1 **Title:** Across-cohort QC analyses of genome-wide association study summary statistics from complex traits

2

3 **Authors:** Guo-Bo Chen[1], Sang Hong Lee[1,2], Matthew R Robinson[1], Maciej Trzaskowski[1], Zhi-Xiang Zhu[3],
4 Thomas W Winkler[4], Felix R Day[5], Damien C Croteau-Chonka[6,7], Andrew R Wood[8], Adam E Locke[9],
5 Zoltán Kutalik[10-12], Ruth J F Loos[13-15], Timothy M Frayling[8], Joel N Hirschhorn[16-19], Jian Yang[1,21], Naomi R
6 Wray[1], The Genetic Investigation of Anthropometric Traits (GIANT) Consortium[20], Peter M Visscher[1,21]

7

8 **Affiliations:**
9 [1] Queensland Brain Institute, The University of Queensland, Brisbane, Queensland, Australia
10 [2] School of Environmental and Rural Science, The University of New England, Armidale, New South
11 Walsh, Australia
12 [3] SPLUS Game, Guangzhou, Guangdong, China
13 [4] Department of Genetic Epidemiology, Institute of Epidemiology and Preventive Medicine, University of
14 Regensburg, Regensburg, Germany
15 [5] Medical Research Council (MRC) Epidemiology Unit, Institute of Metabolic Science, Addenbrooke's
16 Hospital, Cambridge, UK
17 [6] Department of Genetics, University of North Carolina, Chapel Hill, North Carolina, USA
18 [7] Channing Division of Network Medicine, Department of Medicine, Brigham and Women's Hospital and
19 Harvard Medical School, Boston, Massachusetts, USA
20 [8] Genetics of Complex Traits, University of Exeter Medical School, University of Exeter, Exeter, UK
21 [9] Department of Biostatistics and Center for Statistical Genetics, University of Michigan, Ann Arbor,
22 Michigan, USA
23 [10] Department of Medical Genetics, University of Lausanne, Lausanne, Switzerland
24 [11] Institute of Social and Preventive Medicine (IUMSP), Centre Hospitalier Universitaire Vaudois (CHUV),
25 Lausanne, Switzerland
26 [12] Swiss Institute of Bioinformatics, Lausanne, Switzerland
27 [13] The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai,
28 New York, New York, USA
29 [14] The Mindich Child Health and Development Institute, Icahn School of Medicine at Mount Sinai, New
30 York, New York, USA
31 [15] The Genetics of Obesity and Related Metabolic Traits Program, Icahn School of Medicine at Mount Sinai,
32 New York, New York, USA
33 [16] Department of Genetics, Harvard Medical School, Boston, Massachusetts, USA
34 [17] Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge,
35 Massachusetts, USA
36 [18] Center for Basic and Translational Obesity Research, Boston Children's Hospital, Boston, Massachusetts,
37 USA
38 [19] Division of Endocrinology, Boston Children's Hospital, Boston, Massachusetts, USA
39 [20] A full list of members is available in the **Supplementary Note**
40 [21] The University of Queensland Diamantina Institute, Translation Research Institute, Brisbane, Queensland,
41 Australia

42

43 **Correspondence should be addressed to**

44 GBC (chen.guobo@foxmail.com) or PMV (peter.visscher@uq.edu.au)

47

48

**Abstract**

49      Genome-wide association studies (GWASs) have been successful in discovering replicable SNP-trait

50      associations for many quantitative traits and common diseases in humans. Typically the effect sizes of SNP

51      alleles are very small and this has led to large genome-wide association meta-analyses (GWAMA) to

52      maximize statistical power. A trend towards ever-larger GWAMA is likely to continue, yet dealing with

53      summary statistics from hundreds of cohorts increases logistical and quality control problems, including

54      unknown sample overlap, and these can lead to both false positive and false negative findings. In this study

55      we propose a new set of metrics and visualization tools for GWAMA, using summary statistics from cohort-

56      level GWASs. We proposed a pair of methods in examining the concordance between demographic

57      information and summary statistics. In method I, we use the population genetics $F_{st}$ statistic to verify the

58      genetic origin of each cohort and their geographic location, and demonstrate using GWAMA data from the

59      GIANT Consortium that geographic locations of cohorts can be recovered and outlier cohorts can be

60      detected. In method II, we conduct principal component analysis based on reported allele frequencies, and is

61      able to recover the ancestral information for each cohort. In addition, we propose a new statistic that uses the

62      reported allelic effect sizes and their standard errors to identify significant sample overlap or heterogeneity

63      between pairs of cohorts. Finally, to quantify unknown sample overlap across all pairs of cohorts we propose

64      a method that uses randomly generated genetic predictors that does not require the sharing of individual-

65      level genotype data and does not breach individual privacy.

66

## Introduction

Genome-wide association studies (GWASs) have been successful in discovering SNP-trait associations for complex traits[1]. To elucidate genetic architecture, which requires maximized statistical power for discovery of risk alleles of small effect, large genome-wide association meta-analyses (GWAMA) are tending towards ever-larger scale that may contain data from hundreds of cohorts. At the individual cohort level, GWAS analysis is often based on various genotyping chips and conducted with different protocols, such as different software tools and reference populations for imputation, inclusion of study specific covariates and association analyses using different methods and software. Although solid quality control analysis pipelines of GWAMA exist[2], these analyses focus on quality control (QC) for each cohort independently. With ever-increasing sizes of GWAMA there is a need for additional QC that goes beyond the cohort-by-cohort genotype-level analysis performed to date.

In this study, we propose a new set of QC metrics for GWAMA. In contrast to previous QC metrics, our approach explores the genetic and QC context of the all cohorts in GWAMA together rather than by treating them one at a time. These metrics include

(i) a genome-wide comparison of allele frequency differences across cohorts or against a common reference population

(ii) principal component analysis for reported allele frequencies

(iii) a pairwise cohort statistic that uses allele frequency or effect size concordance to detect the proportion of sample overlap or heterogeneity

(vi) an easy to implement analysis to pinpoint each between-cohort overlapping sample that does not require the sharing of individual-level genotype data.

All these applications assume that there is a central analysis hub where summary statistic data from GWAS are uploaded for each cohort. In addition, these metrics reveal information of interest other than merely QC.

## Materials and Methods

**Overview of materials**

**Cohort-level summary statistics.** The GWAS height GWAS summary statistics were provided by the GIANT Consortium and were from 82 cohorts (174 separate files due to different ways a cohort was split into different sexes, different disease statuses) representing a total of 253,288 individuals, and nearly 2.5 million autosome SNPs imputed to the HapMap2 reference[3]. The Metabochip summary statistics were for body mass index (BMI) from 43 cohorts (120 files due to different ways a cohort was split into different sexes, different disease statuses) representing a total of 103,047 samples from multiple ethnicities with about 200,000 SNPs genotyped on customised chips[4,5]. For convenience, we consider each file a cohort. All the summary statistics have already been cleaned using established protocols for GWAS meta-analysis[2].

3

**1000 Genomes Project samples.** 1000 Genomes Project (1KG) reference samples[6] were used as the reference samples for calculating $F_{st}$. When assessing the global-level $F_{st}$ measures, Yoruba represent African samples (YRI, 108 individuals), Han Chinese in Beijing represent East Asian samples (CHB, 103 individuals), and Utah Residents with Northern and Western European Ancestry represent European samples (CEU, 99 individuals) were employed as the reference panels. For calculating within-Europe $F_{st}$, CEU, Finnish (FIN, 99 individuals), and Tuscani (TSI, 107 individuals) were employed to represent northwest, northeast and southern Europeans, respectively. For analyses using a whole European panel, CEU, FIN, TSI, GBR (British, 91 individuals), and IBS (Iberian, 107 individuals) were pooled together as an "averaged" European reference.

**WTCCC GWAS data.** WTCCC GWAS data has 2,934 shared controls for 7 diseases with a total of 14,000 cases[7]. Individual GWAS was conducted for each disease using PLINK[8], and their summary statistics used to estimate $\lambda_{meta}$ (see text below). WTCCC GWAS data were also used for demonstrating pseudo profile score regression (see text below).

**Simulated cohort-level summary statistics.** $M$ independent loci were generated for cohort-level summary statistics. Each locus had allele frequency $p_i$, which was sampled from a uniform distribution ranging from 0.1 to 0.5, and had genetic effect $b_i$, sampled from a standard normal distribution $N(0,1)$. After rescaling, $\Sigma_{i=1}^{M} 2p_i(1-p_i)b_i^2 = h^2$. $p$ and $b$ were treated as true parameters. For a particular cohort with $n$ samples, its $\tilde{p}_i \sim N(p_i, \frac{p_i(1-p_i)}{2n})$, $\tilde{b}_i \sim N(b_i, \frac{1}{2np_i(1-p_i)})$, and the sampling variance for $\tilde{b}_i$ is $\sigma_{b_i}^2 = \frac{1}{2n_ip_i(1-p_i)}$. All cohorts were assumed to share common genetic architecture, and differences were only due to genetic drift, allele frequencies and sampling variance of genetic effects.

**Overview of the methods**

$F_{st}$**-based genetic distance between cohorts.** For a cohort, its $F_{st}$ with reference cohorts, such as CEU, YRI, and CHB, is calculated. Given those three $F_{st}$ values, the coordinate of this cohort can be uniquely projected into the reference equilateral that has CEU, YRI, and CHB at its corners.

**Principal component analysis for cohort-level allele frequencies.** A genetic relationship matrix for cohorts can be constructed based on received allele frequencies. Principal component analysis (PCA) can be implemented on the genetic relationship matrix. The projection of the cohorts into PCA space can reveal the genetic background and relative geographical distance between cohorts.

$\lambda_{meta}$ **for detecting overlapping samples.** In concept, $\lambda_{meta}$ resembles $\lambda_{gc}$, which indicates population stratification for a GWAS[9], but $\lambda_{meta}$ measures the proportion of overlapping samples between a pair of

39     cohorts. Based on reported genetic effects and their sampling variance, $\lambda_{meta}$ can be constructed for a pair of

40     cohorts and follows a chi-square distribution with 1 degree of freedom. $\lambda_{meta}$ will be close to 1 when there

41     is no overlapping samples, smaller than 1 when there are overlapping samples, and greater than 1 when there

42     are heterogeneity between a pair of cohorts. For GWAMA over a single trait across method, we assume

43     heterogeneity is zero.

44

45     **Pseudo profile score regression for pinpointing overlapping samples/relatives.** Pseudo profile score

46     regression (PPSR) provides a framework for pinpointing the overlapping samples/relatives between cohort

47     without sharing genotypes. Each GWAS analyst generates pseudo profile scores (PPS) for each sample on a

48     set of loci, which are chosen by a GWAMA central analyst. If the similarity metric of PPS for a pair of

49     cohorts reaches a similarity threshold, say 1 overlapping samples and 0.5 for first-degree relatives, then

50     overlapping samples/relatives are found. PPSR can have a controlled type-I and type-II error rates in

51     pinpointing overlapping samples, and also can reduce the comprise of privacy. PPSR is an enhanced version

52     of Gencrypt[10], a previous method in pinpointing overlapping samples.

53

54     The technical details of these four methods can be found in the Supplementary notes.

55

56                                            **Results**

57     **Population genetic quality control analysis using $F_{st}$**

58     Allele frequency differentiation among populations reflects population characteristics such as demographic

59     past and geographic locations[11,12]. In GWAMA only summary statistics such as allele frequencies are

60     available to the central analysis hub, and so it is not possible to run principal component analysis for each

61     cohort that requires individual-level data. Therefore it is difficult to quantify genetic distance between

62     cohorts or to a reference in order to identify population outliers. Outlier cohorts can be due to real

63     differences in ethnicity or mistakes in the primary analysis prior to uploading data to the GWAMA analysis

64     hub. Gross differentiation in allele frequencies at specific SNPs between GWAMA cohorts and a reference

65     (such as 1000 Genomes Project, denoted as 1KG)[6] are part of standard QC protocols[2] but checking for more

66     differentiation than expected across the entire genome is not usually part of the QC pipeline. We propose

67     that a genetic distance inferred from $F_{st}$, which reflects genetic distance between pairwise populations, is a

68     useful additional QC statistic to detect cohorts that are population outliers. Using the relationship between

69     $F_{st}$ and principal components[13–15], our $F_{st}$ Cartographer algorithm can be used to estimate the relative

70     genetic distance between cohorts (**Supplementary notes, and Supplementary Fig. 1**).

71

72     We applied the $F_{st}$ metric to the GIANT Consortium body mass index (BMI) Metabochip cohorts (55 male-

73     only cohorts, 55 female-only cohort, and 10 mixed-sex cohorts; for convenience, we called each file a

74     cohort), which were recruited from multiple ethnicities[4], such as Europeans, African Americans in The

175 Atherosclerosis Risk in Communities Study (ARIC) and cohorts from Jamaica (SPT), Pakistan (PROMISE),

176 Philippines (CLHNS) and Seychelles (SEY). For each Metabochip cohort, we sampled 30,000 (see Online

177 method for details) independent markers to calculate $F_{st}$ values with each of three 1KG samples (CEU, CHB,

178 and YRI, respectively). For validation of the method, we also calculated $F_{st}$ values against the 1KG

179 Japanese (JPT, Japanese in Tokyo, Japan), Indian (GIH, Gujarati Indian in Houston, US), Kenyan (LWK,

180 Luhya in Webuye, Kenya) and European samples (IBS, Iberian populations, Spain; FIN, Finnish, Finland;

181 TSI, Toscani, Italy, and GBR, British in England and Scortland, GBR), to see whether the known genetic

182 origins of those cohorts can be recovered.

183

184 According to the origins of the samples, each Metabochip cohort showed a different genetic distance

185 spectrum to the three reference populations (**Fig. 1a**). The JPT and Philippine cohorts had very small genetic

186 distances to CHB, as expected, but large to CEU and YRI; however, the Pakistan cohorts showed much

187 closer genetic distances to CEU than to CHB and YRI, indicating their demographic history. The cohorts

188 sampled from Jamaica, Seychelles, Hawaii, and the African American ARIC cohort had small genetic

189 distances to YRI, but large distances to CHB and CEU. For most European cohorts, as expected, the

190 distances to CEU were very small compared with those to CHB and YRI. Given their relative distances to

191 CEU, CHB, and YRI, using our $F_{st}$ cartographer algorithm (**Supplementary notes, and Supplementary**

192 **Fig. 1**), the cohorts were projected into a two-dimensional space, called $F_{st}$ derived principal components

193 ($F_{PC}$) space, constructed by YRI, CHB, and CEU as the reference populations (**Fig. 1b**). The allocation of

194 the cohorts to the $F_{PC}$ space resembles that of eigenvector 1 against eigenvector 2 in principal component

195 analysis (PCA)[12], and is similar to those observed in PCA using individual-level GWAS data for populations

196 of various ethnicities such as in 1KG samples[6]. Therefore, our method to place cohorts in geographical

197 regions from GWAS summary statistics works well at a global-population scale.

198

199 We next investigated whether our genetic distance method works at a much finer geographic scale. It is

200 known that using individual-level data, principal component analysis can mirror the geographic locations for

201 European samples[11]. Here, we analyzed the 103 GIANT European-ancestry Metabochip cohorts (48 male-

202 only cohorts, 47 female-only cohorts, and 8 mix-sex cohorts) for fine-scale $F_{st}$ genetic distance measure by

203 using the CEU, FIN, and TSI reference populations, which represent northwest, northeast, and southern

204 European populations, respectively. For each of the GIANT European-ancestry Metabochip cohorts, $F_{st}$ was

205 calculated relative to each of these three reference populations and showed concordance with the known

206 origin of the samples (**Fig. 1c**). For example, cohorts from Finland and Estonia were close to FIN but distant

207 to TSI; cohorts from South Europe such as Italy and Greece had small genetic distance to TSI; and cohorts

208 from West European nations had small genetic distance to CEU. Similarly, the projected origin for each

209 European-ancestry Metabochip cohort resembles their geographic location within the European map as

6

210    expected (**Fig. 1d**). Therefore, our QC measure based upon population differentiation also works at a fine

211    scale.

212

213    We next applied the $F_{st}$ genetic distance measures to 174 GIANT height GWAS cohorts (79 male-only

214    cohorts, 76 female-only cohorts, and 19 mixed-sex cohorts; excluding Metabochip data), which were all of

215    European ancestry imputed to the HapMap reference panel[3]. Given the three $F_{st}$ values to CEU, FIN, and

216    TSI (**Fig. 2a**), the geographic origin for each cohort can be inferred as for the GIANT BMI Metabochip data

217    (**Supplementary notes**). The projected coordinates of each GWAS cohort matches its origin very well (**Fig.

218    2b**). For example, a Canadian cohort, the Quebec Family Study (QFS), was closely located to DESIR, a

219    French cohort, consistent with the French genetic heritage of the QFS[16]. In addition, we also observe

220    complexity due to mixed samples from different countries. For example, the DGI/Botnia study had samples

221    recruited from Sweden and Finland, and its inferred geographic location is in between of the Swedish

222    cohorts and Finnish cohorts[17]. We also note that for the MIGEN consortia cohorts, which are from Finland,

223    Sweden, Spain and the US, the same allele frequencies were reported for all their sub-cohorts, and all

224    cohorts were allocated to southern Europe (very closely located to 1KG IBS cohort; **Fig. 2b and

225    Supplementary Fig. 2**). As the allele frequencies, used in QC steps to eliminate low quality loci, were not

226    directly used in estimating genetic effects in the GWAMA, the reported allele frequencies in MIGEN have

227    not impacted on the published GWAMA results[3].

228

229    Next, we show that $F_{st}$ can detect populations that have a different demographic past. Using all 1KG

230    European samples as the reference panel (that is, an "averaged" European reference panel), most cohorts in

231    GIANT had $F_{st} < 0.005$ with this average, which agrees with previously reported results using individual

232    level data from European nations[11]. A few cohorts showed large $F_{st}$, such as the AMISH cohort with

233    $F_{st} = 0.018$, and the North Swedish Population Health Study (NSPHS)[18] with $F_{st} = 0.014$. Consistent with

234    these results, both these populations are known to have been genetically isolated (**Supplementary Fig. 3**).

235

236    **Principal component analysis for allele frequencies**

237    It is well established that given individual-level data principal component analysis (PCA) can reveal the

238    ancestral information for samples[12]. Given the same allele frequencies as used for $F_{st}$-based analysis above,

239    we conducted PCA for allele frequencies, denoted as meta-PCA. In meta-PCA each cohort was analogously

240    considered as an "individual". For example, 120 Metabochip cohorts were considered as a sample of 120

241    "individuals". Although the inferred ancestral information was for each cohort rather than any individuals,

242    implementation of meta-PCA was the same as the conventional PCA (**Supplementary Notes**).

243

244    Meta-PCA was tested with 1KG samples over nearly 1 million SNPs. The cohort-level allele frequencies

245    were calculated first for 26 1KG cohorts, and meta-PCA was conducted. The projected cohorts were

246     consistent to their genetic origin (**Fig. 3**). In contrast, conventional PCA was also conducted on 1KG

247     individual genotypes directly, and the mean coordinates for each cohort was then calculated. As illustrated in

248     Fig. 3, these two techniques resulted in nearly identical projection for 1KG, and the correlation between

249     cohort coordinates remained consistently high for the first eight eigenvectors, $R^2 > 0.8$. It indicated that

250     meta-PCA could reveal genetic background for each cohort as precise as that based on individual-level data.

251

252     We applied meta-PCA to 120 Metabochip cohorts for nearly 34 thousand common SNPs between

253     Metabochip and 1KG variants, with the inclusion of 10 1KG cohorts (East Asian: CHB, JPT; South Asian:

254     GIH; European: CEU, FIN, GBR, IBS, TSI; African: LWK, YRI) as the reference cohorts. Consistent with

255     demographic information, the inferred ancestral information of each cohort agreed well with demographic

256     information. For example, PROMISE (Pakistan) located very close to GIH, CLHNS (Philippines) close to

257     CHB and JPT, ARIC (African American) and SPT (Jamaican) close to YRI and LWK, and the European

258     cohorts close to CEU and FIN (**Fig 4**).

259

260     We also applied meta-PCA to 174 GIANT height GWAS cohorts for nearly 1M SNPs, with the inclusion of

261     10 1KG reference cohorts. At the global-population level, the 174 cohorts were all allocated close to CEU

262     and FIN, consistent with their reported demographic information (**Fig. 5**). For fine-scale inference, we

263     conducted meta-PCA again but with the inclusion of the five European samples. As demonstrated, the

264     resolution of the inferred relative location between European cohorts reflected their real geographical

265     locations, as previously observed using individual-level data[11].

266

267     These results were consistent to what observed from $F_{pc}$ as described in the last section, and also agreed well

268     with demographic information. So, based on the reported allele frequencies, the demographic information

269     could be examined by meta-PCA method.

270

271     **$\lambda_{meta}$ to detect pairwise cohort heterogeneity and sample overlap**

272     For a single cohort GWAS, $\lambda_{GC}$ provides a tool for assessing average trait-SNP associations in GWAS[9], and

273     an value departing from 1 may indicate undesired phenomena such as population stratification. In this study,

274     we use the summary statistics for a pair of cohorts to calculate $\lambda_{meta}$, a metric that examines heterogeneity

275     from the concordance of reported effect sizes and sampling variance. We use 30,000 markers in linkage

276     equilibrium along the genome between a pair of cohorts to estimate $\lambda_{meta}$.

277

278     For a SNP marker (*i*), given its reported estimated effect size ($b_i$) and sampling variance ($\sigma_i^2$) in a pair of

279     cohorts 1 and 2, we can calculate a test statistic $T_i = \frac{(b_{1.i} - b_{2.i})^2}{\sigma_{1.i}^2 + \sigma_{2.i}^2}$, the ratio between the squared difference of

280     their reported effects to the sum of their reported sampling variances. Under the null hypothesis of no

281     overlapping samples/heterogeneity, $T$ follows a chi-square distribution with 1 degree of freedom

282     (**Supplementary notes**). $\lambda_{meta} = \frac{median(T)}{median(\chi_1^2)}$, the ratio between the median of the 30,000 $T$ values and the

283     median of a chi-square statistic with 1 degree of freedom (a value of 0.455), has an expected value of 1 for

284     two independent GWAS summary statistics sets for the same trait. When there is heterogeneity between

285     estimated genetic effects, the expectation is $\lambda_{meta} > 1$, and in contrast $\lambda_{meta} < 1$ if there are overlapping

286     samples. In general, not only overlapping samples but also close relatives present in different cohorts can

287     lead to correlated summary statistics generating $\lambda_{meta} < 1$ (**Supplementary notes**). However, unless the

288     proportion of overlapping relatives is substantial and their phenotypic correlation is high, the correlation of

289     the summary statistics due to the effective number of overlapping samples ($n_o$) is expected to be dominated

290     by the same individuals contributing phenotypic and genetic information to different cohorts

291     (**Supplementary Fig. 4**). Furthermore, if genomic control is applied to adjust the sampling variance[19] then

292     $\lambda_{meta}$ will be reduced relative to its value without genomic control (**Supplementary notes**).

293

294     We estimated $\lambda_{meta}$ from published GWAS summary statistics for a range of traits (other than BMI and

295     height) and were able to find examples of both deflated and inflated $\lambda_{meta}$. First, we tested the $\lambda_{meta}$ on data

296     sets with known overlap. For example, GWAS summary statistics for schizophrenia were available in two

297     phases: the first had 9,394 controls and 12,462 cases[20], and in the next phase about 18,000 Swedish samples

298     were added[21]. Such a substantial overlap sample between these two sets of summary statistics led to the

299     estimated value of $\lambda_{meta}$ as low as 0.257 (**Supplementary Fig. 5**), consistent with this known overlap. In

300     contrast, heterogeneity between data sets (represented by $\lambda_{meta} > 1$), was observed between GWAS

301     summary statistics of rheumatoid arthritis from European and Asian studies[22], for which $\lambda_{meta} = 1.09$

302     (**Supplementary Fig. 6**). In addition, we note that the distribution of the empirical $T$-statistics deviates from

303     expectation at the upper tail of the distribution, suggesting differences in effect size or linkage

304     disequilibrium between these two ancestries.

305

306     Next, we estimated $\lambda_{meta}$ from pairs of cohorts from the 174 GIANT height GWAS cohort[3]. We found no

307     evidence for substantial sample overlap but do observe between-cohort heterogeneity, and technical artifacts.

308     From the 174 GIANT height GWAS (supplied data files)[3], we calculated 15,051 cohort-pairwise $\lambda_{meta}$

309     values, resulting in a bell-shape distribution (**Fig. 6a,b**) with the mean of 1.013 and the empirical standard

310     deviation (S.D.) of 0.022, which was greater than theoretical S.D. of 0.014. The empirical mean and S.D can

311     be used to construct a z-score test for each $\lambda_{meta}$. These results are consistent with a small amount of

312     heterogeneity, which is not unexpected due to variation of actual (unknown) genetic architecture and

313     analysis protocols. However, the mean is close to 1.0 and based upon this QC metric the results are

314     consistent with stringent quality control and data cleaning. The minimum $\lambda_{meta}$ value was around 0.88

315     (between SORBS MEN and SORBS WOMEN, **Fig. 3c**), with $p$-value < 1e-10 (testing for the difference

316     from 1), and the maximum was 1.245 (between SardiNIA and WGHS, **Fig. 6d**), with $p$-value < 1e-10,

317     leading to the most deflated and inflated $\lambda_{meta}$ across GIANT height study cohorts; both were significant

318 after correction for multiple testing. Illustrating $\lambda_{meta}$ (**Fig. 6b**) highlighted that 20 cohorts from the MIGEN

319 consortium showed substantially lower $\lambda_{meta}$ with many other cohorts (right-bottom triangle in **Fig. 6b**)

320 than the average, consistent with over-conservative models for statistical association analyses being used in

321 these cohorts – which may be due to very small sample size (ranging from 36 to 320 for the 20 MIGEN

322 cohorts, with an average sample size of 132). Consistent with this, cohorts from MIGEN also have many of

323 their $\lambda_{GC} < 1$ (**Fig. 7a**). In contrast, the SardiNIA cohort (4,303 samples) showed heterogeneity with nearly

324 all other cohorts (**Fig. 7b**), perhaps due to unknown artifacts or a slightly different genetic architecture for

325 height as result of demographic history[23].

326

327 We investigated the relationship between $\bar{\lambda}_{meta}$ (the mean of all $\lambda_{meta}$ values of a given cohort with each of

328 the other 173 GIANT height cohorts) and $\lambda_{GC}$ among the GIANT height cohorts. If there are no technical

329 issues, such as inflated or deflated sampling variance for the estimated effects, we would expect to see: i) a

330 correlation between $\lambda_{GC}$ and sample size; ii) no correlation between $\bar{\lambda}_{meta}$ and sample size; iii) no

331 correlation between $\bar{\lambda}_{meta}$ and $\lambda_{GC}$ (**Supplementary Fig. S7**). Consistent with a previous study[24], for a

332 polygenic trait such as height $\lambda_{GC}$ of each cohort was related to its sample size (correlation of 0.235, $p =$

333 0.0018). In contrast, the correlation between $\bar{\lambda}_{meta}$ and sample size was of 0.116 ($p = 0.127$) (**Fig 7a,b**).

334 Nevertheless, the correlation between the mean of $\bar{\lambda}_{meta}$ and $\lambda_{GC}$ was 0.836 ($p<10e-16$) for 174 GIANT

335 height cohorts (**Fig 7c**). We note that the 20 MIGEN cohorts had proportionally small $\lambda_{GC}$ and $\bar{\lambda}_{meta}$, with

336 very high correlation between them ($\rho = 0.98$); in contrast, the SardiNIA cohort, which had the largest $\lambda_{GC}$,

337 showed the largest $\bar{\lambda}_{meta}$ ($1.070 \pm 0.049$), standing out as a special case among the GIANT height cohorts.

338 Assuming a polygenic model of $h^2 = 0.5$ over 30,000 independent loci, we simulated 174 cohorts using the

339 actual size samples from the GIANT height cohorts (**Supplementary notes**), and observed an increased

340 correlation ($R^2 = 0.78$) between $\bar{\lambda}_{meta}$ and $\lambda_{GC}$ for simulated cohorts with sample sizes of the MIGEN

341 cohorts (**Fig. 7d**). Other effects, such as inflated/deflated sampling variance of the estimated genetic effects

342 could also lead to correlation between $\bar{\lambda}_{meta}$ and $\lambda_{GC}$ (**Supplementary Fig. S8**). In addition, we constructed

343 a single MIGEN analysis by combining the 20 MIGEN cohorts using an inverse variance weighted meta-

344 analysis[25], and calculated $\lambda_{meta}$ between this combined MIGEN cohort and all 174 cohorts. As expected, the

345 combined MIGEN had $\lambda_{meta} = 0.90 \pm 0.07$ with 20 MIGEN cohorts due to overlapping samples. In

346 contrast, $\lambda_{meta} = 1.01 \pm 0.02$ with 154 other cohorts, was consistent with neither heterogeneity nor sample

347 overlap. Given that the MIGEN (2,340 samples) and SardiNIA (4,303 samples) cohorts contributed less than

348 3% of the total sample size (253,288 samples from the GIANT height GWAS cohorts), any impact of

349 unusual $\lambda_{meta}$ values on the meta-analysis results is very small. Given no heterogeneity between a pair of

350 cohorts, a deflated $\lambda_{meta}$ reflects the effective number of overlapping samples (**Supplementary notes**). For

351 example, the "combined MIGEN" had $\lambda_{meta}$ values proportional to the sample size of each MIGEN cohort

352 (**Fig 7e**).

353

354    The statistical power of detection of overlapping samples is maximized when a pair of cohorts has equal

355    sample size (**Fig. 8a**), or in other words the confidence interval for null hypothesis of no overlapping

356    samples depends on the sample sizes for a pair of cohorts. As a comparison, direct correlation that is

357    estimated between the genetic effects for a pair of cohorts has been proposed to estimate overlapping

358    samples[26,27], but it is confounded with genetic architecture, such as heritability underlying (**Table 1**). When

359    there was heritability, the estimated correlation between genetic effects was biased and leads to incorrect

360    overlapping samples for a pair of cohorts; when there was no heritability, the estimated correlation was

361    correct and agreed well with the one estimated with $\lambda_{meta}$. As existence of heritability is one of the reasons

362    that trigger GWAMA, so $\lambda_{meta}$ is much proper in estimating overlapping samples between cohorts.

363

364    Another parameterization of $\lambda_{meta}$ is to estimate it from differences in allele frequencies between a pair of

365    cohorts instead of differences between estimated effect sizes (**Supplementary notes**). We show that $\lambda_{meta}$

366    constructed on reported allele frequencies from genotyped loci from summary statistics can detect

367    overlapping samples between two cohorts regardless of whether the GWAS is from quantitative traits or

368    case-control data, even for pairs of different traits (**Supplementary notes and Supplementary notes**). For

369    example, 2,934 common controls were shared across the WTCCC 7 diseases[7]. From the 21 pairwise $\lambda_{meta}$,

370    we estimated a mean of the number of overlapping samples, assuming overlapping controls only

371    (**Supplementary notes**), of $\hat{n}_o =$ 2,708 (S.D. = 58.4), which was very close estimate to the actual number of

372    overlapping samples (**Supplementary Fig. 8**). When constructing $\lambda_{meta}$ on the reported genetic effects and

373    their sampling variance, the estimated mean estimate of the number of shared controls was $\hat{n}_o =$ 2,127 (S.D.

374    = 257.7), lower than that estimated from allele frequencies, which is likely due to real genetic heterogeneity

375    between diseases (**Supplementary Fig. 8**). In practice, publically available summary statistics may not

376    include sample specific allele frequencies, but may only be available with reference sample frequencies as a

377    conservative strategy to prevent identification of individuals in a cohort.

378

379    **Detection of overlapping samples using pseudo profile score regression**

380    GWAMAs have grown in sample size and in the number of cohorts that participate, and this trend is likely to

381    continue. The probability that a sample is represented in more than one meta-analysis study is also likely to

382    increase, in particular when very large cohorts such as UK Biobank and 23andMe provide data to multiple

383    studies. While the metric $\lambda_{meta}$ can be transformed to give an estimate of $n_o$ between cohorts for

384    quantitative traits, it cannot give an estimate of overlapping samples in case-control studies due to the ratio

385    of the cases and controls in each study (**Supplementary notes**). Sharing individual genotype data (or

386    imputed genotypes) across the entire study would make it easy to detect identical or near-identical genotype

387    samples (representing real duplicate samples from individuals who participated in the two studies or

388    monozygotic twins). In fact, only a small number of common SNPs is needed to detect sample overlap, and

389    if this is known then individuals could be removed and summary statistics regenerated or the meta-analysis

390    analysis itself can be adapted to correct for potential correlation due to $n_o$[28]. However, in many

391    circumstances, individual cohorts are not permitted to share individual-level data, either by national law or

392    by local ethical review board conditions. To get around this problem, Turchin and Hirshhorn[10] created a

393    software tool, Gencrypt, which utilizes a security protocol known as one-way cryptographic hashes to allow

394    overlapping participants to be identified without sharing individual-level data. To our knowledge, this

395    encryption method has yet to be employed in meta-analysis studies. We propose an alternative approach,

396    pseudo profile score regression (PPSR), which involves sharing of weighted linear combinations of SNP

397    genotypes with the central meta-analysis hub. In essence, multiple random profile scores are generated for

398    each individual in each cohort, using SNP weights supplied by the analysis hub, and the resulting scores are

399    provided back to the analysis hub. PPSR works through three steps (**Supplementary notes and**

400    **Supplementary Fig. 9**), and the purpose of PPSR is to estimate a relationship-like matrix of $n_i \times n_j$

401    dimension for a pair of cohorts, which have $n_i$ and $n_j$ individuals respectively. Each entry of the matrix is

402    filled with genetic similarity for a pair of samples from each of the two cohorts, estimated via the PPSR.

403

404    We use WTCCC data as an illustration to detect 2,934 shared controls between any two of the diseases by

405    PPSR. Among 330K unambiguous SNPs, which are not palindromic (A/T or G/C alleles), we randomly

406    picked $M = 100$, 200, and 500 SNPs, to generate pseudo profile scores. It generated 21 cohort-pair

407    comparisons, leading to the summation for 488,587,090 total individual-pair tests. To have an experiment-

408    wise type I error rate = 0.01, type II error rate = 0.05 (power = 0.95) for detecting overlapping individuals,

409    we needed to generated at least 57 pseudo profile scores (PPS). We generated scores $S = [s_1, s_2, s_3, ..., s_{57}]$,

410    where each $s$ is a vector of $M$ elements, sampled from a standard normal distribution (**Supplementary**

411    **notes**). $S$ is shared across 7 cohorts for generating pseudo-profile scores for each individual. In total 57 PPS

412    were generated for each individual in each cohort. For a pair of cohorts, PPSR was conducted for each

413    possible pair of individuals for any two cohorts over the generated pseudo-profile scores. Once the

414    regression coefficient ($b$) was greater than the threshold, here $b = 0.95$, the pair of individuals was inferred

415    to be having highly similar genotypes, implying that the individual was included in both cohorts

416    (**Supplementary notes**).

417

418    When using 200 and 500 random SNPs, all the known 2,934 shared controls were detected from 21 cohort-

419    pair-wise comparison; when using 100 randomly SNPs, on average 2,931 shared samples were identified,

420    which is more accurate than using $\lambda_{meta}$ constructed using either genetic effects or allele frequencies (**Fig.**

421    **8b**). In addition, for detected overlapping samples, there were no false positives observed – consistent with

422    simulations that show the method was conservative in the controlling type I error rate (**Supplementary**

423    **notes**). For comparison, we also used the Gencrypt to detect overlapping samples using the same set of

424    SNPs as used in PPSR. Although Gencrypt guidelines suggest use of at least 20,000 random SNPs[10],

425    selecting 500 random SNPs in the WTCCC cohorts also provided good accuracy with Gencrypt, and on

426     average about 2,920 (99.6% of the shared controls) overlapping samples were detected, only slightly lower

427     than PPSR. For example, for BP and CAD, Gencrypt detected 2,912 shared controls, but was unable to

428     identify about 20 overlapping controls, due to missing data (on average 1% missing rate). Increasing the

429     number of SNPs when using Gencrypt is likely to overcome the problem of missing data.

430

431     Furthermore, PPSR is able to detect pairs of relatives. For example, between the BD and CAD cohorts, two

432     pairs of apparent first-degree relatives were detected (**Fig. 9a**). In order to find additional first-degree

433     relatives between BD and CAD cohorts, at least 265 PPS were required to have a type I error rate of 0.01

434     and type II error rate of 0.05 (**Supplementary notes**) for a regression coefficient cutoff of 0.45, a threshold

435     for first-degree relatives. As expected, all other individuals that did not show high relatedness did not reach

436     the threshold of 0.45 of the PPS regression coefficient for first-degree relatives (**Fig. 9b**). Gencrypt did not

437     detect any first-degree relatives.

438

439     The speed of PPSR depends on $n_i \times n_j$, the sample sizes for a pair of cohorts, and the number of PPS for

440     each cohort; for the WTCCC data there are 21 cohort-pair comparisons, and each pair took about 20 minutes,

441     on a computer with a 2.3 GHz CPU, given about $5,000 \times 5,000 = 25,000,000$ comparisons. The average

442     sample size of GIANT is about 1,500, and takes about 2 minutes for each pair of cohorts. The two largest

443     datasets are deCODE with 26,790 samples and WGHS with 23,100 samples, and PPSR to detect overlapping

444     samples takes about 8.5 hours. As each pair of individuals is computationally an independent unit, analysis

445     jobs can be parallelized on a cluster. Therefore, even for meta-analyses involving many large cohorts, the

446     computation time is not a limiting factor.

447

448     PPSR for each individual uses very little personal information and can be minimized so that there is very

449     low probability of decoding it. One way to attempt to decode the genotypes from PPS is to reverse the PPSR,

450     so that the individual genotypes can be predicted in the regression (**Supplementary notes**). The individual-

451     level genotypic information that can be recovered by an analyst, who knows the $S$ matrix (the weights for

452     generating PPS), is determined by the ratio between the number of markers ($M$) that generated PPS and the

453     number of PPS ($K$). Therefore, inferred information on individual genotypes can be minimized and tailored

454     to any specific ethics requirements. We suggest $\frac{M}{K} > 5 \sim 10$ to protect the privacy with sufficient accuracy

455     (**Fig 9c**). Of note, if a meta-analysis is conducted within a research consortium, the application of PPSR is

456     even safer because the exchange of information is between the consortium analysis hub and each cohort

457     independently.

458

459                                           **Discussion**

460     In this study, we provide a set of metrics for monitoring and improving the quality of large-scale GWAMA

461     based on summary statistics. These tools not only enrich the toolkit to analysts for GWAMA, but also

162    provide informative summary and visualization for readers to understand the experimental design of

163    GWAMA. As far as we know, no GWAMA to date has checked cohort-level outliers based upon population

164    differentiation metrics or utilized estimated allelic effect sizes to identify and quantify sample overlap.

165

166    Using the $F_{st}$ derived genetic distance measure, we can place all cohorts on an inferred geographic map and

167    can easily identify cohorts that are genetic outliers or that have unexpected ancestry. In application, we

168    should note that the $F_{st}$ measure can identify unusual summary information, such as detected in the MIGEN

169    cohorts from GIANT Consortium GWAMAs, in which the same allele frequencies were reported for all

170    cohorts. Meta-PCA can also be used to infer the genetic background of cohorts. The high concordance

171    between $F_{pc}$ and meta-PCA indicates the both methods are robust. In practice, mete-PCA may be much

172    easier to implement when there are many cohorts, such as GIANT height cohorts and Metabochip BMI

173    cohorts, but the coordinates of a cohort may be slightly shifted with inclusion or exclusion of other cohorts.

174

175    There are limitation for both $F_{pc}$ and meta-PCA. Firstly, the inference depends on the choice of reference

176    cohorts. Meta-PCA is further upon the inclusion or exclusion of other cohorts. However, given the

177    application of the data, we believe the impact will not influence the inference of the genetic background of

178    cohorts in meta-analysis. Secondly, various mechanisms can give the identical projection in PCA[14]. The

179    purpose of both methods is to find the discordance between demographic information and genetic

180    information, or outliers. The projection is not attempt to discover the detailed demographic past that shapes a

181    cohort.

182

183    Our third metric $\lambda_{meta}$ provides information on sample overlap and heterogeneity between cohorts by

184    utilizing the estimated allelic effect sizes and their standard errors. In most meta-analyses, the overall $\lambda_{meta}$

185    is likely to be slightly greater than 1 solely due to unknown heterogeneity, slight as observed, in generating

186    the phenotype and genotype data that cannot be accounted for by QC. The observed mean of $\lambda_{meta}$ for the

187    GIANT height GWAMA was 1.03 but with more variation than expected by chance. The strong correlation

188    between $\lambda_{GC}$ and $\lambda_{meta}$ indicated the reported sampling of the reported data were systematically driven by

189    analysis protocols. For cohorts with $\lambda_{GC} < 1$ and $\lambda_{meta} < 1$, it is likely that the GWAS modeling strategy

190    employed for GWAS in the cohort was too conservative, for example MIGEN cohorts might have on

191    average too small sample size for each cohort. Conversely, for cohorts with $\lambda_{GC} > 1$ and $\lambda_{meta} > 1$ results

192    are too heterogeneous, perhaps reflecting systematically smaller sampling variances of the reported genetic

193    effects. As the GWAMA often uses inverse-variance-weighted meta-analysis[25], such cohorts may lead to

194    incorrect weights to the different cohorts in the meta-analysis, suggesting that the statistical analysis in meta-

195    analyses can be improved by applying better weighting factors.

196

14

497     It is well-recognised that overlapping samples may inflate the type-I error rate of GWAMA and therefore

498     lead to false positives. Although post-hoc correction of the test statistic is possible[26–28], stringent quality

499     control ruling out overlapping samples makes the whole analysis easier and lowers the risk of false positives.

500     A better solution would be to rule out shared samples at the start, for pairs of cohorts that show deflated

501     $\lambda_{meta}$, and we propose PPSR to accomplish this.

502

503     In summary, to maximize the inference from multi-cohort GWAMA, accurate cohort-level information on

504     allele frequencies, estimated effect sizes, and their sampling variance can be exploited to perform additional

505     measures that are likely to lead to reduction in the number of false positives and increasing statistical power

506     for gene discovery. All methods proposed are implemented in freely available software GEAR.

507

516

517     **Author contributions:**

518     GBC and PMV designed the study. GBC, PMV and SHL derived the analytical results. GBC performed all

519     analysis. CGB and ZXZ developed the software. GBC and PMV wrote the first draft of the paper. MRR, JY,

520     NW discussed results and methods, and provided comments that improved earlier versions of the manuscript.

521     Other authors provided cohort-level summary statistics and contributed to improving the study and

522     manuscript.

523

524     **Competing financial interests:**

525     The authors declare no completing financial interests.

526

527     **Web resources:**

528     GEAR (GEnetic Analysis Repository): http://www.complextraitgenomics.com/

529     PGC: http://www.med.unc.edu/pgc/results

530     1000 Genomes Project: http://www.1000genomes.org/

531

## References

632    1    Visscher PM, Brown M a, McCarthy MI, Yang J. Five years of GWAS discovery. *Am J Hum Genet*
633         2012; **90**: 7–24.
634    2    Winkler TW, Day FR, Croteau-Chonka DC, Wood AR, Locke AE, Mägi R *et al.* Quality control and
635         conduct of genome-wide association meta-analyses. *Nat Protoc* 2014; **9**: 1192–212.
637    3    Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S *et al.* Defining the role of common
638         variation in the genomic and biological architecture of adult human height. *Nat Genet* 2014; **46**:
639         1173–1186.
640    4    Locke AE, Kahali B, Berndt SI, Justice AE, Pers TH, Day FR *et al.* Genetic studies of body mass
641         index yield new insights for obesity biology. *Nature* 2015; **518**: 197–206.
642    5    Voight BF, Kang HM, Ding J, Palmer CD, Sidore C, Chines PS *et al.* The Metabochip, a custom
643         genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits.
644         *PLOS Genet* 2012; **8**: e1002793.
645    6    The 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092
646         human genomes. *Nature* 2012; **491**: 56–65.
647    7    The Wellcome Trust Case Control Consortium. Genome-wide association study of 14,000 cases
648         of seven common diseases and 3,000 shared controls. *Nature* 2007; **447**: 661–678.
649    8    Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira M a R, Bender D *et al.* PLINK: a tool set for
650         whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007; **81**:
651         559–75.
652    9    Devlin B, Roeder K. Genomic control for association studies. *Biometrics* 1999; **55**: 997–1004.
653    10   Turchin MC, Hirschhorn JN. Gencrypt: one-way cryptographic hashes to detect overlapping
654         individuals across samples. *Bioinformatics* 2012; **28**: 886–8.
655    11   Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A *et al.* Genes mirror geography
656         within Europe. *Nature* 2008; **456**: 98–101.
657    12   Cavalli-Sforza LL, Menozzi P, Piazza A. *The History and Geography of Human Genes*. Princeton
658         University Press, 1996.
659    13   Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLOS Genet* 2006; **2**:
660         e190.
661    14   McVean G. A genealogical interpretation of principal components analysis. *PLOS Genet* 2009; **5**:
662         e1000686.
663    15   Bryc K, Bryc W, Silverstein JW. Separation of the largest eigenvalues in eigenanalysis of
664         genotype data from discrete subpopulations. *Theor Popul Biol* 2013; **89**: 34–43.
665    16   Chaput J-P, Pérusse L, Després J-P, Tremblay A, Bouchard C. Findings from the Quebec Family
666         Study on the Etiology of Obesity: Genetics and Environmental Highlights. *Curr Obes Rep* 2014; **3**:
667         54–66.
668    17   Diabetes Genetics Initiatives. Genome-wide association analysis identifies loci for type 2
669         diabetes and triglyceride levels. *Science (80- )* 2007; **316**: 1331–1336.
670    18   Igl W, Johansson A, Gyllensten U. The Northern Swedish Population Health Study (NSPHS)--a
671         paradigmatic study in a rural population combining community health and basic research. *Rural
672         Remote Health* 2010; **11**: 1363.
673    19   de Bakker PIW, Ferreira M a R, Jia X, Neale BM, Raychaudhuri S, Voight BF. Practical aspects of
674         imputation-driven meta-analysis of genome-wide association studies. *Hum Mol Genet* 2008; **17**:
675         R122–8.
676    20   Ripke S, Sanders AR, Kendler KS, Levinson DF, Sklar P, Holmans P a *et al.* Genome-wide
677         association study identifies five new schizophrenia loci. *Nat Genet* 2011; **43**: 969–76.
678    21   Ripke S, O'Dushlaine C, Chambert K, Moran JL, Kähler AK, Akterin S *et al.* Genome-wide
679         association analysis identifies 13 new risk loci for schizophrenia. *Nat Genet* 2013; **45**: 1150–9.
680    22   Okada Y, Wu D, Trynka G, Raj T, Terao C, Ikari K *et al.* Genetics of rheumatoid arthritis
681         contributes to biology and drug discovery. *Nature* 2014; **506**: 376–381.
682    23   Calò C, Melis A, Vona G, Piras I. Sardinian Population (Italy): a Genetic Review. *Int J Mod
683         Anthropol* 2010; **1**: 39–64.
684

585 24 Yang J, Weedon MN, Purcell S, Lettre G, Estrada K, Willer CJ *et al.* Genomic inflation factors
586 under polygenic inheritance. *Eur J Hum Genet* 2011; **19**: 807–12.
587 25 Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association
588 scans. *Bioinformatics* 2010; **26**: 2190–1.
589 26 Bolormaa S, Pryce JE, Reverter A, Zhang Y, Barendse W, Kemper K *et al.* A multi-trait, meta-
590 analysis for detecting pleiotropic polymorphisms for stature, fatness and reproduction in beef
591 cattle. *PLOS Genet* 2014; **10**: e1004198.
592 27 Zhu X, Feng T, Tayo BO, Liang J, Young JH, Franceschini N *et al.* Meta-analysis of Correlated
593 Traits via Summary Statistics from GWASs with an Application in Hypertension. *Am J Hum Genet*
594 2015; **96**: 21–36.
595 28 Lin D-Y, Sullivan PF. Meta-analysis of genome-wide association studies with overlapping
596 subjects. *Am J Hum Genet* 2009; **85**: 862–72.
597

598

**Figures & Tables**

**Figure 1 Recovery of cohort-level genetic background and inference of their geographic locations for GIANT BMI Metabochip cohorts using the $F_{st}$ derived genetic distance measure.**

**Figure 2 Using the genetic distance spectrum to infer the geographic origins for GIANT height GWAS cohorts.**

**Figure 3 Comparison between Meta-PCA and genotype PCA on 1KG.**

**Figure 4 Recovery of cohort-level genetic background for GIANT BMI Metabochip cohorts using meta-PCA.**

**Figure 5 The recovery of cohort-level genetic background using meta-PCA analysis for GWAS height cohorts.**

**Figure 6 $\lambda_{meta}$ for the GIANT height GWAS cohorts.**

**Figure 7 $\bar{\lambda}_{meta}$ and $\lambda_{gc}$ for GIANT height GWAS cohorts.**

**Figure 8 Pseudo profile score regression for the WTCCC 7 diseases.**

**Figure 9 PPSR coefficients for identifying shared controls/relatives between WTCCC BD and CAD cohorts.**

**Table 1 The estimated correlation for a pair of cohorts via their summary statistics**

**Figure 1 Recovery of cohort-level genetic background and inference of their geographic locations for GIANT BMI Metabochip cohorts using the $F_{st}$ derived genetic distance measure.** (**a**) Genetic distance spectrum for all Metabochip cohorts to CEU, CHB and YRI. See Supplementary notes for more details. The origins of the cohorts are denoted on the horizontal axis. (**b**) Projection for Metabochip cohort into $F_{PC}$ space defined by YRI, CHB, and CEU reference populations. The x- and y-axis represent relative distances derived from the genetic distance spectrum. Three dashed lines, blue for CEU, green for CHB, and red for YRI, partitioned the whole $F_{PC}$ space to three genealogical subspaces. (**c**) The genetic distance spectrum for Metabochip European cohorts to CEU – Northwest Europeans, FIN – Northeast European, and TSI – Southern Europeans. The nationality of the cohorts are denoted on the horizontal axis. (**d**) The projection for Metabochip European cohorts to the $F_{PC}$ space defined by CEU, FIN, and TSI reference populations. The whole space is further partitioned into three subspaces, CEU-TSI genealogical subspace (red and blue dashed lines), FIN-TSI genealogical subspace (green-blue dashed lines), and CEU-FIN genealogical subspace (red-green dashed lines), respectively. The open circles represent the mean of inferred geographic locations for the cohorts from the same country. Cohort/country codes: AF, African; AU, Australia; DE, Germany; EE, Estonia; EU, European Nations; FI, Finland; FIN, Fins in 1000 Genomes Project (1KG); FR, France; GBR, British in 1KG; GIB, Gujarati Indian in 1KG; GR, Greece; Hawaii, Hawaii in USA; IBS, Iberian Population in Spain in 1KG; IT, Italy; JM, Jamaica; JPT, Japanese in 1KG; LWK, Luhya in 1KG; NO, Norway; PH, the Philippines; PK, Pakistan; SC, Seychelles; SCT, Scotland; SE, Sweden; TSI, Tuscany in 1KG; UK, United Kingdom; US, United States of America.
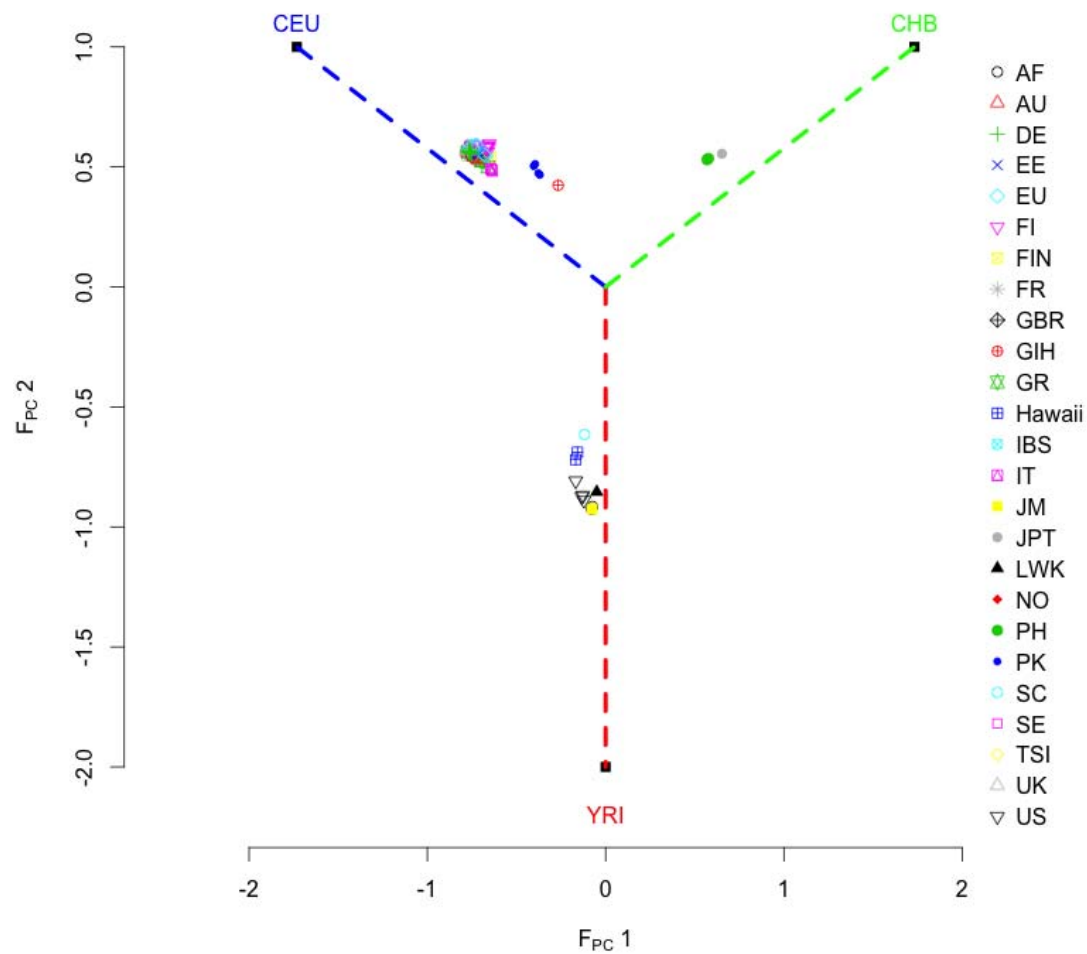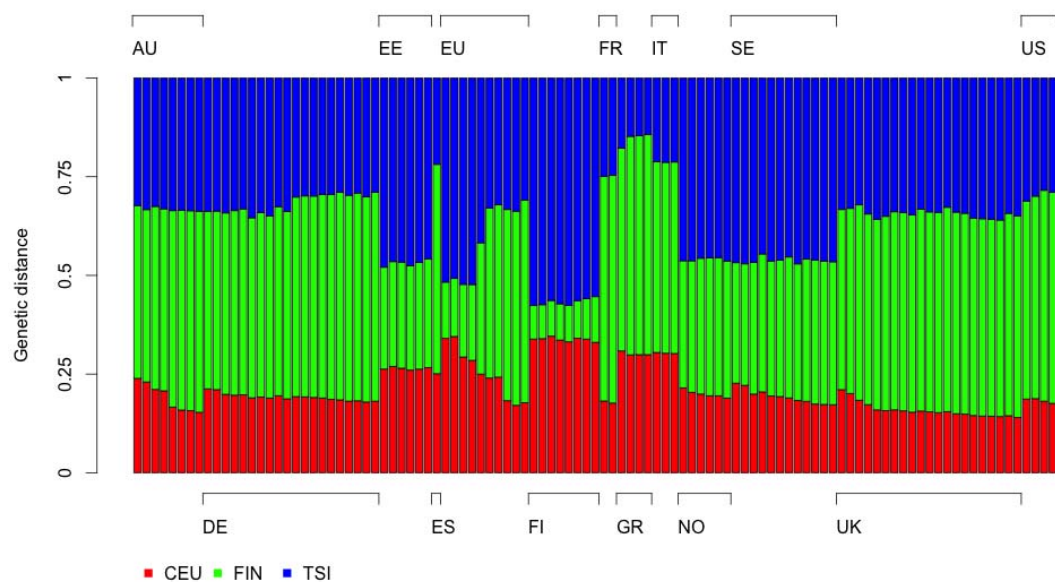
19

546    a)



547

548

549    b)



550

551

652    c)
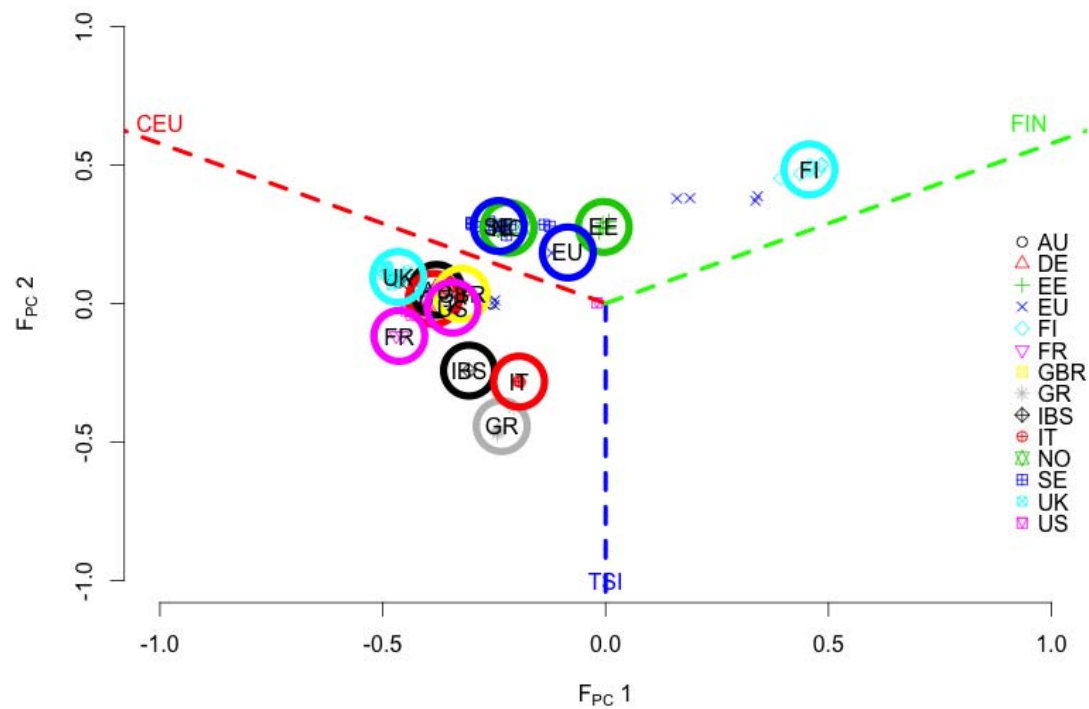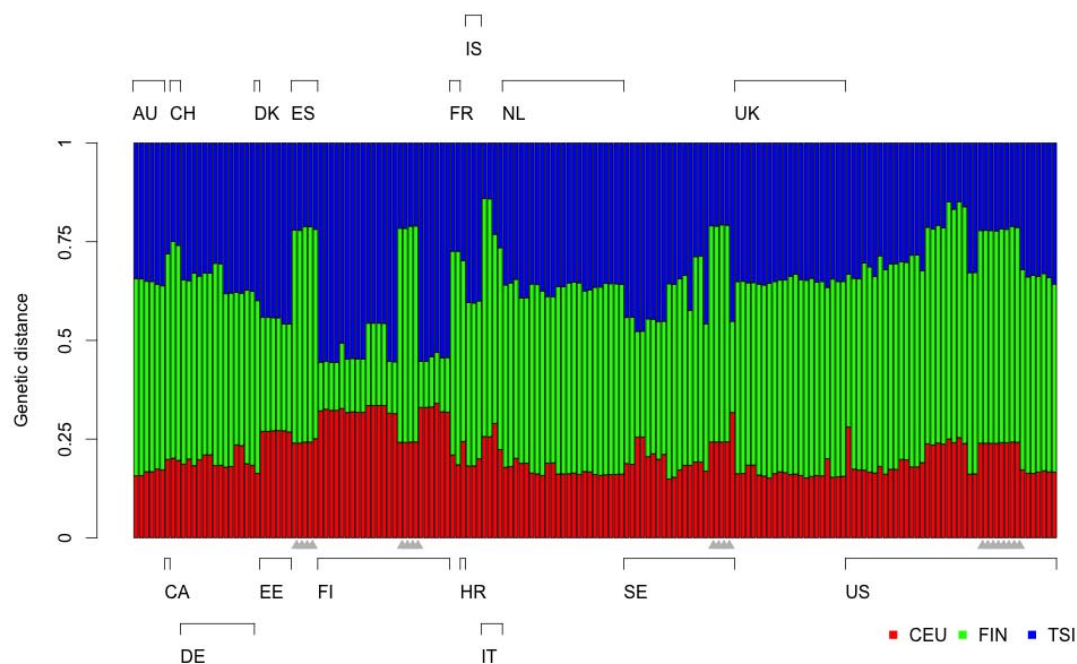


653

654

655    d)



656

657

558 **Figure 2 Using the genetic distance spectrum to infer the geographic origins for GIANT height GWAS**

559 **cohorts. (a)** Each cohort has three $F_{st}$ values by comparing with CEU, FIN, and TSI reference samples. The

560 height of each bar represents its relative genetic distance to these three reference populations. The

561 nationalities of the cohorts were denoted along the horizontal axis. The grey triangles along the x-axis

562 indicate MIGEN cohorts. **(b)** Given the three $F_{st}$ values, the location of each cohort can be mapped. The

563 whole space was partitioned into three subspaces, CEU-TSI genealogical subspace (red and blue dashed

564 lines), FIN-TSI genealogical subspace (green and blue dashed lines), and CEU-FIN genealogical subspace

565 (red and green dashed lines). DGI (in the blue box) had samples from the Botnia study. Across MIGEN

566 cohorts (denoted as red triangles in the red box), the same allele frequencies (likely calculated from a South

567 European cohort) were presented for each cohort. Cohort/country codes: AU, Australia; CA, Canada; CH,

568 Switzerland; DE, Germany; DK, Denmark; EE, Estonia; ES, Iberian Population in Spain in 1KG; FI,

569 Finland; FR, France; GR, Greece; IT, Italy; IS, Iceland; NL, Netherlands; SE, Sweden; UK, United

570 Kingdom; US, United States of America.

571

24

572   a)



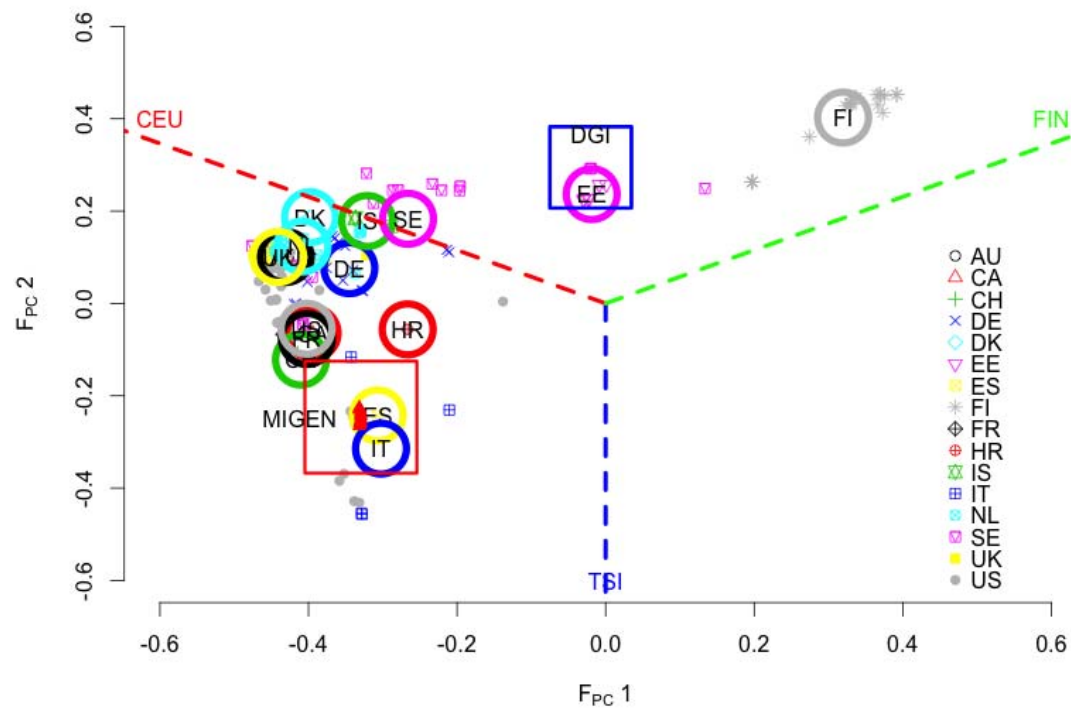573

574

575    b)



576

577

**Figure 3 Comparison between Meta-PCA and genotype PCA on 1KG.** Top left panel is the projection of cohorts based on cohort-level allele frequency for 1KG samples on the first two eigenvectors. Bottom left panel is conventional PCA based on individual genotypes on the first two eigenvectors. Top right panel is the projection by taking the mean of the 1KG individuals within each cohort. Bottom right panel is the correlation, measured in , between meta-PCA and genotype PCA for the first twenty eigenvectors.

**Figure 4 Recovery of cohort-level genetic background for GIANT BMI Metabochip cohorts using meta-PCA.** The x-axis and y-axis represent the first two eigenvectors from meta-PCA. In meta-PCA, Metabochip cohorts could be classified into African ancestry (AFR), European ancestry (EAS), East Asian Ancestry (EAS), and South Asian Ancestry (SAS). The 1KG cohorts, yellow open circles, were added for comparison.
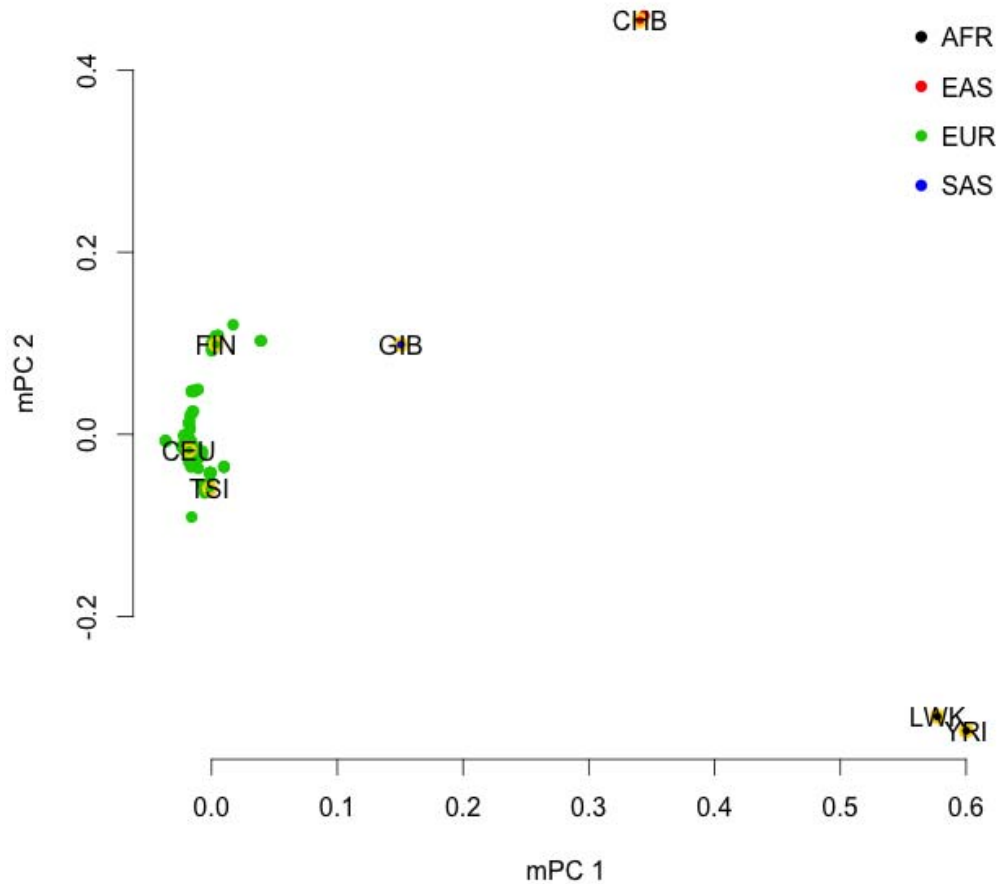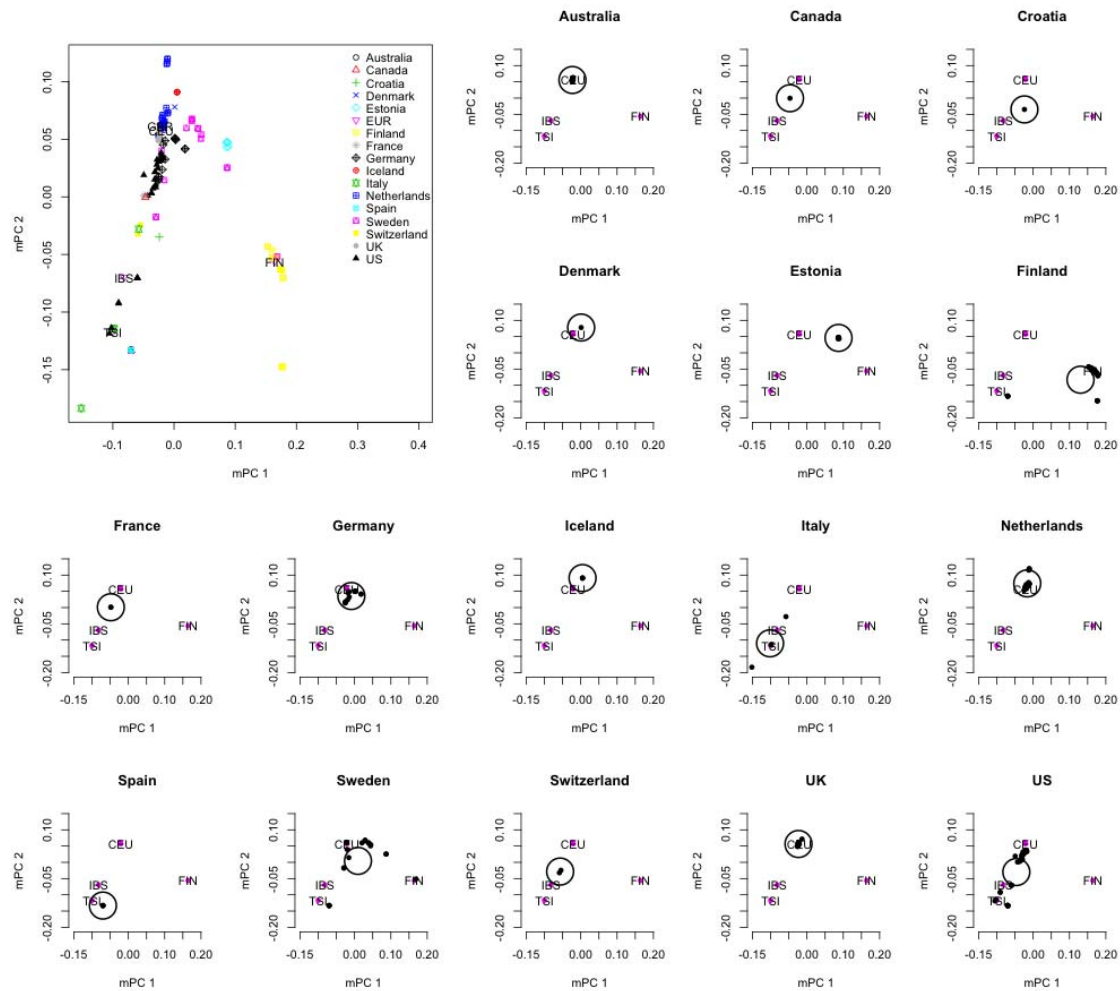
593 **Figure 5 The recovery of cohort-level genetic background using meta-PCA analysis for GWAS height**
594 **cohorts.** The x-axis and y-axis represent the firs two eigenvectors inferred from meta-PCA. a) The genetic
595 background inferred with the inclusion of 10 1KG reference populations. b) The genetic background and
596 relative geographic location for 174 GIANT height cohorts. The large plot on top left was an overview of
597 174 cohorts, and the rest of plots were classified by the reported demographic information of cohorts. Within
598 each country-level plot, the small black points represent one cohort, and the large open circle the mean
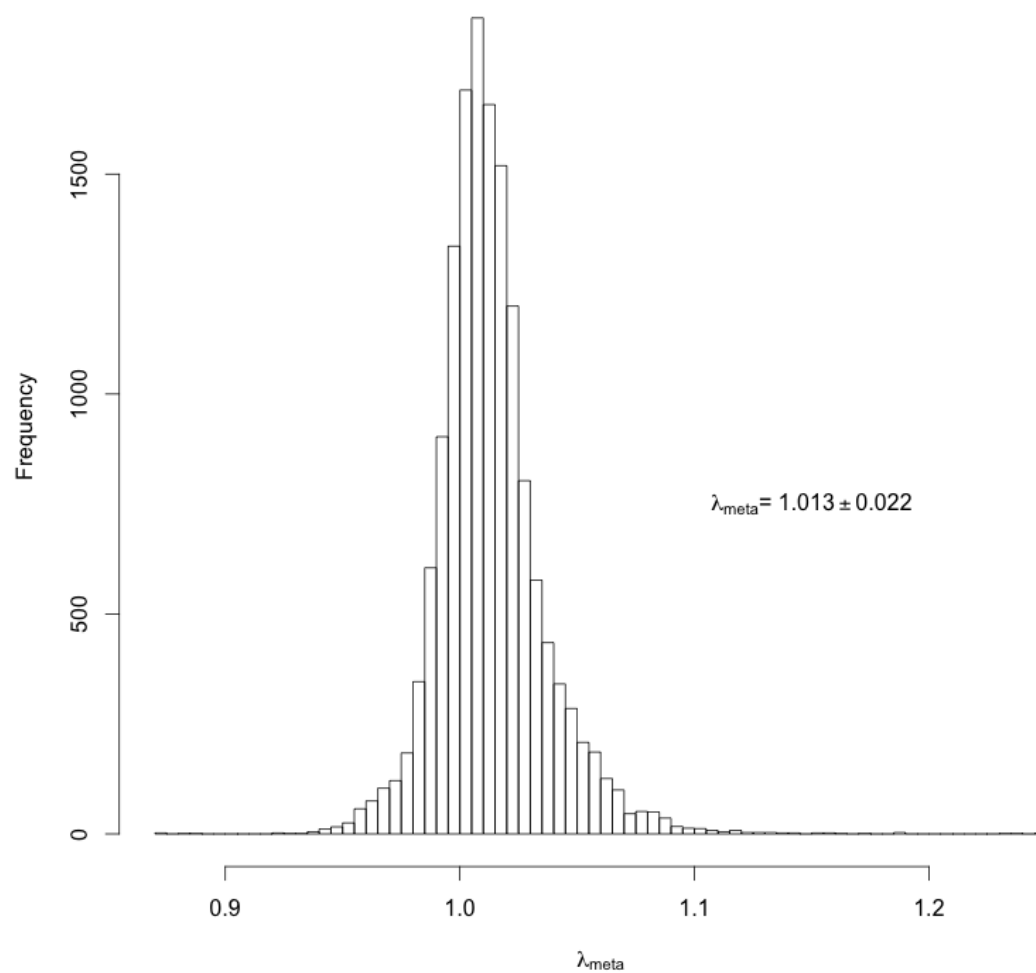599 coordinates for those cohorts from the same country.
600
601 **a**



602
603
604 **b**
605

'06
'07
'08

'09 **Figure 6 $\lambda_{meta}$ for the GIANT height GWAS cohorts.** Given 174 cohorts, there are 15,051 $\lambda_{meta}$ values,

'10 which provide the overview of the quality control of the summary statistics.

'11 **(a)** The distribution of $\lambda_{meta}$ from 174 cohorts/files used in the GIANT height meta-analysis. The overall

'12 mean of 15,051 $\lambda_{meta}$ is 1.013, and standard deviation is 0.022. **(b)** The heat map for $\lambda_{meta}$. Cohorts

'13 showed heterogeneity ($\lambda_{meta} > 1$) are illustrated on left-top triangle, and homogeneity ($\lambda_{meta} < 1$) on right-

'14 bottom triangle. **(c)** Illustration for homogeneity between two cohorts (SORBS MEN & WOMEN), $\lambda_{meta} =$

'15 0.876. **(d)** Illustration of SARDINIA & WGHS, this pair of cohorts has $\lambda_{meta} = 1.245$. The grey band

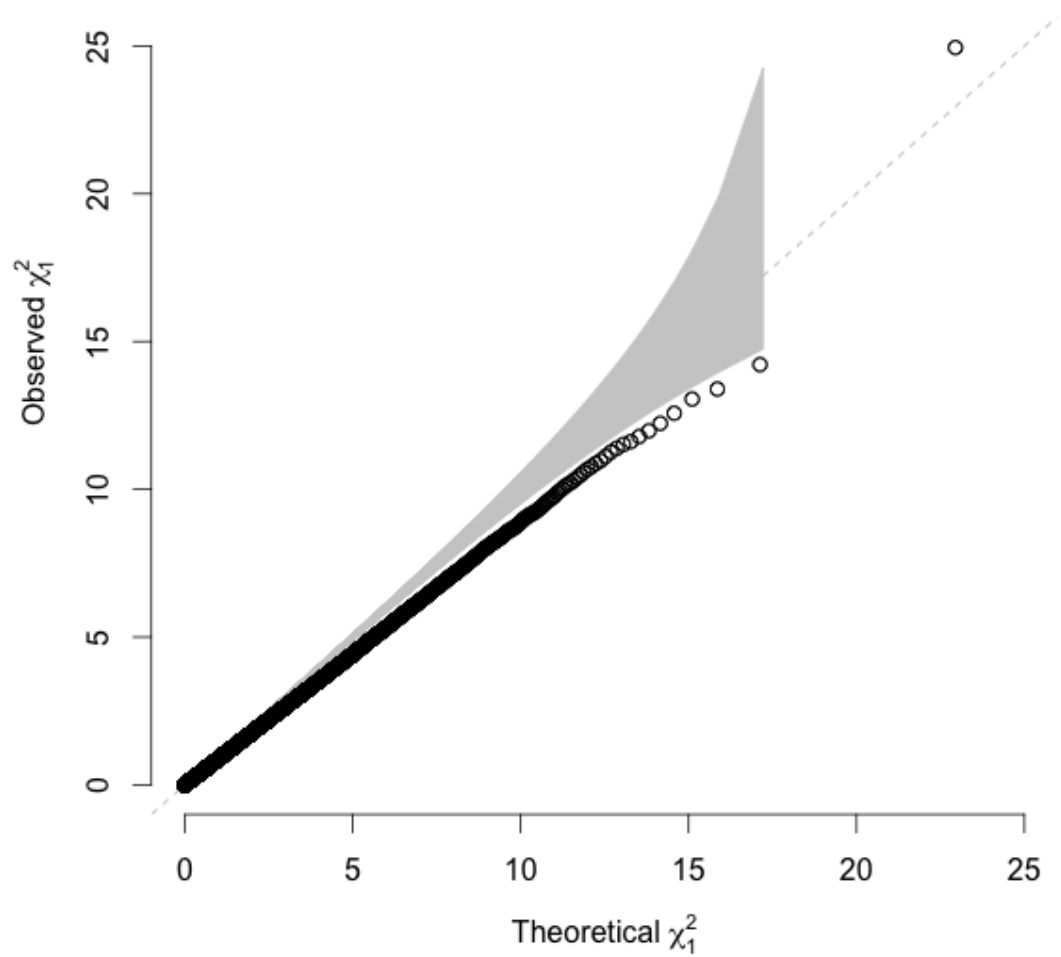'16 represents 95% confidence interval for $\lambda_{meta}$.

'17

31

'18    a)
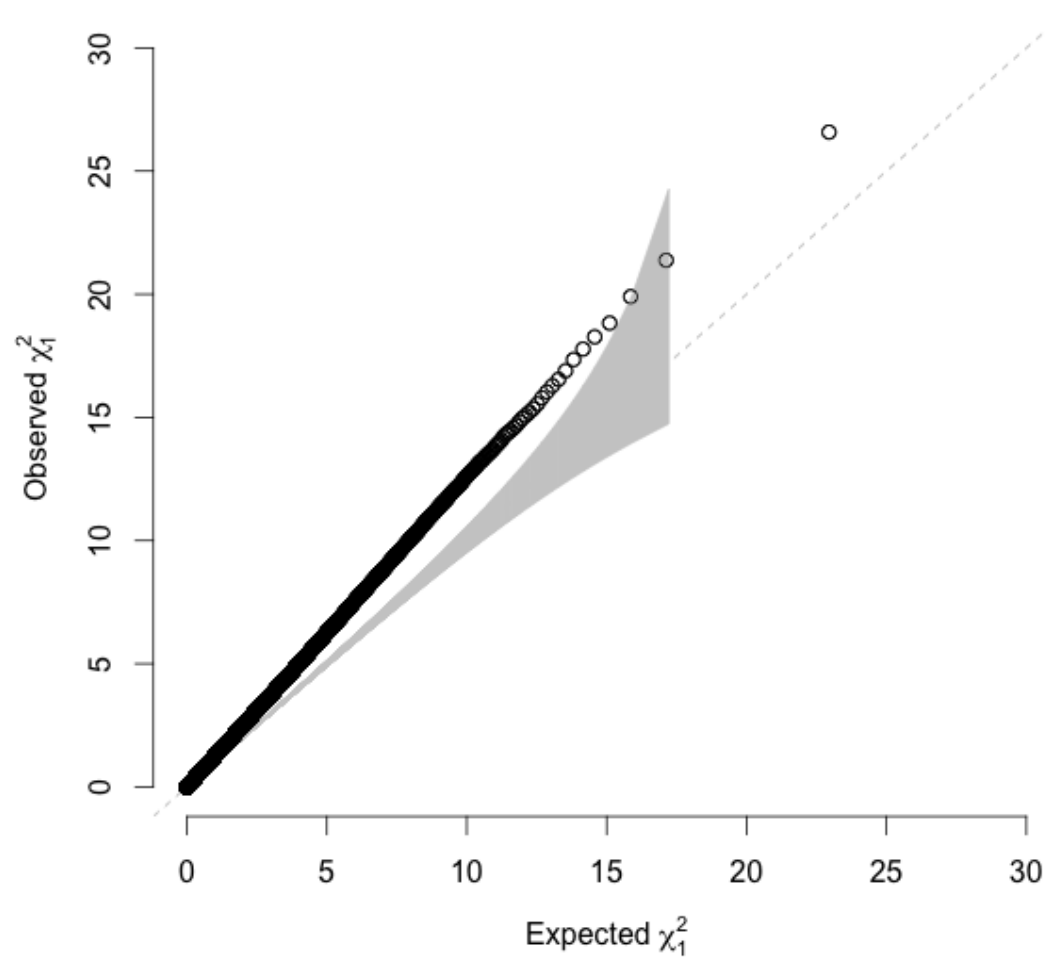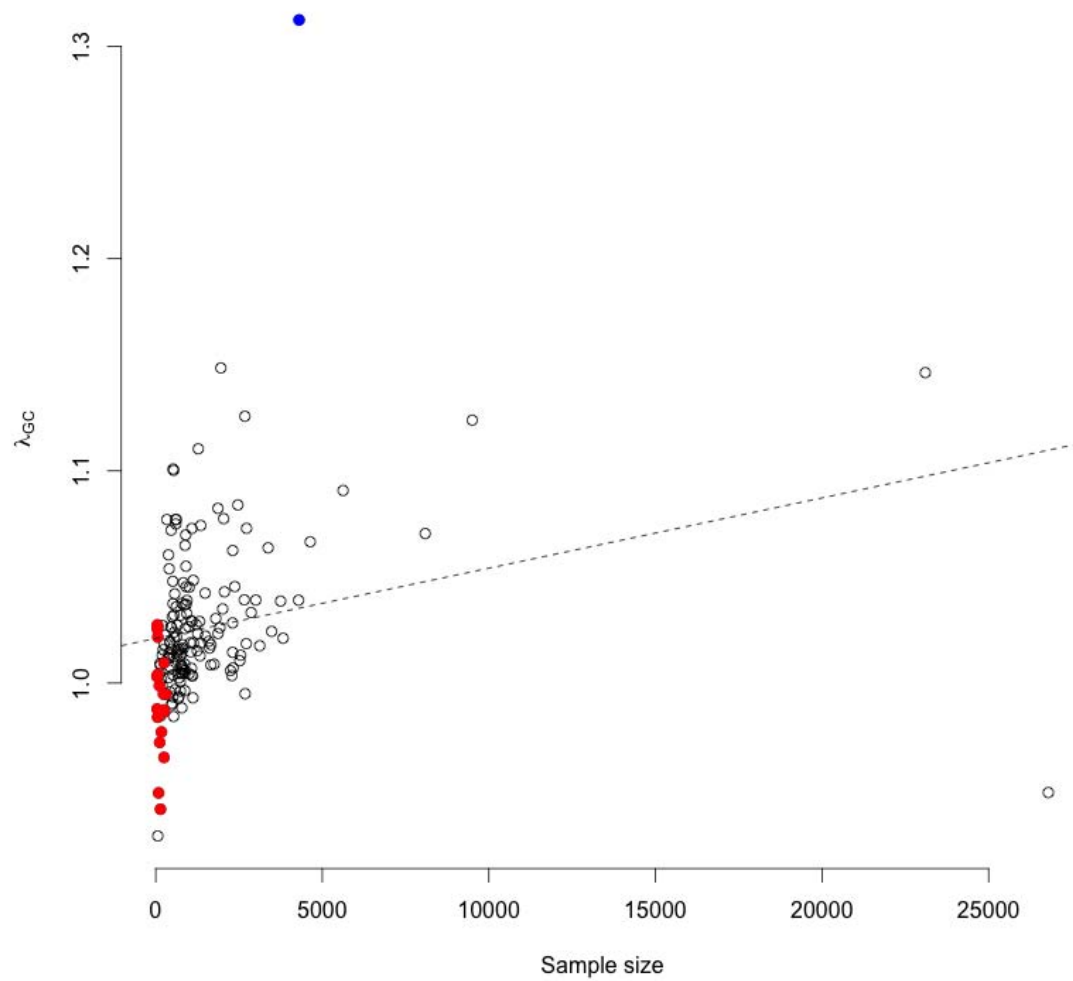


'19

'20

'21    b)



'22

'23

'24　　c)



'25

'26

'27    d)



'28

'29

'30 **Figure 7 $\bar{\lambda}_{meta}$ and $\lambda_{gc}$ for GIANT height GWAS cohorts. (a)** Sample size of each cohort against $\lambda_{GC}$.

'31 The linear regression is presented as a dashed line, $\lambda_{GC} = 1.021 + 0.0000033N$, and $R^2 = 0.013$. **(b)** Sample

'32 size of each cohort against $\bar{\lambda}_{meta}$, which was the mean of a cohort's $\lambda_{meta}$ over all other cohorts. The linear

'33 regression is presented as a dashed line, $\bar{\lambda}_{meta} = 1.012 + 0.00000055N$ (N = reported sample size), and

'34 $R^2 = 0.055$. **(c)** $\lambda_{GC}$ against $\bar{\lambda}_{meta}$ for each cohort, showing a strong correlation, $R^2 = 0.70$. The black dash

'35 line indicates the regression slope for all 174 pairs: $\bar{\lambda}_{meta}=0.7251+0.281\lambda_{GC}+e$. The red dashed line

'36 indicates the regression slope for 20 pairs of MIGEN cohorts: $\bar{\lambda}_{meta}=0.369+0.631\lambda_{GC}+e$. The side of each

'37 circle is proportional to sampling size on logarithm scale. **(d)** Small sample size leads to a correlation

'38 between $\bar{\lambda}_{meta}$ and $\lambda_{GC}$ using 174 GIANT height GWAS sample size. 30,000 independent loci, minor allele

'39 frequency ranged from 0.1~0.5, were simulated, and $h^2 = 0.5$. The red dashed line indicates the regression

'40 slope for 20 simulated MIGEN cohorts, $\bar{\lambda}_{meta} = 0.488 + 0.510\lambda_{GC}+e$ ($R^2 = 0.78$). The side of each circle

'41 is proportional to sampling size on logarithm scale. **(e)** $\lambda_{meta}$ for whole MIGEN to 174 cohorts. 20 MIGEN

'42 files were combined together to make "whole MIGEN" via meta-analysis, and the summary statistics were

'43 used to calculate $\lambda_{meta}$ with 174 cohorts using 30,000 independent loci. As MIGEN cohorts were part of

'44 "whole MIGEN", their $\lambda_{meta}$ were in general below 1. The dashed line is the mean of $\lambda_{meta}$ of the "whole

'45 MIGEN". The subplot (red box) shows a strong correlation of 0.93 between $\lambda_{meta}$ (for "whole MIGEN" vs

'46 each MIGEN cohort), and sample size of each MIGEN cohort.

'47

36
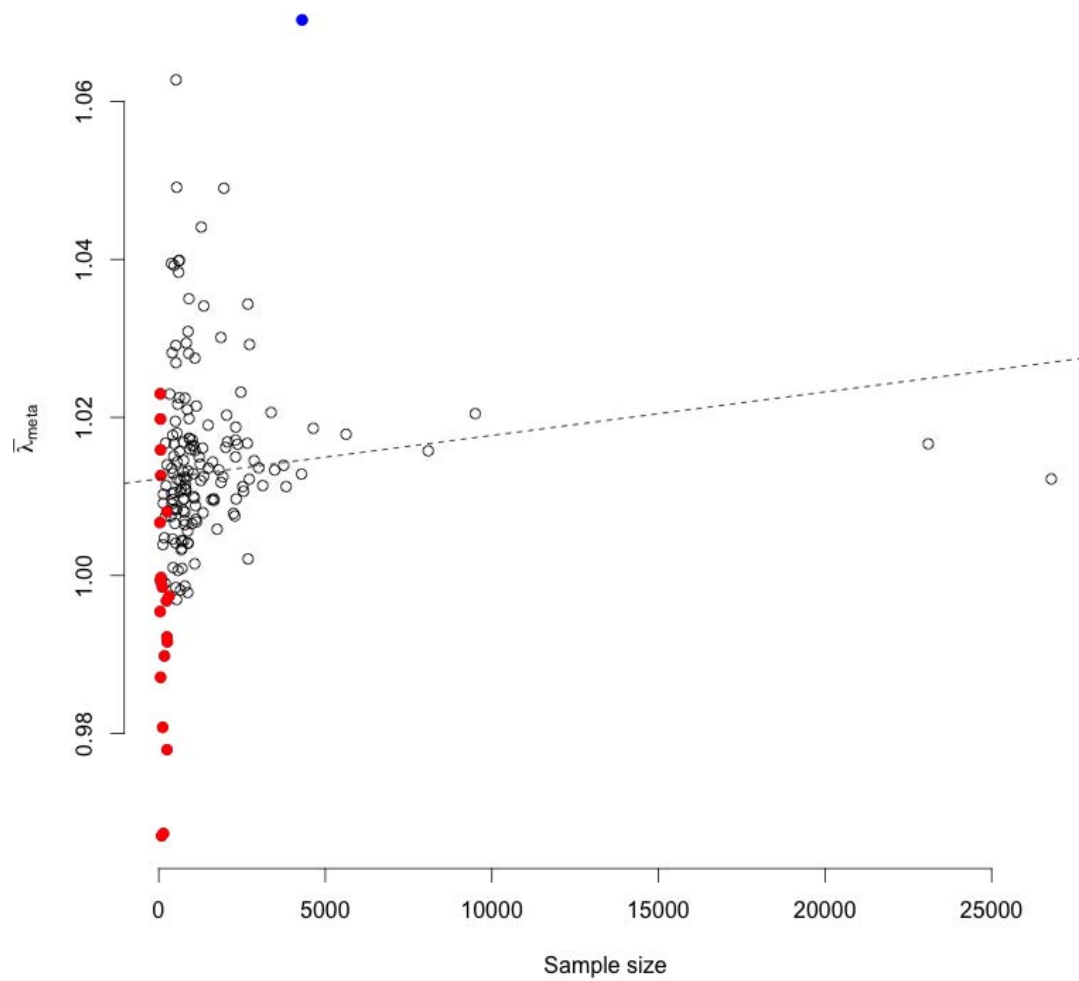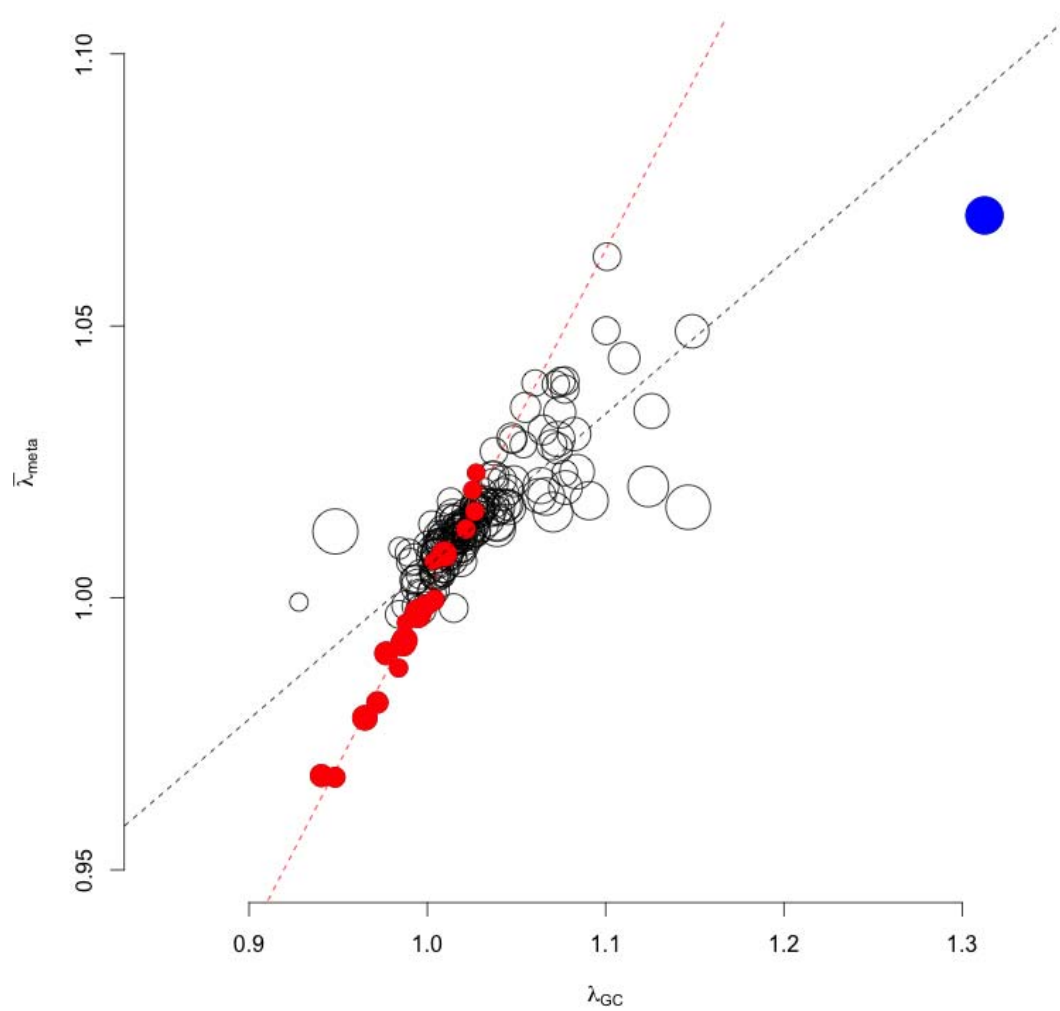
'48    a)



'49

'50

'51    b)


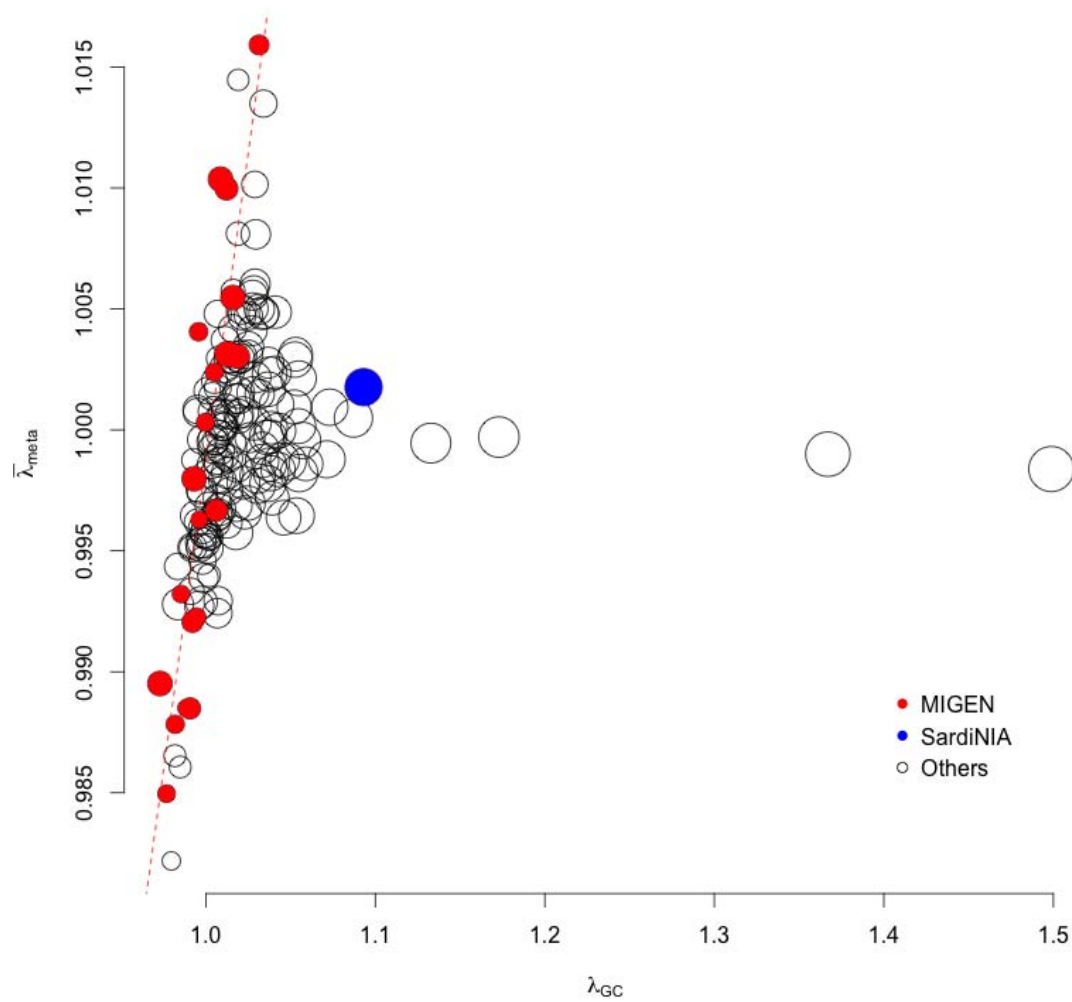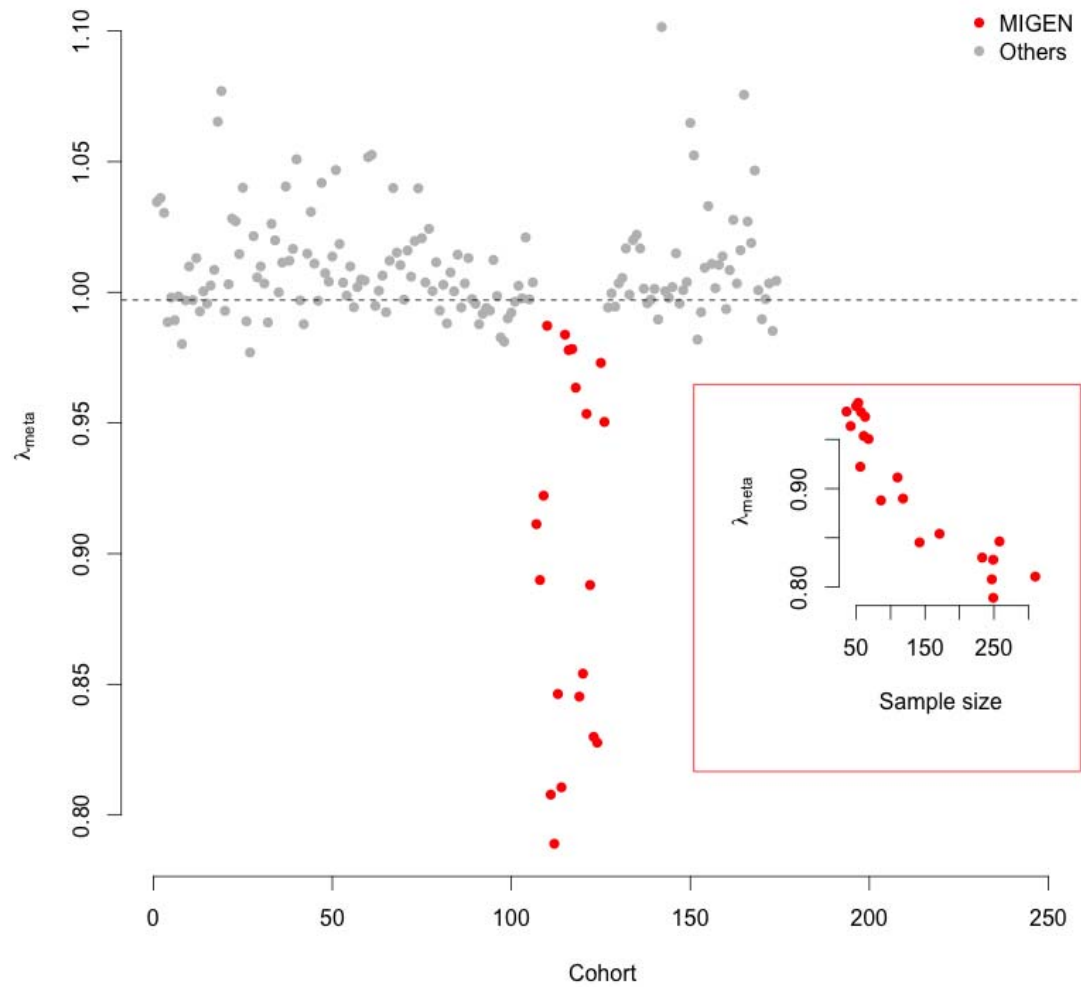
'52

'53

'54    c)



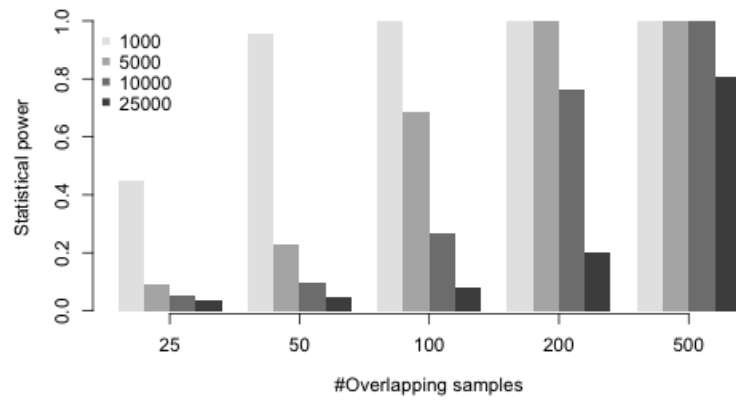'55

'56

'57    d)



'58

'59

′60   e)



′61

′62

'63 **Figure 8 Pseudo profile score regression for the WTCCC 7 diseases.** a) Statistical power for detecting

'64 overlapping samples between a pair of cohorts given type I error rate of 0.05. Top panel: The y-axis

'65 represents statistical power, and the x-axis the number of overlapping samples. Cohort 1 has 1,000, 5,000,

'66 10,000, or 25,000 samples, and cohort 2 has 1,000 samples. The two cohorts have 25, 50, 100, 200, and 500

'67 overlapping samples. Bottom panel: the corresponding 95% confidence interval is given for each scenario in

'68 the top panel. The statistical power is maximized when the two cohorts have the same sample size. b) Each

'69 cluster represents a pair of cohorts as denoted on the x-axis. Within each cluster, from left to right, the

'70 detected overlapping controls using $\lambda_{meta}$ based either on effect size estimates or minor allele frequency

'71 (MAF), PPRS using 100, 200, and 500 markers. WTCCC cohort codes: BD for bipolar disorder, CAD for

'72 coronary artery disease, CD for Crohn's disease, HT for hypertension, RA for rheumatoid arthritis, T1D for

'73 type 1 diabetes, T2D for type 2 diabetes.
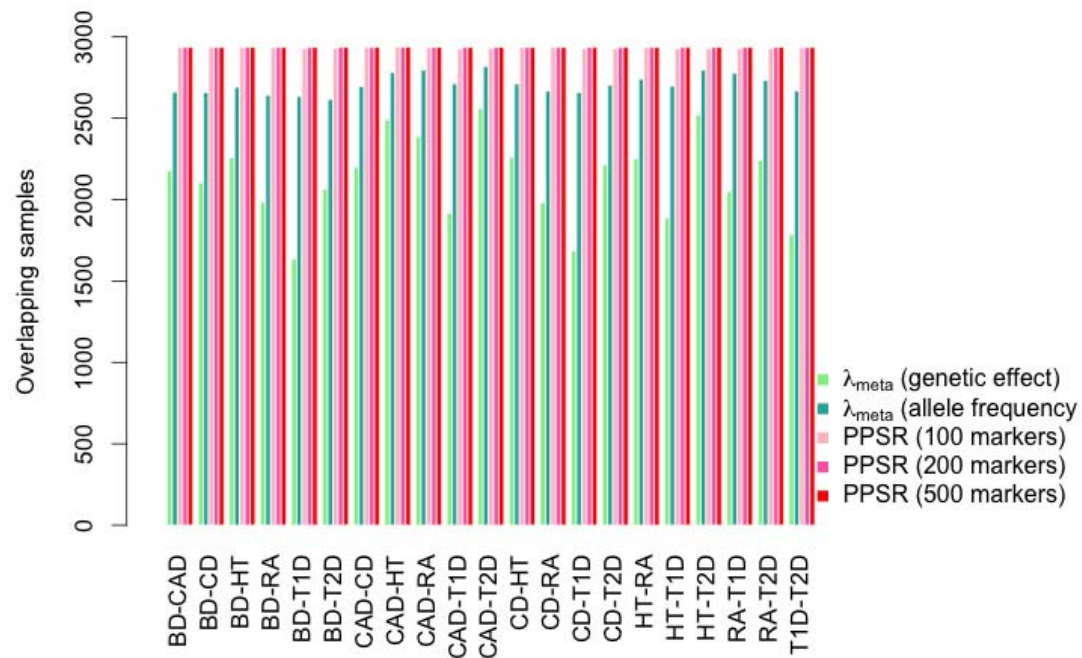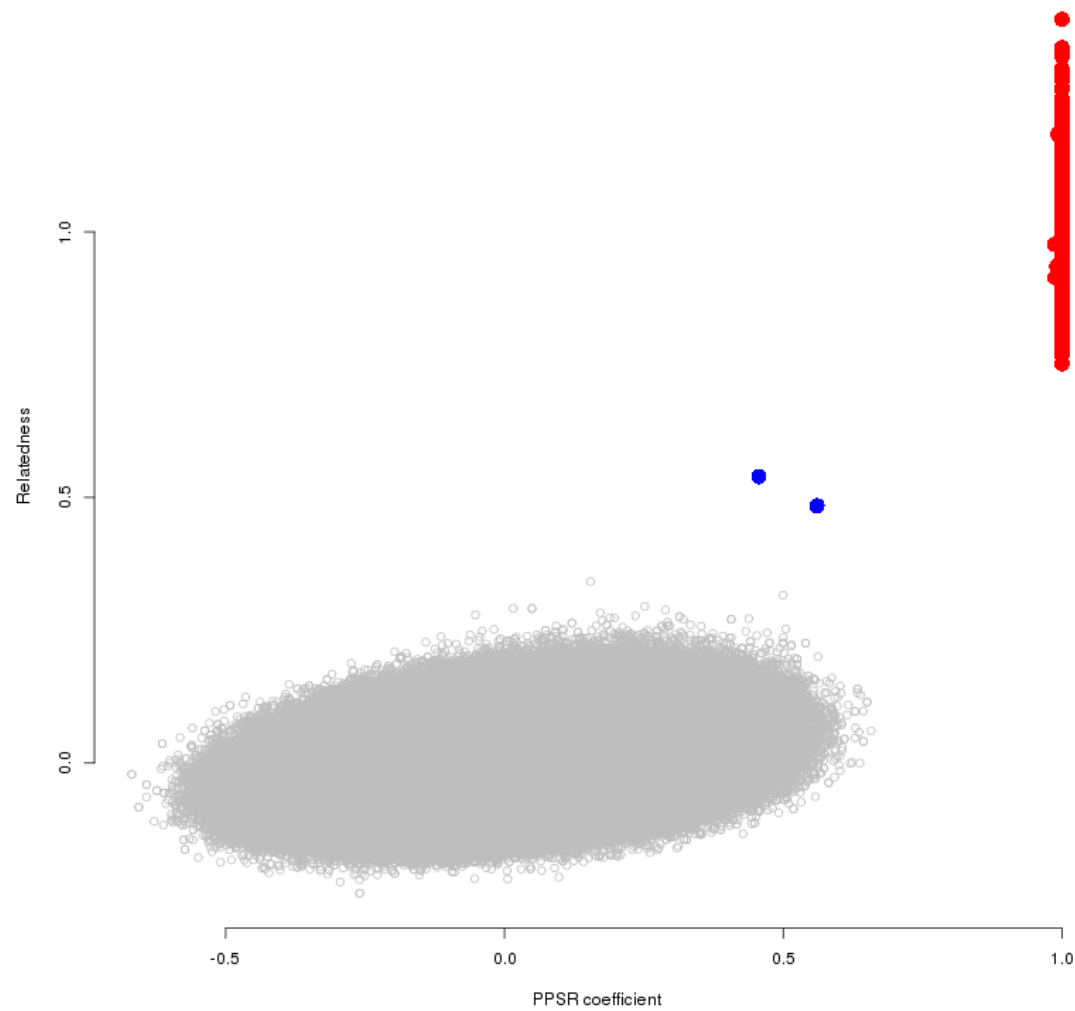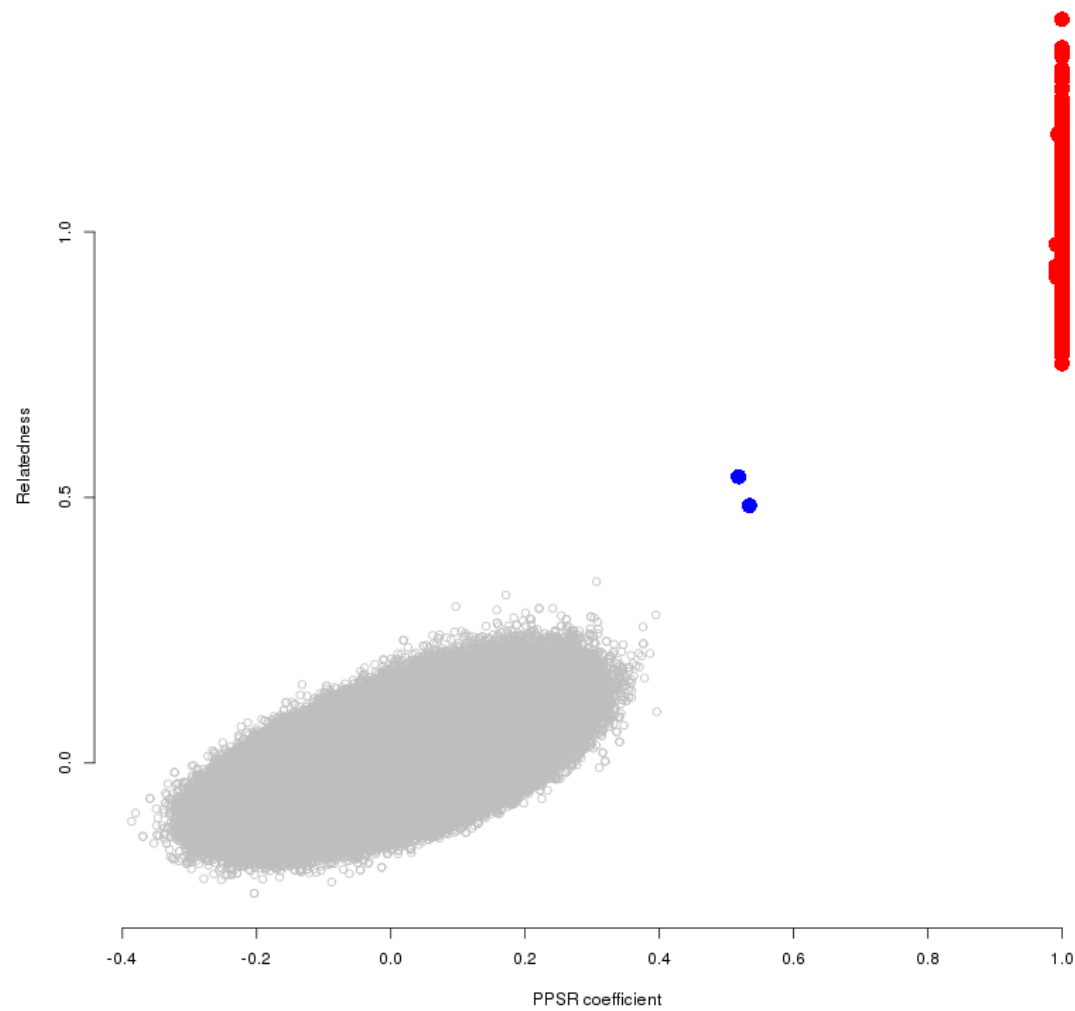
'74

'75    a)



'76

'77

'78   b)



'79

'80

'81　**Figure 9 PPSR coefficients for identifying shared controls/relatives between WTCCC BD and CAD**

'82　**cohorts. (a)** Illustration for regression coefficients between WTCCC BD and CAD from 57 pseudo profile

'83　scores (PPS) generated from 500 markers. The x-axis is the PPSR regression coefficients and y-axis is real

'84　genetic relatedness (as calculated from individual level genotype data). The red points are the shared

'85　controls between two cohorts, and blue points are first-degree relatives. **(b)** The PPS regression coefficients

'86　for detecting overlapping first-degree relatives using 286 PPS generated from 500 markers. **(c)** Decoding

'87　genotypes from the PPS. Given the set of profile scores, one may run a GWAS-like analysis to infer the

'88　genotypes. The ratio between the number of markers ($M$) and number of pseudo profile scores ($K$)

'89　determines the potential discovery of individual-level information. The higher the ratio and, the higher the

'90　allele frequency, the less information can be recovered. From left to right, the profile scores generated using

'91　different number of markers. The y-axis is a $R^2$ metric representing the accuracy between the inferred

'92　genotypes and the real genotypes. From left to right panels 100, 200, 500, and 1000 SNPs were used to

'93　generate 10, 20, 50, and 1000 profiles scores. In each cluster, the three bars are inferred accuracy using

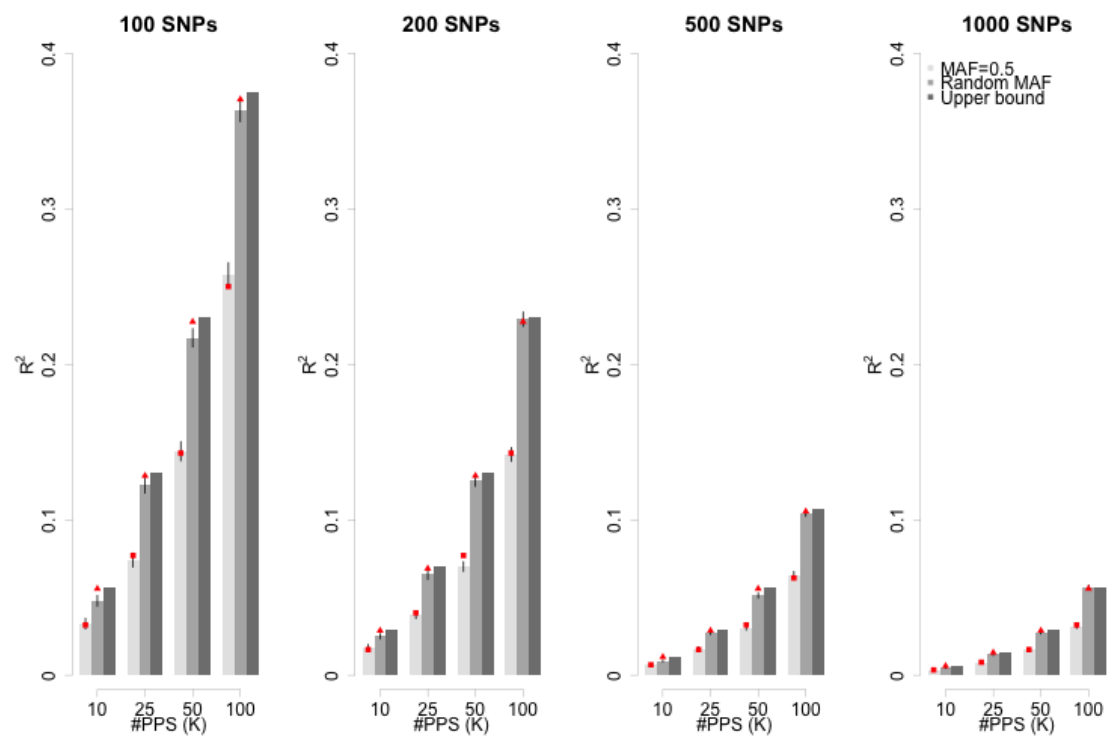'94　different MAF spectrum alleles, given with the SE of the mean.

'95

45

'96  a)



'97

'98

'99    b)



300

301

47

802  c)



803
804

**Table 1 The estimated correlation for a pair of cohorts via their summary statistics**

| $n_1$ | $n_2$ | $n_{1,2}$ | $h^2$ | $M$ | $Q$ | $\gamma_{1,2} = \dfrac{n_{1,2}}{\sqrt{n_1 n_2}}$ | $\hat{\rho}_{1,2} \pm$ SD | $\hat{\gamma}_{1,2} \pm$ SD |
|---|---|---|---|---|---|---|---|---|
| 1,000 | 1,000 | 100 | 0.25 | 30,000 | 1,000 | 0.1 | 0.1072±0.0064 | 0.101±0.0093 |
| 1,000 | 2,000 | 100 | 0.25 | 30,000 | 1,000 | 0.0707 | 0.0814±0.0054 | 0.0709±0.0088 |
| 1,000 | 5,000 | 100 | 0.25 | 30,000 | 1,000 | 0.0447 | 0.0615±0.0055 | 0.0425±0.0096 |
| 1,000 | 10,000 | 100 | 0.25 | 30,000 | 1,000 | 0.0316 | 0.0556±0.0063 | 0.0325±0.0099 |
| | | | | | | | | |
| 1,000 | 1,000 | 1 | 0.25 | 30,000 | 1,000 | 0.001 | 0.0092±0.0056 | 0.0017±0.0093 |
| 1,000 | 2,000 | 1 | 0.25 | 30,000 | 1,000 | 0.0007 | 0.0126±0.0053 | 0.0006±0.0079 |
| 1,000 | 5,000 | 1 | 0.25 | 30,000 | 1,000 | 0.000447 | 0.0189±0.0060 | 0.0016±0.0090 |
| 1,000 | 10,000 | 1 | 0.25 | 30,000 | 1,000 | 0.000316 | 0.0259±0.0059 | 0.0008±0.0092 |
| | | | | | | | | |
| 1,000 | 1,000 | 100 | 0 | 30,000 | 1,000 | 0.1 | 0.0996±0.0052 | 0.094±0.0085 |
| 1,000 | 2,000 | 100 | 0 | 30,000 | 1,000 | 0.0707 | 0.0704±0.0048 | 0.0712±0.0097 |
| 1,000 | 5,000 | 100 | 0 | 30,000 | 1,000 | 0.0447 | 0.0453±0.0057 | 0.0441±0.0090 |
| 1,000 | 10,000 | 100 | 0 | 30,000 | 1,000 | 0.0316 | 0.0335±0.0057 | 0.0325±0.0079 |

$^*$Q is the number of QTLs among M simulated loci. We also tried $Q = 100$, the results were nearly identical.

$\gamma_{1,2}$ represents the true correlation due to overlapping samples

$\hat{\rho}_{1,2}$ represents the estimated correlation estimated via the method proposed by Bolormaa et al[26], and Zhu et al[27]

$\hat{\gamma}_{1,2}$ represents the estimated correlation estimate via $\lambda_{meta}$, $\hat{\gamma}_{1,2} = \dfrac{1-\hat{\lambda}_{meta}}{\frac{2\sqrt{n_1 n_2}}{n_1+n_2}}$.

810 **Table of contents**

815

816

## Method I: $F_{st}$ derived genetic distance

$F_{st}$ is a measure of genetic differentiation between populations. It is usually estimated using individual-level genotype data from multiple samples in two or more populations[1]. Here, we calculate $F_{st}$ using summary data on allele frequencies, which implicitly assumes Hardy-Weinberg equilibrium genotype frequencies within populations. We use summary statistic calculated $F_{st}$ as a metric for quality control for each cohort. If the allele frequencies reported for a cohort depart genome-wide from its expectation based on known ancestry due to technical artifacts, then we may observe an unexpected $F_{st}$ value when comparing to a reference panel of know ancestry.

We calculate $F_{st}$ between each cohort and a reference panel, choosing the appropriate reference sample depending on the purpose of the analysis. For the inference of global-level diversity, we chose YRI, CHB, and CEU as the reference panels. For the inference of within-Europe diversity, we chose CEU, FIN, and TSI as the reference panels. As the different allele frequencies across three samples reflected the real diversity among these reference panels, we did not apply any exclusion criteria on the reference allele frequency. Nevertheless, as GIANT height GWAS samples were imputed to the HapMap panel, the majority of SNPs matched to the 1KG reference samples comprised common SNPs. After ranking the calculated $F_{st}$ in ascending order for all matched SNPs, we sampled 30,000 $F_{st}$ evenly along the ordered $F_{st}$. These 30,000 markers are quasi-independent and evenly distributed across the genomes. The mean of the 30,000 $F_{st}$ was employed to represent the $F_{st}$ measure between a cohort and a reference panel. The sampled 30,000 markers may differ from one pair of cohorts to another pair, but as tested resample 30,000 markers caused ignorable changes of the mean of $F_{st}$. Another reason we chose 30,000 markers is that there are around 30,000 quasi-independent markers for GWAS data as observed in empirical data and expected from theory[2,3].

In this study, $F_{st}$ is calculated from the allele frequencies estimated from cohorts, provided as summary statistics. $F_{st}$ is treated as a data statistic for measuring allele frequency differentiation. In general the interpretation of $F_{st}$ can vary with context[4].

$$F_{st} = \frac{\frac{r}{(r-1)\sum_{i=1}^{r} n_i}[\sum_{i=1}^{r} n_i(p_i - \bar{p})^2]}{\bar{p}(1-\bar{p})} \text{ (Equation 1)}$$

with $p_i$ the estimated reference allele frequency in population $i$ from a sample of $n_i$ alleles, $\bar{p}$ is the weighted average frequency in the entire sample, and $r$ is the number of populations.

51

850　　Here, we only compared each cohort to the 1KG reference panel, so $r = 2$ and the equation

851　　becomes

852　　$$F_{st} = \frac{\frac{2}{n_{1,2}}[\sum_{i=1}^{2} n_i(p_i - \bar{p})^2]}{\bar{p}(1-\bar{p})} \qquad \textbf{(Equation 2)}$$

853　　in which $n_{1,2} = n_1 + n_2$, and $\bar{p} = \frac{n_1}{n_{1,2}}p_1 + \frac{n_2}{n_{1,2}}p_2$ is the mean allele frequency. Alternative

854　　estimators for $F_{st}$ are possible, and a comprehensive comparison of different $F_{st}$ estimators

855　　was recently reported[5].

856

857　　If the two cohorts are not that different in terms of their allele frequencies, for example, the

858　　cohorts from European nations, $p_i \approx p_j \approx \bar{p}$,

859　　$$E(F_{st}) \approx \frac{1}{n_{i,j}} + \frac{2n_i n_j}{n_{i,j}^2}\frac{[E(p_i) - E(p_j)]^2}{\bar{p}(1-\bar{p})} \qquad \textbf{(Equation 3)}$$

860　　At the right side of the equation, the first term represents the sampling variance for allele

861　　frequency for a pair of cohorts, and the second term represents the allele frequency difference

862　　due to divergence from a common ancestor. The estimated $F_{st}$ is influenced by sample size,

863　　and $F_{st} \geq \frac{1}{n_{i,j}}$, which is the sampling variance of $F_{st}$ for a pair of cohorts[1]. As each 1KG

864　　reference population has a sample size around 100, there is no disproportionate impact of

865　　sample size in calculating $F_{st}$.

866

867　　$F_{st}$ **Cartographer algorithm.** The purpose of using the $F_{st}$ Cartographer algorithm is to find

868　　the coordinates of a cohort given its $F_{st}$ to the reference populations. The algorithm can be

869　　expressed in Cartesian geometry. Given three reference populations, a target cohort has three

870　　$F_{st}$ measures, $F_1$, $F_2$, and $F_3$, respectively. Given a Cartesian coordinate system, the

871　　coordinate for these three reference populations are $(a_1, b_1)$, $(a_2, b_2)$, and $(a_3, b_3)$,

872　　respectively. The algorithm tries to find the coordinates $(x, y)_{E_{i,j}}$ on each the edge $(E_{i,j})$ that

873　　connects reference populations $i$ and $j$

874　　$$(x, y)_{E_{i,j}} = [(a_j - a_i)\frac{F_i}{F_i + F_j} + a_i, (b_j - b_i)\frac{F_i}{F_i + F_j} + b_i] \qquad \textbf{(Equation 4)}$$

875　　The coordinates of the gravity of triangle, $(x, y)_G$, that connects $E_{1,2}$, $E_{1,3}$, and $E_{2,3}$ are

876　　$$(x, y)_G = (\frac{x_{E_{1,2}} + x_{E_{1,3}} + x_{E_{2,3}}}{3}, \frac{y_{E_{1,2}} + y_{E_{1,3}} + y_{E_{2,3}}}{3}) \qquad \textbf{(Equation 5)}$$

877

878　　**Inference of cohort origins at the global level.** To assess genetic background, for each

879　　cohort we calculated its $F_{st}$ values using CEU, CHB, and YRI as the reference panel,

880　　respectively. We denote these three $F_{st}$ values as $F_{CEU}$, $F_{CHB}$, and $F_{YRI}$. These values reflect

881 genetic distances between a cohort and the reference panels - the greater the value the further

882 the genetic distance. We developed an algorithm called $F_{st}$ cartographer, which can map a

883 cohort to global genetic variation as previously observed using individual level data from

884 principal component analysis[6]. The steps in the algorithm are as follows (**Supplementary**

885 **Fig. 1**):

886

887 Create the coordinates for the reference samples. Without loss of generality, these three

888 reference populations form an equilateral triangle, and we set the length of each edge to

889 unity. For example, the coordinates CEU, CHB, and YRI are $(-\sqrt{3}, 1)$, $(\sqrt{3}, 1)$, and $(0, -2)$,

890 respectively, and connecting the coordinates of the three reference populations formed an

891 equilateral triangle – the reference space. The gravity of this equilateral triangle is the

892 origin of the Cartesian space. The choice for the coordinates for the reference population is

893 arbitrary.

894

895 **Step 1 Create a cohort triangle using Equation 4.** Finding a point the distances of that to

896 both ends, which represent two populations, is proportional to the ratio of the $F_{st}$ values of

897 the cohort to these two reference populations. Similarly, find the points on the other two

898 edges. For example, Finland Twin Cohort (FTC) had $F_{CEU} = 0.0102$, $F_{YRI} = 0.153$, and

899 $F_{CHB} = 0.099$. On the CEU-YRI edge, a point split the length to 0.0102:0.153, was

900 $(-1.72, 0.98)_{E_{1,2}}$; On the CEU-CHB edge into 0.0102:0.099, was $(-1.70, 1)_{E_{1,3}}$; and on the

901 YRI-CHB edge into 0.153:0.099, was $(1.05, -0.18)_{E_{2,3}}$. Connecting the three coordinates

902 created a "FTC" triangle inside the reference triangle.

903

904 **Step 2 Find the gravity of the cohort triangle using Equation 5.** The gravity of the "FTC"

905 triangle had its coordinates of $(-0.79, 0.60)_{G_{FTC}}$, which is inferred as the geographic

906 coordinates for FTC in $F_{PC}$ space. It had relative distances of 1.03, 2.55, and 2.72 to CEU,

907 CHB, and YRI, respectively. The shorter the distance, the closer the genetic background is.

908

909 **Step 3** Repeat Steps 1, and 2 until the gravity of each cohort is found.

910

911 Plots of the coordinates for each cohort will show the relative distance of each cohort to the

912 reference samples. If a cohort has equal distances to three reference populations, its gravity

913 will be close to the origin of the reference triangle.

914

915

916 ## Method II: Principal component analysis for cohort-level allele
917 ## frequencies

918 PCA has been widely used in genetics[7] and recently proposed for controlling population

919 stratification for GWAS[8,9]. We provide a new method that uses cohort-level allele

920 frequencies, often provided as summary statistics in meta-analysis. We call the new method

921 as meta-PCA.

922

923 Meta-PCA is based on a $G = (C + K) \times M$ matrix, which includes $K$ reference populations

924 and $C$ cohorts of question on $M$ markers. In $G$, the $m^{th}$ column represents the reported

925 reference allele frequencies for the $m^{th}$ marker for (K+C) cohorts. The kernel correlation

926 matrix for PCA is constructed on $\Sigma = G_s \times G_s^T$, in which $G_s$ is the standardization for $G$ for

927 each locus (on each column of $G$). Compared with individual-level data PCA, in the context

928 of meta-PCA each cohort can be viewed as an individual in the conventional sense. Given $\Sigma$

929 matrix, all the implementation is the same as the individual-data PCA.

930

931 There are efforts in establishing genetic interpretation for PCA [8,10–12]. The interpretation of

932 meta-PCA could be approached by $F_{pc}$ as described in the last section.

933

55

## Method III: The detection of overlapping samples with $\lambda_{meta}$

**Inference of cohort origins at the within-Europe level.** To assess genetic background, for each cohort we calculated its $F_{st}$ values using CEU, FIN, and TSI as the reference panel, with coordinates $(-\sqrt{3}, 1)$, $(\sqrt{3}, 1)$, and $(0, -2)$, respectively. For FTC, it had $F_{st}$ values of 0.0102, 0.0052, and 0.0157, to CEU, FIN, and TSI, respectively. Using the $F_{st}$ Cartographer algorithm, the gravity of the FTC triangle had its coordinates of $(0.274, 0.361)_{G_{FTC}}$. It had relative distances of 2.10, 1.59, and 2.42, to CEU, FIN, and TSI, respectively.

**Genealogical subspace.** Furthermore, we partition the $F_{PC}$ space into three subspaces. For example, given coordinates of $(-\sqrt{3}, 1)$, $(\sqrt{3}, 1)$, and $(0, -2)$, for CEU, FIN, and TSI, respectively, connecting the origin and the coordinates for any two reference populations created a subspace, which is defined as a genealogical subspace. We had three genealogical subspaces: CEU-FIN genealogical subspace, CEU-TSI genealogical subspace, and FIN-TSI genealogical subspace, respectively. If a cohort is located inside a subspace, it indicates that this cohort may be derived from these two reference populations that creates the genealogical subspace.

For European cohorts, the coordinates calculated from $F_{st}$ Cartographer algorithm mirror the origins of geographic locations of the cohorts, similar, but less refined, to what has been observed in previous studies using individual level data for European samples[13,14].

**Effective number of overlapping samples $(n_o)$.** If a pair of cohorts has overlapping samples, it leads to a correlation of the estimated genetic effects for each locus. In the recent literature, two kinds of correlation due to overlapping samples were introduced. The first one was defined by directly calculating correlation between all estimated test statistics, $r = cor(Z_1, Z_2)$, in which $Z$ is a vector of $M$ matched loci between two cohorts [15,16]. The second one was defined on the correlation for single locus given overlapping samples, as introduced by Lin and Sullivan[17]. We used the second definition, and then extended the correlation due to any relatives, a generalization of Lin and Sullivan.

For a pair of cohorts of sample sizes $n_1$ and $n_2$ $(n_1 \geq n_2)$, for $M$ matched loci which have GWAS summary statistics, for example additive effects and their standard errors. For the $m^{th}$ locus, estimated association effect sizes are $b_{1.m}$ and $b_{2.m}$ with sampling variance $\sigma_{b_{1.m}}^2$ and $\sigma_{b_{2.m}}^2$, respectively. $b_1$ is assumed to be drawn from a normal distribution $N(b_{1.m}, \sigma_{b_{1.m}}^2)$, and

56

968     $b_2 \sim N(b_{2.m}, \sigma^2_{b_{2.m}})$. In cohort 1, $w_{1|k} = \frac{n_{12|k}}{n_1}$ is proportion of samples with a $k^{th}$-degree

969     relatives in cohort 2 with $n_{12|k}$ the number of relatives of kth degree relatives shared between

970     the samples; the phenotypic variance is assumed to be the same across the cohorts for a

971     quantitative trait. For a locus, the genetic effect is estimated by linear regression $y_1 = a +$

972     $b_1 x + e$ in cohort 1 (the index for the locus is dropped for convenience). If the sampling

973     variance of a locus is assumed to be the same for any subset of samples

$$b_1 = \Sigma^K_{k=0} w_{1|k} \frac{\left[ E(y_{1|k} x_k) - E(y_{1|k}) E(x_k) \right]}{var(x_m)} = \Sigma^K_{k=0} w_{1|k} b_{1|k}$$

974     The standard error of $b_m$ is $\sigma_{b_1} = \sqrt{\frac{(1-h^2_{b_1})\sigma^2_{y_1}}{n_1}} \approx \sqrt{\frac{\sigma^2_{y_1}}{n_1}}$, in which $h^2_{b_1}$ is the proportion of

975     phenotypic variance explained by the locus and $\sigma^2_{y_1}$ is the phenotypic variance of the trait.

976     The sampling variance for $\sigma_{b_{1.k}} = \sqrt{\frac{\sigma^2_{y_1}}{w_{1|k} n_1}}$. This decomposition of the genetic effect can be

977     applied to cohort 2. Consequently, the covariance between $b_1$ and $b_2$ for the locus is

$$cov(b_1, b_2) = cov(\Sigma^K_{k=0} w_{1|k} b_{1.k}, \Sigma^K_{k=0} w_{2|k} b_{2.k}) = \Sigma^K_{k=0} w_{1|k} w_{2|k} cov(b_{1|k}, b_{2|k})$$

978     in which $cov(b_{1|k}, b_{2|k}) = \rho_k \theta_k \sigma_{b_1|k} \sigma_{b_2|k}$ is the covariance between the genetic effects

979     estimated in two cohorts due to the $k$-degree relatives. $\rho_k$ is the phenotypic correlation for the

980     $k$-degree relatives, and $\theta_k$ is the genetic relatedness for the $k$-degree relatives. $\theta_k = \left(\frac{1}{2}\right)^k$ is

981     the coefficient of identity for descent. For duplicated samples, $\rho_0 = h^2 + \rho_{e|0}$, in which $h^2$ is

982     the heritability, and $\rho_{e|0}$, the environmental correlation to be close 1 for overlapping samples;

983     for other relatives ($k \geq 1$), $\rho_k \approx \theta_k h^2$.

984

985     **Correlation between the estimated genetic effects.** The covariance can be generalized as

986     $cov(b_1, b_2) = \Sigma^K_{k=0} \sqrt{w_{1|k} w_{2|k}} \rho_k \theta_k \sqrt{\frac{\sigma^2_{y_1} \sigma^2_{y_2}}{n_1 \ n_2}}$. After adjustment by the sampling variance, the

987     correlation between $b_1$ and $b_2$ is

988     $\rho_{b_1,b_2} = \frac{cov(b_1,b_2)}{\sigma_{b_1}\sigma_{b_2}} = \Sigma^K_{k=0} \rho_k \theta_k \sqrt{w_{1|k} w_{2|k}} = \frac{\Sigma^K_{k=0}\rho_k \theta_k n_{12|k}}{\sqrt{n_1 n_2}} = \frac{n_o}{\sqrt{n_1 n_2}}$ **(Equation 6)**

989     in which $n_o = \Sigma^K_{k=0} \rho_k \theta_k n_{12|k}$, is the effective number of overlapping samples averaged over

990     all relative pairs that are across the two cohorts. As the variance explained by each locus is

991     small, and after further weighted by $\theta_k$, the contribution from overlapping relative is small.

992     When ignoring the first and higher degree relatives $n_o$ equals the contribution from

993     overlapping samples. This is consistent with the results from Lin and Sullivan[17], who

57

994    considered overlapping samples only. So, the correlation at any single locus is largely

995    determined by the overlapping samples ($n_{12|0}$) for summary statistics.

996    $$\rho_{b_1,b_2} = \frac{n_o}{\sqrt{n_1 n_2}} \approx \frac{n_{12|0}}{\sqrt{n_1 n_2}}$$    **(Equation 7)**

997    So, in the text hereafter, $n_e$ indicates overlapping samples only, otherwise specified.

998

999    **Correlation for case-control studies.** The theory above is based on a quantitative trait, but it

1000   holds approximately true for case-control studies if a locus is from the null distribution of no

1001   association with the disease. Given $n_{12.ctrl}$ overlapping controls and $n_{12.cs}$ overlapping cases,

1002   for a locus associated with disease its correlation of the regression coefficient is $\rho_{b_1,b_2} =$

1003   $\frac{n_{12.ctrl}\sqrt{R_1 R_2} + n_{12.cs}\frac{1}{\sqrt{R_1 R_2}}}{\sqrt{n_1 n_2}}$ as indicated by Lin and Sullivan[17], in which $R_i$ is the ratio between

1004   cases and controls in the $i^{th}$ cohort. When it is balanced case-control design – $R. = 1$,

1005   $\rho_{b_1,b_2} = \frac{n_{12.ctrl} + n_{12.cs}}{\sqrt{n_1 n_2}} = \frac{n_o}{\sqrt{n_1 n_2}}$ resembles the correlation for quantitative traits. However, it

1006   should be noticed that for case control data, $n_e$ is confounded with the number of overlapping

1007   cases and controls.

1008
1009   **Theory for $\lambda_{meta}$.** For the summary statistics between a pair of cohorts for the $m^{th}$ locus, we

1010   can construct a statistic

1011   $$T_m = \frac{(b_{1.m} - b_{2.m})^2}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2} = \left[\frac{(b_{1.m} - b_{2.m})^2}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{1,2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}}\right] \times \left[\frac{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{1,2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2}\right]$$ **(Equation 8)**

1012   in which $\rho_{1,2}$ is the correlation between $b_{1,m}$ and $b_{2,m}$.

1013   $$E(T_m) = \left\{\frac{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{1,2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{1,2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}} + \frac{[E(b_{1.m}) - E(b_{2.m})]^2}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{b_1,b_2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}}\right\}\left\{\frac{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2} - \right.$$

1014   $\rho 1,2 2\sigma b1.m\sigma b2.m\sigma b1.m2+\sigma b2.m2=(1+H)(1-\rho 1,2\kappa)$ **(Equation 9)**

1015   in which $H = \frac{[E(b_{1.m}) - E(b_{2.m})]^2}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2 - 2\rho_{1,2}\sigma_{b_{1.m}}\sigma_{b_{2.m}}}$, $\kappa = \frac{2\sigma_{b_{1.m}}\sigma_{b_{2.m}}}{\sigma_{b_{1.m}}^2 + \sigma_{b_{2.m}}^2} = \frac{2\sqrt{n_1 n_2}}{n_1 + n_2}$, and $\rho_{1,2} = \frac{n_o}{\sqrt{n_1 n_2}}$, as

1016   defined in Equation 7, is the correlation for this locus due to overlapping samples between

1017   this pair of cohorts. Of note, $\rho_{1,2}$ is same for each locus regardless of a null locus or a locus

1018   associated to genetic effects. For convenience, the subscript $b$ was dropped in the text

1019   hereafter.

1020

1021   Under the null hypothesis of no heterogeneity ($H = 0$) and no correlation ($\rho_{1,2} = 0$), $T_0 \sim \chi_1^2$,

1022   a standard 1-degree-of-freedom chi-square distribution. $\rho_{1,2} = \frac{n_o}{\sqrt{n_1 n_2}}$, in which $n_o$ is the

1023    effective number of overlapping samples. Of note, since the majority of markers are likely

1024    sampled from the null distribution or have very small effect sizes, we can approximate

1025    $E(b_{1.m}) = 0$ and $E(b_{2.m}) = 0$, and therefore $H \approx 0$ for most marker pairs between a pair of

1026    cohorts. For the $m^{th}$ marker that is in linkage disequilibrium with causal variants, $E(b_{1.m}) =$

1027    $\sum_{j=1}^{J_1} \beta_{1.j}\ell_{1.j}$, in which $J_1$ is the number of causal variants in linkage disequilibrium with the

1028    $m^{th}$ marker for cohort 1, $\beta_{1.j}$ is the $j^{th}$ causal variants in linkage disequilibrium with the $m^{th}$

1029    marker, and $\ell_{1.j}$ is the LD correlation between the $m^{th}$ marker and the $j^{th}$ causal variant[18].

1030    Similarly for $E(b_{2.m}) = \sum_{j=1}^{J_2} \beta_{2.j}\ell_{2.j}$. If the cohorts are from the same ethnicity, the

1031    difference in the LD correlation can be ignored, for example for samples from cohorts with

1032    European ancestry. So, under a polygenic model $H$ is expected to be zero, or close to zero.

1033

1034    The $T$ statistic is calculated for each matched SNP between a pair of cohorts. After ordering

1035    all $T$ values, we evenly sample 30,000 independent markers from the order statistic of all $T$

1036    values. Each pair of cohorts may sample $T$ values based on 30,000 markers different from

1037    another pair of cohorts.

1038    $\lambda_{meta} = \frac{median(T)}{median(\chi_1^2)} = 1 - \frac{2n_o}{n_1+n_2}$ **(Equation 10)**

1039    in which $median(\chi_1^2) = 0.4549$. Under the null hypothesis of no heterogeneity and

1040    overlapping samples ($n_o = 0$), plotting the ordered $T$ against its corresponding quartiles

1041    from $\chi_1^2$, will be along the diagonal, leading to $\lambda_{meta} = 1$. Heterogeneity between two

1042    cohorts, equivalent to a "negative" number of overlapping samples, will drive $\lambda_{meta} > 1$, and

1043    overlapping samples will make $\lambda_{meta} < 1$. The distribution of $\lambda_{meta}$ can be assessed via the

1044    beta distribution, and $\lambda_{meta}$ follows asymptotically a normal distribution $N(1,0.0136)$ given

1045    30,000 independent markers.

1046

1047    **Factors that influence $\lambda_{meta}$.** A number of factors will influence the $\lambda_{meta}$. 1) Sample

1048    overlap, including close relatives across cohorts, reduces the value of $\lambda_{meta}$ (**Supplementary**

1049    **Fig. 2**) Conservative modeling, such as inclusion of covariates in the association model that

1050    are genetically correlated with the phenotype or the 'genomic control' approach (adjusting

1051    the sampling variance with $\lambda_{GC}$, $z = \frac{b}{\sqrt{\lambda_{GC}}\sigma}$), will inflate the sampling variance, and deflate

1052    $\lambda_{meta}$. 3) Genetic heterogeneity, which can be caused by differences in genetic architecture

1053    or methodological difference, will inflate $\lambda_{meta}$. 4) As characterized by Equation 10, the

59

1054    lower bound (cohort 2 is completed included in cohort 1, given $n_1 > n_2 = n_o$) of $\lambda_{meta}$ is

1055    $1 - \frac{2}{\frac{n_1}{n_2}+1}$, upon the ratio of the samples sizes of the two cohorts.

1056

1057    **Estimating overlapping samples.** As shown in Equation 10, $\lambda_{meta}$ is a linear function of $n_o$,

1058    hence the statistical power to detect overlapping samples is equivalent to asking how $\lambda_{meta}$

1059    departs from the null distribution. Assuming $H = 0$, the overlapping samples can be

1060    estimated as $\hat{n}_o = (1 - \hat{\lambda}_{meta})\frac{(n_1+n_2)}{2}$, and $\sigma_{\hat{n}_o} = \frac{n_1+n_2}{2} \times 0.0136 \approx 0.0068(n_1 + n_2)$ given

1061    30,000 independent markers. Hence, using summary statistics only the proportion of

1062    overlapping samples can be estimated for quantitative traits. Given the type I error rate of

1063    0.05 ($\alpha = 0.05$), the statistical power for detecting $\tilde{n}_o$ overlapping samples between two

1064    cohorts is $p = \Phi^{-1}(T, \tilde{n}_o, \sigma_{\tilde{n}_o})$, in which $\Phi^{-1}$ represents the accumulation power function of

1065    a normal distribution with the mean of $\tilde{n}_o$ and standard deviation of $\sigma_{\tilde{n}_o}$. The statistical power

1066    is determined by $T$, the threshold for significance, $\tilde{n}_0$, the real overlapping samples, and $\sigma_{\tilde{n}_o}$,

1067    the standard deviation of the null hypothesis that there is no overlapping samples. Without

1068    loss of generality, $T = 1.96\sigma_{\tilde{n}_o} \approx 0.13(n_1 + n_2)$ given $\alpha = 0.05$. The 95% confidence

1069    interval is $[-0.13(n_1 + n_2), 0.13(n_1 + n_2)]$. The statistical power is maximized when

1070    $n_1 = n_2$, i.e. when a pair of cohorts has the same sample size.

1071

1072    For case-control studies, as $\hat{n}_o = n_{12.ctrl}\sqrt{R_1 R_2} + n_{12.cs}\frac{1}{\sqrt{R_1 R_2}}$, the estimate cannot

1073    distinguish between overlapping cases and overlapping controls; when $R_1 = 1$ and $R_2 = 1$

1074    (balanced case-control design for both cohorts), $\hat{n}_o = n_{12.ctrl} + n_{12.cs}$, indicating the overall

1075    overlapping samples between two cohorts, summed across cases and controls. If we know

1076    that only controls (cases) were shared between two cohorts, then $\hat{n}_o = n_{12.ctrl}\sqrt{R_1 R_2}$

1077    ($\hat{n}_o = n_{12.cs}\frac{1}{\sqrt{R_1 R_2}}$), so then an estimate of $n_o$ indicates the number of overlapping controls

1078    (cases). Therefore, quantifying overlapping samples for case-control studies is more difficult

1079    than that for quantitative traits.

1080

1081

## Method IV: Pseudo profile score regression (PPSR)

1082
1083 **PPSR** resembles the previously proposed Gencrypt method[19], but PPSR is more powerful in

1084 detecting various degree of relatives and more robust to missing data and imputation errors.

1085 For each individual, the PPS can be generated as below

1086 $A_i = S \times G_i$ **(Equation 11)**

1087 in which $A_i$ is the PPS for the $i^{th}$ individual, $S$ is a $K \times M$ score matrix, and $G_i$ is vector for

1088 the genotypes for the chosen $M$ loci.

1089 In detail,

$$\begin{bmatrix} a_{i1} \\ a_{i2} \\ \vdots \\ a_{iK} \end{bmatrix} = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1M} \\ s_{21} & s_{22} & \cdots & s_{2M} \\ \vdots & \vdots & \ddots & \vdots \\ s_{K1} & s_{K2} & \cdots & s_{KM} \end{bmatrix} \begin{bmatrix} g_{i1} \\ g_{i2} \\ \vdots \\ g_{iM} \end{bmatrix}$$

1090

1091 in which $a_{ik}$ is the $k^{th}$ profile score for the $i^{th}$ individual, $s_{km}$ is the additive effect at the

1092 $m^{th}$ locus ($m$ from 1 to $M$) for the $k^{th}$ profile score, and $g_{im}$ is the standardized genotype at

1093 the $k^{th}$ locus for the $i^{th}$ individual. Each $s$, the pseudo genetic effect, follows a standard

1094 normal distribution $N(0,1)$; each pseudo genetic effect is independent to another. For each

1095 PPS, $var(a_i) = \Sigma_{m=1}^{M} var(g_{im}s_{k.}) = \Sigma_{m=1}^{M} g_{im}^2 var(s_{k.}) = M$, in which $s_{k.}$ is the $k^{th}$ column

1096 for the $S$ matrix, and on average each locus explains $\frac{1}{M}$ of the variation. For an individual a

1097 pair of PPS, say $a_{l1}$ and $a_{l1}$, has $cov(a_{i1}, a_{i2}) = \sum_{m=1}^{M} g_{im}^2 cov(s_{l_1 m}, s_{l_2 m}) = 0$.

1098
1099 Each PPS can be seen as a trait with $h^2 = 1$ because it does not have any sampling variance.

1100 For a pair of individuals, individual $i$ and individual $j$, when both $A_i$ and $A_j$ have been

1101 standardized, their covariance for the $k^{th}$ PPS $cov(a_{i1}, a_{j1}) = \theta h^2$, in which $\theta$ is the

1102 relatedness scores in terms of identity by state[20]. Depending on the relatedness between a pair

1103 of individuals, $\theta = 1$ for monozygous twins or to a duplicated sample, $\theta = 0.5$ for first-

1104 degree relatives such as parent and offspring or full sibs. In general, for $r^{th}$-degree of

1105 relatives, $E(\theta_r) = 0.5^r$.

1106
1107 The theory presented above provides a theoretical basis for detecting overlapping samples

1108 using PPS other than sharing individual level genotypes. Assuming that each individual has

1109 $K$ independent PPS ($A_i$ having $K$ elements), for individual $i$ and $j$, we can regress $A_i$ on $A_j$,

1110 $A_i = \mu + bA_j + e_{ij}$ **(Equation 12)**

1111   in which $\mu$ is the grand mean, $b$ is the regression coefficient, and $e_{ij}$ is the residual. $E(b) =$

1112   $\frac{cov(A_i, A_j)}{var(A_j)} = \theta_r$. $E(b) = 0$ if individual $i$ is not correlated with individual $j$, $E(b) = 0.5$ for

1113   first-degree relatives, and $E(b) = 1$ if individual $i$ and $j$ are genetically same, say an

1114   overlapping sample or the homozygous twins. The sampling variance of $b$ is $\sigma_b^2 =$

1115   $\frac{\sigma_{A_i}^2 - \sigma_{A_j}^2 \theta_r^2}{\sigma_{A_j}^2 K} = \frac{1 - \theta_r^2}{K}$. Under the null distribution for no related or overlapping samples,

1116   $b \sim N(0, \frac{1}{K})$. The residual $e_{ij}$ accounts the discordant genotypes, including missing genotypes

1117   and genotyping or imputation errors. For current GWAS data, after quality control, the

1118   discordant rate is often smaller than 1%.

1119
1120   If now we have $C$ cohorts for which the individual genotypes of which cannot be disclosed to

1121   the central analysis hub, overlap between cohorts can be identified if PPS are supplied. By

1122   regressing their PPS to each other the overlapping individuals could be detected if $b \approx \theta_r$.

1123   Assuming there are $N_c$ samples in each cohort, a total of $N = \sum_{c_1=1}^{C} \sum_{c_2 > c1}^{C} N_{c_1} \times N_{c_2}$

1124   regressions need to be carried out as defined in Equation 12. If we want to control the

1125   experiment-wise type I error rate $\alpha$ under the null hypothesis and type II error rate $\beta$ (with

1126   power= $1 - \beta$) for $b = \theta_r$, the required number of pseudo profile scores for each individual

1127   is

1128   $K \geq \left( \frac{z_{(1-\beta)} \sqrt{1-b^2} + z_{(1-\alpha)}}{b} \right)^2$     **(Equation 13)**

1129   in which $z_{(1-\beta)}$ and $z_{(1-\alpha)}$ are $z$ scores under the given $p$-values at the subscripts. To

1130   accommodate technical errors, such as missing genotypes and genotype error, a cutoff of 0.95

1131   for $b$ is adopted for detection of overlapping samples, and $0.4 \sim 0.45$ for detecting first-degree

1132   relative.

1133
1134   The standardization of genotypes can either use the allele frequency from each cohort, or

1135   from a reference sample. Throughout the study, we used the allele frequency calculated from

1136   WTCCC bipolar disorder cohort as the reference, and using it as an approximation to

1137   standardize genotypes for all cohorts in comparison.

1138
1139   **Workflow for PPSR.** Given the statistical method for detecting overlapping samples as

1140   described above, the whole workflow for detecting can be split into three steps

1141   (**Supplementary Fig. 9**).

62

1142

1143 **In step 1, the required type I and type II error rates are defined and from that the**

1144 **required number of pseudo profiles to be generated.** The GWAMA central analyst selects

1145 consensus SNP markers across cohorts, and determines additive effects matrix $S$ that will be

1146 used to generate pseudo profile scores for each cohort. In order to avoid strand issues, the loci

1147 having palindromic loci (A/T alleles or G/C alleles) are excluded.

1148
1149 **In step 2, each cohort generates PPSR for each individual with the set of consensus**

1150 **markers and the marker weights received from the GWAMA coordinator.** After

1151 generate the PPS, they send them back to the coordinator. This will be a file that contains N

1152 rows and K columns with pseudo-profile scores.

1153
1154 **In step 3, the coordinator runs PPSR for each sample in a cohort on each PPS generated**

1155 **for another cohort.** The final product of running PPSR is to generate a $n_i \times n_j$ matrix for a

1156 pair of cohorts, which have $n_i$ and $n_j$ samples respectively. For each pair of individuals in

1157 comparison, we take the one from cohort $i$ as the response variable and from cohort $j$ as the

1158 predictor variable in PPSR. In principle, swapping the response variable and the predictor

1159 variable do not affect the performance of PPSR. Each entry, the regression coefficient of

1160 PPSR, in the $n_i \times n_j$ matrix represents genetic similarity for these pair of individuals in

1161 comparison. Once the regression coefficients are above the threshold, it indicates there are

1162 samples duplicated. The central analyst can then request each cohort that is implicated in

1163 containing samples that are also in other cohorts to drop those samples, without revealing

1164 where the duplication occurred.

1165

1166 **Privacy issues when using PPSR.** As the exchange of the PPS is within a meta-analysis

1167 facility, it is not as vulnerable as that of releasing the GWAS summary to the public domain

1168 as discussed in previous studies[21–23]. However, as PPS are generated from genotypes, it is

1169 worth to consider whether the PPS will reveal individual genotype information, or can be

1170 decoded from PPS. As a demonstration for the principle-of-proof, we consider to reverse

1171 Equation 11 to estimate genotypes. We consider the case where the additive effect matrix in

1172 Equation 11 is known, otherwise it is nearly impossible to recover genotype information.

1173 Given the workflow of PPSR, the analysts who coordinate the meta-analysis know the

1174 additive effect matrix, $S$ in Equation 11, and receive PPS from each cohort have the

1175 information to decode genotypes that are employed to generate PPS.

1176

1177   After reversing Equation 11, using the standard regression method, the genotype in each

1178   locus can be estimated as

1179   $A_i = \mu + g_{im} \times s_{.m} + e$          **(Equation 14)**

1180   In detail,

$$\begin{bmatrix} a_{i1} \\ a_{i2} \\ \vdots \\ a_{iK} \end{bmatrix} = \mu + g_m \begin{bmatrix} s_{im} \\ s_{im} \\ \vdots \\ s_{im} \end{bmatrix} + e$$

1181   in which $s_{.m}$ is the $m^{th}$ column in the additive effects matrix in Equation 11. Although

1182   $E(g_{im}) = g_{im}$, which is an unbiased estimate of the genotype, its sampling variance is

1183   $\sigma_{g_{im}} = \sqrt{\frac{\Sigma_{m=1}^{M}[1-(1-p_m)^2]\sigma_{s.m}^2}{K}}$. The sampling variance can be further written as $\sigma_{g_{im}} =$

1184   $\sqrt{\frac{M}{K}E(\mathcal{P}_m)}$ because $\sigma_{s.m}^2 = 1$ and $[1 - (1 - p_m)^2]$ is denoted as $\mathcal{P}_m$. The greater the ratio

1185   between $\frac{M}{K}$ and $E(\mathcal{P}_m)$, the larger the sampling variance, and consequently the lower

1186   probability to construct the real genotype.

1187

1188   Without loss of generality, the accuracy of the estimated $\hat{g}$, a continuous variable, and $g$, a

1189   discrete variable with values of 2, 1, and 0, can be measure using the squared correlation

1190   $(R^2)^{24}$,

1191   $R^2 = \frac{E[\sigma_g^2]}{E[\sigma_g^2] + \frac{M}{K}E[g^2]}$          **(Equation 15)**

1192   in which $E(g^2)$ and $E(\sigma_g^2)$ are:

$$E(g^2) = \tilde{p}_{AA}x_{AA}^2 + \tilde{p}_{Aa}x_{Aa}^2 + \tilde{p}_{aa}x_{aa}^2$$

$$E(\sigma_g^2) = \tilde{p}_{AA}(x)(x_{AA} - 2p)^2 + \tilde{p}_{Aa}(x_{Aa} - 2p)^2 + \tilde{p}_{aa}(x_{aa} - 2p)^2$$

1193   $x_{AA} = 2, x_{Aa} = 1$, and $x_{aa} = 0$ if $A$ is the reference allele, and $\tilde{p}_{AA}, \tilde{p}_{Aa}$, and $\tilde{p}_{aa}$ are

1194   weighted frequency given the distribution of $g$. $f = \tilde{p}_{AA} + 0.5\tilde{p}_{Aa}$.

1195

1196   When the reference allele frequency follows a uniform distribution between $(a_1, a_2)$,

1197   assuming that the loci follow Hardy-Weinberg proportions, $p_{AA} = p^2$, $p_{Aa} = 2pq$, and

1198   $p_{aa} = q^2$, in which $p$ follows a uniform distribution between $a_1$ and $a_2$ and $q = 1 - p$.

$$p_{AA} = \int_{a_1}^{a_2} p^2 = \frac{1}{3}p^3\big|_{a_1}^{a_2} = \frac{1}{3}(a_2^3 - a_1^3)$$

64

$$p_{Aa} = \int_{a_1}^{a_2} 2pq = \left(p^2 - \frac{2}{3}p^3\right)\Big|_{a_1}^{a_2} = (a_2^2 - a_1^2) - \frac{2}{3}(a_2^3 - a_1^3)$$

$$p_{aa} = \int_{1-a_2}^{1-a_1} q^2 = \frac{1}{3}q^3\Big|_{1-a_2}^{1-a_1} = \frac{1}{3}[(1-a_1)^3 - (1-a_2)^3]$$

1199    and $\tilde{p}_{AA} = \frac{p_{AA}}{p_{AA}+p_{Aa}+p_{aa}}, \tilde{p}_{Aa} = \frac{p_{Aa}}{p_{AA}+p_{Aa}+p_{aa}}$, and $\tilde{p}_{aa} = \frac{p_{aa}}{p_{AA}+p_{Aa}+p_{aa}}$.

1200

1201    If the reference allele frequency follows a uniform distribution between (0, 0.5), $R^2 =$

1202    $\frac{\frac{5}{12}}{\frac{5}{12}+\frac{2M}{3K}} = \frac{5}{5+8\frac{M}{K}}$.

1203

1204    Given $M$ loci with MAF of 0.5, the expected frequencies for $AA$, $Aa$, and $aa$ are $\tilde{p}_{AA} = 0.25$,

1205    $\tilde{p}_{Aa} = 0.5, \tilde{p}_{aa} = 0.25$, and $f = 0.5$. $E(g^2) = 1.5$, and $E(\sigma_g^2) = 0.5$. Plugging them in to

1206    the Equation 13 leads to $R^2 = \frac{0.5}{0.5+1.5\frac{M}{K}} = \frac{1}{1+3\frac{M}{K}}$.

1207    Equation 13 can be rewritten as $R^2 = \frac{1}{1+\varphi\frac{M}{K}}$, in which $\varphi = 3$ if MAF is 0.5, and $\varphi = 1.6$ if

1208    MAF in nearly from a uniform distribution. From Equation 13, it is easy to calculate the ratio

1209    between the number of markers and the number of PPS given a controlled $R^2$,

1210    $\frac{M}{K} \geq \frac{1-R^2}{\varphi R^2}$       **(Equation 16)**

1211

1212    For uniform distribution of MAF, if $R^2 \leq 0.1$ is set as the threshold, $\frac{M}{K} \geq 5.4$; if $R^2 \leq 0.05$,

1213    $\frac{M}{K} \geq 11.4$, and if $R^2 \leq 0.01$, $\frac{M}{K} \geq 59.4$. In general, the higher the ratio between $M$ and $K$, the

1214    less information can be inferred. We suggest $\frac{M}{K} \geq 5{\sim}10$ may be sufficient.