

1 **Title:** EigenGWAS: finding loci under selection through genome-wide association studies of
2 eigenvectors in structured populations
3

4 **Authors:** Guo-Bo Chen¹, Sang Hong Lee^{1,2}, Zhi-Xiang Zhu³, Beben Benyamin¹, Matthew R.
5 Robinson¹
6

7 **Affiliation:**

8 ¹Queensland Brain Institute, The University of Queensland, Brisbane, QLD 4072, Australia;

9 ²School of Environmental and Rural Science, The University of New England, Armidale,
10 NSW 2351, Australia; ³SPLUS Game, Guangzhou, Guangdong 510665, China
11

12 **Running Title:** EigenGWAS for selection
13

14 **Keywords:** GWAS, Eigenvector, selection, population structure
15

16 **Correspondence should be addressed to**

17 GBC (chen.guobo@foxmail.com)

18 Queensland Brain Institute

19 The University of Queensland

20 Brisbane, QLD 4072, Australia
21

22 MRR (m.robinson11@uq.edu.au)

23 Queensland Brain Institute

24 The University of Queensland

25 Brisbane, QLD 4072, Australia
26
27

28

Abstract

29 We apply the statistical framework for genome-wide association studies (GWAS) to
30 eigenvector decomposition (EigenGWAS), which is commonly used in population genetics to
31 characterise the structure of genetic data. The approach does not require discrete sub-
32 populations and thus it can be utilized in any genetic data where the underlying population
33 structure is unknown, or where the interest is assessing divergence along a gradient. Through
34 theory and simulation study we show that our approach can identify regions under selection
35 along gradients of ancestry. In real data, we confirm this by demonstrating *LCT* to be under
36 selection between HapMap CEU-TSI cohorts, and validated this selection signal across
37 European countries in the POPRES samples. *HERC2* was also found to be differentiated
38 between both the CEU-TSI cohort and within the POPRES sample, reflecting the likely
39 anthropological differences in skin and hair colour between northern and southern European
40 populations. Controlling for population stratification is of great importance in any
41 quantitative genetic study and our approach also provides a simple, fast, and accurate way of
42 predicting principal components in independent samples. With ever increasing sample sizes
43 across many fields, this approach is likely to be greatly utilized to gain individual-level
44 eigenvectors avoiding the computational challenges associated with conducting singular
45 value decomposition in large datasets. We have developed freely available software to
46 facilitate the application of the methods.

47

48

Introduction

49 In population genetics, eigenvectors have been routinely used to quantify genetic
50 differentiation across populations and to infer demographic history (Cavalli-Sforza *et al.*,
51 1996; Novembre *et al.*, 2008; Reich *et al.*, 2009). More recently, eigenvectors are commonly
52 used as covariates in genome-wide association studies (GWAS) to adjust for population
53 stratification (Price *et al.*, 2006). Eigenvectors are usually estimated for each individual
54 (individual-level eigenvectors, involving the inversion of a $N \times N$ matrix, where N is sample
55 size). Theoretical studies have suggested that individual-level primary eigenvectors are
56 measures of population differentiation reflecting F_{st} among subpopulations (Patterson *et al.*,
57 2006; McVean, 2009; Bryc *et al.*, 2013) and can be interpreted as the divergence of
58 individuals from their most recent common ancestor. Eigenvectors can also be estimated for
59 each SNP (SNP-level eigenvectors, which involve inversion of a $M \times M$ matrix, M is the
60 number of SNPs) and these SNP-level eigenvectors can be interpreted as F_{st} metrics of each
61 SNP (Weir, 1996). SNP-level eigenvectors from a reference population are useful for
62 revealing the population structure of independent samples (Zhu *et al.*, 2008) as they can be
63 used to project, or predict, the eigenvector values of individuals. However, due to high-
64 dimensional nature of GWAS data (commonly expressed as $M \gg N$), direct estimation of
65 SNP-level eigenvectors is nearly impossible when using millions of single nucleotide
66 polymorphisms (SNPs).

67

68 Singular value decomposition (SVD) enables SNP-level eigenvalues to be obtained in a
69 computationally efficient manner for any set of genotype data (Chen *et al.*, 2013), however, it
70 is not possible to determine the SNPs that contribute most to the leading eigenvector, or to
71 test whether specific SNPs are differentiated along the genetic gradient described by the
72 eigenvector. Here, we propose an alternative simple, fast approach for the estimation of SNP-
73 level eigenvectors. By using individual-level eigenvectors as phenotypes in a linear
74 regression, we demonstrate that the regression coefficients generated by single-SNP
75 regression are equivalent to SVD SNP effects as proposed by Chen *et al.* (Chen *et al.*, 2013).
76 As the single-SNP regression resembles the popular single-marker GWAS method, as
77 implemented in PLINK (Purcell *et al.*, 2007), we call this method EigenGWAS. We show
78 that the EigenGWAS framework represents an alternative way for identifying regions under
79 selection along gradients of ancestry.

80

81

82

Results

83 Properties of the estimating SNP effects for eigenvectors

84 We applied EigenGWAS to the HapMap cohort, a known structured population. Eigenvectors
85 were estimated via principal component analysis based on the \mathbf{A} matrix using all 919,133
86 SNPs. We conducted EigenGWAS for HapMap, using E_k , the k^{th} eigenvector, as the
87 phenotype and investigated the performance of EigenGWAS from E_1 to E_{10} . From E_1 to E_{10} ,
88 we found 546,716 significant signals (231,677 quasi-independent signals after clumping) on
89 E_1 and gradually reduced to 236 (163 after clumping) selection signals on E_{10} (**Fig. 1**). The
90 large number of genome-wide significant loci are likely because HapMap3 was comprised of
91 samples from different ethnicities, and these loci can be interpreted as ancestry informative
92 marker (AIM). For each E_k , its associated eigenvalue was highly correlated with the λ_{GC} , the
93 genomic inflation factor that is commonly used in adjusting population stratification for
94 GWAS (Devlin and Roeder, 1999), resulted from its EigenGWAS. The top five eigenvalues
95 associated to HapMap samples were 100.14, 47.66, 7.168, 5.92, and 4.40, and the
96 corresponding λ_{GC} of EigenGWAS were 103.72, 44.69, 6.47, 5.17, and 3.96, respectively
97 (**Table 1**). The large eigenvalues observed were consistent with previous theory that the
98 magnitude of eigenvalues indicating structured population (Patterson *et al.*, 2006). The
99 connection between λ_{GC} and eigenvalues, provides a straightforward interpretation: a large
100 λ_{GC} indicates underlying population structure (Devlin and Roeder, 1999). Therefore,
101 correction for λ_{GC} will filter out signals due to population stratification, allowing loci under
102 selection to be identified. These observations agreed well with our theory (see Methods &
103 Materials).

104

105 We demonstrate theoretically that for EigenGWAS, the estimated SNP effects using single-
106 marker GWAS are equivalent to the estimates from BLUP, and the correlation between the
107 estimates from these two methods was very high (greater than 0.98 on average) (**Fig. 2**), even
108 in HapMap samples that consist of a mix of ethnicities where the \mathbf{A} matrix is non-zero for
109 off-diagonal elements (**Supplementary Fig. 1**). This confirms that our EigenGWAS
110 approach provides an accurate representation of the SNP effects on eigenvalues.

111

112 We also conducted EigenGWAS on the POPRES samples, from which we selected 2,466
113 European samples. On E_1 , there were 10,885 (3,004 quasi-independent signals after

114 clumping) genome-wide significant signals, and reduced to 1,639 (90 after clumping) on E_{10}
115 (**Table 1**). As in the HapMap sample, we observed a concordance between eigenvalues and
116 λ_{GC} in POPRES. The top five eigenvalues were 5.104, 2.207, 2.157, 2.077, and 1.971, with
117 their associated EigenGWAS λ_{GC} were 5.005, 1.929, 1.910, 1.464, and 1.866, respectively
118 (**Table 1**), indicating population structure. The genetic relationship matrix (GRM) estimated
119 from the POPRES data resembled a diagonal matrix, which had off-diagonal elements close
120 to zero, suggesting that POPRES is a more homogenous samples as compared to HapMap
121 (**Supplementary Fig. 1**). Correlations between the estimates from EigenGWAS and BLUP
122 were high, with an average of greater than 0.999 from E_1 to E_{10} (**Supplementary Fig. 2**),
123 close to one as expected.

124

125 The chi-square statistics of the estimated SNP effects on eigenvectors from EigenGWAS
126 were correlated with F_{st} for each SNP, consistent with previous established relationship
127 between eigenvectors and F_{st} (Patterson *et al.*, 2006; McVean, 2009). Using naïve threshold
128 of $E_k > 0$, 2,466 POPRES samples were divided into nearly two even groups, which would
129 be served as two subgroups in calculating F_{st} . $E_1 > 0$ split the POPRES samples into North
130 and South Europe; samples from UK, Ireland, Germany, Austria, and Australia were in one
131 group, and samples from Italy, Spain, and Portugal were in the other group; samples from
132 Switzerland and France were nearly evenly split into two groups. F_{st} for each SNP was
133 consequently calculated based on these two groups. For every eigenvector until E_{10} , we
134 observed strong correlations between F_{st} and the chi-square test statistics for EigenGWAS
135 signals (**Fig. 3**), and the averaged correlation was 0.925 (S.D., 0.067). For example, the
136 correlation was 0.89 (p -value $<1e-16$) between chi-square test statistics and F_{st} for E_1 in
137 POPRES (**Supplementary Table 1**). This correlation is consistent with our theory, where F_{st}
138 has a strong linear relationship with its EigenGWAS chi-square test statistic.

139

140 We also validated our results in the simulation scheme I, in which there was neither selection
141 nor population stratification. Given 2,000 simulated samples, each of which had 500,000
142 unlinked SNPs, the EigenGWAS showed few GWAS signals (2 genome-wide significant
143 signals on E_1 , (**Supplementary Fig. 4**). After splitting the samples into 2 groups depending
144 on $E_i > 0$, the correlation between chi-square test statistics and F_{st} is about 0.67 from E_1 to
145 E_{10} (**Supplementary Fig. 5**). As expected, λ_{GC} ranged from around 1.124 to 1.130, with a
146 mean of 1.124 for EigenGWAS on the top 10 eigenvectors, indicating little population

147 stratification for the simulated data. Furthermore, we also validated the theory in the
148 simulation scheme II, in which there was population stratification. We wanted to know
149 whether the adjustment of the test statistic with the greatest eigenvalue could render the
150 distribution of the test statistics immune to population stratification. Given various sample
151 sizes for two subdivisions, after the adjustment for the test statistic with the largest
152 eigenvalue, the test statistic followed the null distribution, which was a chi-square
153 distribution of 1 degree of freedom (**Supplementary Fig. 6**), indicating a well control of
154 population stratification after correction. The statistical power of EigenGWAS was also
155 evaluated. As demonstrated, the power of EigenGWAS in detecting a locus under selection
156 was determined by the ratio between the specific F_{st} of a locus and the averaged population
157 stratification in the sample (**Supplementary Fig. 7**).

158

159 **Using EigenGWAS to identify loci under selection in structured populations**

160 We propose EigenGWAS as a method of finding loci differentiated among populations, or
161 across a gradient of ancestry. Intuitively, every EigenGWAS hit is an AIM, which differ in
162 allele frequency along an eigenvector due to genetic drift or selection. A locus under
163 selection should be more differentiated across populations than genetic drift can bring out. Thus,
164 correction for λ_{GC} , controls for background population structure, providing a test of whether
165 an AIM shows greater allelic differentiation than expected under the process of genetic drift.

166

167 We pooled together CEU (112 individuals) and TSI (88 individuals), which represent
168 Northwestern and Southern European populations in HapMap. EigenGWAS was conducted
169 on E_1 , which partitioned CEU and TSI into two groups accurately using $E_1 > 0$ as threshold
170 (**Supplementary Fig. 8**). We corrected for λ_{GC} , which was 1.723, for CEU&TSI. Adjustment
171 for λ_{GC} significantly reduced population stratification (**Supplementary Fig. 9**), and was
172 consequently possible to filter out the baseline difference between these two cohorts. After
173 correction, we found evidence of selection at the lactose persistence locus, *LCT* (p -
174 value=1.21e-20). Due to hitchhiking effect, the region near *LCT* also showed divergent allele
175 frequencies. For example, the *DARS* gene, 0.15M away from *LCT*, was also significantly
176 associated with E_1 (p -value=1.51e-23). *HERC2* was slightly below genome-wide
177 significance level (p -value=8.22e-08), indicating that anthropological difference reflected
178 geographic locations of two cohorts but not under selection as strong as *LCT*.

179

180 We then conducted EigenGWAS in the POPRES sample by treating E_1 as a quantitative trait,
181 and calculated the approximate F_{st} for each SNP given two groups split by the threshold of
182 $E_1 > 0$ (**Supplementary Fig. 10**). Given 643,995 SNPs, the genome-wide threshold was p -
183 value $< 7.76e-08$ for the significance level of $\alpha = 0.05$. $\lambda_{GC} = 5.00$, which indicated
184 substantial population stratification as expected for POPRES. Correcting for λ_{GC}
185 systematically reduced the EigenGWAS χ^2 test statistics (**Supplementary Fig. 11**), and we
186 replicated the significance of *LCT* (p -value=1.23e-22) and *DARS* (p -value=8.99e-22) (**Table**
187 **2**), suggesting selection at these regions. *HERC2* was also replicated with p -value 8.15e-09,
188 and with F_{st} of 0.041.

189

190 **Prediction accuracy for projected eigenvector**

191 We investigated three aspects of EigenGWAS prediction: 1) the number of loci needed to
192 achieve high accuracy for the projected eigenvectors; 2) the required sample size of the
193 training set; 3) the importance of matching the population structure between the training and
194 the test sets.

195

196 Using the POPRES samples, we split 5% (125 individuals), 10% (250 individuals), 20% (500
197 individuals), 30% (750 individuals), 40% (1000 individuals), and 50% (1250 individuals) of
198 the sample as the training set, and used the remainder of the samples as the test set.
199 Eigenvectors were estimated using all markers in each training set. As predicted by our
200 theory (Eq 7), the prediction accuracy of the projected eigenvector was consistent with
201 $R^2 = \frac{1}{1 + \frac{N_e}{M}}$ in which $N_e = 1,000$ for E_1 empirically. If only 100 and 1,000 random SNPs were
202 sampled as predictors, the expected maximal $R^2 = 0.091$ and 0.5, respectively and accuracy
203 reached almost 1 if more than 100,000 SNPs were sampled. In agreement with our theory
204 (**Fig. 6**), if the number of predictors were too small the prediction accuracy was poor, with
205 prediction accuracy increasing with the addition of more markers for E_1 . When the sample
206 size of the discovery was 1,000 or above, maximal prediction accuracy was achieved, as
207 predicted in our theory. Therefore, a discovery with a sample size greater than 1,000 should
208 be sufficient to predict the first eigenvector of an independent set, provided that population
209 structure is the same across the discovery and prediction samples (**Fig. 6**). In contrast, the
210 prediction accuracy for prediction eigenvectors decreased (**Fig. 6**) quickly for eigenvectors
211 other than E_1 . For example, the prediction accuracy for E_2 was below $R^2 < 0.2$ and $R^2 <$
212 0.15 for E_3 . For $E_4 \sim E_{10}$, the prediction accuracy dropped down to nearly zero. This is

213 consistent with the top 2~3 eigenvectors explaining the majority of variation (McVean,
214 2009), if the training and the test sets had their population structure matched.

215

216 If EigenGWAS SNPs of low p -value were likely to be AIMs, we would hypothesise that
217 AIM markers would be more efficient in giving high accuracy for the predicted eigenvectors
218 (**Fig. 6**). For E_1 , the prediction accuracy reached 1 more quickly by using markers selected by
219 p -value thresholds. The prediction accuracy for projected E_2 was dependent upon the
220 threshold. For projected E_2 given a 50:50 split of POPRES sample, applying the threshold of
221 p -value $< 1e-6$ (927 SNPs), $R^2 = 0.136$, as high as using all markers. For other projected
222 eigenvectors, the pattern of accuracy did not change much after applying p -value thresholds
223 because in general, the prediction accuracy was low. This indicated that eigenvectors other
224 than the first two eigenvectors capture little replicable population structure in POPRES.

225

226 In practice, the training and the test set may not match perfectly on population structure, and
227 this will likely lead to a reduction in prediction accuracy. To demonstrate this, we split the
228 POPRES samples into two sets: pooling Swiss (991 samples) and French (96 samples)
229 samples into one group (SF), and the rest of the samples into the other group (NSF). We used
230 SF as the training and the NSF as the testing. As SF was almost an average of North
231 European and South European gene flow, making a less stratified population, its
232 EigenGWAS effects would be consequently small and less “heritable”. When using all SNPs
233 effects estimated from SF set, the observed prediction accuracy for NSF set was $R^2 = 0.33$
234 and 0.005 for E_1 and E_2 , respectively. These results indicate that a matched training and test
235 set is important for prediction accuracy of the projected eigenvectors.

236

237 Ancestry information may still be elucidated well even if the training set and the test set do
238 not match well in their population structure. Using HapMap3 as the training set, we also tried
239 to infer the ancestry of the Puerto Rican cohort (PUR, 105 individuals) and Pakistani cohorts
240 (PIL, 95 individuals) from 1000 Genomes project (The 1000 Genomes Project Consortium,
241 2012). In chromosome 1, 74,500 common SNs were found between HapMap3 and 1000
242 Genomes project. As illustrated, using only 74,500 common markers between HapMap3 and
243 1000 Genome projects SNPs on chromosome 1, it projected Eigenvectors accurately revealed
244 the demographic history of Puerto Rican cohort, an admixture of African and European gene
245 flows, and Pakistan cohort, an admixture of Asian and European gene flows (**Fig. 7**).

246

247 As a negative control, we replicated the prediction study for simulated data used in the
248 previous section. The simulated data was split to two equal sample size. As there was no
249 population structure in the simulated data, the prediction accuracy was poor, $R^2 < 0.01$ from
250 E_1 to E_{10} . This demonstrates that prediction can be used to validate whether population
251 structure exists within a genotype sample.

252

253 We concluded that to achieve high prediction accuracy of projected eigenvectors for
254 independent samples, there are several conditions to be met: 1) the training set should
255 harbour sufficient population stratification; 2) the sample size of the training should be
256 sufficiently large; 3) the test sets should be as concordant as possible in its population
257 structure; 4) when there is no real population structure, the prediction accuracy is very low
258 close to zero; 5) depending on the population, high prediction was largely achievable for the
259 projected E_1 .

260

261

262

Discussion

263 Eigenvectors have been routinely employed in population genetics, and various approaches
264 have been proposed to offer interpretation and efficient algorithms (Patterson *et al.*, 2006;
265 Rokhlin *et al.*, 2009; McVean, 2009; Chen *et al.*, 2013; Galinsky *et al.*, 2015). In this study,
266 we created a GWAS framework for studying and validating population structure, and offer an
267 interpretation of eigenvectors within this framework. The EigenGWAS framework (least
268 square) identifies ancestry informative markers and loci under selection across gradients of
269 ancestry.

270

271 We integrated SVD, BLUP, and single-marker regression into a unified framework for the
272 estimation of SNP-level eigenvectors. SVD is a special case of BLUP when heritability is of
273 1 for the trait and the target phenotype is an eigenvector. Furthermore, the BLUP is
274 equivalent to the commonly used GWAS method for estimating SNP effects. As
275 demonstrated, the correlation between BLUP and GWAS is almost 1 for the estimated SNP
276 effects. EigenGWAS offers an alternative way in estimating F_{st} that can replace conventional
277 F_{st} when population labels are unknown, populations are admixed, or differentiation occurs

278 across a gradient. As demonstrated for CEU&TSI samples, EigenGWAS brings out nearly
279 identical estimation of F_{st} compared with conventional estimation.

280

281 Different from conventional GWAS, which requires conventional phenotypes, the proposed
282 EigenGWAS provides a novel method for finding loci under selection based on eigenvectors,
283 which are generated from the genotype data itself. An EigenGWAS hit may reflect the
284 consequence of process and thus additional evidence is needed to differentiate selection from
285 drift. *LCT* is a known locus under selection, which differs in its allele frequency as indicated
286 by F_{st} statistic between Northern and Southern Europeans (Bersaglieri *et al.*, 2004). We
287 replicated the significance of *LCT* in CEU&TSI samples and POPRES European samples.
288 *DARS* has been found in association with hypomyelination with brainstem and spinal cord
289 involvement and leg spasticity (Taft *et al.*, 2013). In addition, we also found *HERC2* locus
290 independently, which may indicate the existence of anthropological difference in certain
291 characters, such as hair, skin, or eyes color across European nations (Voight *et al.*, 2006;
292 Visser *et al.*, 2012).

293

294 Although by definition selection and genetic drift are different biological processes, both lead
295 to allele frequency differentiation across populations and often difficult to tear them apart. In
296 this study, with and without adjustment for λ_{GC} from EigenGWAS offers a straightforward
297 way to filter out population stratification. For example, with adjustment for λ_{GC} , *LCT* and
298 *DARS* were still significant in both EigenGWAS, while *HERC2* was only significant in
299 POPRES. If adjustment for λ_{GC} removed the average genetic drift since the most recent
300 common ancestor for the whole sample, it might indicate that *HERC2* reflected the
301 anthropological difference between subsamples but not under selection as strong as that for
302 *LCT*. Nevertheless, *LCT* was differentiated due to selection that was on top of genetic drift,
303 and for *DARS*, it might be significant due to hitchhiking effect. So, *LCT*, *DARS*, and *HERC2*
304 were significant in EigenGWAS for different mechanisms.

305

306 In EigenGWAS application, it provides a clear scenario that λ_{GC} is necessary if genetic
307 drift/population stratification should be filtered out. It has been debated whether correction
308 for λ_{GC} is necessary for GWAS (Devlin and Risch, 1995; Yang, Weedon, *et al.*, 2011). If the
309 inflation is due to population stratification, as initially λ_{GC} introduced, it seems necessary to
310 control for it. In contrast, if it is due to polygenic genetic architecture, then correction for λ_{GC}

311 will be a overkilling for GWAS signals. Interestingly, Patterson et al (Patterson *et al.*, 2006)
312 found that the top eigenvalues reflect population stratification, and in our study we found λ_{GC}
313 from EigenGWAS was numerically so similar to its corresponding eigenvalues. It in another
314 aspect indicates λ_{GC} captures population stratification. So, in concept and implementation, the
315 correction for λ_{GC} is technically reasonable. Of note, Galinsky et al also proposed a similar
316 procedure to filter out population stratification in a study similar to ours (Galinsky *et al.*,
317 2015), but we believe our framework is much easier to understand and implement in practice.
318

319 Once we have EigenGWAS SNP effects estimated, it is straightforward to project those
320 effects onto an independent sample. The prediction of population structure was to that of
321 recent studies (Chen *et al.*, 2013). We found that the prediction accuracy for the top
322 eigenvector could be as high as almost 1. Given a training set of about 1,000 samples, the
323 prediction accuracy could be very high if there were a reasonable number of common
324 markers in the order of 100,000. This number, which needs to be available in both reference
325 set and the target set, is achievable. Further investigation may be needed to check whether
326 this number of markers is related to effective number or markers after correction for linkage
327 disequilibrium for GWAS data. When the population structure of the test sample resembles
328 the training sample, high accuracy will be achieved for the leading projected eigenvectors.
329 Therefore, this approach is likely to be extremely beneficial for extremely large samples,
330 such as UK Biobank samples and 23andMe, both of which have more than half million
331 samples where direct eigenvector analysis may be infeasible. Our results suggest that
332 sampling about 1,000 individuals from the whole sample as the training set and subsequently
333 project EigenGWAS SNP effects to the reminding samples will be sufficient to reach a
334 reasonable high resolution of the population structure.

335

336 Many improvements to the inference of ancestry using projected eigenvectors have been
337 suggested (Chen *et al.*, 2013). As the concordance of population structure between the
338 training and test sets is often unknown (population structure, upon from genetic or social-
339 cultural perspectives, its definition can be difficult or controversial), improvement of the
340 inference of ancestry may or may not be achieved dependent upon the scale of the precision
341 required for a sample. However, for classification of samples at ethnicity level, projected
342 eigenvectors are likely to have high accuracy, as demonstrated in the Puerto Rican cohort and

343 the Pakistani cohort. Therefore, when identifying ethnic outliers, using projected eigenvectors
344 from HapMap is likely to be sufficient in practice.

345

346 Eigenvector analysis of GWAS data is an important well utilized data technique, and here we
347 show that its interpretation depends on many factors, such as proportion of different
348 subpopulations, and F_{st} between subpopulations. Our EigenGWAS approach provides
349 intuitive interpretation of population structure, enabling ancestry informative markers (AIM)
350 to be identified, and potentially loci under selection to be identified. To facilitate the use of
351 projected eigenvectors, we provide estimated SNP effects from HapMap samples and
352 POPRES and software that can largely reduce the logistics involved in conventional way in
353 generating eigenvectors, such as reference allele match, and strand flips.

354

355 **Methods and Materials**

356 **HapMap3 samples.** HapMap3 samples were collected globally to represent genetic diversity
357 of human population (Altshuler *et al.*, 2010). HapMap3 contains representative samples from
358 many continents: CEU and TSI represent population from north and south Europe, CHB and
359 JPT from East Asia, and CHD Chinese collected in Denver, Colorado. Loci with palindrome
360 alleles (A/T alleles, or G/C alleles) were excluded, and 919,133 HapMap3 SNPs were used
361 for the analysis.

362

363 **1000 Genomes project.** 1000 Genomes project samples were used as a prediction set for
364 projecting eigenvectors (The 1000 Genomes Project Consortium, 2012). We selected the
365 Puerto Rico cohort (PUR, 105 samples) and the Pakistan cohort (Punjabi from Lahore,
366 Pakistan, 95 samples) for analysis.

367

368 **POPRES samples.** POPRES (Nelson *et al.*, 2008) is a reference population for over 6,000
369 samples from Asian, African, and European nations. In this study, we selected 2,466
370 European descendants. The POPRES genotype sample was imputed to a 1000 Genomes
371 reference panel (The 1000 Genomes Project Consortium, 2012). Imputation for the POPRES
372 was performed in two stages. First, the target data was haplotyped using HAPI-UR (Williams
373 *et al.*, 2012). Second, Impute2 was used to impute the haplotypes to the 1000 genomes
374 reference panel (Howie *et al.*, 2011). We then selected SNPs which were present across all
375 datasets at an imputation information score of >0.8 . A full imputation procedure is described

376 at <https://github.com/CNSGenomics/impute-pipe>. After quality control and removing loci
377 with palindromic alleles (A/T alleles, or G/C alleles) 643,995 SNPs for POPRES remained.
378 In addition, we also conducted the analysis using non-imputed 234,127 common markers
379 between POPRES and HapMap3. As the results were between these two datasets were very
380 similar, this report focused on the results from 643,995 SNPs, which were more informative.

381

382 **Simulation scheme I: null model without population structure.** 2,000 unrelated samples
383 with 500,000 biallelic markers, which were in linkage equilibrium to each other, were
384 simulated. The minor allele frequencies ranged from 0.01~0.5, and Hardy-Weinberg
385 equilibrium was assumed for each locus. All individuals were simulated from a homogeneous
386 population, with no population stratification. In order to calculate F_{st} at each locus, we
387 divided the sample into sub-populations based upon eigenvectors that were estimated from a
388 genetic relationship matrix calculated using all 500,000 markers (see below).

389

390 **Simulation scheme II: null model with population structure.** In general, this simulation
391 scheme was followed Price et al (Price *et al.*, 2006). 2,000 unrelated samples with 10,000
392 biallelic markers, which were in linkage equilibrium to each other, were generated. For each
393 marker, its ancestral allele frequency was sampled from a uniform distribution between 0.05
394 to 0.95, and its frequency in a subpopulation was sampled from Beta distribution with
395 parameters $p \frac{1-F_{st}}{F_{st}}$ and $(1-p) \frac{1-F_{st}}{F_{st}}$. The Beta distribution had mean of p and sampling
396 variance of $p(1-p)F_{st}$. Once the allele frequency for a subpopulation over a locus was
397 determined as p_s , individuals were generated from a binomial distribution $Binomial(2, p_s)$. It
398 agreed with the quantity that measures the genetic distance between a pair of subpopulations
399 (Cavalli-Sforza *et al.*, 1996).

400

401 **Calculating individual-level eigenvectors**

402 We assume that there is a reference sample consisting of N unrelated individuals and M
403 markers. $X_i = (x_{i1}, x_{i2}, \dots, x_{iM})^T$, is a vector of the i^{th} individual's genotypes along M loci,
404 with x the number of the reference alleles. An $N \times N$ genetic relatedness (correlation) matrix
405 **A** (matrix in bold font) for each pair of individuals is defined as $A_{ij} = \frac{1}{M} \sum_{l=1}^M \frac{(x_{il}-2f_l)(x_{jl}-2f_l)}{2f_l(1-f_l)}$,
406 in which f_l is the frequency of the reference allele. The principal component analysis (PCA)

407 is then implemented on the A matrix (Price *et al.*, 2006), generating \mathbb{E} , which is an $N \times K$
408 ($K \leq N$) matrix, in which E_k is the eigenvector corresponding to the k^{th} largest eigenvector.

409

410 **Unified framework for BLUP, SVD, and EigenGWAS**

411 Theoretically, PCA can also be implemented on a $M \times M$ matrix, but this is often infeasible
412 because the $M \times M$ matrix is very large. However, for individual i , eigenvector k can also be
413 written as:

$$414 E_{k,i} = \beta_k X_i^T \quad \text{(Equation 1)}$$

415 in which β_k is a $M \times 1$ SNP-level vector of the SNP effects on E_k , and x_i is the genotype of
416 the i^{th} individual across M loci. In the text below, we denote individual-level eigenvector as
417 eigenvector ($N \times 1$ vector), and SNP-level eigenvector ($M \times 1$) as SNP effects.

418

419 We review three possible methods to estimate β given eigenvectors. The first method is best
420 linear prediction (BLUP), which is commonly used in animal breeding and recently has been
421 introduced to human genetics for prediction (Henderson, 1975; Goddard *et al.*, 2009). The
422 second method is to convert an individual-level eigenvector to SNP-level eigenvector using
423 SVD, as proposed by Chen *et al.* (Chen *et al.*, 2013). The third method is the approach
424 outlined here, EigenGWAS, which is a single-marker regression, as commonly used in
425 GWAS analysis.

426

427 **Method 1 and 2: BLUP and SVD**

428 For a quantitative trait, $y = \mu + \beta X + e$, in which y is the phenotype, μ is the grand mean, β
429 is the vector for additive effects, X is the genotype matrix, and e is the residual. Without loss
430 of generality, the BLUP equation can be expressed as:

$$431 \hat{\beta} = \tilde{X}^T V^{-1} y \quad \text{(Equation 2)}$$

432 in which $\hat{\beta}$ is the estimates of the SNP effects, \tilde{X} is the standardized genotype matrix, V is the
433 variance covariance with $V = \sigma_A^2 A + (\sigma_y^2 - \sigma_A^2) I$, and y is the trait of interest (Henderson,
434 1975). Replacing y with individual-level eigenvector (E_k), Eq 2 can be written as

$$435 \hat{\beta}_k = \tilde{X}^T A^{-1} E_k \quad \text{(Equation 3)}$$

436 in which β_k is the BLUP estimate of the SNP effects, E_k is the k^{th} eigenvector estimated
437 from the reference sample., The V matrix can be replaced with A because the eigenvector has
438 no residual error (i.e. $h^2=1$). This method has also been proposed as an equivalent computing
439 algorithm for genomic predictions (Maier *et al.*, 2015).

440

441 In addition, the connection between PCA and SVD can be established through the
442 transformation between the $N \times N$ matrix to the $M \times M$ matrix (McVean, 2009). Let
443 $\mathbf{A} = \mathbf{PDP}^{-1}$, in which \mathbf{D} is a $N \times N$ diagonal matrix with λ_k , \mathbf{P} is $N \times N$ matrix with the
444 eigenvectors. $\mathbf{B} = \mathbf{X}^T(\mathbf{PDP}^{-1})^{-1}\mathbf{P} = \mathbf{X}^T\mathbf{PD}^{-1}$, in which \mathbf{B} is $M \times N$ matrix. This is
445 equivalent to the equation used in Chen et al (Chen *et al.*, 2013) where $\mathbf{B}^T = \mathbf{D}^{-1}(\mathbf{X}^T\mathbf{P})^T$.
446 Thus, eigenvector transformation can be viewed as a special case of BLUP in which the
447 heritability is 1 (Eq 3). However, under SVD another analysis step is then required to
448 evaluate the significance of the estimated SNP effect. In an EigenGWAS framework an
449 empirical *p-value* is produced when estimating the regression coefficient.

450

451 **Method 3: estimating SNP effects on eigenvectors with EigenGWAS**

452 Given the realized genetic relationship matrix \mathbf{A} , for unrelated homogeneous (i.i.d.) samples,
453 $E(A_{ij}) = 0$ ($i \neq j$), and consequently $E(\mathbf{A}) = \mathbf{I}$, an **identity matrix**. Due to sampling
454 variance of the genetic relationship matrix \mathbf{A} , the off diagonal is a number slightly different
455 from zero even for unrelated samples (Chen, 2014). If we replace the matrix with its
456 mathematical expectation – the identity matrix, Equation 3 can be further reduced to $\beta_k =$
457 $\tilde{\mathbf{X}}^T E_k$, which is equivalent to single-marker regression $E_k = a + bx + e$, as implemented in
458 PLINK (Purcell *et al.*, 2007). Furthermore, standardization for \mathbf{X} is not required because it
459 will not affect *p-value*. Thus, SNP effects can be estimated using the single-marker
460 regression, which is computationally much easier in practice and is implemented in many
461 software packages. Each SNP effect, $\hat{\beta}_{k,m}$, is estimated independently, and the *p-value* of
462 each marker can be estimated, which requires additional steps in BLUP and SVD.

463

464 We summarise the properties and their transformation of SVD, BLUP, and EigenGWAS as
465 below:

- 466 1) E_k is determined by the \mathbf{A} matrix, or in another words, it is determined by the
467 genotypes completely. If we consider each E_k is the trait of interest – a quantitative
468 trait, its heritability is 1.
- 469 2) $h^2 = 1$. SVD and BLUP are both computational tool in converting a vector from
470 $N \times N$ matrix to a $M \times M$ matrix. SVD is a special case to BLUP when $h^2 = 1$ for
471 BLUP.

472 3) $h^2 = 1$ and $E(\mathbf{A}) = \mathbf{I}$. When these two conditions are set, BLUP is further reduced to
 473 single-marker association studies, which is EigenGWAS as suggested in this study.

474

475 Recently, in an independent work Galinsky et al (Galinsky *et al.*, 2015) introduced an
 476 approximation to find the proper scaling for SNP effects (“SNP weight” in Galinsky’s
 477 terminology) estimated from SVD, in order to produce accurate p -values. In our EigenGWAS
 478 framework, p -values for individual-level SNP eigenvector are automatically generated. In
 479 practice, it is conceptually easier to conduct EigenGWAS on eigenvectors than to conduct
 480 BLUP/SVD. Also, if computational speed is of concern, EigenGWAS can be easily
 481 parallelized for each chromosome, each region, or even each locus.

482

483 **Interpretation for EigenGWAS**

484 We can write a linear regression model $E_k = a + \beta x + e$, in which both E_k and x is
 485 standardized. Assuming that a sample has two subdivisions, which have sample size n_1 and

486 n_2 , $\beta = \frac{2\sqrt{w(1-w)}(p_1-p_2)}{\sqrt{2\bar{p}\bar{q}}}$, and the sampling variance for β is $\sigma_\beta^2 = \frac{\sigma_e^2}{n\sigma_x^2} = \frac{1}{n}$. A chi-square test

487 for β is

$$488 E(\chi_1^2) = 4nw(1-w)F_{st}^N \text{ (Equation 4)}$$

489 in which $F_{st}^N = \frac{(p_1-p_2)^2}{2\bar{p}\bar{q}}$ is Nei’s estimator of genetic difference for a biallelic locus (Nei,
 490 1973).

491

492 In principal component analysis, the proportion of the variance explained by the largest

493 eigenvalue is equal to F_{st}^W (McVean, 2009), in which $F_{st}^W = \frac{2\sum_{i=1}^2 w_i [(p_i - \bar{p})^2]}{\bar{p}\bar{q}}$ for a pair of

494 subpopulations as defined in Weir (Weir, 1996). So $\lambda_1 \approx \bar{F}_{st}^W \times n$, in which \bar{F}_{st}^W characterizes

495 the average divergence for a pair of subpopulations. When the test statistic, Eq 4, is adjusted

496 by the largest eigenvalue λ_1 , an equivalent technique in GWAS for the correction of

497 population stratification, $E(\chi_{1,\lambda_1}^2) = \frac{4nw(1-w)F_{st}^N}{\lambda_1} = 4w(1-w)\frac{F_{st}^N}{\bar{F}_{st}^W}$. For a population with a

498 pair of subdivisions $4w(1-w)F_{st}^N = F_{st}^W$. So

$$499 E(\chi_{1,\lambda_1}^2) = \frac{F_{st}^W}{\bar{F}_{st}^W} \text{ (Equation 5)}$$

500 after the adjustment of the largest eigenvalue, the test statistic immune of population

501 stratification, at least for a divergent sample.

502

503 For a locus under selection, which should have a greater F_{st} than \bar{F}_{st} the background
 504 divergence. So the statistical power for detecting whether a locus is under selection is
 505 determined by the strength of selection, which can be defined as the ratio between F_{st}
 506 of a particular locus and \bar{F}_{st} the average divergent in the sample. It is analogous to
 507 consider a chi-square test with non-centrality parameter (NCP), $NCP = \frac{F_{st}^W}{\bar{F}_{st}^W} - 1$.

508 Otherwise specified, in this study F_{st} is referred to the one defined in Weir (Weir, 1996).
 509

510 **Validation and prediction for population structure**

511 Once β_k is estimated, it is straightforward to get genealogical profile for an independent
 512 target sample. In general, it is equivalent to genomic prediction, and the theory for prediction
 513 can be applied (Daetwyler *et al.*, 2008; Dudbridge, 2013). The predicted genealogical score
 514 can be generated as

$$515 \quad \tilde{E}_k = \hat{\beta}_{k.m} \mathbf{X} \quad \text{(Equation 6)}$$

516 in which E_k is the predicted k^{th} eigenvector, $\hat{\beta}_{k.m}$ is the estimated SNP effects, and \mathbf{X} is the
 517 genotype for the target sample. We focus on the correlation between the predicted
 518 eigenvectors and the direct eigenvectors, and thus it does not matter whether \mathbf{X} or $\tilde{\mathbf{X}}$ is used.
 519

520 In contrast to conventional prediction studies, which focus on a metric phenotype of interest,
 521 prediction of population structure is focussed on a “latent” variable. This latent variable is the
 522 genetic structure of population, which is shaped by allele frequency and linkage
 523 disequilibrium of markers. Thus, expectations of prediction accuracy differ from what has
 524 been established for conventional prediction (Daetwyler *et al.*, 2008; Dudbridge, 2013)

525 $R^2 = h^2 \left(\frac{h^2}{h^2 + \frac{M}{N_e}} \right) < h^2 \ll 1$. We therefore assess prediction of accuracy for E_1 across
 526 markers, when using different prediction thresholding (Purcell *et al.*, 2009).
 527

528 Here we proposed an equation for prediction accuracy, especially for E_1

$$529 \quad R^2 = \frac{\left(h^2 + \frac{M}{N_e} \right)^2}{h^2 \left(1 + 2 \frac{M}{N_e} \right) + \frac{M}{N_e} \left(1 + \frac{M}{N_e} \right)} \approx \frac{1}{1 + N_e/M} \quad \text{(Equation 7)}$$

530 when there is no heritability, the predictor can be simplified to $R^2 = \frac{1}{1 + \frac{N_e}{M}}$, meaning that as
 531 the number of markers increases prediction accuracy should rapidly reach 1. Here the h^2 is
 532 interpreted as the genetic difference in the source population, or real ancestry informative

533 markers. For a homogeneous population, the genetic difference is large due to genetic drift,
534 and $h^2 \approx 0$.

535

536 For this study, the genetic relationship matrix (**A** matrix), principal component analysis, and
537 BLUP estimation were conducted using GCTA software (Yang, Lee, *et al.*, 2011). Single-
538 marker GWAS was conducted using PLINK (Purcell *et al.*, 2007), or GEAR
539 (<https://github.com/gc5k/GEAR/wiki/EigenGWAS>;
540 <https://github.com/gc5k/GEAR/wiki/ProPC>).

541

542 **Web resource and data availability**

543 GEAR is available at <http://cnsgenomics.com/>

544 GCTA is available at <http://cnsgenomics.com/>

545 [PLINK is available at http://pngu.mgh.harvard.edu/~purcell/plink/index.shtml](http://pngu.mgh.harvard.edu/~purcell/plink/index.shtml)

546 [1000 Genomes Project: http://www.1000genomes.org/](http://www.1000genomes.org/)

547

548 **Acknowledgements**

549 This research was funded by ARC (DE130100614 to SHL), NHMRC (APP1080157 to SHL,
550 APP1084417 and APP1079583 to BB, and APP1050218 to MRR), and GBC was supported
551 by IAP P7/43-BeMGI from the Belgian Science Policy Office Interuniversity Attraction
552 Poles (BELSPO-IAP) program. We thank Peter M. Visscher for discussion, helpful
553 comments, and for proposing the name EigenGWAS. Robert Maier assisted with ggplot, and
554 Alex Holloway helped with Github. We also thank to the Information Technology group, the
555 Queensland Brain Institute. The POPRES dataset were obtained from dbGaP at
556 <http://www.ncbi.nlm.nih.gov/gap> through accession number phs000145.v4.p2.

557

558 **Author contributions:** GBC, SHL, and BB conceived study. GBC, SHL, BB, and MRR
559 designed the experiment. GBC and SHL developed the theory and methods. BB conducted
560 the quality control for HapMap data, and MRR conducted quality control for POPRES data.
561 GBC performed the analyses of the study. GBC and ZZZ developed GEAR software. GBC,
562 MRR, SHL, and BB wrote the paper.

563

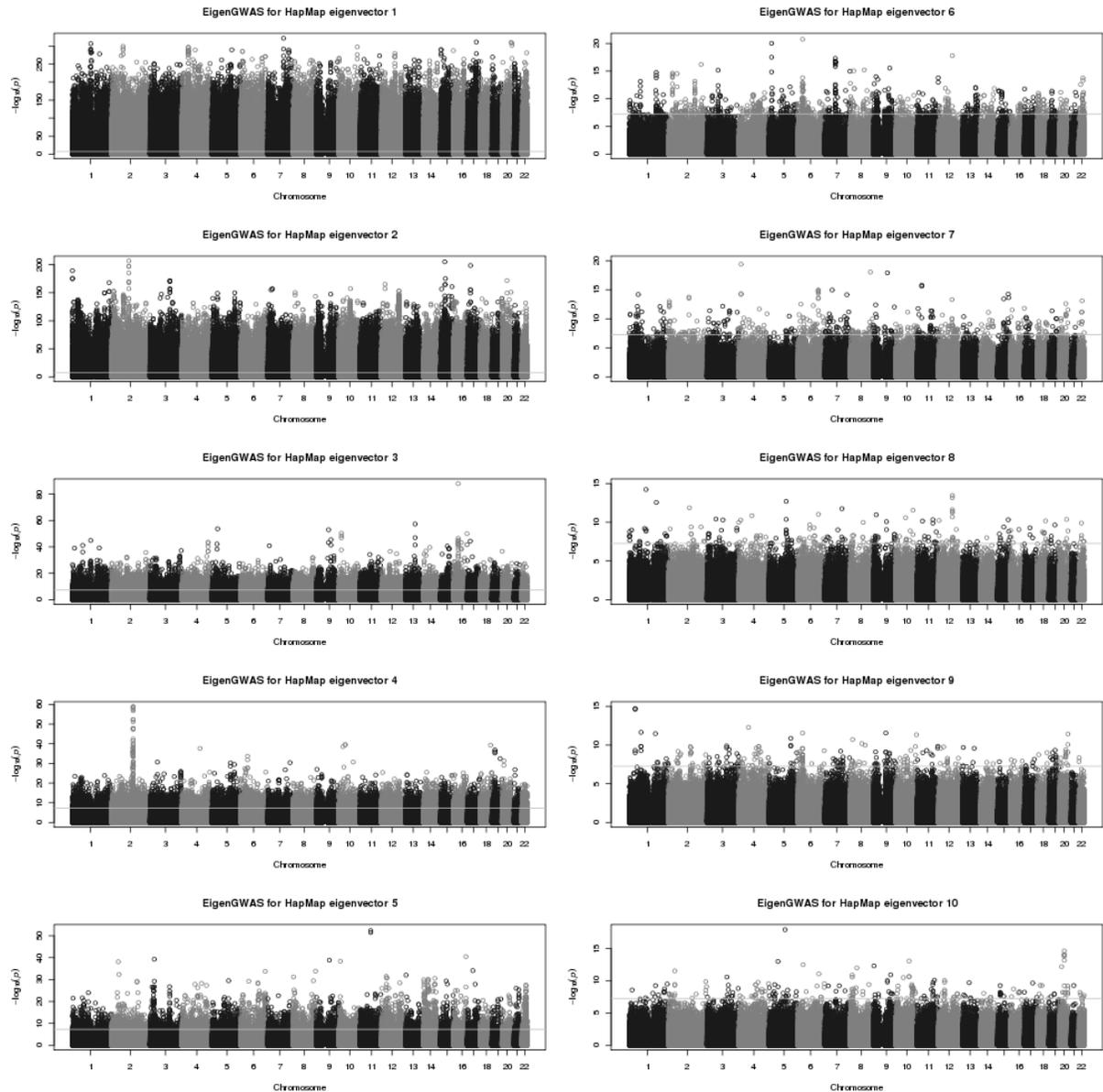
564 **References**

565 Altshuler DM, Gibbs RA, Peltonen L, Dermitzakis E, Schaffner SF, Yu F, *et al.* (2010).
566 Integrating common and rare genetic variation in diverse human populations. *Nature*

- 567 **467**: 52–8.
- 568 Bersaglieri T, Sabeti PC, Patterson N, Vanderploeg T, Schaffner SF, Drake JA, *et al.* (2004).
569 Genetic Signatures of Strong Recent Positive Selection at the Lactase Gene. *Am J Hum*
570 *Genet* **74**: 1111–1120.
- 571 Bryc K, Bryc W, Silverstein JW (2013). Separation of the largest eigenvalues in
572 eigenanalysis of genotype data from discrete subpopulations. *Theor Popul Biol* **89**: 34–
573 43.
- 574 Cavalli-Sforza LL, Menozzi P, Piazza A (1996). *The History and Geography of Human*
575 *Genes*. Princeton University Press.
- 576 Chen G-B (2014). Estimating heritability of complex traits from genome-wide association
577 studies using IBS-based Haseman-Elston regression. *Front Genet* **5**: 107.
- 578 Chen C-Y, Pollack S, Hunter DJ, Hirschhorn JN, Kraft P, Price AL (2013). Improved
579 ancestry inference using weights from external reference panels. *Bioinformatics* **29**:
580 1399–406.
- 581 Daetwyler HD, Villanueva B, Woolliams J a (2008). Accuracy of predicting the genetic risk
582 of disease using a genome-wide approach. *PLoS One* **3**: e3395.
- 583 Devlin B, Risch N (1995). A comparison of linkage disequilibrium measures for fine-scale
584 mapping. *Genomics* **29**: 311–22.
- 585 Devlin B, Roeder K (1999). Genomic control for association studies. *Biometrics* **55**: 997–
586 1004.
- 587 Dudbridge F (2013). Power and predictive accuracy of polygenic risk scores. *PLOS Genet* **9**:
588 e1003348.
- 589 Galinsky KJ, Bhatia G, Loh P, Georgiev S, Mukherjee S, Nick J (2015). Fast principal
590 components analysis reveals independent evolution of ADH1B gene in Europe and East
591 Asia. *bioRxiv*: <http://dx.doi.org/10.1101/018143>.
- 592 Goddard ME, Wray NR, Verbyla K, Visscher PM (2009). Estimating Effects and Making
593 Predictions from Genome-Wide Marker Data. *Stat Sci* **24**: 517–529.
- 594 Henderson CR (1975). Best Linear Unbiased Estimation and Prediction under a Selection
595 Model. *Biometrics* **31**: 423–447.
- 596 Howie B, Marchini J, Stephens M, Chakravarti A (2011). Genotype Imputation with
597 Thousands of Genomes. *G3* **1**: 457–470.
- 598 Maier R, Moser G, Chen G, Ripke S, Group CW, Consortium PG, *et al.* (2015). Joint
599 Analysis of Psychiatric Disorders Increases Accuracy of Risk Prediction for
600 Schizophrenia, Bipolar Disorder, and Major Depressive Disorder. *Am J Hum Genet* **96**:
601 283–294.

- 602 McVean G (2009). A genealogical interpretation of principal components analysis. *PLOS*
603 *Genet* **5**: e1000686.
- 604 Nei M (1973). Analysis of gene diversity in subdivided populations. *Proc Natl Acad Sci U S*
605 *A* **70**: 3321–3.
- 606 Nelson MR, Bryc K, King KS, Indap A, Boyko AR, Novembre J, *et al.* (2008). The
607 Population Reference Sample, POPRES: A Resource for Population, Disease, and
608 Pharmacological Genetics Research. *Am J Hum Genet* **83**: 347–358.
- 609 Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, *et al.* (2008). Genes mirror
610 geography within Europe. *Nature* **456**: 98–101.
- 611 Patterson N, Price AL, Reich D (2006). Population structure and eigenanalysis. *PLOS Genet*
612 **2**: e190.
- 613 Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick N a, Reich D (2006). Principal
614 components analysis corrects for stratification in genome-wide association studies. *Nat*
615 *Genet* **38**: 904–9.
- 616 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira M a R, Bender D, *et al.* (2007).
617 PLINK: a tool set for whole-genome association and population-based linkage analyses.
618 *Am J Hum Genet* **81**: 559–75.
- 619 Purcell SM, Wray NR, Stone JL, Visscher PM, O’Donovan MC, Sullivan PF, *et al.* (2009).
620 Common polygenic variation contributes to risk of schizophrenia and bipolar disorder.
621 *Nature* **460**: 748–752.
- 622 Reich D, Thangaraj K, Patterson N, Price AL, Singh L (2009). Reconstructing Indian
623 population history. *Nature* **461**: 489–94.
- 624 Rokhlin V, Szlam A, Tygert M (2009). A randomized algorithm for principal component
625 analysis. *SIAM J Matrix Anal Appl* **31**: 1100–1124.
- 626 Taft RJ, Vanderver A, Leventer RJ, Damiani S a, Simons C, Grimmond SM, *et al.* (2013).
627 Mutations in DARS cause hypomyelination with brain stem and spinal cord involvement
628 and leg spasticity. *Am J Hum Genet* **92**: 774–80.
- 629 The 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from
630 1,092 human genomes. *Nature* **491**: 56–65.
- 631 Visser M, Kayser M, Palstra R-J (2012). HERC2 rs12913832 modulates human pigmentation
632 by attenuating chromatin-loop formation between a long-range enhancer and the OCA2
633 promoter. *Genome Res* **22**: 446–55.
- 634 Voight BF, Kudaravalli S, Wen X, Pritchard JK (2006). A map of recent positive selection in
635 the human genome. *PLOS Biol* **4**: e72.
- 636 Weir BS (1996). *Genetic data analysis*, 2nd edn. Sinauer Associates, Inc.: Sunderland, MA,
637 USA.

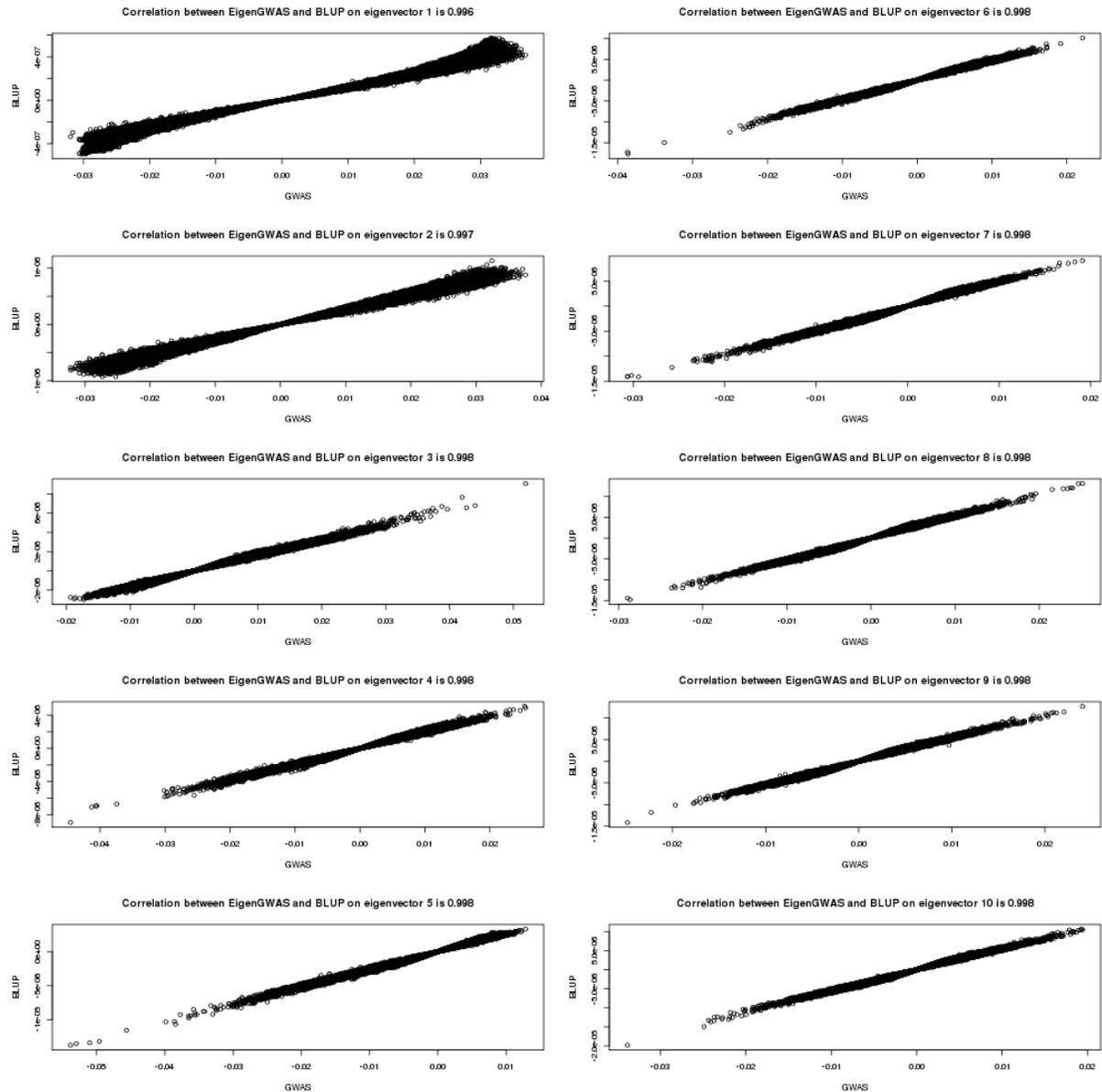
- 638 Williams AL, Patterson N, Glessner J, Hakonarson H, Reich D (2012). Phasing of many
639 thousands of genotyped samples. *Am J Hum Genet* **91**: 238–251.
- 640 Yang J, Lee SH, Goddard ME, Visscher PM (2011). GCTA: a tool for genome-wide complex
641 trait analysis. *Am J Hum Genet* **88**: 76–82.
- 642 Yang J, Weedon MN, Purcell S, Lettre G, Estrada K, Willer CJ, *et al.* (2011). Genomic
643 inflation factors under polygenic inheritance. *Eur J Hum Genet* **19**: 807–12.
- 644 Zhu X, Li S, Cooper RS, Elston RC (2008). A Unified Association Analysis Approach for
645 Family and Unrelated Samples Correcting for Stratification. *Am J Hum Genet* **82**: 352–
646 365.
- 647



648

649 **Figure 1** Manhattan plots for EigenGWAS for top 10 eigenvectors for HapMap. Using E_i as
650 the phenotype, the single-marker association was conducted for nearly 919,133 markers. The
651 left panel illustrates from $E_1 \sim E_5$; the right panel from $E_6 \sim E_{10}$. The horizontal lines indicate
652 genome-wide significant after Bonferroni correction.

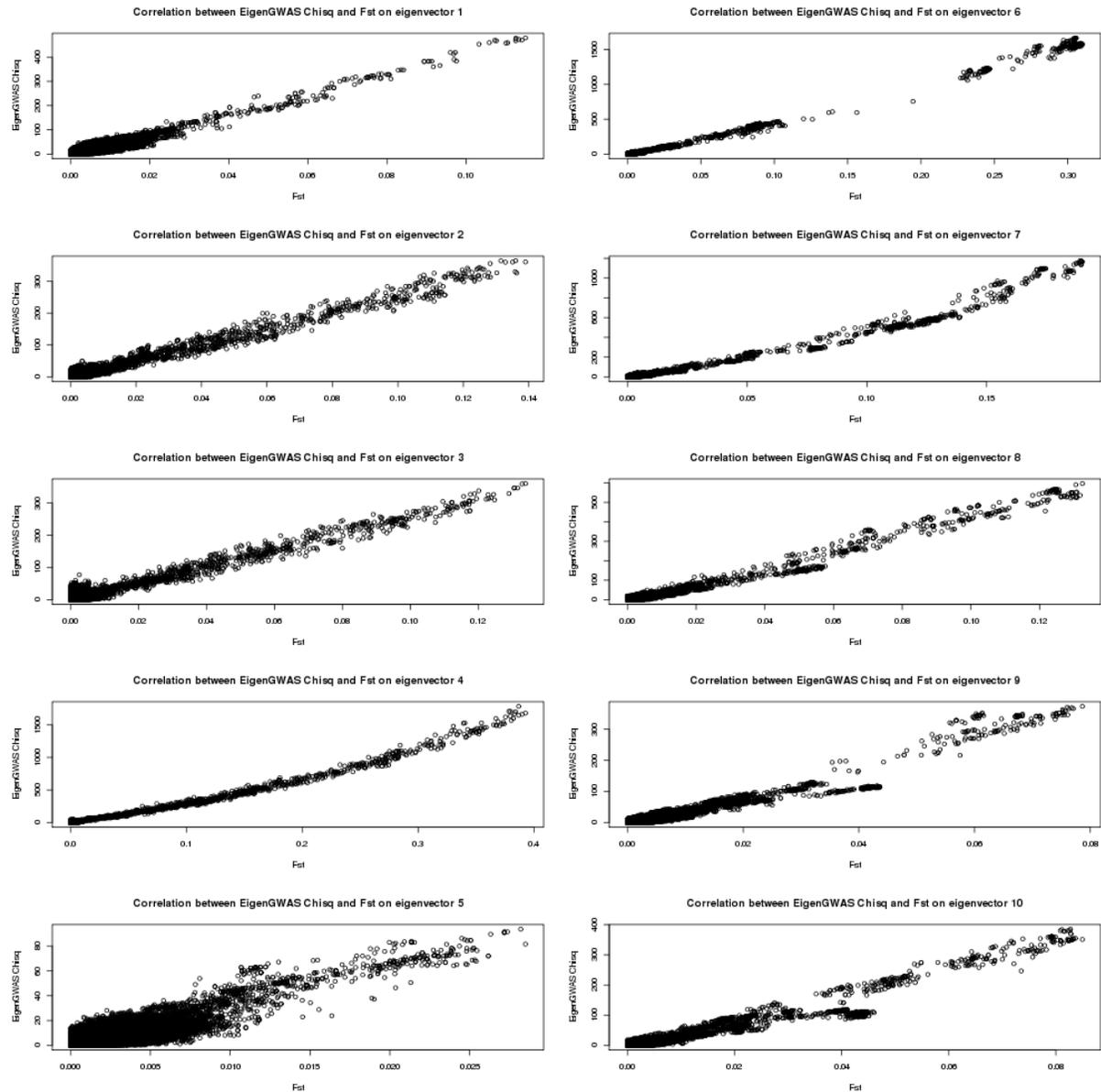
653



654

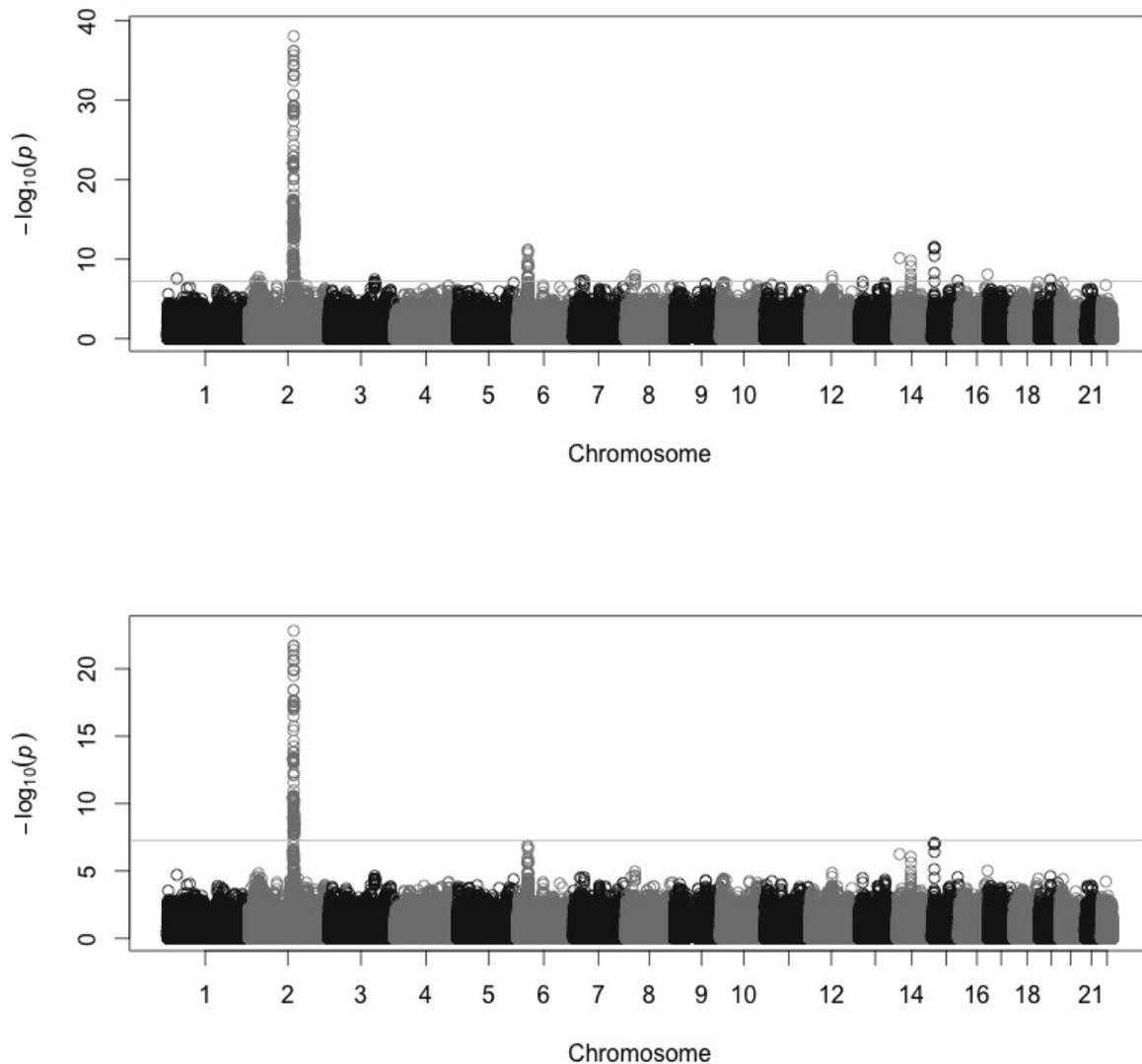
655 **Figure 2** Correlation for the SNP effects estimated using EigenGWAS and BLUP for
656 HapMap3. The x-axis represents EigenGWAS estimation for SNP effects, and the y-axis
657 represents BLUP estimation for SNP effects. The left panel illustrates from $E_1 \sim E_5$; the right
658 panel from $E_6 \sim E_{10}$. As illustrated, the correlation is nearly 1.

659



660
661
662
663
664

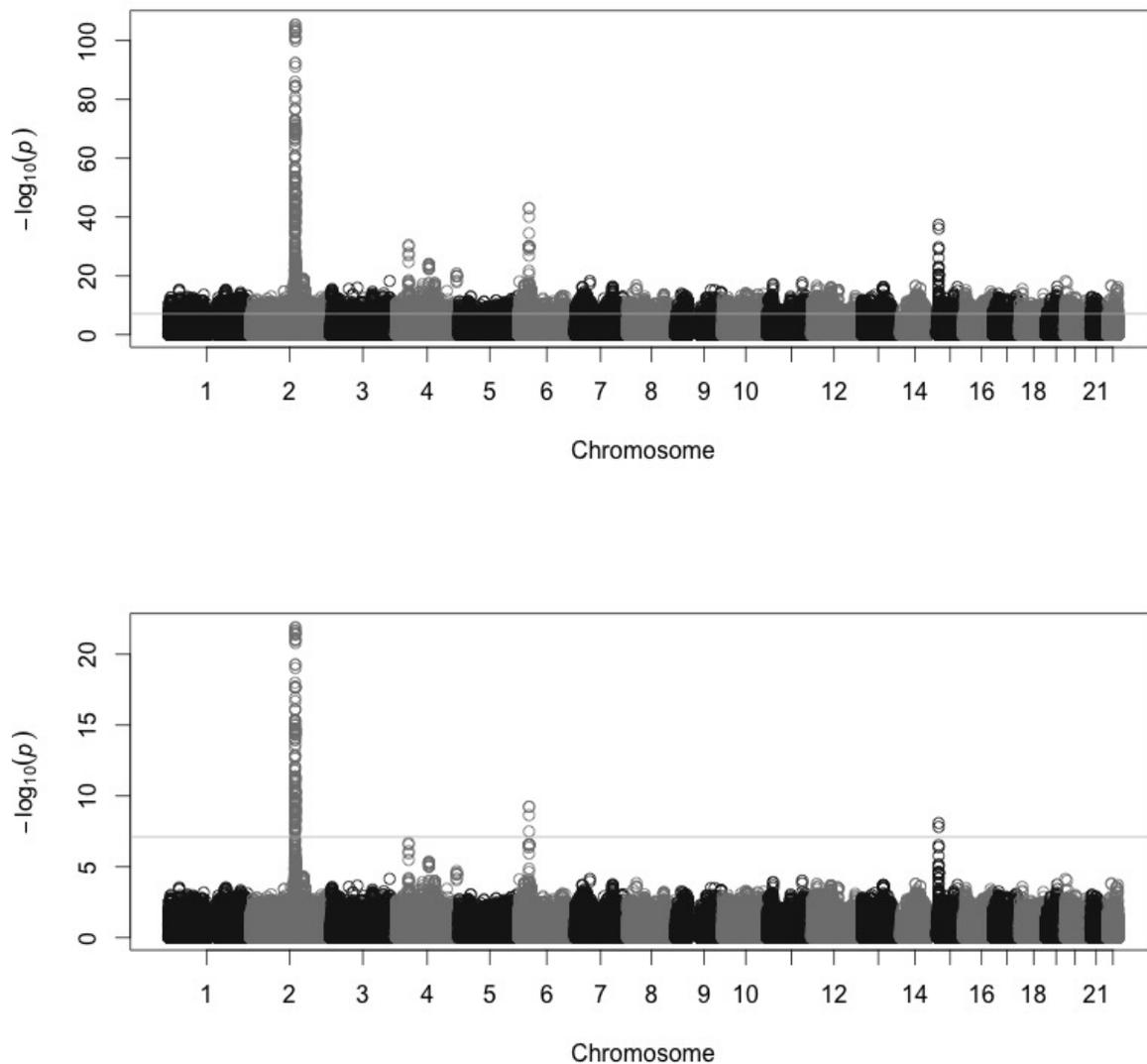
Figure 3 The correlation between F_{st} and χ^2_1 for EigenGWAS SNP effects for POPRES. For each eigenvector, upon $E_i > 0$ or $E_i \leq 0$, the samples were POPRES samples were split into two groups, upon which F_{st} was calculated for each locus.



665

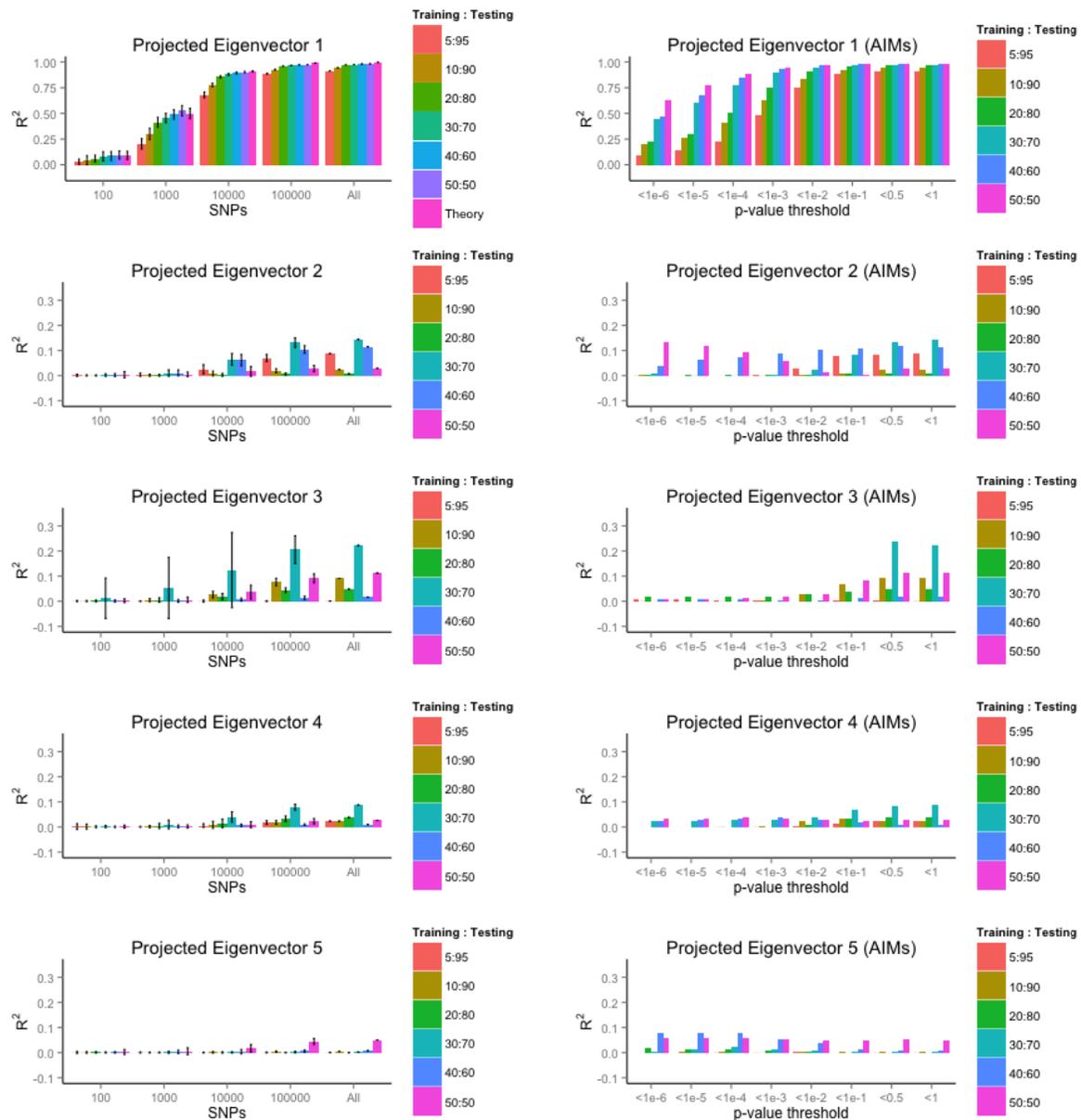
666 **Figure 4** EigenGWAS for CEU (112 samples) & TSI (88 samples) from HapMap. a)
667 Manhattan plot for EigenGWAS on E_1 without correction for λ_{GC} . When there was no
668 correction, on chromosome 2 found *LCT*, chromosome 6 *MICA* (HMC region), chromosome
669 14 *HIF1A*, and chromosome 15 *HERC2*. The line in the middle was for genome-wide
670 significant level at $\alpha = 0.05$ given multiple correction. b) Manhattan plot for EigenGWAS
671 on E_1 with λ_{GC} correction, and *LCT* was still significant, and *HERC2* slightly below whole
672 genome-wide significance level. The genome-wide significance threshold was p -values =
673 $5.44e-08$ for $\alpha = 0.05$.

674



675
676
677
678
679
680

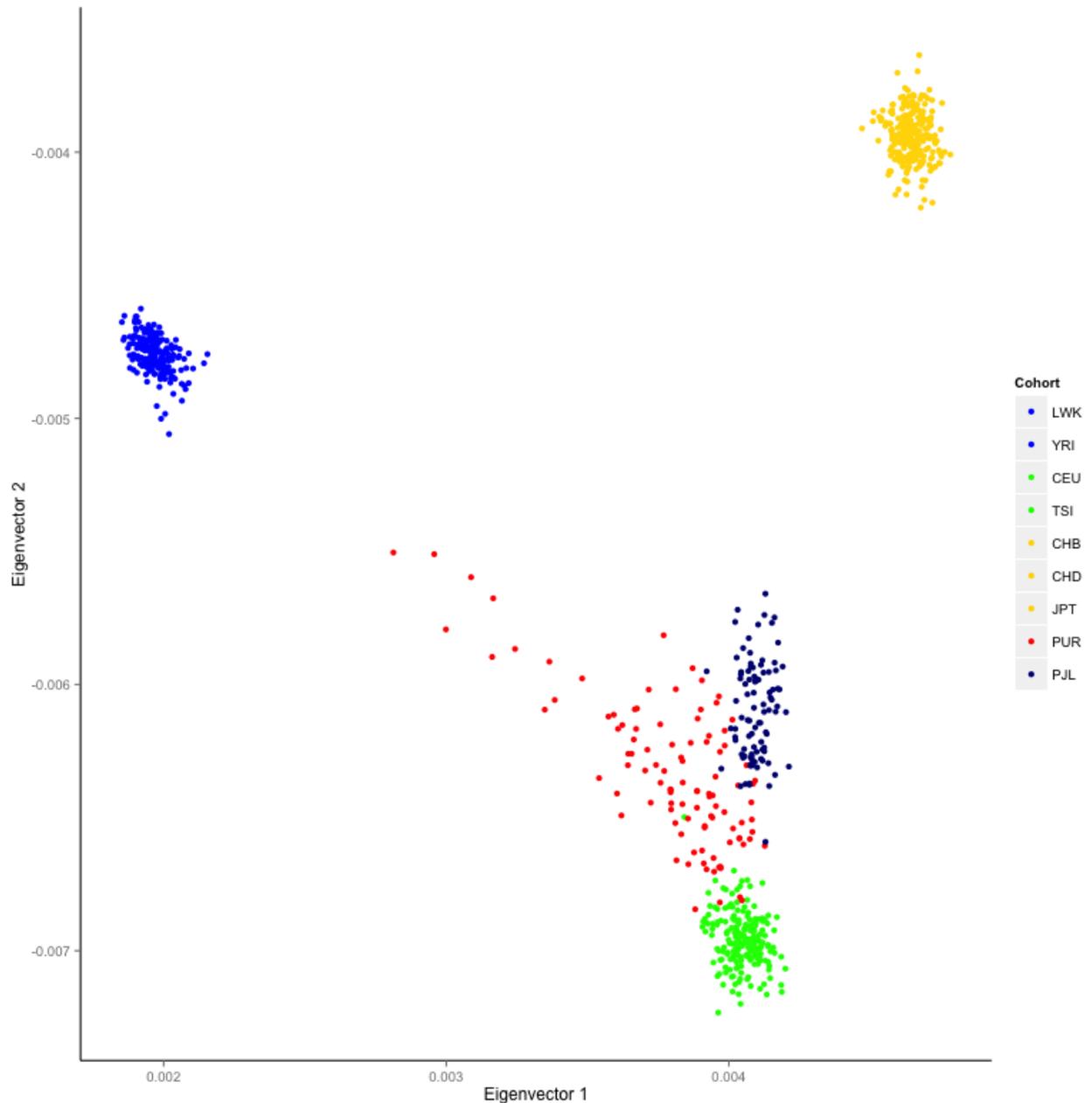
Figure 5 EigenGWAS for POPRES samples on eigenvector 1. a) Manhattan plot for EigenGWAS without correction for λ_{GC} . b) After correction for λ_{GC} , on Chromosome 2 found *LCT*, chromosome 6 *SLC44A4*, and chromosome 15 *HERC2*. The genome-wide significance level was p -values = $7.76e-08$ given $\alpha = 0.05$.



681

682 **Figure 6** Prediction accuracy of the projected eigenvectors for POPRES samples. Given
 683 2,466 POPRES samples, the data were split to 5%:95%, 10%:90%, 20%:80%, 30%:70%,
 684 40%:60%, and 50%:50, as training and test set. The left columns represent prediction
 685 accuracy (R^2) using randomly selected numbers (100, 1,000, 10,000, 100,000, all) of
 686 markers, the 95% confidence interval were calculated from 30 replication for resampling
 687 given number of markers. In contrast, the right columns represent the predicted accuracy for
 688 8 p-value thresholds ($1e-6$, $1e-5$, $1e-4$, $1e-3$, $1e-2$, $1e-1$, 0.5 , and 1) for EigengWAS SNPs.

689



690

691 **Figure 7 Projected eigenvectors for Puerto Rican cohort (PUR) and Pakistan cohort**
692 **(PJL) in 1000 Genome project.** The training set was HapMap3 samples build on 919,133
693 SNPs. The eigenvectors 1 and 2 for were generated based on the 74,500 common SNPs on
694 chromosome 1. PUR showed an admixture of African and European gene flows, and PJL
695 Asian and European gene flows.
696

Table 1 GWAS signals for on eigenvectors for HapMap and POPRES

Eigenvector (E_i)	HapMap				POPRES			
	Eigen value	GWAS λ_{GC}	#GWAS hits	#After clumping	Eigen value	GWAS λ_{GC}	#GWAS hits	#After clumping
1	100.135	103.715	546,716	231,677	5.104	5.005	10,885	3,004
2	47.658	44.686	382,867	161,022	2.207	1.929	1,254	289
3	7.168	6.471	33,317	15,344	2.157	1.910	1,201	340
4	5.923	5.173	21,935	12,401	2.077	1.464	1,353	331
5	4.402	3.964	9,554	4,727	1.971	1.866	781	76
6	2.449	1.982	1,113	567	1.871	1.295	1,162	111
7	2.285	1.986	593	389	1.843	1.337	1,239	130
8	2.107	1.742	236	171	1.818	1.486	1,259	152
9	2.056	1.729	268	174	1.807	1.503	1,701	113
10	2.0217	1.661	236	163	1.798	1.492	1,639	90

Notes: HapMap has 988 samples, and 919,133 SNPs; its GWAS hits were those had p -values $< 5.44e-08$ given $\alpha = 0.05$. POPRES has 2,466 European samples, and 643,995 SNPs; its GWAS hits were those had p -values $< 7.76e-08$ given $\alpha = 0.05$.

After clumping, the reported numbers were quasi-independent GWAS hits. Within 250K bp and linkage disequilibrium of $r^2 > 0.5$ only the most significant GWAS hit was counted as a GWAS hit (see PLINK --clump default option).

λ_{GC} was calculated as the ratio between the median of observed χ^2 from EigenGWAS to the median of χ^2 value, which is 0.455.

Table 2 Gene discovery using EigenGWAS

Gene	Lead SNP	Position	Allele	p -value (λ_{GC})	MAF (TSI:CEU)	F_{st}	Annotation
CEU & TSI samples							
<i>LCT</i>	rs6719488	2:135817629	G/T	6.68e-34 (1.21e-20)	0.733:0.206	0.558	Lactose persistent locus
<i>DARS</i>	rs13404551	2:135964425	C/T	8.18e-39 (1.51e-23)	0.756:0.206	0.604	Genetic hitchhiking due to <i>LCT</i> .
<i>MICA</i>	rs2256175	6:31412672	T/C	8.94e-10 (2.60e-6)	0.665:0.360	0.183	MHC class I polypeptide-related sequence A
<i>HIF1A</i>	rs2256205	14:61670944	A/G	1.51e-10 (8.86e-7)	0.464:0.179	0.192	HIF-1A thus plays an essential role in embryonic vascularization, tumor angiogenesis and pathophysiology of ischemic disease.
<i>HERC2</i>	rs8039195	15:26189679	C/T	2.75e-12 (8.22e-08)	0.403:0.122	0.212	Genetic variations in this gene are associated with skin/hair/eye pigmentation variability
POPRES European samples							
					Southern Europeans : Northern Europeans		
<i>LCT</i>	rs3754686	2:135817629	T/C	3.30e-106 (1.23e-22)	0.514:0.279	0.110	
<i>DARS</i>	rs13404551	2:135964425	C/T	6.32e-102 (8.99e-22)	0.518:0.293	0.106	
<i>SLC4A4</i>	rs605203	6:31819235	C/A	8.94e-44 (5.77e-10)	0.214:0.343	0.040	Defects in this gene can cause sialidosis, a lysosomal storage disease
<i>HERC2</i>	rs1667394	15:26189679	C/T	3.90e-38 (8.15e-09)	0.276:0.173	0.041	

Notes: The p -value cutoff for CEU&TSI was $5.44e-08$ (919,133 SNPs), for POPRES was $7.76e-08$ (643,995 SNPs) at genome-wide significance level of $\alpha = 0.05$. $\lambda_{GC} = 1.725$ for CEU&TSI, and $\lambda_{GC} = 5.00$ for POPRES.

F_{st} is calculated by partitioning the sample into two groups upon $E_1 > 0$. For TSI&CEU set, partitioning on E_1 perfectly separated TSI (88 samples) and CEU (112 samples). For POPRES, partitioning on E_1 separated southern European population (1,092 samples) and northern European population (1,374 samples).