) LocPerm. Relative risks vary from 2 to 4 on the x-axis.

756

757 **Figure S4. Histogram of powers for methods with a correct type I error rate for the large size**
758 **Worldwide sample at the level $\alpha = 0.001$.** (A) Principal components. (B) Linear Mixed Models.
759 (C) LocPerm. Relative risks vary from 2 to 4 on the x-axis.

760

761 **Figure S5. Histogram of adjusted powers of the correction methods the small size sample at the**
762 **level $\alpha = 0.001$.** (A) Scenarios with 50 cases from Southern-Europe. (B) Scenarios with 50 cases

763 from the whole Europe. (C) Scenarios with 50 cases from the Worldwide sample. The relative risk is
764 fixed at 4.

765 **Figure legends**

766

767 **Figure 1. Graphical representation of the European sample.** (A) PCA plots of the 4,887 samples
768 comprising the 3,104 samples from our in-house cohort HGID and 1000 genomes (1KG) individuals
769 including African (AFR), Ad Mixed American (AMR), East-Asian (EAS), European (EUR) and
770 South-Asian (SAS). Common variants were used to produce these plots. The European reference
771 individual is singled out. (B) 1,523 individuals with European ancestry selected. The dashed vertical
772 lines correspond to empirical separations between Northern (n=127 including 1KG FIN and HGID
773 samples), Middle (n=651 including 1KG CEU and GBR and HGID samples), and South European
774 ancestry (n=745 including 1KG TSI and IBS and HGID samples).

775

776 **Figure 2. Graphical representation of the Worldwide sample.** (A) PCA plots of the 4,887 samples
777 comprising the 3,104 samples from our in-house cohort HGID and 1000 genomes (1KG) individuals
778 including African (AFR), Ad Mixed American (AMR), East-Asian (EAS), European (EUR) and
779 South-Asian (SAS). Common variants were used to produce these plots. Reference individuals are
780 singled out. (B) The selected 1,967 individuals with European (n=700), Middle-Eastern (n=543),
781 North-African (n=359) and South-Asian (n=365) ancestries are colored. The remaining individuals
782 are left in grey.

783

784 **Figure 3. Histogram of powers for methods with a correct type I error rate for the large size**
785 **European sample (n=1,523) at the level $\alpha = 0.001$.** (A) Without stratification. (B) With moderate
786 stratification. (C) With high stratification. Relative risks considered vary from 2 to 4 on the x-axis.

787

788

789 **Figure 4. Histogram of powers for methods with a correct type I error rate for the large size**
790 **Worldwide sample (n=1,967) at the level $\alpha = 0.001$.** (A) Without stratification. (B) With moderate
791 stratification. (C) With high stratification. Relative risks considered vary from 2 to 4 on the x-axis.

792

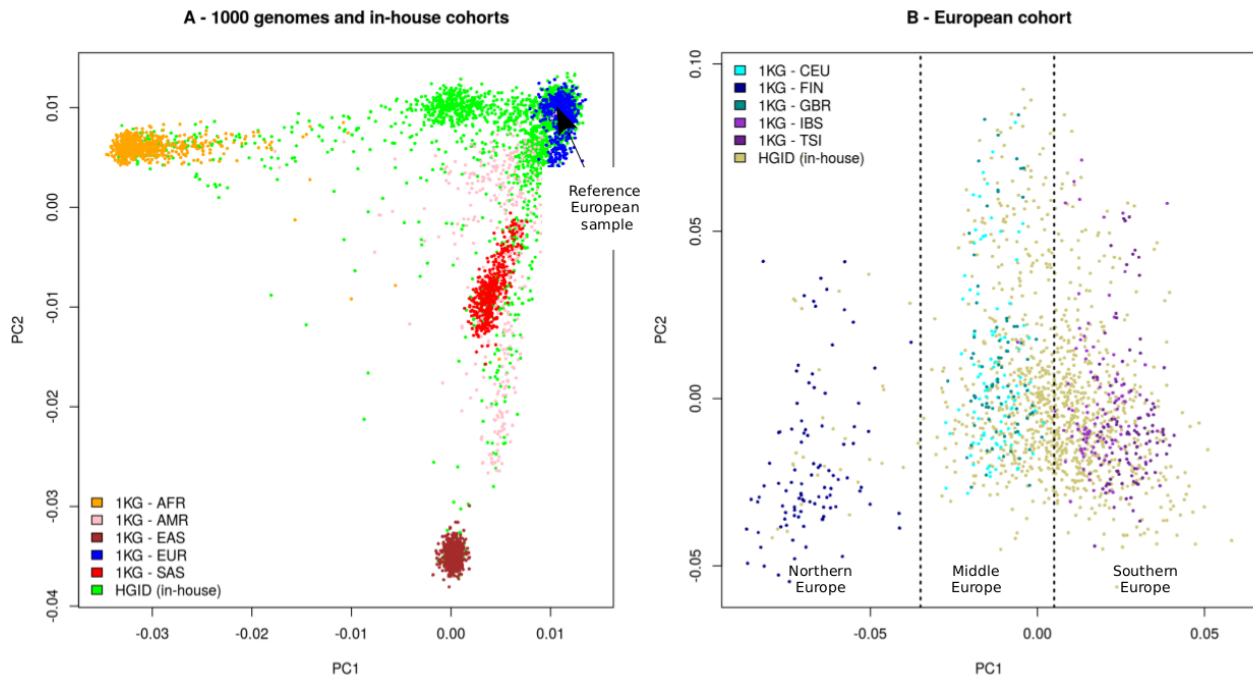
793

794 **Figure 5. Power for methods with a correct type I error rate under H_0 for the small size sample**
795 **at the level $\alpha = 0.001$.** (A) Scenarios with 50 cases from Southern-Europe. (B) Scenarios with 50
796 cases from the whole Europe. (C) Scenarios with 50 cases from the Worldwide sample. The relative
797 risk is fixed at 4.

799 **Figures**

800

801 **Figure 1**

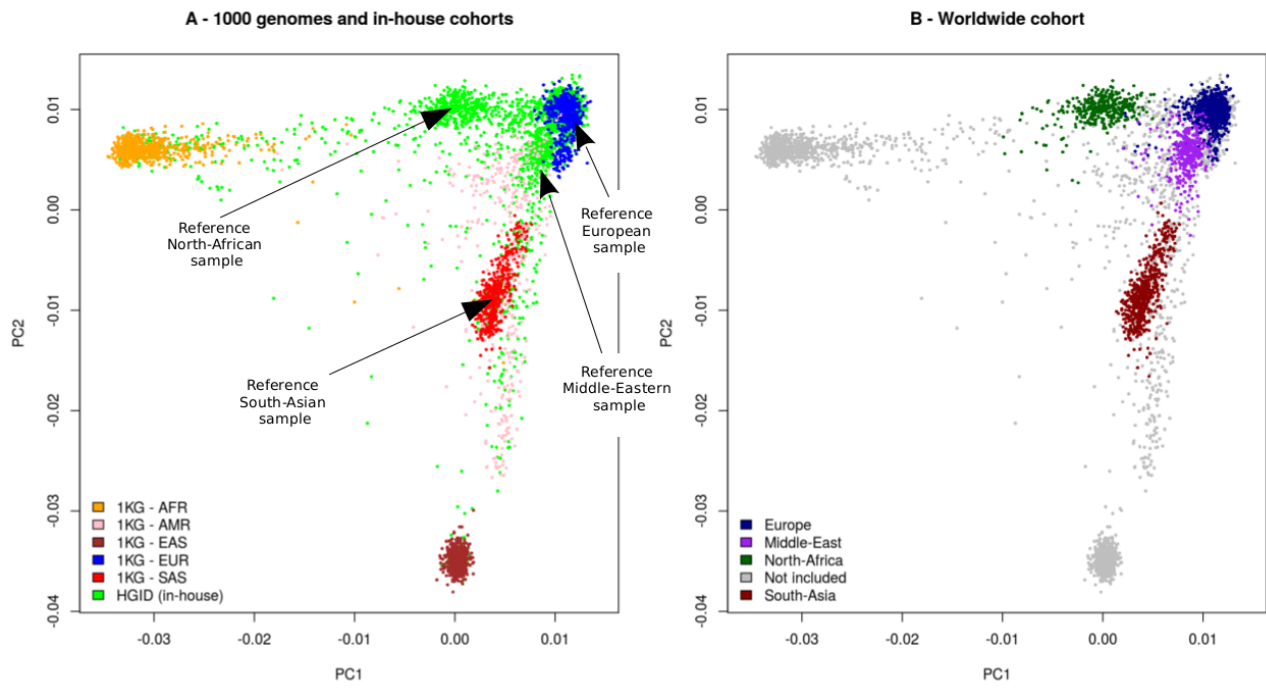


802

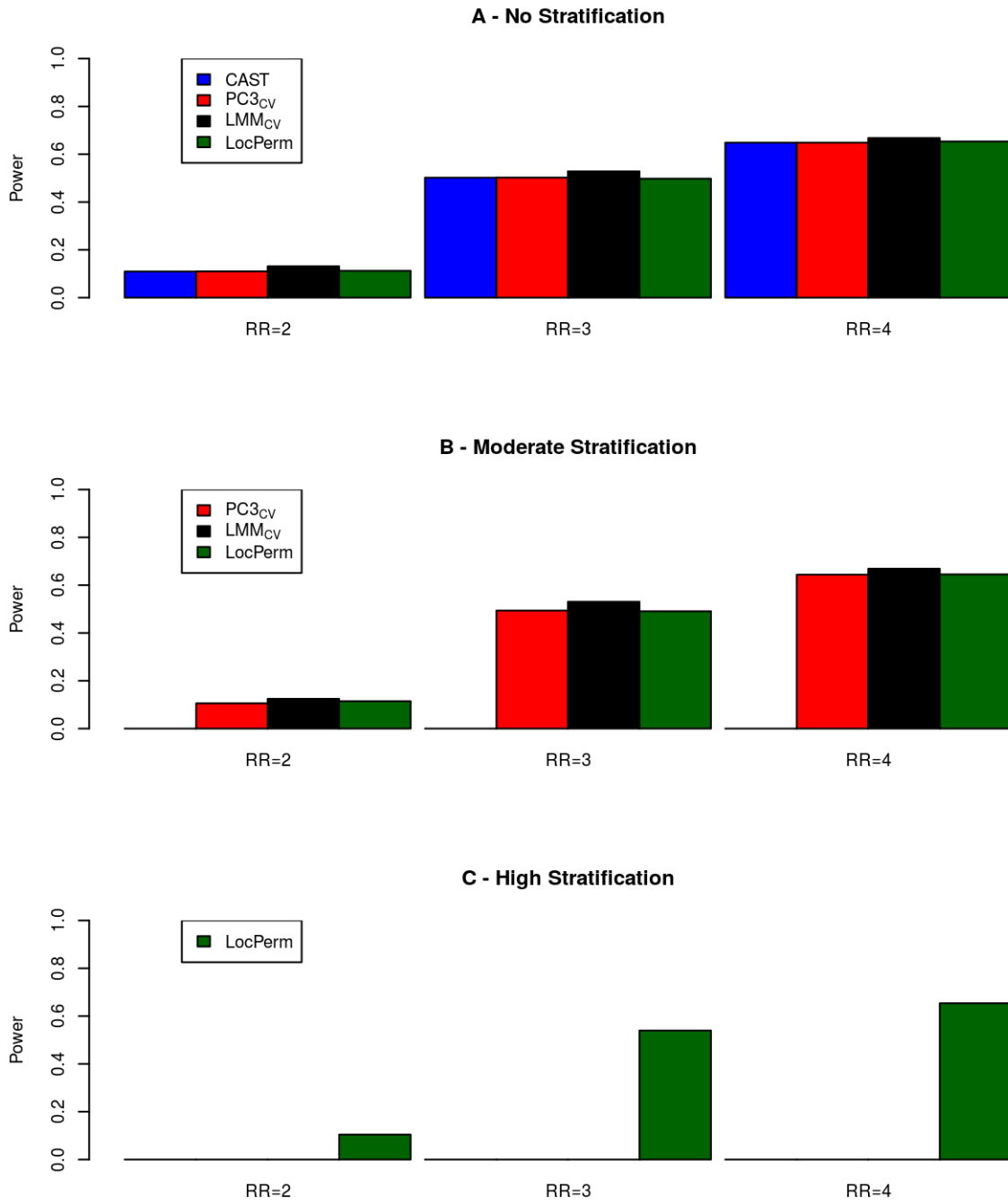
803

804 **Figure 2**

805



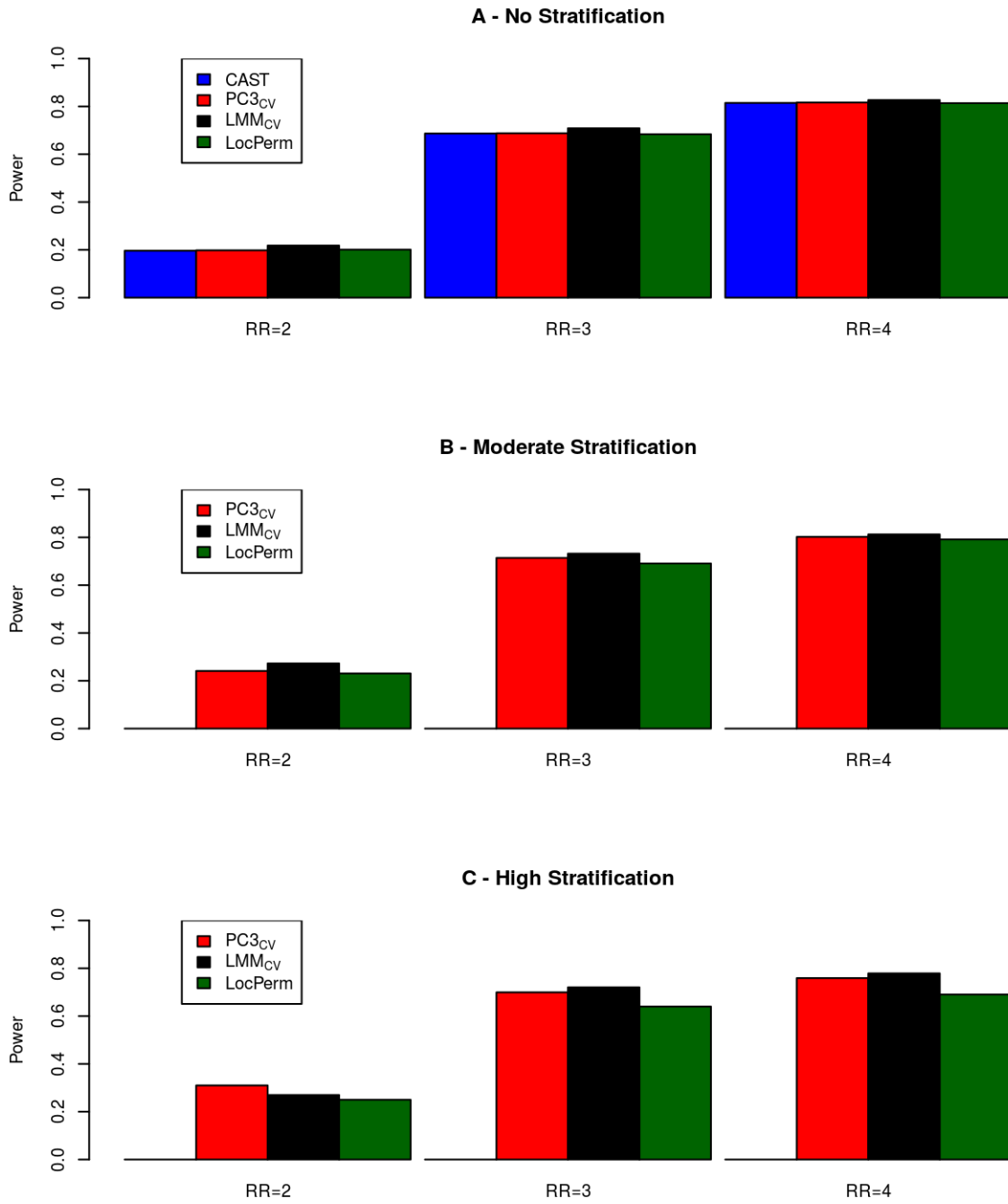
806 **Figure 3**



807

808

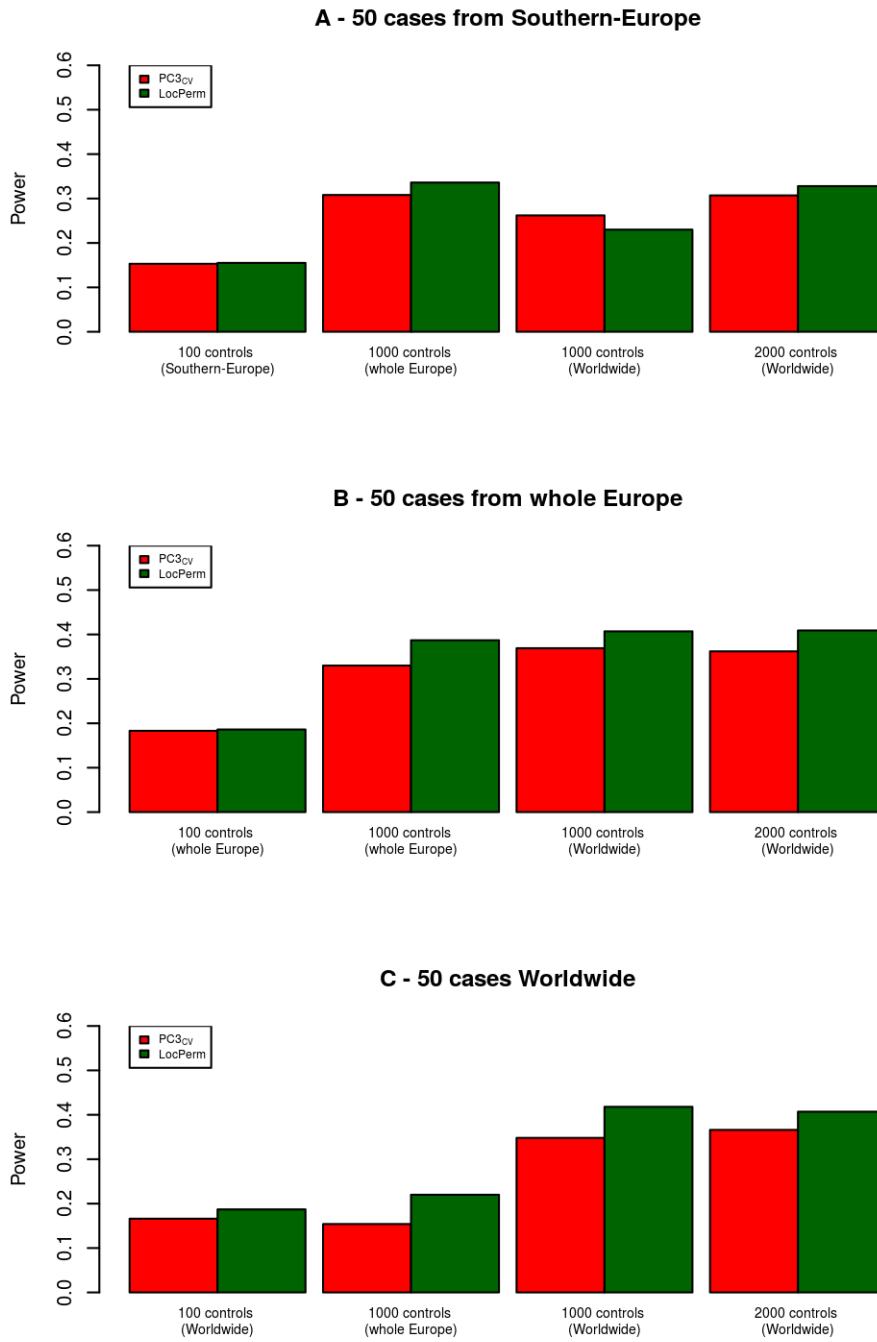
809 **Figure 4**



810

811

812 **Figure 5**



813

814