

1 **Boosting the power of rare variant association studies by imputation**
2 **using large-scale sequencing population**

3 **Running head:** rare variant imputation and genome-wide association study

4

5 Jinglan Dai, MS^{1#}; Yixin Zhang, MS^{1#}; Zaiming Li, MS¹; Hongru Li, MS¹; Sha Du,
6 PhD^{1,2}; Dongfang You, PhD^{1,3}; Ruyang Zhang, PhD^{1,4}; Yang Zhao, PhD^{1,4}; Zhonghua
7 Liu⁵, PhD; David C. Christiani, MD^{6,7}; Feng Chen, PhD^{1,2,3*}; Sipeng Shen, PhD^{1,2,4*}

8

9 ¹Department of Biostatistics, Center for Global Health, School of Public Health,
10 Nanjing Medical University, Nanjing 211166, China

11 ²Jiangsu Key Lab of Cancer Biomarkers, Prevention and Treatment, Jiangsu
12 Collaborative Innovation Center for Cancer Personalized
13 Medicine, Nanjing Medical University, 211166, Nanjing, China

14 ³China International Cooperation Center of Environment and Human Health, Nanjing
15 Medical University

16 ⁴Key Laboratory of Biomedical Big Data of Nanjing Medical University, Nanjing
17 211166, China

18 ⁵Department of Biostatistics, Columbia University, New York, NY, USA

19 ⁶Department of Environmental Health, Harvard T.H. Chan School of Public Health,
20 Harvard University, Boston, MA 02115, USA

21 ⁷Pulmonary and Critical Care Division, Massachusetts General Hospital, Department
22 of Medicine, Harvard Medical School, Boston, MA 02114, USA

23 [#]These authors contributed equally to this study.

24

25 ***Send correspondence to:**

26 Dr. Sipeng Shen, SPH Building Room 406, 101 Longmian Avenue, Nanjing, Jiangsu

27 211166, China. Email: sshen@njmu.edu.cn

28 Dr. Feng Chen, SPH Building Room 412, 101 Longmian Avenue, Nanjing, Jiangsu

29 211166, China. Email: fengchen@njmu.edu.cn

30

31

32 **Abstract**

33 Rare variants can explain part of the heritability of complex traits that are ignored by
34 conventional GWASs. The emergence of large-scale population sequencing data
35 provides opportunities to study rare variants. However, few studies systematically
36 evaluate the extent to which imputation using sequencing data can improve the power
37 of rare variant association studies. Using whole genome sequencing (WGS) data (n =
38 150,119) as the ground truth, we described the landscape and evaluated the
39 consistency of rare variants in SNP array (n = 488,377) imputed from TOPMed or
40 HRC+UK10K in the UK Biobank, respectively. The TOPMed imputation covered
41 more rare variants, and its imputation quality could reach 0.5 for even extremely rare
42 variants. TOPMed-imputed data was closer to WGS in all MAC intervals for three
43 ethnicities (average Cramer's $V > 0.75$). Furthermore, association tests were performed
44 on 30 quantitative and 15 binary traits. Compared to WGS data, the identified rare
45 variants in TOPMed-imputed data increased 27.71% for quantitative traits, while it
46 could be improved by ~10-fold for binary traits. In gene-based analysis, the signals in
47 TOPMed-imputed data increased 111.45% for quantitative traits, and it identified 15
48 genes in total, while WGS only found 6 genes for binary traits. Finally, we
49 harmonized SNP array and WGS data for lung cancer and epithelial ovarian cancer.
50 More variants and genes could be identified than from WGS data alone, such as
51 *BRCA1*, *BRCA2*, and *CHRNA5*. Our findings highlighted that incorporating rare
52 variants imputed from large-scale sequencing populations could greatly boost the
53 power of GWAS.

54

55 **Keywords:** rare variant; genotype imputation; whole-genome sequencing; association
56 analysis

57

58 **Introduction**

59 Conventional genome-wide association studies (GWASs) are well-powered to detect
60 thousands of common variants associated with human traits and diseases [1-3].
61 However, GWASs underrepresent rare variants due to the limitations of SNP array [4].
62 Rare variants have larger effects and behave differently from common variants to
63 explain a fraction of human traits or disease heritability [5, 6]. With the emergence of
64 next-generation sequencing (NGS) [e.g., whole genome sequencing (WGS)] based on
65 biobank-level populations, massive rare variants [minor allele frequency (MAF) <0.01]
66 or ultra-rare variants (MAF <0.0001) could be accurately captured [7]. For example,
67 several studies aimed at identifying novel germline variants in cancer or rare diseases
68 have benefited from NGS [8, 9].

69 However, given that the application of deep WGS is limited by cost for large-scale
70 populations [10], SNP array is a cost-effective and major approach for human
71 genomics exploration thus far, such as UCLA ATLAS [11], China Kadoorie Biobank
72 [12], FinnGen [13], and various case-control GWASs [14]. Even if there are
73 high-quality population-based sequencing data, the incidence number of complex
74 disease cases is often low and the power for identifying rare variants is insufficient
75 compared to case-control studies based on SNP array. For example, the UK Biobank
76 (UKB) has whole-exome sequencing (WES) data with a large sample size ($n \approx$
77 450,000) [15, 16], but most cancer cases have a sample size of less than 10,000.
78 Nevertheless, unlike sequencing data, SNP array-based genotyping can only represent
79 a fraction of the genetic variation in the genome. Thus, it is a more appropriate
80 method to power association analysis to detect rare variants that were not genotyped
81 to the SNP array by imputing from the external high-quality large-scale sequencing
82 panels [10, 17], such as the UK10K, Haplotype Reference Consortium (HRC) and
83 Trans-Omics for Precision Medicine (TOPMed). In TOPMed, 290 million (~97%)
84 variants had an MAF $<1\%$, which might be strong candidates for association analyses

85 [18, 19]. However, few studies have focused on evaluating the accuracy of rare
86 variant imputation based on sequencing data panel and the extent to which it can
87 enhance the power of rare variant exploration.

88 Here, we present a comparative study that leverages the WGS data of 150,119
89 individuals from the UKB as the ground truth to evaluate the coverage and accuracy
90 of rare variant imputation using HRC+UK10K and TOPMed. We investigated the
91 association tests of WGS data and imputed variants with 30 biochemistry biomarkers
92 and 15 complex diseases, respectively. Finally, we harmonized WGS data and diverse
93 SNP array data to reveal the rare variant signals in lung cancer and epithelial ovarian
94 cancer.

95

96

97 **Results**

98 *Landscape of rare variant imputation results*

99 We analyzed WGS data from 150,119 individuals and genotype imputation data
100 generated from the HRC+UK10K and TOPMed reference panels (n=488,377) from
101 the UKB, respectively (Figure 1). The intersection of variants between WGS and
102 imputed data is shown in Figure 2A. TOPMed imputed a substantially larger number
103 of genetic variants that could detect approximately 26% single nucleotide variants
104 (SNVs) of WGS data, which was higher than the 10.7% detected in HRC+UK10K.
105 The TOPMed-imputed data could detect 22 million singleton and doubleton variants,
106 which was far more than HRC+UK10K-imputed data but was still fewer than 333
107 million variants in WGS data. For ultra-rare variants, the TOPMed-imputed variants
108 could reach 66.6% of WGS and were ~3.4-fold of that in HRC+UK10K. For rare and
109 common variants, the number of imputed variants was close to WGS (Figure 2B and
110 Table S1). For protein-coding variants, TOPMed-imputed data detected 44,440
111 loss-of-function (LoF) variants, 509,759 synonymous variants and 880,589 missense

112 variants, which was at least four-fold greater than that in the HRC+UK10K
113 imputation (Figure 2B).

114 To deeply understand the performance of different imputed data compared to WGS,
115 all variants were divided into five intervals according to the minor allele count (MAC).
116 As expected, we observed an increase in the fraction of variants as MAC increased.
117 When $MAC \in (10, 20]$, TOPMed-imputed data had already detected ~65% variants of
118 WGS data, HRC+UK10K-imputed data could only detect ~20% (Figure 2C and Table
119 S2).

120 For the imputed data, post-imputation quality control (QC) is generally performed
121 based on INFO/RSQ filtering. The number of different INFO/RSQ variants in the
122 imputed data is presented in Figure 2D. TOPMed-imputed data achieved substantially
123 higher coverage of rare variants than HRC+UK10K, and as the MAC increased, the
124 proportion of high-quality TOPMed-imputed variants with higher INFO/RSQ
125 increased. Even for extremely rare variants ($MAC \in (0, 5)$), the average INFO/RSQ of
126 imputed data could also reach 0.5 (Table S3).

127 *Evaluation of genotype consistency between WGS and imputed data*

128 We used WGS data as the ground truth to conduct correlation analysis on the two
129 imputed datasets to evaluate the consistency of rare variants in three ethnicities
130 (White, Asian and African) (Table S4). TOPMed had a better imputation performance
131 that was closer to WGS than HRC+UK10K in all MAC intervals, and the average
132 Cramer's V imputed by TOPMed was above 0.75 in three ethnicities (Figure 3A and
133 Table S5). Even for $MAC \in (0, 5)$, the Cramer's V of TOPMed-imputed data exceeded
134 0.6 in all ethnicities. The difference in Cramer's V of various intervals between
135 TOPMed-imputed data and HRC+UK10K-imputed data reached a maximum in
136 extremely rare variants ($MAC \in (0, 5)$), which was 0.33 for White, 0.25 for Asian and
137 0.17 for African (Figure 3A and Table S5).

138 We also described the relationship between INFO/RSQ and Cramer's V (Figure 3B).
139 With the increase in INFO/RSQ bins, Cramer's V also increased in all three
140 ethnicities, indicating that the variants with high imputation quality could be closer to
141 the WGS data to some extent. Moreover, for the same INFO/RSQ bin, using different
142 reference panels would obtain different Cramer's V, and the consistency of
143 TOPMed-imputed data was higher than that of HRC+UK10K in each bin. Even when
144 the imputation quality was not particularly good ($\text{INFO/RSQ} \in (0.3, 0.4]$), the
145 TOPMed-imputed data had an advantage over that from HRC+UK10K which showed
146 a stable consistency (all Cramer's $V > 0.5$ in three ethnicities) (Table S6). Therefore,
147 the rare variants imputed from TOPMed were applicable for subsequent association
148 analysis.

149 ***Rare variant association analysis for quantitative biochemistry*** 150 ***biomarkers***

151 We performed association tests between rare variants ($\text{MAF} < 0.01$) and 30
152 biochemistry biomarkers in WGS ($n=150$ k) and imputed data ($n=150$ k and 480 k) of
153 European descent. In the single-variant tests of 16 traits, compared to WGS data, the
154 two imputed data ($n=480$ k) could detect more significant rare variants;
155 HRC+UK10K-imputed data could improve 4.7% while the TOPMed-imputed data
156 could improve 27.71% (Figure 4A-4B). However, the number of significant rare
157 variants found in WGS data was more than that of imputed data in 12 traits, which
158 might be due to the extremely large number of sequencing sites in WGS. Meanwhile,
159 we noticed that the imputed data of $n=150$ k had a weaker ability to detect rare
160 variants than WGS (Figure S2). On average, the number of significant rare variants in
161 the TOPMed-imputed data was 661 and in the HRC+UK10K-imputed data was 450,
162 which were less than 2,479 in WGS data (Table S7).

163 We also evaluated the gene-based tests results to capture the effects of rare variants.
164 On average, 22 genes were found in the TOPMed-imputed data and 18 genes were

165 found in the HRC+UK10K-imputed data (n=480 k), while the WGS data identified
166 only 7 genes for all traits (Table S8). The average improvement ratio of the
167 HRC+UK10K-imputed data was 55.81%, and that of the TOPMed-imputed data
168 reached 111.45% (Figure 4B). When comparing the WGS data with the imputed data
169 of n=150 k, the number of significant genes in the TOPMed-imputed data decreased
170 by 27.10%, while that in the HRC+UK10K-imputed data decreased by 33.66%
171 (Figure S2).

172 Based on these association tests of the significant variants, we depicted the
173 relationship between the chi-square statistics of the association analysis in WGS data
174 and two imputed data (n=150 k). The test statistics in different imputed data both
175 presented a strong relationship with that in WGS data, which meant their association
176 test statistics were highly correlated (on average, Pearson's $r = \square 0.97$ for
177 TOPMed-imputed data and 0.95 for HRC+UK10K-imputed data) (Figure 4C and
178 Table S9). Therefore, the rare variant-trait associations using imputed genotypes were
179 robust.

180 ***Rare variant association analysis for complex diseases***

181 We also conducted single-variant tests and gene-based tests for 15 complex diseases,
182 including ten chronic diseases and five cancers (Table S10). In single-variant tests, the
183 number of significant rare variants in the imputed data (n=480 k) was approximately
184 10-fold that in WGS data, owing to its larger sample sizes with sufficient disease
185 cases (Figure S1 and Table S13). On average, 22 significant rare variants were
186 identified in TOPMed-imputed data, and 24 were identified in HRC+UK10K-imputed
187 data. However, only 2 variants could be found in the WGS data on average (Figure 5
188 and Table S11). We also found that for the imputed data with a sample size of 150 k,
189 its ability to detect rare variants was a slightly weaker than that of WGS (Figure S3).
190 Therefore, it could be expected that using large-scale SNP array-based data imputed
191 from sequencing panels would greatly improve the ability of rare variant findings

192 compared to WGS data.

193 In the gene-based tests, the number of significant genes in the imputed data was
194 higher than that in the WGS data, and the performance of the TOPMed-imputed data
195 was better. HRC+UK10K-imputed data could find 14 significant genes in total, and
196 the TOPMed-imputed data could find 15 genes, while the WGS data only found six
197 genes owing to low cases (Figure 5 and Table S12, S14). In addition, we also noticed
198 that some classical genes identified in the two imputed datasets were roughly
199 consistent but undetected in WGS. For example, *IL33* has been certified to be
200 associated with asthma and a rare mutation in *IL33* could decrease the risk of asthma
201 [20, 21]. For prostate cancer, the gene *CHEK2* which was found only in
202 TOPMed-imputed data has been confirmed to be a risk gene containing rare variants
203 [22]. Most genes found in the TOPMed-imputed data had supporting association or
204 biological evidence from previous studies (Table S21).

205 Based on the results above, we concluded that incorporating SNP array imputed from
206 large-scale sequencing populations could enhance the ability to detect rare
207 variant-trait associations.

208 ***Powering the rare variant hits by harmonizing WGS and SNP array***

209 Furthermore, we leveraged UKB WGS and several SNP array data imputed by the
210 TOPMed reference panel to perform association tests for lung cancer (LC) and
211 epithelial ovarian cancer (EOC) (Figure 6A, 6D). Compared with limited cases in
212 WGS ($n_{LC} = 1,494$; $n_{EOC} = 628$), the SNP array could provide sufficient cases ($n_{LC} =$
213 $23,168$; $n_{EOC} = 23,822$) that were 16-fold and 38-fold for LC and EOC respectively.

214 In single-variant tests, no rare variants could pass the genome-wide significance level
215 in WGS. However, when combined with SNP-array data, we detected 12 SNVs and
216 22 SNVs, mapping to six and four independent genetic regions, respectively (Figure
217 6B, 6E, Table S17-S18). In gene-based tests, no genes could pass $P < 2.5 \times 10^{-6}$ in WGS.

218 In the meta-analysis, we identified *BRCA2* and *CHRNA5* in lung cancer, while

219 *BRCA2* and *BRCA1* were significant in EOC (Figure 6B, 6E, Table S19-S20). Using
220 different *P* value thresholds, the WGS+SNP-array strategy could detect more variants,
221 especially under stringent thresholds (e.g., 10^{-7} in single-variant and 10^{-5} in gene-based)
222 (Figure 6C, 6F).

223

224

225 **Discussion**

226 It is a common practice to directly apply sequencing data or imputed genotype data to
227 GWAS. Limited by the cost of NGS technology, biobank-level sequencing data are
228 still in the minority [23], resulting in the difficulty of using rare variant analysis which
229 provides important genomic architectures. Large-scale GWASs tend to use imputed
230 genotype data, but few studies systematically compare the power of imputed data for
231 rare variants. Therefore, evaluating the consistency and association power
232 performance of rare variants imputed from sequencing data can provide better
233 guidance and reference for rare variant exploration in complex traits.

234 In this study, we comprehensively conducted a series of comparative analyses for rare
235 variant imputation using SNP array data from the UKB based on TOPMed or
236 HRC+UK10K against the WGS data of 150,119 individuals as the ground truth. We
237 described the rare variant imputation results, and further carried out the correlation
238 analysis in three ethnicities (White, Asian, and African) to evaluate how close rare
239 variant imputation could be to WGS data. We further investigated association analysis
240 on 30 biochemistry markers and 15 complex diseases to explain the ability of imputed
241 rare variants to improve the association analysis at both single-variant and gene-based
242 levels.

243 Our summarized description shows that the TOPMed panel could impute more
244 high-quality rare variants than HRC+UK10K. Although the imputed variants were far
245 inferior to WGS in singleton and doubleton variants, the total number of rare variant

246 imputed by TOPMed also had a certain scale. Even for the extremely rare variants
247 (MAC<5), the average imputation quality of TOPMed-imputed variants was
248 acceptable (INFO/RSQ≥0.5). Although not all rare variants can be estimated reliably,
249 they could provide a crucial supplement in addition to WGS.

250 For the three ethnicities, we observed that rare variants imputed from TOPMed were
251 closer to WGS than those imputed from HRC+UK10K in each MAC interval. In
252 Europeans, the consistency of rare variants with MAC≤20 in HRC+UK10K-imputed
253 data is poor, which is not recommended for subsequent application. In addition, the
254 overall correlation between TOPMed-imputed data and WGS was slightly stronger in
255 Africans than in the other two ethnicities, which may be due to the low sample size of
256 Africans (n<1,000) in WGS, and unbalanced ethnic sample sizes affect the stability of
257 Cramer's V. Nevertheless, the TOPMed-imputed data still had a stable imputation
258 performance that was closer to WGS data, indicating its reliability to be applied in
259 rare variant association studies.

260 In the association analysis of quantitative and binary traits, although the ability of
261 both imputed data to identify significant rare variants was weaker than that of WGS
262 data in the same population (n = 150k), we could still identify a few rare variants or
263 gene sets. When the sample size of the imputed data increased to n = 480 k, the ability
264 to identify significant rare variant could improve, but this improvement was slightly
265 different in distinct types of traits. For quantitative traits, the power of WGS data for
266 finding significant rare variants was basically sufficient, and the two imputed data had
267 limited improvement compared to WGS data. However, for binary traits, the number
268 of disease cases in WGS data was far from sufficient to ensure optimistic power,
269 causing fewer significant rare variants to be found. To improve the power, external
270 data with additional disease cases must be supplemented. Our study indeed showed
271 that larger SNP array data imputed from sequencing data could greatly improve the
272 ability to find significant rare variants.

273 Although the natural population cohorts are large with adequate sample sizes, the
274 disease cases are generally insufficient, which is known as a case-control imbalance
275 problem and leads to low statistical power to identify rare variants [24]. Here, we
276 harmonize WGS population and multiple case-control design GWASs and
277 successfully identify classical rare variants and gene sets in two cancers, including the
278 well-known *BRCA2* [25], *BRCA1* [26], and *CHRNA5* [27], which could not be
279 detected in UKB WGS. Therefore, integrating SNP-array data through imputation is
280 of great significance for discovering rare variants.

281 There are several strengths in our study. First, we comprehensively described the
282 landscape of the imputed rare variants from different sequencing panels in terms of
283 variant amount, coverage, and imputed quality. Second, we leveraged WGS data as
284 the ground truth to perform correlation analysis by different ethnicities, and analyzed
285 their consistency with imputed rare variants. We demonstrated the feasibility of using
286 imputed rare variants for association analysis. Third, we conducted various
287 association tests on 45 traits of the UKB and harmonized six GWAS case-control
288 datasets totally. Through both single-variant and gene-based tests, we quantitatively
289 explained the improvement of SNP array-based data imputed from sequencing data in
290 rare variant association studies.

291 The results we presented here also have some limitations. First, our genotype
292 consistency evaluation was only conducted on the UKB SNP array (~50,000 UK
293 BiLEVE Axiom array and ~450,000 UK Biobank Axiom array) that did not consider
294 other array types, which might influence the results. Second, we currently impute the
295 genotype in TOPMed's imputation online server (97,256 reference samples) due to
296 computing power limitations, but it is convenient and feasible for most genomic
297 studies. Theoretically, the results of imputation in larger populations (e.g., UKB
298 whole WGS population) should be more accurate [28], which needs further studies to
299 confirm. Third, we attempted to integrate the SNP-array data and sequencing data in
300 two exemplary cancers. The effect of integrating other various diseases in large-scale

301 population cohorts ($n \geq 100,000$) needs further evaluation.

302 In conclusion, incorporating rare variants imputed from large-scale sequencing
303 populations can greatly enhance the power of GWAS research. Our study shows that
304 combining multi-source genomic data with a sufficient number of cases can
305 accurately identify a wider range of rare variants, allowing us to take full advantage of
306 different types of genetic data and to gain a deeper understanding of the causal links
307 between rare genetic variants in human complex traits and diseases.

308

309

310 **Methods**

311 *UK Biobank population and phenotypic data collection*

312 The UK Biobank (UKB) is a population-based prospective cohort of individuals aged
313 40–69 years, enrolled between 2006 and 2010 [29]. The work described herein was
314 approved by the UK Biobank under applications no. 92675 and 83445. All phenotypic
315 data were accessed in July 2022.

316 Blood biochemistry data were collected from Category 17518. The UK Biobank
317 embarked on a project to measure a wide range of biochemical markers in biological
318 samples collected at baseline (2006-2010) in all 500,000 participants. All 30
319 biochemistry biomarkers were included as quantitative traits.

320 Health-related outcomes were ascertained via individual record linkage to national
321 cancer and mortality registries and hospital in-patient encounters. Cancer diagnoses
322 were coded by International Classification of Diseases version 10 (ICD-10) codes
323 from data fields 41270 (Diagnoses - ICD10), 41202 (Diagnoses - main ICD10), and
324 40001 (primary cause of death: ICD10). Individuals with at least one recorded
325 incident diagnosis were defined as cases. We included 15 common chronic diseases or
326 cancers as binary traits, including insulin-dependent diabetes mellitus,

327 non-insulin-dependent diabetes mellitus, obesity, depressive episode, hypertension,
328 chronic ischemic heart disease, heart failure, COPD, asthma, cholelithiasis, bladder
329 carcinoma, breast carcinoma, non-Hodgkin lymphoma, prostate cancer, and
330 Melanoma.

331 ***UK Biobank genomics data collection***

332 We collected whole-genome sequencing (WGS) of 150,119 people in the UK Biobank
333 (data field 23352) [30], which were sequenced to an average coverage of 32.5× (at
334 least 23.5× per individual) using Illumina NovaSeq sequencing machines at deCODE
335 Genetics (90,667 individuals) and the Wellcome Trust Sanger Institute (59,452
336 individuals). Sequence reads were mapped to the human reference genome GRCh38
337 using BWA. SNPs and short indels were jointly called over all individuals using
338 GraphTyper (v2.7.1) [31], which provided more accurate genotype calls. This
339 constitutes a set of high-quality variants, including 585,040,410 single-nucleotide
340 variants (SNVs) and 58,707,036 indels.

341 We also collected imputed genotype data based on two mainstream reference panels.
342 The first is the Haplotype Reference Consortium (HRC) and UK10K haplotype
343 resource, which increases the number of testable variants over 100-fold to ~96 million
344 variants (data field 22828). The genetic data was imputed using two different
345 reference panels. The HRC panel (64,976 haplotypes) was used wherever possible,
346 but for SNPs not in that reference panel the UK10K + 1000 Genomes panel (12,570
347 haplotypes) was used [32]. The raw positions are in GRCh37 coordinates and then
348 lifted over to GRCh38.

349 The second is imputation from genotype using the TOPMed R2 panel (97,256 deeply
350 sequenced genomes), performed by the TOPMed Informatics Research Center [33].
351 After phasing the UK Biobank genetic data (carried out on 81 chromosomal chunks
352 using Eagle v.2.4), the phased data were converted from GRCh37 to GRCh38 using
353 LiftOver. Imputation was performed using Minimac4 v1.0.2. Imputation was

354 performed in 1 Mb chunks and merged back together by chromosome.

355 ***Quality control of imputed SNP array data***

356 If an imputed variant on N samples has an imputation quality metric scored at α , it
357 implies that the statistical power of association tests is approximately equivalent to αN
358 perfectly observed genotype data [34]. To perform GWAS on the UK Biobank data
359 with ~480,000 samples, it is typical to use variants with imputation quality higher
360 than 0.3, equivalent to ~150,000 perfectly observed samples (WGS sample size). Thus,
361 markers with poor imputation quality were not retained in the subsequent analysis
362 (excluding the Minimac4 imputation quality metric $RSQ < 0.3$ or IMPUTE
363 imputation quality metric $INFO < 0.3$).

364 ***Genotype consistency evaluation***

365 In two imputed datasets (TOPMed and HRC+UK10K) and WGS data, we divided the
366 population into three ethnicities --- White, Asian, African based on ethnic data field
367 21000. To evaluate the consistency between the imputed genotype and WGS data, we
368 matched the imputed data (n=480k) with the WGS data (n=150k) according to the
369 individual ID of each ethnicity.

370 Then, we used PLINK2.0 to calculate the number of minor allele count (MAC) and
371 minor allele frequency (MAF) of genetic variants in each ethnicity. Meanwhile, based
372 on MAC and MAF, genetic variants were classified into five categories, namely (0,5],
373 (5,10], (10,20], (20,50], >50 and $MAF < 0.01$ (the interval was abbreviated as > 50 in
374 the subsequent analysis). We chose Cramer's V as the metric to evaluate the
375 consistency of imputed variants and sequencing variants:

$$\text{Cramer's V} = \sqrt{\chi^2 / nm}$$

376 where χ^2 represents the chi-square value of the current contingency table, n represents
377 the sample size, and m represents the smaller value of the two degrees of freedom (r-1)

378 or (c-1) of the two variables.

379 ***Single-variant and gene-based association tests***

380 Single-variant and gene-based association analyses were performed using REGENIE
381 v3.2.6 [35]. It is a machine-learning based method to fit a whole-genome regression
382 model for quantitative and binary phenotypes. Quantitative traits were rank-based
383 inverse normal transformed. Saddlepoint approximation (SPA) was used to handle
384 extreme case-control imbalance of binary traits. All association tests were performed
385 in the population of European descent (ethnicity = “White”) only.

386 The single-variant association tests included rare variants with $MAF < 0.01$. The
387 variants were functionally annotated using Variant Effect Predictor (VEP) software
388 [36]. The genome-wide significant threshold of single-variant tests was set to $P < 5 \times$
389 10^{-8} .

390 For gene-based analysis, we included rare and ultra-rare variants with $MAF < 0.01$.
391 Three genetic models were considered: loss of function (LoF), LoF+missense and
392 LoF+missense+synonymous. Of all the combinations, we reported the association
393 results with the lowest P value to collectively capture a wide range of genetic
394 architectures [5]. The genome-wide significance threshold of gene-based tests was set
395 to $P < 2.5 \times 10^{-6}$.

396 In the association analyses of quantitative traits, we adjusted the covariates including
397 age, sex, and the top ten principal components (PCs). In the association analyses of
398 binary traits, we adjusted the covariates including age, sex (excluding sex-specific
399 diseases), body mass index (BMI), smoking status (binary), drinking status (binary),
400 and the top ten principal components (PCs).

401 ***SNP array data collection for lung cancer and ovarian cancer***

402 We also collected large-scale SNP array data of lung cancer and epithelial ovarian
403 cancer to harmonize them with UKB WGS to improve the power.

404 For lung cancer, we investigated the association of imputed genetic variants with lung
405 cancer in three additional cohorts from the Prostate, Lung, Colorectal, and Ovarian
406 (PLCO) cancer screening trial [37], the International Lung Cancer OncoArray
407 Consortium (ILCCO-OncoArray) [38], and the Transdisciplinary Research in Cancer
408 of the Lung (TRICL) research team [39] (Table S15).

409 For ovarian cancer, we collected data in three additional cohorts from The Follow-up
410 of Ovarian Cancer Genetic Association and Interaction Studies (FOCI)-OncoArray
411 [40], FOCI-Exome Chip [41], and Consortium of Investigators of Modifiers of
412 BRCA1/2 (CIMBA) [42] (Table S16).

413 Samples excluded were those who lacked disease status, were second-degree relatives
414 or closer having identity by descent (IBD) > 0.2 or had low-quality DNA (call rate $<$
415 95%), or sex inconsistency, or were non-European. SNPs were removed if they met of
416 the following criteria: (1) sex chromosome, (2) MAF < 0.05 , (3) call rate $< 95\%$, and
417 (4) Hardy-Weinberg equilibrium (HWE) test $P < 1.00 \times 10^{-7}$ in controls or $P <$
418 1.00×10^{-12} in cases.

419 All genotype data were imputed on the TOPMed online imputation server. Poorly
420 imputed SNVs with imputation quality score $R^2 < 0.3$ and SNVs on sex chromosomes
421 were excluded from the analyses. The effect sizes and 95% confidence interval (CI) of
422 genes were estimated by burden tests and then summarized by fix-effects
423 meta-analysis.

424

425

426

427 **Declarations**

428 **Consent for publication**

429 All authors have reviewed and approved this manuscript.

430 **Data availability**

431 UK Biobank data is available from: <https://www.ukbiobank.ac.uk/>.

432 ILCCO-Oncoarray data is available from:
433 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-ibin/study.cgi?study_id=phs001273.v3.p2
434 [p2](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-ibin/study.cgi?study_id=phs001273.v3.p2)

435 TRICL data is available from:
436 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001681.v1.p1
437 [p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001681.v1.p1)

438 PLCO data is available from:
439 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi-study_id=phs001286.v2.p2
440 [p2](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi-study_id=phs001286.v2.p2).

441 FOCl-Exome Chip data is available from:
442 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001131.v1.p1
443 [p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001131.v1.p1)

444 FOCl-OncoArray Chip: data is available from:
445 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001882.v1.p1&phv=435464&phd=&pha=&pht=10492&phvf=&phdf=&phaf=&phtf=&dssp=1&consent=&temp=1
446 [p1&phv=435464&phd=&pha=&pht=10492&phvf=&phdf=&phaf=&phtf=&dssp=1&consent=&temp=1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001882.v1.p1&phv=435464&phd=&pha=&pht=10492&phvf=&phdf=&phaf=&phtf=&dssp=1&consent=&temp=1)
447 [consent=&temp=1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001882.v1.p1&phv=435464&phd=&pha=&pht=10492&phvf=&phdf=&phaf=&phtf=&dssp=1&consent=&temp=1)

448 CIMBA: data is available from:
449 https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001321.v1.p1
450 [p1](https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs001321.v1.p1)

451 **Online resources**

452 TOPMED imputation server: <https://imputation.biodatacatalyst.nhlbi.nih.gov/>.

453 gnomAD: <http://gnomad.broadinstitute.org/>.

454 dnSNP: <https://www.ncbi.nlm.nih.gov/snp/>.

455 **Code availability**

456 The codes that support our findings are available from the corresponding author by a
457 request.

458 **Conflicts of interest statements/Financial Disclosure statement**

459 The authors report no conflicts of interest.

460 **Funding**

461 National Natural Science Foundation of China (82373685 and 82373685 to S.S.,
462 82220108002 to F.C., 82103946 and 82173620 to Y.Z.), Natural Science Foundation
463 of the Jiangsu Higher Education Institutions of China (21KJB330004 to S.S.), and US
464 NIH (NCI) grant #U01CA209414 to DCC.

465 **Author contributions**

466 SS and FC contributed to the study design. SS and JD contributed to data collection.

467 JD performed statistical analyses and interpretation. JD and SS drafted the manuscript.

468 HL, SD, DY, YZ, RZ, ZL, and DC revised the final manuscript. All authors approved

469 the final version of the manuscript.

470 **Acknowledgement**

471 We want to acknowledge the participants and investigators of UK Biobank, TOPMed,

472 PLCO, ILCCO-OncoArray, TRICL, FOCI, and CIMBA.

473

474

475 Reference

- 476 1. Chiou J, Geusz RJ, Okino ML, Han JY, Miller M, Melton R, Beebe E, Benaglio P, Huang S,
477 Korgaonkar K *et al*: **Interpreting type 1 diabetes risk with genetics and single-cell**
478 **epigenomics**. *Nature* 2021, **594**(7863):398-402.
- 479 2. Bouatia-Naji N, Bonnefond A, Cavalcanti-Proenca C, Sparso T, Holmkvist J, Marchand M,
480 Delplanque J, Lobbens S, Rocheleau G, Durand E *et al*: **A variant near MTNR1B is**
481 **associated with increased fasting plasma glucose levels and type 2 diabetes risk**. *Nat*
482 *Genet* 2009, **41**(1):89-94.
- 483 3. Lyssenko V, Nagorny CL, Erdos MR, Wierup N, Jonsson A, Spegel P, Bugliani M, Saxena R,
484 Fex M, Pulizzi N *et al*: **Common variant in MTNR1B associated with increased risk of**
485 **type 2 diabetes and impaired early insulin secretion**. *Nat Genet* 2009, **41**(1):82-88.
- 486 4. Kaisinger LR, Kentistou KA, Stankovic S, Gardner EJ, Day FR, Zhao Y, Mörseburg A, Carnie
487 CJ, Zagnoli-Vieira G, Puddu F *et al*: **Large-scale exome sequence analysis identifies sex-**
488 **and age-specific determinants of obesity**. *Cell genomics* 2023, **3**(8):100362.
- 489 5. Wang Q, Dhindsa RS, Carss K, Harper AR, Nag A, Tachmazidou I, Vitsios D, Deevi SVV,
490 Mackay A, Muthas D *et al*: **Rare variant contribution to human disease in 281,104 UK**
491 **Biobank exomes**. *Nature* 2021, **597**(7877):527-532.
- 492 6. Karczewski KJ, Solomonson M, Chao KR, Goodrich JK, Tiao G, Lu W, Riley-Gillis BM, Tsai
493 EA, Kim HI, Zheng X *et al*: **Systematic single-variant and gene-based association testing**
494 **of thousands of phenotypes in 394,841 UK Biobank exomes**. *Cell Genom* 2022,
495 **2**(9):100168.
- 496 7. Sazonovs A, Barrett JC: **Rare-Variant Studies to Complement Genome-Wide Association**
497 **Studies**. *Annu Rev Genomics Hum Genet* 2018, **19**:97-112.
- 498 8. Shen S, Li Z, Jiang Y, Duan W, Li H, Du S, Esteller M, Shen H, Hu Z, Zhao Y *et al*: **A**
499 **Large-Scale Exome-Wide Association Study Identifies Novel Germline Mutations in**
500 **Lung Cancer**. *Am J Respir Crit Care Med* 2023, **208**(3):280-289.
- 501 9. Greene D, Pirri D, Frudd K, Sackey E, Al-Owain M, Giese APJ, Ramzan K, Riaz S,
502 Yamanaka I, Boeckx N *et al*: **Genetic association analysis of 77,539 genomes reveals rare**
503 **disease etiologies**. *Nature Medicine* 2023, **29**(3):679-688.
- 504 10. Tam V, Patel N, Turcotte M, Bosse Y, Pare G, Meyre D: **Benefits and limitations of**
505 **genome-wide association studies**. *Nat Rev Genet* 2019, **20**(8):467-484.
- 506 11. Johnson R, Ding Y, Bhattacharya A, Knyazev S, Chiu A, Lajonchere C, Geschwind DH,
507 Pasaniuc B: **The UCLA ATLAS Community Health Initiative: Promoting precision**
508 **health research in a diverse biobank**. *Cell genomics* 2023, **3**(1):100243.
- 509 12. Walters RG, Millwood IY, Lin K, Schmidt Valle D, McDonnell P, Hacker A, Avery D, Edris A,
510 Fry H, Cai N *et al*: **Genotyping and population characteristics of the China Kadoorie**
511 **Biobank**. *Cell genomics* 2023, **3**(8):100361.
- 512 13. Sun BB, Kurki MI, Foley CN, Mechakra A, Chen CY, Marshall E, Wilk JB, Chahine M,
513 Chevalier P, Christé G *et al*: **Genetic associations of protein-coding variants in human**
514 **disease**. *Nature* 2022, **603**(7899):95-102.
- 515 14. Abdellaoui A, Yengo L, Verweij KJH, Visscher PM: **15 years of GWAS discovery: Realizing**

- 516 **the promise.** *Am J Hum Genet*, **110**(2):179-194.
- 517 15. Backman JD, Li AH, Marcketta A, Sun D, Mbatchou J, Kessler MD, Benner C, Liu D, Locke
518 AE, Balasubramanian S *et al*: **Exome sequencing and analysis of 454,787 UK Biobank**
519 **participants.** *Nature* 2021, **599**(7886):628-634.
- 520 16. Sun BB, Kurki MI, Foley CN, Mechakra A, Chen CY, Marshall E, Wilk JB, Biogen Biobank T,
521 Chahine M, Chevalier P *et al*: **Genetic associations of protein-coding variants in human**
522 **disease.** *Nature* 2022, **603**(7899):95-102.
- 523 17. Barton AR, Sherman MA, Mukamel RE, Loh PR: **Whole-exome imputation within UK**
524 **Biobank powers rare coding variant association and fine-mapping analyses.** *Nat Genet*
525 2021, **53**(8):1260-1269.
- 526 18. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, Kang HM,
527 Fuchsberger C, Danecek P, Sharp K *et al*: **A reference panel of 64,976 haplotypes for**
528 **genotype imputation.** *Nat Genet* 2016, **48**(10):1279-1283.
- 529 19. Wuttke M, Konig E, Katsara MA, Kirsten H, Farahani SK, Teumer A, Li Y, Lang M, Gocmen
530 B, Pattaro C *et al*: **Imputation-powered whole-exome analysis identifies genes associated**
531 **with kidney function and disease in the UK Biobank.** *Nat Commun* 2023, **14**(1):1287.
- 532 20. Grotenboer NS, Ketelaar ME, Koppelman GH, Nawijn MC: **Decoding asthma: translating**
533 **genetic variation in IL33 and IL1RL1 into disease pathophysiology.** *J Allergy Clin*
534 *Immunol* 2013, **131**(3):856-865.
- 535 21. Smith D, Helgason H, Sulem P, Bjornsdottir US, Lim AC, Sveinbjornsson G, Hasegawa H,
536 Brown M, Ketchem RR, Gavalda M *et al*: **A rare IL33 loss-of-function mutation reduces**
537 **blood eosinophil counts and protects from asthma.** *PLoS Genet* 2017, **13**(3):e1006659.
- 538 22. Southey MC, Goldgar DE, Winqvist R, Pylkas K, Couch F, Tischkowitz M, Foulkes WD,
539 Dennis J, Michailidou K, van Rensburg EJ *et al*: **PALB2, CHEK2 and ATM rare variants**
540 **and cancer risk: data from COGS.** *J Med Genet* 2016, **53**(12):800-811.
- 541 23. Zhou W, Kanai M, Wu K-HH, Rasheed H, Tsuo K, Hirbo JB, Wang Y, Bhattacharya A, Zhao
542 H, Namba S *et al*: **Global Biobank Meta-analysis Initiative: Powering genetic discovery ac**
543 **ross human disease.** *Cell genomics*, **2**(10):100192.
- 544 24. Zhou W, Nielsen JB, Fritsche LG, Dey R, Gabrielsen ME, Wolford BN, LeFaive J, VandeHaar
545 P, Gagliano SA, Gifford A *et al*: **Efficiently controlling for case-control imbalance and**
546 **sample relatedness in large-scale genetic association studies.** *Nature genetics*,
547 **50**(9):1335-1341.
- 548 25. Wang Y, McKay JD, Rafnar T, Wang Z, Timofeeva MN, Broderick P, Zong X, Laplana M,
549 Wei Y, Han Y *et al*: **Rare variants of large effect in BRCA2 and CHEK2 affect risk of**
550 **lung cancer.** *Nature genetics*, **46**(7):736-741.
- 551 26. Hyman DM, Spriggs DR: **Unwrapping the implications of BRCA1 and BRCA2 mutations**
552 **in ovarian cancer.** *JAMA*, **307**(4):408-410.
- 553 27. Olfson E, Saccone NL, Johnson EO, Chen LS, Culverhouse R, Doheny K, Foltz SM, Fox L,
554 Gogarten SM, Hartz S *et al*: **Rare, low frequency and common coding variants in**
555 **CHRNA5 and their contribution to nicotine dependence in European and African**
556 **Americans.** *Mol Psychiatry*, **21**(5):601-607.
- 557 28. Hofmeister RJ, Ribeiro DM, Rubinacci S, Delaneau O: **Accurate rare variant phasing of**

- 558 **whole-genome and whole-exome sequencing data in the UK Biobank. *Nature genetics*,
559 **55(7):1243-1249.****
- 560 29. Szustakowski JD, Balasubramanian S, Kvikstad E, Khalid S, Bronson PG, Sasson A, Wong E,
561 Liu D, Wade Davis J, Haefliger C *et al*: **Advancing human genetics research and drug**
562 **discovery through exome sequencing of the UK Biobank.** *Nature genetics*, **53(7):942-948.**
- 563 30. Halldorsson BV, Eggertsson HP, Moore KHS, Hauswedell H, Eiriksson O, Ulfarsson MO,
564 Palsson G, Hardarson MT, Oddsson A, Jensson BO *et al*: **The sequences of 150,119 genomes**
565 **in the UK Biobank.** *Nature* 2022, **607(7920):732-740.**
- 566 31. Eggertsson HP, Jonsson H, Kristmundsdottir S, Hjartarson E, Kehr B, Masson G, Zink F,
567 Hjorleifsson KE, Jonasdottir A, Jonasdottir A *et al*: **GraphTyper enables population-scale**
568 **genotyping using pangenome graphs.** *Nat Genet* 2017, **49(11):1654-1660.**
- 569 32. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, Motyer A, Vukcevic D,
570 Delaneau O, O'Connell J *et al*: **Genome-wide genetic data on ~500,000 UK Biobank**
571 **participants.** *bioRxiv* 2017:166298.
- 572 33. Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, Taliun SAG, Corvelo A,
573 Gogarten SM, Kang HM *et al*: **Sequencing of 53,831 diverse genomes from the NHLBI**
574 **TOPMed Program.** *Nature* 2021, **590(7845):290-299.**
- 575 34. Marchini J, Howie B: **Genotype imputation for genome-wide association studies.** *Nat Rev*
576 *Genet* 2010, **11(7):499-511.**
- 577 35. Mbatchou J, Barnard L, Backman J, Marcketta A, Kosmicki JA, Ziyatdinov A, Benner C,
578 O'Dushlaine C, Barber M, Boutkov B *et al*: **Computationally efficient whole-genome**
579 **regression for quantitative and binary traits.** *Nat Genet* 2021, **53(7):1097-1103.**
- 580 36. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, Flicek P, Cunningham F:
581 **The Ensembl Variant Effect Predictor.** *Genome biology* 2016, **17(1):122.**
- 582 37. Tammemagi MC, Katki HA, Hocking WG, Church TR, Caporaso N, Kvale PA, Chaturvedi
583 AK, Silvestri GA, Riley TL, Commins J *et al*: **Selection criteria for lung-cancer screening.**
584 *N Engl J Med* 2013, **368(8):728-736.**
- 585 38. McKay JD, Hung RJ, Han Y, Zong X, Carreras-Torres R, Christiani DC, Caporaso NE,
586 Johansson M, Xiao X, Li Y *et al*: **Large-scale association analysis identifies new lung**
587 **cancer susceptibility loci and heterogeneity in genetic susceptibility across histological**
588 **subtypes.** *Nat Genet* 2017, **49(7):1126-1132.**
- 589 39. Wang Y, McKay JD, Rafnar T, Wang Z, Timofeeva MN, Broderick P, Zong X, Laplana M,
590 Wei Y, Han Y *et al*: **Rare variants of large effect in BRCA2 and CHEK2 affect risk of**
591 **lung cancer.** *Nat Genet* 2014, **46(7):736-741.**
- 592 40. Phelan CM, Kuchenbaecker KB, Tyrer JP, Kar SP, Lawrenson K, Winham SJ, Dennis J, Pirie
593 A, Riggan MJ, Chornokur G *et al*: **Identification of 12 new susceptibility loci for different**
594 **histotypes of epithelial ovarian cancer.** *Nat Genet* 2017, **49(5):680-691.**
- 595 41. Permut JB, Pirie A, Ann Chen Y, Lin HY, Reid BM, Chen Z, Monteiro A, Dennis J,
596 Mendoza-Fandino G, Group AS *et al*: **Exome genotyping arrays to identify rare and low**
597 **frequency variants associated with epithelial ovarian cancer risk.** *Hum Mol Genet* 2016,
598 **25(16):3600-3612.**
- 599 42. Chenevix-Trench G, Milne RL, Antoniou AC, Couch FJ, Easton DF, Goldgar DE, Cimbá: **An**

600 **international initiative to identify genetic modifiers of cancer risk in BRCA1 and BRCA2**
601 **mutation carriers: the Consortium of Investigators of Modifiers of BRCA1 and BRCA2**
602 **(CIMBA). *Breast Cancer Res* 2007, **9**(2):104.**
603
604

605 **Figure legends**

606 **Figure 1. Study workflow**

607 The schematic of the analytical pipeline summarizes the main steps for conducting a
608 comparative study using WGS as the ground truth. We first compared the rare variant
609 imputation for two imputed data (TOPMed-imputed and HRC+UK10K imputed) and
610 WGS data. Then genotype consistency evaluation was investigated on three
611 ethnicities. The rare variant association analysis of 45 traits was carried out after
612 above evaluation in different data. Finally, we harmonized the WGS and SNP-array
613 data for the two cancers.

614 **Figure 2. Landscape of rare variant imputation results**

615 (A) Venn diagram showing the intersection of variants between TOPMed-imputed
616 data, HRC+UK10K-imputed data and WGS data.

617 (B) Bar chart showing the number of five distinct variant types for different data (left).

618 **Variant type description :

619 Singleton: The minor allele occurs only once in the entire sample

620 Doubleton: The minor allele occurs twice in the entire sample

621 Ultra-rare variant: $MAF < 0.0001$ except singleton and doubleton

622 Rare variant: MAF in $0.0001-0.01$

623 Common variant: $MAF > 0.01$

624 Bar chart showing the number of three distinct variant annotations for different
625 data (right).

626 (C) Line chart of variant coverage proportions in different imputed data. The x-axis
627 represents the MAC intervals, and the y-axis represents the proportion of variants
628 that imputed data could detect related to WGS data.

629 (D) The number of variants imputed by different panels and their distribution across
630 INFO/RSQ bins.

631 **Figure 3. Consistency evaluation of the imputed datasets in different ethnicities**

632 (A) Correlation of three ethnicities between imputed data and WGS data in different
633 MAC intervals.

634 (B) Box plot of the relationship between the INFO/RSQ of imputed data and the
635 Cramer's V.

636 **Figure 4. Rare variant association analysis results for biochemistry biomarkers**

637 (A) The number of additional significant rare variants ($MAF < 0.01$) found in the
638 imputed data compared to the WGS data for 30 quantitative traits.

639 (B) The average improvement ratio compared to WGS data for single-variant tests and
640 gene-based tests.

641 (C) Pairwise Pearson correlations between chi-square statistics produced by
642 association tests using imputed data ($n=150$ k) and WGS data respectively.

643 **Figure 5. Rare variant association analysis results for complex diseases**

644 Multiple-trait Manhattan plot of single-variant tests and gene-based tests for 15
645 diseases in different data. The x-axis labels each disease, and the y-axis shows $-\log_{10}P$.
646 The red dotted line represents the significance filtering threshold of the P value ($P <$
647 5×10^{-8} for single-variant tests, $P < 2.5 \times 10^{-6}$ for gene-based tests), and the gray dotted
648 line represents the suggestive filtering threshold of the P value ($P < 5 \times 10^{-6}$ for
649 single-variant tests)

650 **Figure 6. Harmonizing WGS and SNP array to perform association tests for lung**
651 **cancer (left panel) and epithelial ovarian cancer (right panel)**

652 (A-C) Results of lung cancer. (D-F) Results of epithelial ovarian cancer. (A, D) Circos
653 plots of the single-variant and gene-based association results using UKB WGS data or
654 WGS+SNP array data. (B, E) Rare variants and genes identified in each dataset. The

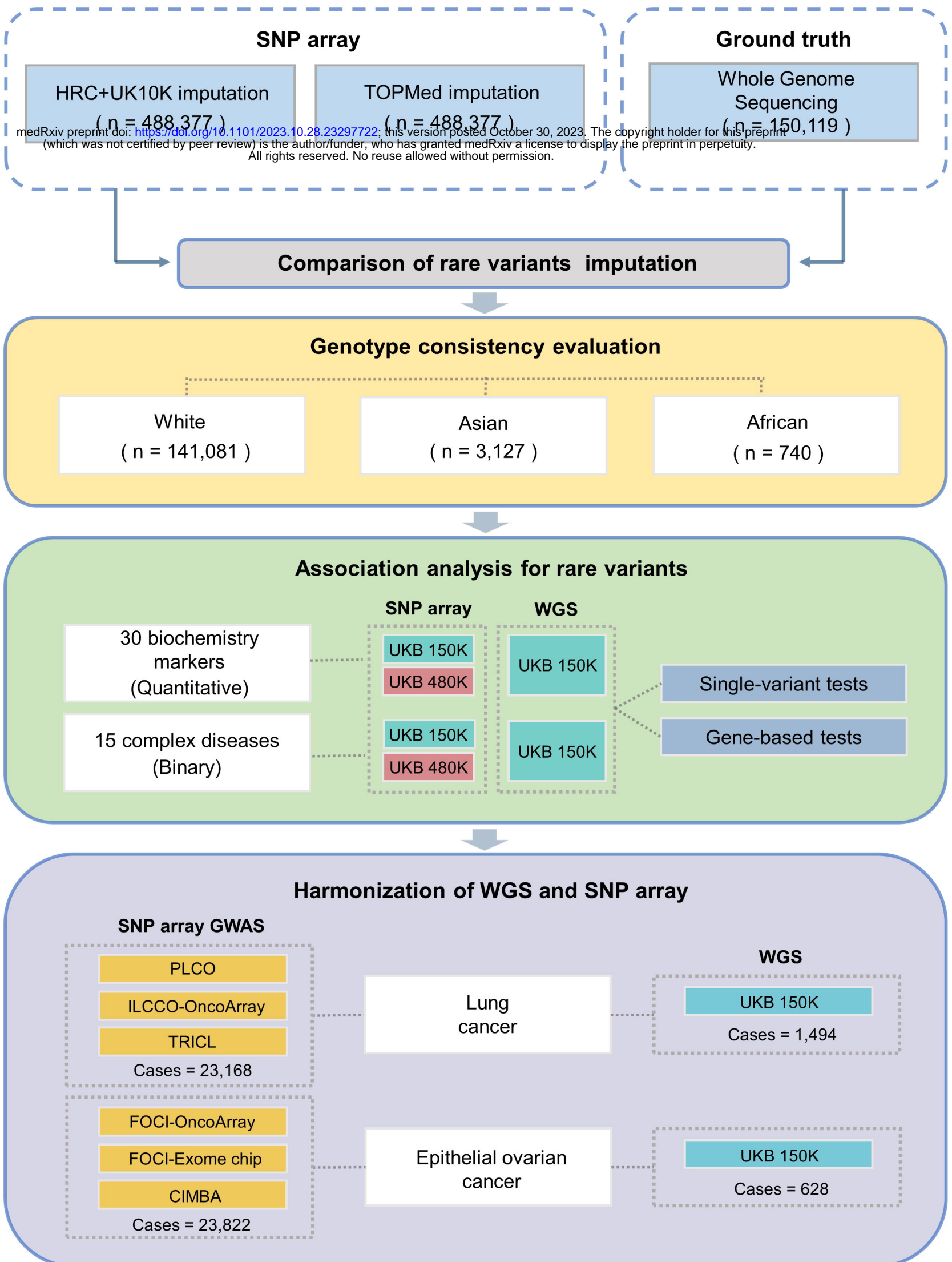
655 blocks are marked if the P values reach nominal significance ($P < 0.05$). (C, F)

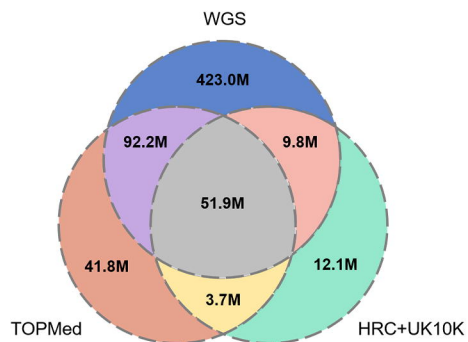
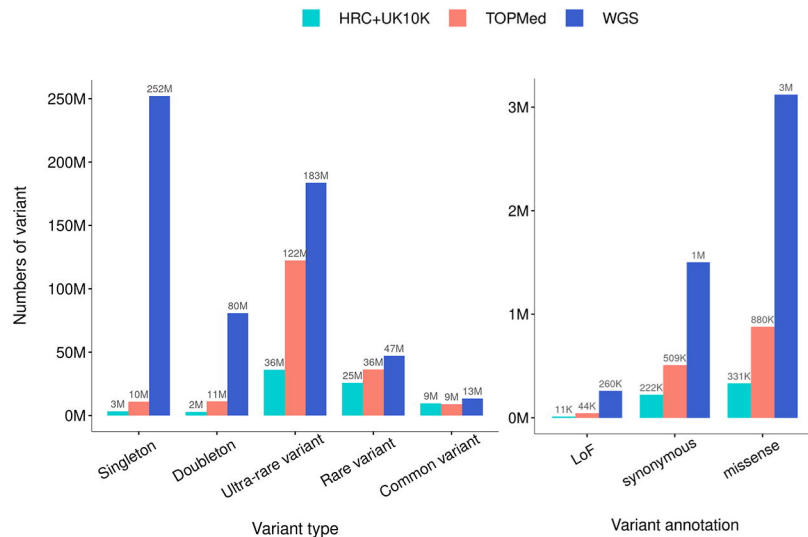
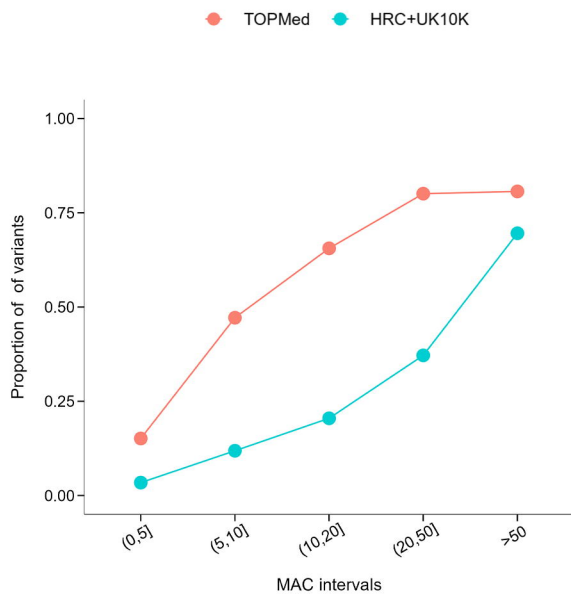
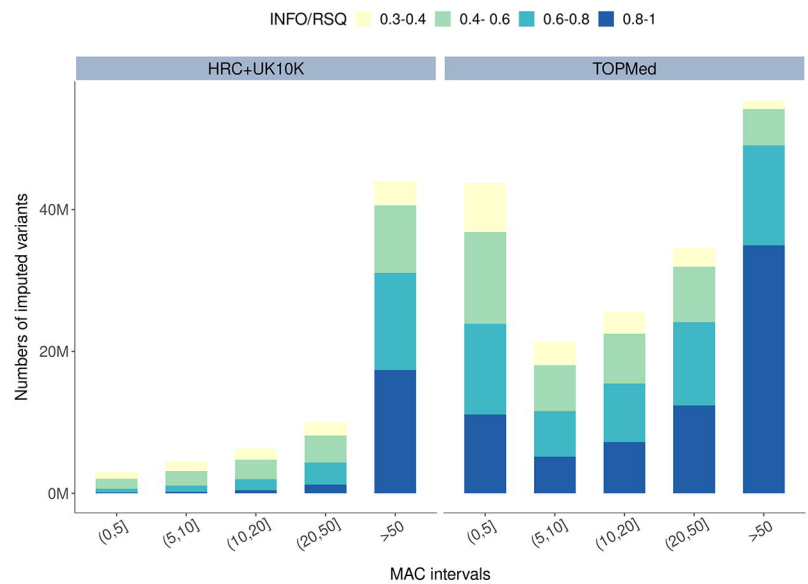
656 Comparison of the identified signals in the single-variant and gene-based association

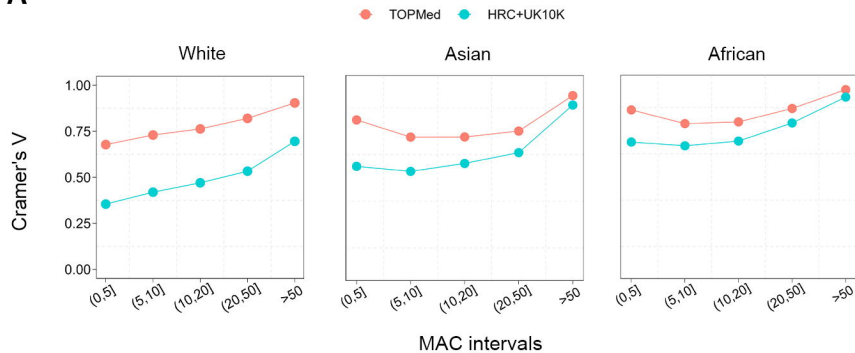
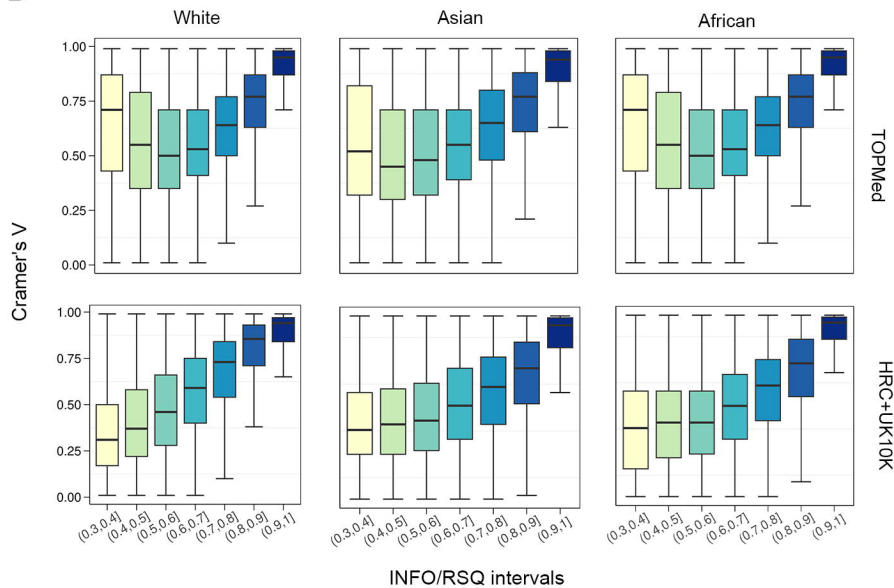
657 tests under different P value thresholds.

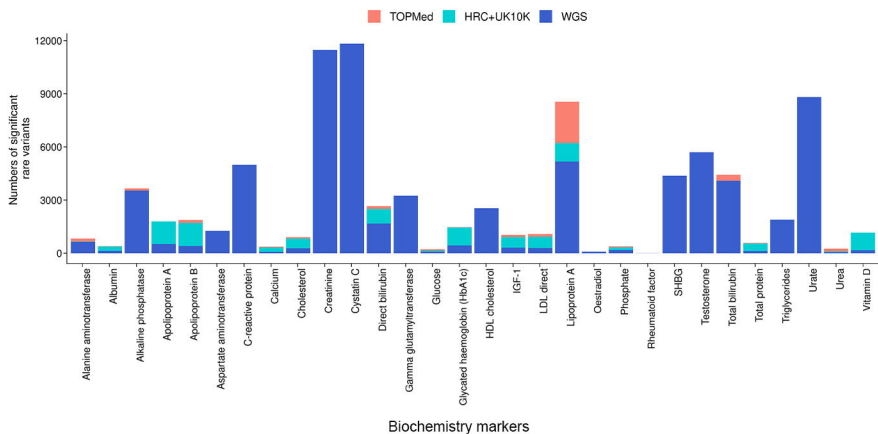
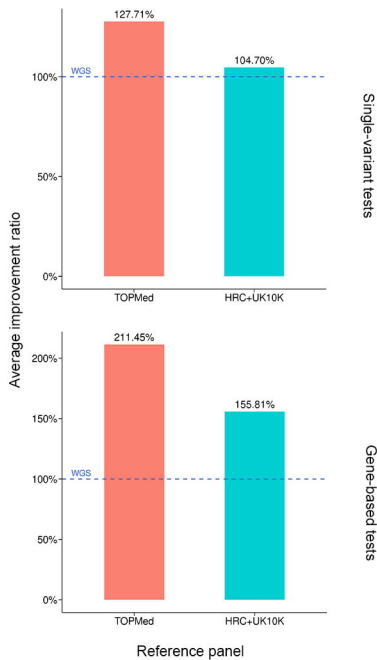
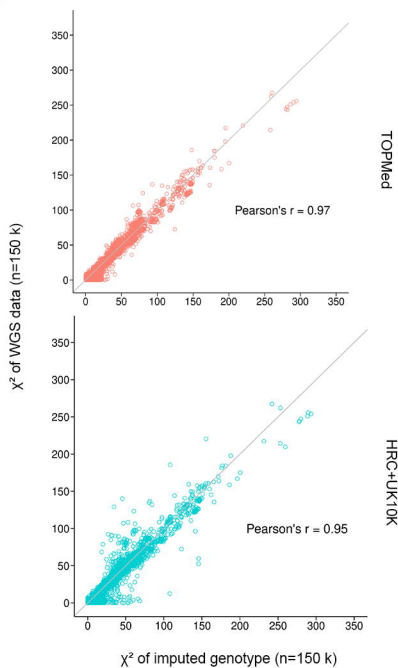
658

659

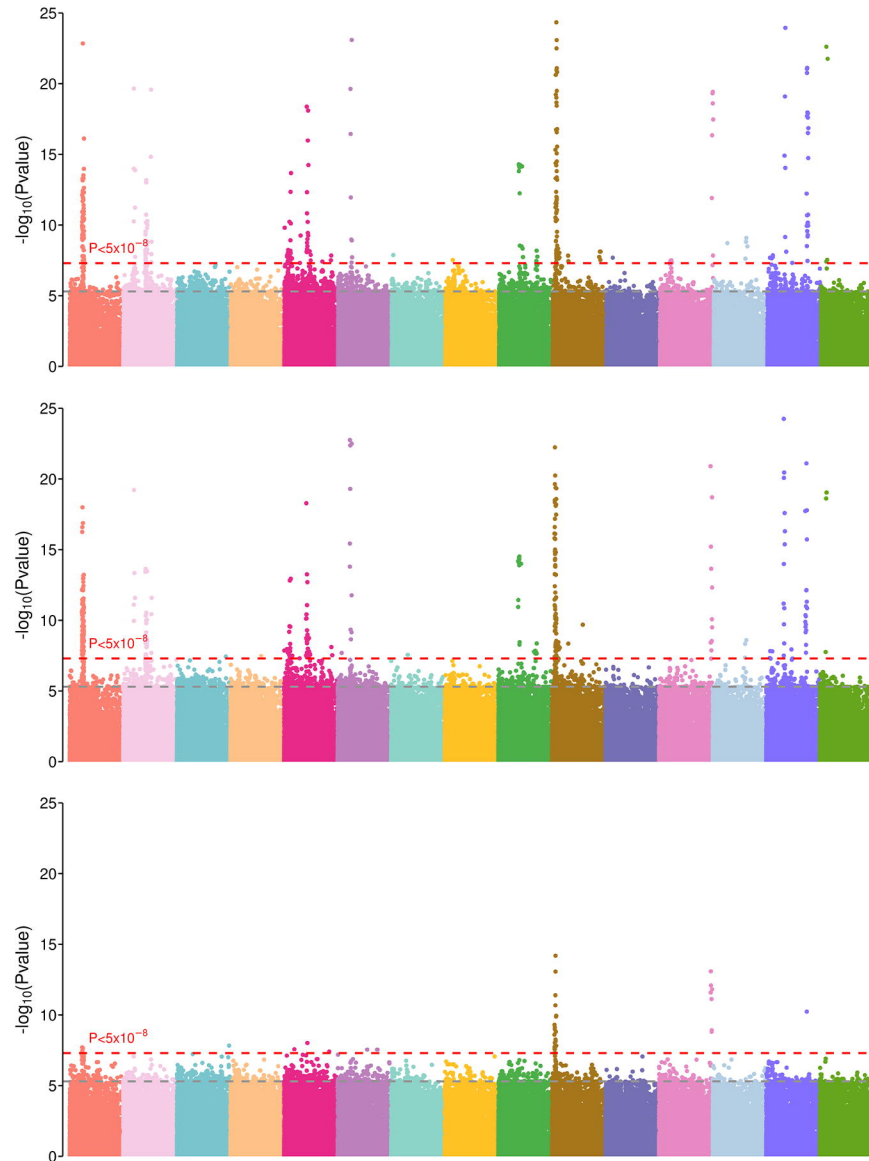


A**B****C****D**

A**B**

A**B****C**

Single-variant tests



Gene-based tests

