# isGWAS: ultra-high-throughput, scalable and equitable inference of genetic associations with disease

Christopher N Foley[1,2*], Zhana Kuncheva[1,2], Riccardo E. Marioni[3], Heiko Runz[4], Benjamin B Sun[4*].

[1]Optima Partners, Edinburgh, UK.

[2]Bayes Centre, University of Edinburgh, Edinburgh, UK.

[3]Centre for Genomic and Experimental Medicine, Institute of Genetics and Cancer, University of Edinburgh, Edinburgh, UK.

[4]Translational Sciences, Biogen Inc., Cambridge, MA, US.

*Correspondence:

Dr Christopher N Foley

Email: chris.foley@optimapartners.co.uk or chris.neal.foley@gmail.com

Dr Benjamin B Sun

Email: benjamin.sun@biogen.com or bbsun92@outlook.com

# Abstract

2   Genome-wide association studies (GWAS) have proven a powerful tool for human geneticists

3   to generate biological insights or hypotheses for drug discovery. Nevertheless, a dependency

4   on sensitive individual-level data together with ever-increasing cohort sample sizes, numbers

5   of variants and phenotypes studied put a strain on existing algorithms, limiting the GWAS

6   approach from maximising potential. Here we present in-silico GWAS (isGWAS), a uniquely

7   scalable algorithm to infer regression parameters in case-control GWAS from cohort-level

8   summary data. For any sample size, isGWAS computes a variant-disease association

9   parameter in ~1 millisecond, or ~11m variants in UK-Biobank within ~4 minutes (~1500-fold

10  faster than state-of-the-art). Extensive simulations and empirical tests demonstrate that

11  isGWAS results are highly comparable to traditional regression-based approaches. We further

12  introduce a heuristic re-sampling algorithm, leapfrog re-sampler (LRS), to extrapolate

13  association results to semi-virtually enlarged cohorts. Owing to significant computational

14  gains we anticipate a broad use of isGWAS and LRS which are customizable on a web

15  interface.

# Main

17  Genome wide association studies (GWAS) have been immensely successful in unravelling the

18  genetic contribution to human disease. Cost-effective genotyping and large biobank cohorts

19  now make it possible to routinely conduct GWAS for tens of millions of variants in hundreds

20  of thousands of individuals across thousands of phenotypes[1]. With the advent of population-

21  scale whole genome sequencing and expansion of GWAS to research participants of non-

22  European ancestries, these numbers can be expected to increase by another magnitude over the

23  next few years[2], [3].

24

25  Current GWAS approaches that compute variant-disease associations in a regression

26  framework, such as PLINK[4], fastGWA[5], BOLT-LMM[6], SAIGE[7] and REGENIE[8],

27  require access to and input from ever increasing individual-level data (ILD). The efforts of

28  individual-level GWAS sample collection, genotyping and data analysis tend to grow as a

29  polynomial function of sample size[7], [8]. Moreover, the exchange of ILD between researchers

30  is non-trivial due partly to data size but increasingly to strict – but essential - data protection

31  regulations, which can limit the scope of collaborative analyses and biological insights

32  gained[9]–[12]. Finally, the substantial computational and financial burden of running massive-

33  scale GWAS, especially for binary disease outcomes, is exacerbating inequity between

34  researchers, typically favouring already well-equipped institutions. There is therefore a pressing

35  need for innovative approaches that help attenuate the increasing resource and financial

36  inequities for conducting contemporary GWAS and to help decide where limited resources

37  should best be allocated.

38

39  Here we present in-silico GWAS (isGWAS), a biobank-scalable and computationally highly

40  efficient algorithm to infer genetic regression parameters in case-control GWAS from just four

41  broadly ascertained cohort-level summary parameters: the counts of cases and controls within

42  a cohort, as well as case and control minor allele frequencies (MAFs). isGWAS is highly

43  parallelisable, exceeding efficiencies of current GWAS analysis tools by several orders of

44  magnitude. Furthermore, we demonstrate that isGWAS yield association summary statistics

45  highly comparable to traditional ILD regression-based approaches through extensive

46  simulations and empirical tests in UK Biobank[13], Biobank Japan[14] and the Psychiatric

47  Genomics Consortium cohort[15]. Owing to the sizeable computational gains, we introduce a

48  heuristic re-sampling algorithm, called the leapfrog re-sampler (LRS), which can confidently

49    extrapolate GWAS results to larger sample sizes, both at a locus or genome-wide scale. Our

50    underlying methodology also leads to several desirable high-utility properties. We release a

51    web tool available to the wider public to conduct customized isGWAS at www.optima-

52    isgwas.com.

# Results

## Genome-wide association testing from sufficient statistics

isGWAS assumes disease-variant associations can be evaluated via a logistic-link function and, similar to widely used methods[7], [8], uses a Firth adjusted maximum likelihood procedure and Newton-Raphson solver to estimate genetic effects, standard errors and association *p*-value[7], [8], [16]. isGWAS' notable advance is based on the insight that the Newton-Raphson procedure can be simplified so that: (a) elements of the Fisher information matrix and score function vector are collapsed by taking expectation over the empirical or *a priori* distribution of a genetic variant; and (b) sufficient statistics – a specific type of summary data - are used as input variables in the score function (see **Methods** for details). We provide several options to initialise the Newton-Raphson algorithm[17], [18] that improve computational performance and reduce analysis time (**Supplementary Information**). In brief, let $y_i$ denote disease status for the $i$-th individual and $g_{ij,M}$ denote the $j$-th genetic variant under model $M$ (e.g., additive, recessive or dominant). The sufficient statistic triple used by isGWAS is:

$$\left\{ T_{1j} = \sum_{i=1}^{N_j} y_i = N_j^*, \qquad T_{2j} = \sum_{i=1}^{N_j} y_i\, g_{ij,M}, \qquad T_{3j} = \sum_{i=1}^{N_j} g_{ij,M} \right\},$$

where $T_{1j}$ is the total number of cases for variant $j$, $T_{2j}$ is the covariance between the outcome $y$ and genotype $g$ for variant $j$ under model $M$, and $T_{3j}$ is the minor allele count for variant $j$ under model $M$. For each variant, data can be provided as either: the sufficient statistic triple $\{T_{1j}, T_{2j}, T_{3j}\}$ plus sample size $N_j$ (necessary for computing standard errors) or separately, on assuming Hardy-Weinberg equilibrium (HWE), as $\{N_j, N_j^*, MAF_{j,M}, MAF_{j,M}^*\}$. Default GWAS analyses assume HWE, making input data widely available[13]–[15], [19]–[21] for researchers to perform isGWAS, replicate or further expand on classical GWAS analyses (**Methods**). If MAFs for cases and controls are supplied, isGWAS will automatically convert to the pair

76 $\{MAF_{j,M}, MAF^*_{j,M}\}$ (**Methods**). After convergence, which is guaranteed for most scenarios by a

77 re-initialization approach (empirically all scenarios converged using isGWAS-Firth), the

78 estimated genetic effect parameter $\hat{\beta}_{j,M}$ and standard error $SE(\hat{\beta}_{j,M})$ are used to construct Wald

79 *p*-values (**Methods**). Additional options include a sample-level likelihood ratio-test or *p*-values

80 computed using sandwich-robust standard errors (**Supplementary Methods**). A simplified

81 illustration highlighting differences and computational advantages of isGWAS against ILD-

82 based genetic association analyses is summarised in **Figure 1**.

83

84 **isGWAS reliably identifies genetic associations across cohorts and diseases**

85 We benchmarked isGWAS in real-data settings and performed simulation studies to compare

86 isGWAS performance and results relative to several existing individual-level data (ILD)-based

87 approaches. Our assessments broadly fall into two categories: (1) methods which require ILD,

88 i.e., REGENIE[8] , logistic and Firth corrected regression[16], and (2) approaches which do not

89 require ILD directly, i.e., the logistic ad-hoc estimator[17] and Fisher's Exact Test (FET)[22].

90 We note FET was successfully leveraged for efficient large-scale GWAS analyses recently[23].

91 Using data from UK Biobank (UKB), we first assessed isGWAS performance against the

92 popular ILD based regression approach REGENIE[8] by deploying both methods for analyses

93 of seven diseases some of which were previously used for establishing GWAS methodology[5],

94 [7], [8]. Second, we evaluated isGWAS's ability to replicate 309 significantly associated

95 variants from the Biobank Japan (BBJ) meta-analysis of 30 diseases[24] using only the

96 published sample-summaries, i.e., without access to ILD, per variant and disease pair. Lastly,

97 by considering multiple iterations of nested schizophrenia meta-analysis GWAS[25]–[27] we

98 assessed isGWAS's ability to accurately predict genomic regions and significant novel

99 associations. The isGWAS sample size of the meta-analysis from 2014 was virtually expanded

100    to match the numbers from the larger 2022 GWAS whilst holding constant the allele frequency

101    and prevalence information in 2014 cohort.

102

103    *Performance against ILD based regression GWAS in UK Biobank*

104    We compared isGWAS results and computational performance against the state-of-the-art ILD-

105    based method REGENIE for seven diseases in UKB: asthma (IC10:J45); atherosclerosis

106    (IC10:I25); colon cancer (IC10:C18); hypertension (IC10:I10); glaucoma (IC10:H40); stroke

107    (IC10:I63); and thyroid gland cancer (IC10:C73). Case-control ratios varied from 1:2 in

108    hypertension to 1:669 in thyroid gland cancer across the diseases (**Supplementary Table 1a**),

109    allowing for the review of performance in near balanced to highly imbalanced case-control

110    settings. To help attenuate the possible influence of confounders when deploying isGWAS,

111    particularly sample relatedness and population structure, we describe and apply additional data

112    quality control (QC) steps before computing the required sample-level sufficient statistics

113    (**Methods**). After additional QC, the total sample size analysed was ~335,000 individuals per

114    disease (**Supplementary Table 1**). This sample was used to perform and contrast analyses in

115    both isGWAS and REGENIE. We review approaches which leverage additional insight from

116    any removed samples in the **Discussion**. For each of the seven diseases, we applied isGWAS

117    (no covariates) and two-step REGENIE with Firth correction (including the covariates age, sex

118    and ethnicity principal components) to ~11 million autosomal variants. Results are presented

119    in **Table 1**, **Figures 2-3, Supplementary Figures 1-15, Supplementary Tables 1-9 and**

120    **Supplementary Files 1-2**.

121

122    Across all seven traits tested we observed close to perfect consistency between REGENIE and

123    isGWAS association results, as illustrated in the mirrored Manhattan and *p-p* plots for asthma

124    (**Figure 2a**) and other diseases (**Supplementary Figures 1-6**). Concordance between

125    REGENIE and isGWAS is further validated by benchmarking accuracy (**Table 1**) and Pearson

126    correlations between estimated $p$-values $(cor(p_{isG}, p_R)) > 0.94$ at $log_{10}$ scale

127    (**Supplementary Tables 2-3**). The results are consistent for varying prevalence levels

128    (**Supplementary Figures 1-6**) and are not affected by covariate adjustment (**Supplementary**

129    **Tables 2-3**). The consistency translated to the regional locus level. This is exemplified by a

130    locus zoom plot of the *FLG2* gene region for asthma (**Figure 2b**) where isGWAS not only

131    nominated the identical GWAS lead variants but also largely recapitulated the overall

132    association pattern identified through REGENIE. This observation is consistent across the lead

133    independent loci from the asthma GWAS (**Supplementary Figure 7**) and translates to all other

134    diseases studied (**Supplementary Figure 8**). For a comprehensive numerical comparison of

135    association results, we took REGENIE derived *p*-values as the ground truth, retaining all SNPs

136    with $p < 0.01$ and setting the true positive threshold as $p < 5 \times 10^{-8}$ (excluding stroke

137    (**Supplementary Figure 4**) which did not yield any significant associations). We computed the

138    accuracy, false positive (FPR), true positive (TPR) and false discovery rates (FDR) of isGWAS

139    (**Figure 3a** and **Table 1**). Accuracy of isGWAS was $\geq 99.98\%$ for each disease, highlighting

140    excellent overall correspondence between methods. The FPR was low, i.e., $FPR \lesssim 10^{-5}$, and

141    TPR was generally good at $> 88\%$ - excluding hypertension which had a $TPR = 0.63$. The

142    FDR was below $\leq 5\%$ for each disease, revealing that the positive predictive value of isGWAS

143    was greater than $95\%$.

144

145    Importantly, isGWAS and REGENIE results differed for two broad categories: (i) the

146    estimation of genetic effect sizes; and (ii) computational performance. When non-confounding

147    covariates are excluded $(\beta_{no\_cov})$ or included $(\beta_{cov})$ in a model, previous and extensive

148    investigations of effect size estimates in logistic regression deduce that $|\beta_{no\_cov}| \leq |\beta_{cov}|$, i.e.,

149    regression estimates are smaller in magnitude when excluding covariates but the null-

150 hypothesis of no association is maintained[28], [29]. Overall, we replicate these results in our

151 analyses. isGWAS computed effect sizes are smaller in absolute value, but largely concordant

152 with covariate-adjusted REGENIE. Moreover, we fail to reject the null hypothesis for the same

153 variants almost always between methods - suggesting that the isGWAS QC helps attenuate

154 issues of population confounding (**Figures 2d-e, Supplementary Figures 1d-e – 6d-e, 9, 10**

155 **and 11**). An investigation of the performance of isGWAS without removing related individuals

156 highlights potential expansion of isGWAS beyond the recommended QC (**Methods**,

157 **Supplementary Information**), but further investigation – possibly leveraging the re-sampling

158 potential of isGWAS - is required on the reliability of isGWAS in family-based cohorts and

159 ethnically diverse populations (**Supplementary Figures 12-13, Supplementary Tables 5-7**).

160 In our full-QC analyses, all estimated effects between isGWAS and REGENIE were observed

161 to be in the same direction, and the correlation between estimates was on average

162 $cor(\beta_{isG}, \beta_{REGENIE}) \approx 0.7$ (**Supplementary Tables 1-3**). The relative drop in the correlation

163 between effect estimates ($\approx 0.7$) and p-values ($\gtrsim 0.94$) is anticipated[28] and can be explained

164 on noting that, across all diseases more precise effect estimates (i.e., those with smaller standard

165 errors) have stronger concordance between approaches (**Figure 2e and Supplementary**

166 **Figures 1d-6d**). Overall, we found that at least 98% of isGWAS and REGENIE confidence

167 intervals (CI) overlap, (**Supplementary Table 4**). When effect estimates are viewed as a

168 function of MAF, the absolute value of REGENIE-derived estimates seemingly increase (along

169 with standard errors) as MAF decreases across all scenarios. This contrasts with isGWAS where

170 the relationship between MAF and effect size is less clear: fewer variants with low MAF are

171 associated with relatively larger effect sizes. However, the correspondingly narrower standard

172 errors guarantee the same significance p-values as REGENIE. The isGWAS derived

173 distribution of effect sizes is consistent with the hypothesis of a flattened heritability

174 distribution under negative selection[30]. Genomic inflation computed from isGWAS results

175 across all analyses was on average $\approx 1.07$ and ranged between $(0.94, 1.26)$ which was similar

176 to REGENIE with average $\approx 1.1$ and range $(1.01, 1.3)$ (**Supplementary Figure 14** and

177 **Supplementary Table 8**). We deploy isGWAS with genotype imputation in our primary

178 analyses and, as secondary sensitivity analyses, without imputation. Our investigation reveals

179 some surprising results. Imputation occasionally led to changes in MAF between cases and

180 controls such that estimated genetic effects switched sign (i.e., effect direction) relative to

181 results computed from non-imputed data (**Supplementary Figure 15, Supplementary Table**

182 **3, Supplementary Files 2-3**). The approach might be used to efficiently flag ambiguous

183 significant results in analyses that are the result of the missing values imputation strategy

184 (mean-imputed in the case of REGENIE).

185

186 Finally, the computational gains of isGWAS relative to REGENIE Step 2 are striking: a full

187 genome-wide association assessment for each disease took approximately 4 minutes using

188 isGWAS and, on average over different prevalence, this is around 1,300 times faster than a like-

189 for-like assessment using REGENIE Step-2 (**Figure 7, Supplementary Table 9, and**

190 **Supplementary File 4**).

191

192 *Replicating significant associations in Biobank Japan analyses*

193 Using only publicly available summary information from Biobank Japan (BBJ), i.e., without

194 access to ILD, we looked to compare and replicate BBJ GWAS results across 42 diseases[24].

195 We considered 309 variants that were identified in [24] as genome-wide significant ($p <$

196 $5 \times 10^{-8}$) across 30 of the 42 diseases. Our results reveal very close alignment between

197 isGWAS computed associations and those of [24] - correlation between p-values at $log_{10}$ scale

198 was $cor(p_{isG}, p_{BBJ}) = 0.98$ with 92.2% of isGWAS computed genetic effects within the 95%

199 CI of the original study (**Supplementary Figure 16**). Using the published study-level

200    results[24]   as the ground truth, isGWAS demonstrated good sensitivity and specificity

201    (**Supplementary Figure 17**). We alternatively assessed performance when setting more

202    stringent significance thresholds - returning near identical conclusions when classifying

203    variants at $p < 9.58 \times 10^{-9}$ (used in the original publication). Results for X-chromosome

204    variants in males and females were similarly concordant (results not presented).

205

206    **isGWAS model validation using simulations**

207    We generated simulated datasets to assess performance of isGWAS - with and without Firth

208    correction - against a variety of classical methods which either: (a) do not require ILD, the

209    logistic ad-hoc estimator[17] and Fisher's Exact Test[22]; or (b) require ILD, logistic and Firth

210    corrected regression[7], [8], [16]. We perform two simulation studies (**Figure 4**,

211    **Supplementary Figure 18, Supplementary Information**). isGWAS-Firth outperformed all

212    other approaches in terms of either computational cost or robustness of results over the range

213    of scenarios considered. It is well documented that computational performance is reduced when

214    using Firth's bias correction in ILD regression analyses[7], [8], [16], we discover, however,

215    that no-ILD isGWAS-Firth regression has significantly improved performance relative to

216    uncorrected isGWAS (**Supplementary Table 21**). As anticipated[23], when disease prevalence

217    is rare (i.e., $\pi \leq 0.01$) parameter estimates computed using non-Firth corrected ILD regression

218    were unreliable. The MSE and distribution of parameters estimated via ILD logistic regression

219    were often orders of magnitude poorer than other methods (**Figure 4a-c**). **Figure 4f-h**

220    highlights the chronological evolution of no-ILD *p*-value estimates, from Sasieni's logistic ad-

221    hoc estimator (1997)[17], Fisher's Exact Test (1922)[22] to isGWAS-Firth, illustrating

222    improvements in estimation via successive approaches. See **Supplementary Information** for

223    detailed results review.

224

## Leapfrog re-sampling: using isGWAS to extrapolate variant association results to larger sample sizes

227    When ILD are available, the computational benefits of isGWAS make it possible to deploy

228    resampling approaches to estimate empirical effect sizes, *p*-values and corresponding

229    confidence intervals, previously considered computationally daunting in GWAS[31], [32]. We

230    extend the idea by introducing a heuristic leapfrog re-sampling (LRS) algorithm to help forecast

231    future results in larger hypothetical GWAS sample sizes (**Methods**). The LRS is summarised

232    in three key steps: (1) specify a target sample size along with the number and size of sub-

233    samples to be generated; (2) (leapfrog-step) compute sufficient statistics in the sub-samples and

234    re-scale the estimated number of cases and controls to match the larger target sample size; and

235    (3) deploy isGWAS in each leapfrog sample to recover a distribution of association *p*-values

236    over the collection of sub-samples. In our testing of the LRS, we use the median *p*-value as a

237    generally robust estimate of a target p-value (weighted or distribution-based summaries can

238    alternatively be considered). Thus, the LRS leverages variation in both genotype and disease

239    status between individuals in the current sample to help predict updates of parameters after

240    adding new samples. Despite perceived similarities, traditional GWA power calculators[33]

241    and the isGWAS-LRS are different. isGWAS-LRS does not require input of case-control ratios,

242    heritability (i.e., beta estimates) or type-I error rates. Instead, multiple regression analyses are

243    combined to forecast and test parameter estimates in expanding sample sizes.

244    We run the leapfrog re-sampler in both simulation and real-data settings, informed by the seven

245    studied diseases in UKB (**Methods and Supplementary Information**). We evaluate

246    performance over a range of initialisations, starting from a 10% increase to a maximum of 100%

247    (i.e., 2-fold) increase in GWAS sample size relative to the current actual UKB sample size.

248    Results are presented in **Figure 5**, **Supplementary Tables 10-11**. As is standard, we assume a

249    true positive association of p<5×10−8 in the target sample. Our results in simulated scenarios

250    (**Figure 5f**) reveal that: when doubling sample size from $N_{current} = 276,204$ to a maximum

251    $N_{target} = 552,408$ , the accuracy and TPR progressively dropped for subsequent increases in

252    the target sample size, but values for each measurement were typically ≥80% across the range.

253    Our real world LRS analyses of UK Biobank data replicate and further elucidate performance

254    across the six of the seven diseases (**Figure 5a-e**). Using a subsample size of $N_{current} \approx$

255    135,000 we increased target sample size up to $N_{target} \approx 270,000$, taken as the maximum

256    observed sample size we could benchmark against. For all choices of target sample size, and

257    across each disease, we observe high accuracy rates (≥95%). However, the TPR was sensitive

258    to disease prevalence, reducing monotonically as the target sample size increased. Broadly,

259    TPR remained reasonable ($\gtrsim$ 60%) up to a 2-fold increase in sample size, except for the very

260    rare (case-control ratio of 1:669) thyroid gland cancer. This is due to fewer significant variants

261    being included in the assessment as a result of lower percentage of heritability explained, which

262    can artificially reduce the TPR for each new locus with relatively high odds ratios. Naturally,

263    TPR reduces as a function of decreasing disease prevalence, as re-sampling from fewer cases

264    can increase the variability in MAFs and thus isGWAS forecasting. We note that our theoretical

265    sub-sampling approach had better predictive capabilities, owing to the prevalence preserving

266    sampling strategy taken (**Methods and Supplementary Figure 19**).

267    We also assessed isGWAS's ability to extrapolate results when ILD were not available, using

268    a highly constrained version of the leapfrog re-sampler (**Supplementary Information**). In this

269    scenario, MAF and disease prevalence per variant are fixed, computed from the maximum

270    current sample (i.e., without sub-sampling), and the number of cases and controls are

271    proportionately increased to match the target sample size. We did this for two GWAS of

272    schizophrenia: (a) 2014 analyses with up to $N = 77,096$ (cases $= 33,640$, controls $=$

273    43456) European ancestry individuals[25]; and (b) the larger (and future) 2022 analyses with

274    up to $N = 130,644$ (cases $= 53386$, controls $= 77258$) European ancestry individuals[27].

275    We treat the 2014 study as the current sample size and the 2022 sample size as the target future

276    state, which we benchmark predictive performance against. The studies were selected because,

277    for each variant $j$, necessary data to run isGWAS, i.e., $\{N_j, N_j^*, MAF_j, MAF_j^*\}$, were made

278    publicly available. Note these data are pooled estimates, computed across all European cohorts.

279    Despite not accessing ILD, our results reveal reasonable concordance between isGWAS 2014

280    extrapolated results and the published analyses of 2022 (**Figure 6**, **Supplementary Figures 20-**

281    **21, Supplementary Tables 13-16**). Like our Biobank Japan analyses, we also used a more

282    stringent significance threshold ($p < 10^{-10}$) to help attenuate false positives, observing

283    improved overall performance by recovering a good TPR $\geq 70\%$ (**Supplementary Table 13**).

284    We do not report FDR as these cannot be accurately computed when filtering results based on

285    a p-value inclusion/exclusion threshold. Of the overlapping 608 clumped variants considered,

286    isGWAS-LRS identified 136 associations that were not yet deemed GWAS significant (i.e,.

287    $p > 5e - 8$) in the 2014 study but later identified as significant in the 2022 study. Moreover,

288    of the 436 significant associations predicted by isGWAS, 75% overlap with observed

289    significant associations in 2022. isGWAS predicted an additional 74 associations as significant

290    that were not significant in 2022 – of those 52 were near the significance threshold with $p <$

291    $9e - 07$. There were 121 variants not correctly predicted by the 2014 cohort. This could be due

292    to increased ethnical and relatedness heterogeneity in the 2022 cohort that was not present in

293    the 2014 analysis.

294

295    **Computational performance and convergence details**

296    isGWAS is an iterative algorithm whose convergence (i.e., ability to estimate model

297    parameters) depends on several tuning parameters (**Methods**). Using default parameter settings,

298    isGWAS-Firth converged in all real-data and simulated scenarios tested (**Figure 4e and**

14

299 **Supplementary Tables 17, 19-20**). Convergence was achieved in around 0.001 seconds per

300 variant (**Supplementary Table 21**) on a 2.4 GHz 8-Core Intel Core i9 processor. The non-Firth

301 corrected isGWAS algorithm may require more iterations, particularly for diseases with lower

302 prevalence (e.g., case:control ratio of 1:94 and lower) which included scenarios where

303 convergence was not achieved (**Figure 4e, Supplementary Tables 17-20, Supplementary**

304 **Information**).

305

306 When distributed over 32 CPU cores on a high-performance cluster, Firth-corrected isGWAS

307 analysed a single disease from UK Biobank across ~11 million SNPs and for ~335,000

308 individuals in ~4 minutes (**Figure 7, Supplementary Tables 9 and 18**). This means that

309 isGWAS-Firth can perform around 1,500 disease GWAS for every one GWAS performed using

310 an alternative methodology. The same analysis with a small number of CPU cores was

311 completed in tens of minutes using isGWAS-Firth (**Figure 7**). Further computational gains at

312 larger sample sizes will likely be achieved as ILD methods can scale poorly with sample size,

313 whereas isGWAS has near fixed computational cost at any size. As isGWAS currently

314 computes associations for each variant independently, additional improvements such as

315 parallelisation are possible. Full details are available in **Supplementary File 4.**

316

## Discussion

318 In this study, we developed isGWAS, an efficient, biobank-scalable method for genetic

319 association testing which can: (a) compute regression parameters and test for a variant-disease

320 association in real-time (i.e., approximately one millisecond) for any sample size; (b) bypass

321 the need to run large-scale GWAS using high-performance computing facilities owing to ultra-

322 low system resource demands (i.e., runtime and memory); and (c) infer GWAS results from

323 virtually enlarged sample sizes using a novel re-sampling procedure. The isGWAS algorithm

324 design allows analyses to be run without the need to hold or access individual-level data (ILD)

325 directly, thereby providing a single methodological framework to utilise a wide range of data

326 sources such as published summary-level data from biobanks and repositories.

327

328 isGWAS draws inspiration from classical methodologies to overcome significant

329 computational bottlenecks associated with massive-scale analyses. The practical simplicity and

330 quick runtime of classical approaches have seen them deployed in a recent large-scale

331 analysis[23]. Rather than using ILD, as contemporary GWAS regression analyses do, isGWAS

332 distils the required input data down to sufficient statistics – a low-dimensional summary of ILD

333 that captures all necessary information required to compute a genetic-disease association model

334 parameter. In combination with modifications to the Newton-Raphson procedure, used to

335 estimate model parameters in a logistic regression, our use of sufficient statistics dramatically

336 reduces the computational time for disease association testing relative to existing methods.

337 Achieving up to a 1,500-fold improvement in computational runtime, when benchmarked

338 against a state-of-the-art GWAS tool, isGWAS reduced time to genome-wide insight from

339 several days down to ~4 minutes. Thereby unlocking potential for massive scale exploration of

340 genetic-disease associations in real-time and making feasible the routine assessment of

341 thousands of disease endpoints and studies. Computational bottlenecks associated with existing

342 GWAS methodologies are fast approaching. Analyses of resources such as UK Biobank WGS

343 data, the emerging massive cohorts of the Global Biobank Initiative[1], and Our Future

344 Health[34] are expected to push current GWAS tools to their system resource limits with

345 significant associated time-to-insight penalties. Conversely, with no computational sensitivity

346 to sample size, as the number of variants assessed and sample sizes continue to increase, the

347 relative savings and benefits of isGWAS can be expected to grow.

348

349    To attenuate possible issues of confounding and population stratification, we propose that

350    additional QC-steps are performed before computing sufficient statistics for isGWAS. In our

351    analyses, these steps reduced our UK Biobank sample size from ~408k individuals (used in the

352    original testing of REGENIE[8]) down to a more homogeneous sample of ~335k individuals

353    used to generate and compare results from isGWAS and REGENIE. Our results reveal an often-

354    striking concordance between approaches genome-wide as well as at the regional locus level.

355    The reduction in sample size was compensated by the isGWAS leapfrog re-sampler (LRS),

356    which we demonstrate efficiently helped extrapolate GWA results onto larger sample sizes (up

357    to 2-times). While we note sensitivity of an LRS extrapolation to disease prevalence, across the

358    range considered the TPR and FPR were well calibrated to at least a 1.5-fold increase in sample

359    size. The LRS might therefore be leveraged to aid GWAS cohort design, for example to

360    quantify the potential benefit of sampling more participants with a disease of interest against

361    cost. In our analyses of the PGC Schizophrenia cohort, we deployed a highly restricted (i.e., no

362    ILD) version of the LRS: forecasting results from a smaller 2014 cohort[25] onto a sample size

363    that matched a future 2022 study[27]. Despite no guarantees of sufficiency, isGWAS LRS

364    identified 75% of significant variants that were later identified in the larger 2022 cohort (almost

365    double the size) while maintaining a low FDR. Unlike extrapolation via the ILD leapfrog re-

366    sampler, this naïve extrapolation does not account for differences in the MAF of cases and

367    controls between 2014 and 2022 data. Regardless, the above findings highlight potential for

368    isGWAS to furnish reasonable forecasts of future results without accessing ILD directly. Our

369    recent predictions from FinnGen consortium data[18], [35] provide confidence that the

370    isGWAS algorithm is applicable also to multi-ethnic GWAS through analyzing each ethnicity

371    separately and combining results in a meta-analysis, as it is common practice[36].

372

17

373    Beyond, isGWAS can be applied to help address routinely asked questions about future

374    scenarios and evaluate enrichment contribution of biobanks to disease-specific associations[35]

375    or to protein-specific variant associations[18], particularly in the rare spectrum. isGWAS-Firth

376    provides a timely, rapid regression-based analysis of common, rare and ultra-rare variants.

377    Unlike ILD-based analyses, where Firth's correction significantly increases computational

378    time[7], [8], there is no computational penalty when using Firth's correction in the isGWAS

379    framework - in fact, we observe improved computational performance. The advantage of

380    considerable improvements in computational runtime is that it allows for the introduction of

381    forecasting, re-sampling and other non-parametric techniques -  the LRS being one example.

382    These might widen robust association testing strategies as well as provide new avenues to tackle

383    confounding or population sub-structure. For now, we envisage the possibility that the wider

384    human genetics community routinely compute and make available the sufficient statistics, i.e.

385    MAF in the cases and the cohort, and the corresponding sample sizes per variant, toward a

386    publicly available, privacy compliant, data asset. In addition to avoiding the need for expensive

387    high-performance computing facilities and memory intensive data storage, the data asset might

388    enhance meta-analyses and biological insight, improve equitable access, and enable faster

389    collaborations between teams and help bridge financial and resource gaps between institutions

390    and research groups internationally.

391

392

393

394

# Methods

## Disease SNP association model

Let $S_M \in \{1,2\}$ be the maximum number of copies of the effect allele for an individual under model $M \in \{A, D, R\}$, where $A$ denotes an additive model, $R$ a recessive model and $D$ a dominant model, i.e.,

$$S_M = \begin{cases} 2, & M = A, \\ 1, & M = D \cap R. \end{cases}$$

Furthermore let,

$$MAF_{j,M} = MAF_j \mid M$$

$$MAF_{j,M}^* = MAF_j^* \mid M$$

where, for a given model $M$, $MAF_{j,M}$ is the minor allele frequency for variant $j$ in the sample, ancestry, or population and $MAF_{j,M}^*$ the minor allele frequency in the cases. We let $Y$ denote disease status and $G_j$ the $j$th genotype in the sample. For convenience we write $G_{j,M} = G_j|M$. It is assumed that the outcome model for $Y$, conditional on $G_j$, is given by:

$$\mathbb{E}[Y \mid G_j, M] = h^{-1}\left(\alpha_{j,M} + \beta_{j,M} G_{j,M}\right), \qquad j = 1,2,\dots,Q,$$

where, conditional on model $M$, the pair $\{\alpha_{j,M} \; \beta_{j,M}\}$ denote the intercept and genotype effect and $h$ is a function linking the outcome to genotype $G_{j,M}$ for all $j = 1,2,\dots,Q$ genetic variants considered. In deriving the isGWAS estimation procedure we assume that $h$ is the logit function, i.e.,

$$\pi_{Y|G_{j,M}} = P\left(Y = 1 \mid G_j, M\right) = \frac{e^{\left(\alpha_{j,M} + \beta_{j,M} G_{j,M}\right)}}{1 + e^{\left(\alpha_{j,M} + \beta_{j,M} G_{j,M}\right)}}.$$

The isGWAS methodology can, however, be broadened to other link functions and outcome types. We take $\left\{\hat{\alpha}_{j,M}, \hat{\beta}_{j,M}, \hat{\sigma}_{\alpha_{j,M}}, \hat{\sigma}_{\beta_{j,M}}\right\}$ to be sample based estimates of the intercept $\alpha_{j,M}$ and coefficient $\beta_{j,M}$, and their associated standard errors $\left\{\hat{\sigma}_{\alpha_{j,M}}, \hat{\sigma}_{\beta_{j,M}}\right\}$. We allow the genetic effect

19

416 $\beta_{j,M}$ and genotype $G_j$ (or an observation thereof $g_j$) to be analyzed on either (i) the standardized

417 scale or (ii) non standardized scale. To accommodate this, we introduce the variable $s_m^*$, so that:

418
$$\{\beta_{j,s_m^*}, g_{j,s_m^*}\} = \left\{(1 + (\sigma_{g_i} - 1)s_m^*)\beta_j, \left. {}^{g_j}\middle/ {(1 + (\sigma_{g_i} - 1)s_m^*)}\right.\right\}.$$

419 Hence, when $s_m^* = 0$ analyses are performed on the non-standardized scale and $s_m^* = 1$ on the

420 standardized scale. Default analyses assume the genetic effect is assessed on the non-

421 standardized scale, i.e., $s_m^* = 0$. Note that, while p-values are generally invariant to the choice

422 of effect scale $s_m^*$, betas and standard errors are dependent on the specification of $s_m^*$.

423

## Sample-Level Newton-Raphson (SaLN-R) algorithm

425 Here we detail the isGWAS procedure for computing summary statistics

426 $\left\{\hat{\alpha}_{j,M}, \hat{\beta}_{j,M,s_m^*}, \hat{\sigma}_{\alpha_{j,M}}, \hat{\sigma}_{\beta_{j,M,s_m^*}}\right\}$ using only four data points,

427
$$\left\{N_j, \sum_{i=1}^{N_j} y_i = N_j^*, \sum_{i=1}^{N_j} y_i g_{ij,M}, \sum_{i=1}^{N_j} g_{ij,M} = n_{j1} + 2n_{j2}I(M = A)\right\}$$

428 or, as we show, the quadruple

429
$$\{N_j, N_j^*, MAF_{j,M}, MAF_{j,M}^*\},$$

430 where $N_j$ denotes the study or population sample size and $N_j^*$ the number of cases in the sample,

431 see ad-hoc estimator (Supplementary Information) for definitions of $n_j.$. Note that we have

432 allowed the sample size $N$ and number of cases $N^*$ to vary by genotype $j$, this is useful when

433 emulating results from GWAS. This is because the number of individuals analyzed in GWAS

434 can vary by genotype owing to (e.g.,) quality of imputation or available data per variant and

435 participant in a study. Ideally the sample size and number of cases would not vary by genotype

436 and when using isGWAS to forecast GWAS results, users do not need not vary $N_j$, i.e,.

437
$$N_j = N \quad \text{and} \quad N_j^* = N^*, \qquad j = 1,2,\dots,Q.$$

20

438    Given a vector of observed data $\{\boldsymbol{y}, \boldsymbol{g}_{j,s_m^*}\}$, where $\boldsymbol{y} = \{y_1, y_2, \ldots, y_{N_j}\}$ and $\boldsymbol{g}_{j,s_m^*} =$

439    $\{g_{1j}, g_{2j}, \ldots, g_{N_j}\}$, estimates of model parameters are typically derived by maximizing the log-

440    likelihood function

$$L(\boldsymbol{\beta}_{j,M,s_m^*}) = \sum_{i=1}^{N_j} \log P(y_i | \boldsymbol{\beta}_{j,M,s_m^*}, g_{ij,M,s_m^*}),$$

441    which is equivalent to identifying parameter values $\boldsymbol{\beta}_{j,M,s_m^*} = \{\alpha_{j,M}, \beta_{j,M,s_m^*}\}$ which satisfy:

$$\frac{\partial L(\boldsymbol{\beta}_{j,M,s_m^*})}{\partial \boldsymbol{\beta}_{j,M,s_m^*}} = V(\boldsymbol{\beta}_{j,M,s_m^*}, I_F) = 0,$$

444    where $V(\boldsymbol{\beta}_{j,M,s_m^*}, I_F)$ denotes the logistic score function, i.e.,

$$V(\boldsymbol{\beta}_{j,M,s_m^*}, I_F) = \tilde{\boldsymbol{g}}_{j,M,s_m^*}^T (\boldsymbol{y} - \pi_{y|\tilde{g}_{j,M,s_m^*}}) + I_F K(\boldsymbol{\beta}_{j,M,s_m^*}) = 0,$$

446    with $I_F$ denoting an indicator function used to highlight that a Firth modified version of the

447    score function has been used. For ease of mathematical presentation initially, we detail the Firth

448    adjusted SaLN-R algorithm later, i.e., we set $I_F = 0$ in this section. Additionally, to improve

449    succinctness of notation, we drop the use of the parameter $s_m^*$ - reintroducing where necessary

450    - and set $\boldsymbol{\beta}_{j,M} = \{\alpha_{j,M}, \beta_{j,M}\}$ and $\tilde{\boldsymbol{g}}_{j,M} = (1, g_{j,M})$ above, so that $\boldsymbol{\beta}_{j,M}\tilde{\boldsymbol{g}}_{j,M}^T = \alpha_{j,M} + \beta_{j,M}g_{j,M}$.

451    We compute candidate solutions to by expanding $V(\boldsymbol{\beta}_{j,M})$ as a Taylor series about a value

452    $\boldsymbol{\beta}_{j0,M}$ and up to second order, i.e., using the Newton-Raphson (N-R) method:

$$V(\boldsymbol{\beta}_{j,M}) = V(\boldsymbol{\beta}_{j0,M}) + \frac{\partial V(\boldsymbol{\beta}_{j,M})}{\partial \boldsymbol{\beta}_{j,M}}\bigg|_{\boldsymbol{\beta}_{j0,M}} (\boldsymbol{\beta}_{j,M} - \boldsymbol{\beta}_{j0,M}) + \mathcal{O}\left((\boldsymbol{\beta}_{j,M} - \boldsymbol{\beta}_{j0,M})^2\right),$$

454    Which is re-written as

$$\boldsymbol{\beta}_{j,M} = \boldsymbol{\beta}_{j0,M} + \mathcal{I}^{-1}(\boldsymbol{\beta}_{j0,M})V(\boldsymbol{\beta}_{j0,M}) + \mathcal{O}\left((\boldsymbol{\beta}_{j,M} - \boldsymbol{\beta}_{j0,M})^2\right)$$

456    and generalized into an N-R iterative algorithm:

457    $\boldsymbol{\beta}_{j(k+1),M,s_m^*} = \boldsymbol{\beta}_{jk,M,s_m^*} + \mathcal{I}^{-1}(\boldsymbol{\beta}_{jk,M,s_m^*})V(\boldsymbol{\beta}_{jk,M,s_m^*}) + \mathcal{O}\left((\boldsymbol{\beta}_{j(k+1),M,s_m^*} - \boldsymbol{\beta}_{jk,M,s_m^*})^2\right), \quad k = 0,1,\ldots,K,$

458     where we have re-introduced $s_m^*$ to highlight that the algorithm is dependent on the choice of

459     effect scale. The variable $\mathcal{J}^{-1}$ denotes the inverse Fisher Information matrix, where

460

$$\mathcal{J}\left(\boldsymbol{\beta}_{jk,M}\right) = \left. -\frac{\partial V\left(\boldsymbol{\beta}_{j,M}\right)}{\partial \boldsymbol{\beta}_{j,M}} \right|_{\boldsymbol{\beta}_{j,M}=\boldsymbol{\beta}_{jk,M}}$$

461

$$= \left. \widetilde{\boldsymbol{g}}_{j,M}^{T} \operatorname{diag}\left(\pi_{y_i|g_{j,M}}\left(1-\pi_{y_i|g_{jk,M}}\right)\right) \widetilde{\boldsymbol{g}}_{j,M} \right|_{\boldsymbol{\beta}_{j,M}=\boldsymbol{\beta}_{jk,M}}$$

462

$$= \left. \begin{pmatrix} \sum_{i=1}^{N_j} \pi_{y_i|g_{ij,M}}\left(1-\pi_{y_i|g_{ij,M}}\right) & \sum_{i=1}^{N_j} \pi_{y_i|g_{ij,M}}\left(1-\pi_{y_i|g_{ij,M}}\right)g_{ij,M} \\ \sum_{i=1}^{N_j} \pi_{y_i|g_{ij,M}}\left(1-\pi_{y_i|g_{ij,M}}\right)g_{ij,M} & \sum_{i=1}^{N_j} \pi_{y_i|g_{ij,M}}\left(1-\pi_{y_i|g_{ij,M}}\right)g_{ij,M}^2 \end{pmatrix} \right|_{\boldsymbol{\beta}_{j,M}=\boldsymbol{\beta}_{jk,M}}$$

463     and the score function is given by

464

$$V\left(\boldsymbol{\beta}_{jk,M}\right) = \left. \widetilde{\boldsymbol{g}}_{j,M}^{T}\left(\boldsymbol{y}-\pi_{y|g_{j,M}}\right) \right|_{\boldsymbol{\beta}_{,Mj}=\boldsymbol{\beta}_{jk,M}}$$

465

$$= \left. \begin{pmatrix} \sum_{i=1}^{N_j}\left(y_i-\pi_{y_i|g_{ij,M}}\right) \\ \sum_{i=1}^{N_j}\left(y_i-\pi_{y_i|g_{ij,M}}\right)g_{ij,M} \end{pmatrix} \right|_{\boldsymbol{\beta}_{j,M}=\boldsymbol{\beta}_{jk,M}}$$

466     Both $\mathcal{J}\left(\boldsymbol{\beta}_{jk,M}\right)$ and $V\left(\boldsymbol{\beta}_{jk,M}\right)$ above require individual-level data to compute their values.

467     isGWAS aims to estimate values for these variables using sample-level information only,

468     thereby avoiding the immediate need for individual data. To achieve this, we approximate both

469     the Fisher Information matrix and the Score function via the pair $\{\mathcal{J}_e\left(\boldsymbol{\beta}_{jk,M}\right), V_e\left(\boldsymbol{\beta}_{jk,M}\right)\}$,

470     where:

471

$$\mathcal{J}_e\left(\boldsymbol{\beta}_{jk,M}\right) = N_j \begin{pmatrix} \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}\left(1-\pi_{y|g_{j,M}}\right); \boldsymbol{\beta}_{jk,M}\right] & \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}\left(1-\pi_{y|g_{j,M}}\right)g_{j,M}; \boldsymbol{\beta}_{jk,M}\right] \\ \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}\left(1-\pi_{y|g_{j,M}}\right)g_{j,M}; \boldsymbol{\beta}_{jk,M}\right] & \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}\left(1-\pi_{y|g_{j,M}}\right)g_{j,M}^2; \boldsymbol{\beta}_{jk,M}\right] \end{pmatrix}$$

472     and

473

$$V_e\left(\boldsymbol{\beta}_{jk,M}\right) = \begin{pmatrix} \sum_{i=1}^{N_j} y_i - N_j\,\mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}; \boldsymbol{\beta}_{jk,M}\right] \\ \sum_{i=1}^{N_j} y_i g_{ij,M} - N_j\,\mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}g_{j,M}; \boldsymbol{\beta}_{jk,M}\right] \end{pmatrix}$$

22

474
$$\approx \begin{pmatrix} N_j^* - N_j\, \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}};\, \boldsymbol{\beta}_{jk,M}\right] \\ N_j^*\, \mathbb{E}_{g_{j,M}}\left[g_{j,M}\mid y=1\right] - N_j\, \mathbb{E}_{g_{j,M}}\left[\pi_{y|g_{j,M}}g_{j,M};\, \boldsymbol{\beta}_{jk,M}\right] \end{pmatrix},\quad N_j^* \gg 1.$$

475

476 with $\mathbb{E}_{g_{j,M}}[\,\cdot\,;\boldsymbol{\beta}_{jk,M}]$ denoting that expectation is taken with respect to $g_{j,M}$ and evaluated at

477 $\boldsymbol{\beta}_{j,M} = \boldsymbol{\beta}_{jk,M}$. Note that $\{\mathcal{I}_{\mathbb{e}}(\boldsymbol{\beta}_{jk,M}), V_{\mathbb{e}}(\boldsymbol{\beta}_{jk,M})\}$ are motivated by switching from empirical, i.e.,

478 sample-based, estimates in $\{\mathcal{I}(\boldsymbol{\beta}_{jk,M}), V(\boldsymbol{\beta}_{jk,M})\}$ to their expected value analogues, which

479 reverses the usual mode of estimation. Sample size $N_j$ is presumed large and thus switching

480 from sample-based to expected values in the N-R algorithm is well motivated. However, when

481 the number of cases $N_j^*$ is 'small', an approximation of $\sum_{i=1}^{N_j} y_i g_{j,M} \approx N_j^* \mathbb{E}_{g_{j,M}}[g_{j,M}\mid y=1]$

482 becomes weaker and we recommend using the statistic $\sum_{i=1}^{N_j} y_i g_{j,M}$. Values for the elements in

483 $\{\mathcal{I}_{\mathbb{e}}(\boldsymbol{\beta}_{jk,M}), V_{\mathbb{e}}(\boldsymbol{\beta}_{jk,M})\}$ are computed via:

484
$$\mathbb{E}_{g_{j,M,s_m^*}}\left[\pi_{y|g_{j,M,s_m^*}}\left(1 - \pi_{y|g_{j,M,s_m^*}}\right)g_{j,M,s_m^*}^c;\, \boldsymbol{\beta}_{jk,M,s_m^*}\right] =$$

485
$$\sum_{l=0}^{S_M} \pi_{y|g_{j,M}=l/w_{s_m^*};\,\boldsymbol{\beta}_{jk,M,s_m^*}}\left(1 - \pi_{y|g_{j,M}=l/w_{s_m^*};\,\boldsymbol{\beta}_{jk,M,s_m^*}}\right)\left(\frac{l}{w_{s_m^*}}\right)^c p\left(g_{j,M} = {}^l/w_{s_m^*}\right)$$

486
$$= e_{j,M}^{(c,k)}$$

487 and

488
$$\mathbb{E}_{g_{j,M,s_m^*}}\left[\pi_{y|g_{j,M,s_m^*}}g_{j,M,s_m^*}^c \mid \boldsymbol{\beta}_{jk,M,s_m^*}\right] =$$

489
$$\sum_{l=0}^{S_M} \pi_{y|g_{j,M}=l/w_{s_m^*};\,\boldsymbol{\beta}_{jk,M,s_m^*}}\left(\frac{l}{w_{s_m^*}}\right)^c p\left(g_{j,M} = {}^l/w_{s_m^*}\right)$$

490
$$= \tilde{e}_{j,M}^{(c,k)},$$

491 where $w_{s_m^*} = \left(1 + (\sigma_{g_i} - 1)s_m^*\right)$ and the superscript and subscript in $e_M^{(c,k)}$, $\tilde{e}_M^{(c,k)}$ are used to

492 highlight that expectation has been taken conditional on k-th iteration $\beta_{jk,M}$ and under

23

493     modelling assumption $M$ (and implicitly effect scale $s_m^*$). Probability mass $p\left(g_{j,M} = {}^{l}/{w_{s_m^*}}\right)$

494     is either defined a-priori or can be approximated empirically, which we detail later. In

495     combination, therefore, it follows that:

496
$$\mathcal{I}_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right) = N_j \begin{pmatrix} e_{j,M}^{(0,k)} & e_{j,M}^{(1,k)} \\ e_{j,M}^{(1,k)} & e_{j,M}^{(2,k)} \end{pmatrix}$$

497     and

498
$$\mathcal{I}_{\mathbb{e}}^{-1}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right) = \frac{1}{N_j\left(e_{j,M}^{(0,k)} e_M^{(2,k)} - \left(e_{j,M}^{(1,k)}\right)^2\right)} \begin{pmatrix} e_{j,M}^{(2,k)} & -e_{j,M}^{(1,k)} \\ -e_{j,M}^{(1,k)} & e_{j,M}^{(0,k)} \end{pmatrix}.$$

499     Following the same process that led to the above, we re-write the sample-level Score function

500     $V_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right)$ as:

501
$$V_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right) = \begin{pmatrix} N_j^* - N_j \, \tilde{e}_{j,M}^{(0,0)} \\ \sum_{i \, : \, y_i=1} g_{ij,M,s_m^*} - N_j \tilde{e}_{j,M}^{(1,0)} \end{pmatrix}$$

502
$$\approx \begin{pmatrix} N_j^* - N_j \, \tilde{e}_{j,M}^{(0,0)} \\ \dfrac{S_M N_j^* MAF_{j,M}^*}{w_{s_m^*}} - N_j \tilde{e}_{j,M}^{(1,0)} \end{pmatrix}, \quad N_j^* \gg 1,$$

503     where we have used the following approximation:

504
$$\sum_{i \, : \, y_i=1} g_{ij,M,s_m^*} \approx N_j^* \mathbb{E}_{g_{j,M}}\left[g_{j,M,s_m^*} \mid y = 1\right] = \frac{S_M N_j^* MAF_{j,M}^*}{w_{s_m^*}}.$$

505     Hence, replacing the pair $\left\{\mathcal{I}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right), V\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right)\right\}$ with the sample-level approximations

506     $\left\{\mathcal{I}_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right), V_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right)\right\}$ we furnish the SaLN-R algorithm:

507
$$\boldsymbol{\beta}_{j(k+1),M,s_m^*} = \boldsymbol{\beta}_{jk,M,s_m^*} + \mathcal{I}_{\mathbb{e}}^{-1}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right) V_{\mathbb{e}}\left(\boldsymbol{\beta}_{jk,M,s_m^*}\right)$$

508
$$= \boldsymbol{\beta}_{jk,M,s_m^*} + \frac{1}{N_j\left(e_{j,M}^{(0,k)} e_{j,M}^{(2,k)} - \left(e_{j,M}^{(1,k)}\right)^2\right)} \begin{pmatrix} e_{j,M}^{(2,k)} & -e_{j,M}^{(1,k)} \\ -e_{j,M}^{(1,k)} & e_{j,M}^{(0,k)} \end{pmatrix} \begin{pmatrix} N_j^* - N_j \, \tilde{e}_{j,M}^{(0,k)} \\ \dfrac{S_M N_j^* MAF_{j,M}^*}{w_{s_m^*}} - N_j \tilde{e}_{j,M}^{(1,k)} \end{pmatrix}.$$

509     The standard error of the updates, $\widehat{\boldsymbol{\sigma}}_{\boldsymbol{\beta}_{j(k+1),M,s_m^*}} = \left\{\widehat{\sigma}_{\alpha_{j(k+1),M,s_m^*}}, \widehat{\sigma}_{\beta_{j(k+1),M,s_m^*}}\right\}$, are given by the

510     diagonal of the inverse Fisher information matrix, i.e.,

511
$$\hat{\sigma}_{\alpha_{j(k+1),M,s_m^*}} = \sqrt{\frac{e_{j,M}^{(2,k)}}{N_j\left(e_{j,M}^{(0,k)}e_{j,M}^{(2,k)}-\left(e_{j,M}^{(1,k)}\right)^2\right)}},$$

512
$$\hat{\sigma}_{\beta_{j(k+1),M,s_m^*}} = \sqrt{\frac{e_{j,M}^{(0,k)}}{N_j\left(e_{j,M}^{(0,k)}e_{j,M}^{(2,k)}-\left(e_{j,M}^{(1,k)}\right)^2\right)}}.$$

513   It can be shown from the above that:

514
$$\boldsymbol{\beta}_{jk,M,s_m^*=1} = \sigma_{g_i}\boldsymbol{\beta}_{jk,M,s_m^*=0},$$

515
$$\hat{\sigma}_{\beta_{j(k+1),M,s_m^*=1}} = \sigma_{g_i}\hat{\sigma}_{\beta_{j(k+1),M,s_m^*=0}}.$$

516   We set $s_m^* = 0$ to compute values for the pair $\left\{\hat{\boldsymbol{\beta}}_{jk,M,0}, \hat{\sigma}_{\beta_{j(k+1),M,0}}\right\}$ and use the above identities

517   to return parameter estimates on the standardized scale $s_m^* = 1$. The data required to run the

518   SaLN-R algorithm are:

519
$$\left\{N_j, N_j^*, \sum_{i=1}^{N_j} y_i g_{ij,M}, p(g_{j,M}=1)\right\} \bigcup \left\{\begin{matrix}\emptyset, & S_M=1,\\ p(g_{j,M}=2), & S_M=2.\end{matrix}\right\}$$

520   We use the approximations

521
$$p(g_{j,M}=1) \approx \frac{n_{j1}}{N_j} \xrightarrow[HWE]{N_j \gg 1} \begin{cases}MAF_{j,M}, & S_M=1,\\ 2MAF_{j,M}(1-MAF_{j,M}), & S_M=2,\end{cases}$$

522
$$p(g_{j,M}=2) \approx \frac{n_{j2}}{N_j} \xrightarrow[HWE]{N_j \gg 1} MAF_{j,M}^2,$$

523   where $\xrightarrow[HWE]{}$ is used to denote under Hardy-Weinberg equilibrium.

524   The SaLN-R algorithm is extended to include Firth's penalty function (see **Supplementary**

525   **Information for more details**):

526
$$\boldsymbol{\beta}_{j(k+1),M,s_m^*} = \boldsymbol{\beta}_{jk,M,s_m^*} + \frac{1}{\left(\left(e_M^{(0,k)} - \frac{\partial_{\alpha_{j,M}}K_e^{(\alpha_{j,M})}}{N_j}\right)\left(e_M^{(2,k)} - \frac{\partial_{\beta_{j,M}}K_e^{(\beta_{j,M})}}{N_j}\right) - \left(e_M^{(1,k)} - \frac{\partial_{\beta_{j,M}}K_e^{(\alpha_{j,M})}}{N_j}\right)^2\right)}$$

527
$$\times \begin{pmatrix} e_M^{(2,k)} - \dfrac{\partial_{\alpha_{j,M}} K_{\circleddash}^{(\alpha_{j,M})}}{N_j} & -e_M^{(1,k)} + \dfrac{\partial_{\beta_{j,M}} K_{\circleddash}^{(\alpha_{j,M})}}{N_j} \\[2em] -e_M^{(1,k)} + \dfrac{\partial_{\beta_{j,M}} K_{\circleddash}^{(\alpha_{j,M})}}{N_j} & e_M^{(0,k)} - \dfrac{\partial_{\beta_{j,M}} K_{\circleddash}^{(\beta_{j,M})}}{N_j} \end{pmatrix} \begin{pmatrix} \pi_j^* - \tilde{e}_M^{(0,k)} + K_{\circleddash}^{(\alpha_{j,M})} / N_j \\[2em] \dfrac{S_M \pi_j^* MAF_{j,M}^*}{w_{s_m^*}} - \tilde{e}_M^{(1,k)} + K_{\circleddash}^{(\beta_{j,M})} / N_j \end{pmatrix},$$

528

529    where

530
$$K_{\circleddash}(\boldsymbol{\beta}_{j,M}) = \frac{1}{2\left(e_{j,M}^{(0,k)} e_{j,M}^{(2,k)} - \left(e_{j,M}^{(1,k)}\right)^2\right)} \begin{pmatrix} d_{\alpha,j,M}^{(0,k)} e_{j,M}^{(2,k)} + d_{\alpha,j,M}^{(2,k)} e_{j,M}^{(0,k)} - 2 d_{\alpha,j,M}^{(1,k)} e_{j,M}^{(1,k)} \\[1em] d_{\beta,j,M}^{(0,k)} e_{j,M}^{(2,k)} + d_{\beta,j,M}^{(2,k)} e_{j,M}^{(0,k)} - 2 d_{\beta,j,M}^{(1,k)} e_{j,M}^{(1,k)} \end{pmatrix}$$

531    and

532
$$d_{\alpha,j,M}^{(c+1,k)} = d_{\beta,j,M}^{(c,k)} = \frac{\partial e_{j,M}^{(c,k)}}{\partial \beta_{j,M}}.$$

533

## isGWAS is computed using sufficient statistics

535    Under Hardy-Weinberg equilibrium, the quadruple $\left\{N_j, N_j^*, MAF_{j,M}, MAF_{j,M}^*\right\}$ are combined to

536    form the global and local (under a wide radius of convergence) sufficient statistics from the

537    logistic model. Consequently, they hold all necessary information to compute regression

538    parameter estimates $\left\{\hat{\alpha}_{j,M}, \hat{\beta}_{j,M,s_m^*}, \hat{\sigma}_{\alpha_{j,M}}, \hat{\sigma}_{\beta_{j,M,s_m^*}}\right\}$ over a broad range of scenarios. Regardless

539    of Hardy-Weinberg being valid or not, we show that the triple $\left\{T_{1j}, T_{2j}, T_{3j}\right\}$,

540
$$\left\{ T_{1j} = \sum_{i=1}^{N_j} y_i = N_j^*, \qquad T_{2j} = \sum_{i=1}^{N_j} y_i \, g_{ij,M,s_m^*}, \qquad T_{3j} = \sum_{i=1}^{N_j} g_{ij,M,s_m^*} \right\},$$

541    are the two global and one local sufficient statistics and these can alternatively be used as input

542    variables in isGWAS. To show this, we write:

543
$$L\left(\boldsymbol{\beta}_{j,M,s_m^*}\right) = \sum_{i=1}^{N_j} \log P\left(y_i \mid \boldsymbol{\beta}_{j,M,s_m^*}, g_{ij,M,s_m^*}\right)$$

26

544
$$= \alpha_{j,M,s_m^*} \sum_{i=1}^{N_j} y_i + \beta_{j,M,s_m^*} \sum_{i=1}^{N_j} y_i\, g_{ij,M,s_m^*}$$

545
$$+ \sum_{i=1}^{N_j} \log\left(1 - P\left(y_i = 1 \,|\, \boldsymbol{\beta}_{j,M,s_m^*}, g_{ij,M,s_m^*}\right)\right)$$

546
$$= \alpha_{j,M,s_m^*} T_{1j} + \beta_{j,M,s_m^*} T_{2j} - \sum_{i=1}^{N_j} \log\left(1 + \exp \alpha_{j,M,s_m^*}\right)$$

547
$$- \sum_{i=1}^{N_j} \log\left(1 + \frac{\exp \alpha_{j,M,s_m^*}}{1 + \exp \alpha_{j,M,s_m^*}}\left(\left(\exp \beta_{j,M,s_m^*}\, g_{ij,M,s_m^*}\right) - 1\right)\right)$$

548
$$= \alpha_{j,M,s_m^*} T_{1j} + \beta_{j,M,s_m^*}\left(T_{2j} - \frac{\exp \alpha_{j,M,s_m^*}}{1 + \exp \alpha_{j,M,s_m^*}} T_{3j}\right) - N_j \log\left(1 - \exp \alpha_{j,M,s_m^*}\right)$$

549
$$+ \mathcal{O}\left(\frac{\exp \alpha_{j,M,s_m^*}}{1 + \exp \alpha_{j,M,s_m^*}}\left(\beta_{j,M,s_m^*} g_{\cdot j,M,s_m^*}\right)^2\right)$$

550
$$= f\left(T_{1j}, T_{2j}, T_{3j}; \alpha_{j,M,s_m^*}, \beta_{j,M,s_m^*}\right) + \mathcal{O}\left(\frac{\exp \alpha_{j,M,s_m^*}}{1 + \exp \alpha_{j,M,s_m^*}}\left(\beta_{j,M,s_m^*} g_{\cdot j,M,s_m^*}\right)^2\right)$$

551 and valid when

552
$$\frac{\exp \alpha_{j,M,s_m^*}}{1 + \exp \alpha_{j,M,s_m^*}}\left|\left(\exp \beta_{j,M,s_m^*}\, g_{\cdot j,M,s_m^*}\right) - 1\right| < 1.$$

553 Hence, the global sufficient statistics are $\{T_{1j}, T_{2j}\}$ and (on assuming random $g_{ij,M,s_m^*}$ as in the

554 SaLN-R algorithm) the locally sufficient statistic is $\{T_{3j}\}$, where:

555
$$T_{1j} = \sum_{i=1}^{N_j} y_i = N_j^*, \qquad T_{2j} = \sum_{i=1}^{N_j} y_i\, g_{ij,M,s_m^*}$$

556 and

557
$$T_{3j} = \sum_{i=1}^{N_j} g_{ij,M,s_m^*} = \begin{cases} n_{j1}, & S_M = 1, \\ n_{j1} + 2n_{j2}, & S_M = 2. \end{cases}$$

558 Under Hardy-Weinberg equilibrium, we can write

559
$$T_{2j} = s_m^* N_j^* MAF_j^* \quad \text{and} \quad T_{3j} = s_m^* N_j MAF_j.$$

27

560

**Leapfrog re-sampler: forecasting results in target sample sizes**

To estimate regression parameters $\{\alpha_{j,M,s_m^*}, \beta_{j,M,s_m^*}\}$ in larger target sample sizes, i.e., $\widehat{N}_j > N_j$, we propose the following strategy:

1. **Specify number $K$, sub-sample $\gamma_1$ and target sample $\gamma_2$ parameters**, where $K \geq 1$, $0 < \gamma_1 < 1$ and $\gamma_2 > 1$.

2. **Generate random sub-samples of individuals of size** $\widetilde{N}_j = \gamma_1 N_j < N_j$. For each of $k = 1, 2, \ldots, K$, generate a random sub-sample $D_{k,\gamma_1} \subset D$, where $\left|D_{k,\gamma_1}\right| = \widetilde{N}_j = \gamma_1 N_j$.

3. **(Leapfrog-step) Compute subsample quadruple and project to target sample size** $\widehat{N}_j$. For each subsample $D_{k,\gamma_1}$, compute values $\left\{\widetilde{N}_{ij}^*, \widehat{MAF}_{kj,M}, \widehat{MAF}^*_{kj,M}\right\}$ and project these on to the target sample size, i.e., $d_{k,\gamma_{1,2}} = \left\{\left(\frac{\gamma_2}{\gamma_1}\right)\widetilde{N}_{kj}, \left(\frac{\gamma_2}{\gamma_1}\right)\widetilde{N}_{ij}^*, \widehat{MAF}_{kj,M}, \widehat{MAF}^*_{kj,M}\right\}$ for sample $D_{k,\gamma_1}$

   - Note that $\left(\frac{\gamma_2}{\gamma_1}\right)\widetilde{N}_{kj} = \widehat{N}_j$, which is the target 'future' sample size.

4. **Deploy isGWAS across all $K$ (projected) quadruples** $d_{k,\gamma_{1,2}}$ and record each estimate of the genetic effects, standard error and p-value $\left\{\hat{\beta}_{k,j,M,s_m^*}, \hat{\sigma}_{\beta_{k,j,M,s_m^*}}, p_{k,j,M,s_m^*}\right\}_{k=1:K}$.

5. **Estimate p-value in target sample size** as a summary point estimate (e.g., median) or range across all $K$ sub-samples,

$$p_{target_{j,M,s_m^*}} = \text{median}\{p_{k,j,M,s_m^*}\}_{k=1:K}.$$

578

### Data Quality Control: preparation of sufficient statistics for isGWAS

580  In order to deploy isGWAS successfully, the sufficient statistics are required to be

581  prepared in a sample where only a single individual (preferably case) from pairs or n-

582  tuples of 3rd, 2nd and 1st degree relatives is retained. Additionally, ethnical outliers must

583  also be removed. In summary, to deploy isGWAS successfully we require either: (a)

584  access to the sufficient statistics computed after duplications of related n-tuples and

585  ethnical outliers are removed; or (b) access to the individual level data, whereupon the

586  sufficient statistics can be prepared as described in (a). We provide a detailed outline of

587  recommended Quality Control for genetic Individual Level Data (ILD) to running

588  successfully isGWAS in **Supplementary Information**.

589

### Application to Biobank data

591  The GWAS results used in the assessment of isGWAS were taken from large-scale analyses of

592  UK Biobank[13], Biobank Japan[14] and the Psychiatric Genomics Consortium[15].

593  The UK Biobank[13] is a large-scale biomedical database and research resource containing

594  in-depth genetic and health information from half a million UK participants. From the full

595  available UK Biobank cohort, we obtain phenotypes for seven different diseases with varying

596  levels of prevalence. These are Hypertension (IC10:I10), Asthma (IC10:J45), Atherosclerosis

597  (IC10:I25), Glaucoma (IC10:H40), Stroke (IC10:I63), Colon Cancer (IC10:C18) and Thyroid

598  Gland Cancer (IC10:C73) patients. From a total cohort of 502,422 participants, we used the

599  following inclusion criteria: white British (Field 22006), non-related (>3rd degree), no

600  patients with difference in reported (Field 31) and genetic (Field 22001) sex, no patients with

601  aneuploidy (Field 22019), no patients with unusual heterozygosity and high missing rates

602    (Field 22027). The ethnicity component is obtained from samples who self-identified as

603    'White British' according to Field 21000 and have very similar genetic ancestry based on a

604    principal components analysis of the genotypes. Retaining one related individual (where we

605    favour the retention of cases) we obtain a working sample size of ~335,000 individuals; the

606    approximate value is owing to small differences in the number of cases between disease

607    phenotypes (**Supplementary Information**). Comparative analysis for these varying

608    populations is reported in the main text.  The prevalence ratios and exact number of cases and

609    controls are provided in **Supplementary Table 1**. The variant based statistics needed for

610    isGWAS were obtained from the imputed UK Biobank dataset. A quality info score>0.9 is

611    applied to the data, and the number of cases and controls per variant and the MAF for variant

612    in cases and controls is based on patients with non-missing genotypes for the variant using

613    software PLINK[4]. Sample-level MAF>0.001 is used as inclusion criteria for the variants to

614    analyse. For each disease, we run isGWAS analysis using default settings under the 'additive'

615    genetic model. In addition, we also perform GWAS analysis using two-step REGENIE[8]

616    applied to all variants with a MAF>0.001 and Genotype Score>0.99. Firth correction was

617    enabled and performed on variants with p-value<0.1. REGENIE was also adjusted for

618    covariate information (age, sex, ancestry). For each disease we provide the following

619    diagnostic plots: 1) mirrored Manhattan plot comparing directly p-values for isGWAS and

620    REGENIE, 2) p-value – p-value plots comparing REGENIE and isGWAS, 3) $\beta - \beta$ plots

621    comparing REGENIE and isGWAS where we have colored the values by a) MAF and b)

622    ratios of computed standard error (SE) between methods, i.e., $log_2(\frac{SE(isGWAS)}{SE(REGENIE)})$ . Across all

623    diseases and variants considered, we compare performance of isGWAS and isGWAS-Firth to

624    REGENIE-Firth.

625    Schizophrenia data from the Psychiatric Genomics Consortium[15] was used to conduct two

626    different large-scale GWAS analysis. The first GWAS analysis was executed with data from

627　77,096 European individuals (33,640 cases, 43,456 controls)[25]. The second GWAS analysis

628　was executed with data from the larger 130,644 European individuals (53,386 cases, 77,258

629　controls)[27]. We used the 2014 dataset to infer the 2022 results. To do this, we refine the

630　significant results from both imputed 2014 and 2022 summary statistics using clumping with

631　PLINK. The European 1000 Genomes Project v3[19] dataset was used as a reference population

632　for the clumping procedure. Twelve strategies for clumping were explored: three were LD $R^2$-

633　based only, the other nine were a combination of clumping by LD block information and p-

634　value thresholding. The refined variants are used to assess the inference capabilities of isGWAS

635　both within each of the two datasets and the enrichment capabilities of isGWAS to infer p-

636　values of the 2022 dataset using the 2014 dataset. For the 2014 dataset, 225 variants were

637　remaining after the clumping. For the 2022 dataset, 451 variants were remaining after the

638　clumping. From those, 54 are overlapping and 608 is the unique set between the two datasets.

639　The Biobank Japan data was used to conduct a large-scale GWAS with 212,453 Japanese

640　individuals across 42 different diseases[24]. We obtained the published significantly associated

641　loci (P < 5e-08) in autosomes from the GWAS findings which amounted to 309 variants across

642　30 different diseases. Similarly, we used the significantly associated X chromosome findings

643　for males and females that amounted to a total of nine significantly associated loci across five

644　diseases, although results are omitted from text. We applied isGWAS to the three different sets

645　of variants using default parameters to assess the performance of isGWAS. To aid association

646　interpretation, we use the following additional statistical tests to assess the accuracy and

647　sensitivity of the isGWAS calculator for the Biobank Japan data. First, a classical ROC curve

648　was produced where the true/false actual value was determined by various p-value thresholds

649　(benchmarked against published Biobank Japan results). The isGWAS calculator is an

650　inferential tool thus this usage of the ROC curve is unconventional, however, it provides us

651　with the opportunity to assess the sensitivity to the choice of thresholds used to correct for

652 multiple testing. These are $10^{-10}$, $10^{-8}$, $5 \times 10^{-8}$, $10^{-7}$, where we have also used

653 $9.58 \times 10^{-9}$ for Biobank Japan as recommended by the authors[24]. AUC values were not

654 obtained as this is not a standard classification problem and they are not interpretable in this

655 context. Second, an adapted ROC curve was produced which accounts for two different

656 thresholds – one more stringent one to determine the true positive rate and one less stringent

657 one to determine the true negative rate. **Supplementary Figure 22** showcases this scenario and

658 highlights the importance of a threshold choice and its impact on a sensitivity analysis. The

659 main aim of isGWAS calculator is to be used as an inferential tool for truly significant or truly

660 non-significant genetic signals. Thus, using two thresholds – one for truly significant and one

661 for truly non-significant – provides us the assess the sensitivity of isGWAS to this scientific

662 question. Third, the obtained $\beta$ values were compared to the true ones by obtaining the

663 percentage of 1) predicted $\beta$ values in the 95% C.I.s of the true $\beta$ values and 2) 95% C.I.s of

664 the predicted $\beta$ values in the 95% C.I.s of the true $\beta$ values.

**Simulation scenarios**

665

666 In the first scenario, for each individual $i$ and iteration index $k$, we randomly generate disease

667 status via $y_{ik} \sim Ber(\pi_{ik}; \alpha_k, \beta)$ with probability of disease $\pi_{ik} = expit(\alpha_k + \beta g_{ik})$ and

668 $g_{ik} \sim Bin(2, MAF_k)$. Minor allele frequency is randomly selected from the set $MAF_k \in$

669 $\{10^{-4}, 5 \times 10^{-4}, 0.01\} \cup \{0.025, 0.05, \ldots, 0.5\}$ and the genetic effect on disease risk is fixed as

670 $\beta = 0.5$. In the second study, we allow the genetic effect to vary, i.e., $\beta \equiv \beta_k$, by fixing disease

671 status per individual and generating genotype data in controls $g_{ik}|y_{ik} = 0 \sim Bin(2, MAF_k)$ or

672 cases $g_{ik}|y_{ik} = 1 \sim Bin(2, MAF_k^*)$, where minor allele frequency in cases is taken as the outer

673 product with the sample minor allele frequency, with a random increase or decrease in

674 frequency (which controls the magnitude and direction of genetic effect), i.e., we introduce the

675 set $MAF_k^* \in MAF_k \otimes (1 \pm MAF_k)$. The parameter $\beta_k$ is then estimated via each of the 5

676 estimators using the vector of simulated data $\{\boldsymbol{y}_k, \boldsymbol{g}_k\}$.

677   We compare isGWAS and isGWAS-Firth against classical logistic and Firth corrected

678   regression[16], [37], [38]. Details for the second scenario are provided alongside full

679   description of the simulation protocol in the **Supplementary Information**.

680

681   Leapfrog re-sampler: simulation and real-data analyses

682   The parameters $\{K, \gamma_1, \gamma_2\}$ in the leapfrog re-ampler are assessed over a variety of values. To

683   attenuate the computational burden of a 3-dimensional grid search, we considered scenarios in

684   which: $K = 100$, $\gamma_1 = {}^1\!/_{\gamma_2}$ and a $\gamma_2$-fold increase in sample size of $\gamma_2 \in$

685   $\{1.1, 1.25, 1.5, \dots, 2.5\}$, i.e., a 10% to 150% increase in sample size. Furthermore, we take our

686   working sample size to be 276,204 individuals, which matches the number of all unrelated

687   individuals in our UKB sample (i.e., on not retaining any member of a related pair – which is

688   therefore fixed between diseases). We used our simulation protocol (**Supplementary**

689   **Information**) to generate synthetic samples and additionally assessed performance across all

690   seven disease datasets in UK Biobank. Variants for assessment were selected after pruning in

691   PLINK[4] was applied to the ~11 million variants with the following parameters: genotype

692   quality>0.99, MAF>0.01, HWE $p < 10e - 15$, 1000 bp windows, 100 variant increments,

693   $R^2 > 0.9$. From the pruned variants, 5% were selected uniformly from variants with $p >$

694   $10^{-6}$ and all variants with $p \leq 10^{-6}$ were retained. Final number of variants progressed for

695   LRS for the seven diseases are provided in **Supplementary Table 10**. For simulated data,

696   data for smaller sub-samples were simulated using full cohort and empirical distributions for

697   MAF and disease prevalence. In our tests of the LRS, we assess the predictive properties of

698   isGWAS on real-life data where the ground truth is either computed from the entire sample or

699   provided in the literature.  Predictions from X-fold increases in sample size are compared

700   using standard accuracy, FDR, FPR and TPR measures based on a putative true significance

701   threshold of $5e - 08$.

**Computational resources**

Real-life analyses were performed using up to 48 virtual CPU cores of a 2.5 GHz Intel Xeon Gold 6240R processor with 64 GB of memory. Simulation analyses were performed using up to 8 virtual CPU cores of a 2.4 GHz Intel Core i9 processor.

*Computational comparison protocol*

We contrast the computational performance of isGWAS and REGENIE (Step-2 only). For clarity, REGENIE Step-1 simplifies the outcome and model by projecting out covariate information, before variant-disease association analyses are performed in Step-2. To directly compare both methods, we performed individual GWA analyses of each of the seven diseases considered in UK Biobank across ~11m variants for ~335,000 individuals. Owing to computational cost of the ILD method, we summarise results from a single GWA analysis per trait. Performance of isGWAS across repeated runs, for varying numbers of SNPs and available CPUs, up to a maximum of 10m variants, is also performed.

# Data availability

The genotype data, phenotype status and allele counts were extracted from UK Biobank[13] to support the findings of this study. The genome-wide association summary data with available allele frequencies and cohort counts that was used to support the findings of this study are available from: Psychiatric Genomics Consortium[15] and Biobank Japan[14].

# Code availability

The tool is available for use on the webportal www.optima-isgwas.com. The isGWAS algorithm is also available on github (https://github.com/cnfoley/isgwas/).

# References

724

725   [1]   W. Zhou *et al.*, 'Global Biobank Meta-analysis Initiative: Powering genetic discovery

726         across human disease', *Cell Genomics*, vol. 2, no. 10, Oct. 2022, doi:

727         10.1016/J.XGEN.2022.100192.

728   [2]   A. Abdellaoui, L. Yengo, K. J. H. Verweij, and P. M. Visscher, '15 years of GWAS

729         discovery: Realizing the promise', *Am. J. Hum. Genet.*, vol. 110, no. 2, pp. 179–194,

730         Feb. 2023, doi: 10.1016/J.AJHG.2022.12.011.

731   [3]   M. Claussnitzer *et al.*, 'A brief history of human disease genetics', *Nature*, vol. 577, no.

732         7789, pp. 179–189, Jan. 2020, doi: 10.1038/S41586-019-1879-7.

733   [4]   C. C. Chang, C. C. Chow, L. C. A. M. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee,

734         'Second-generation PLINK: Rising to the challenge of larger and richer datasets',

735         *GigaScience*, vol. 4, no. 1, Feb. 2015, doi: 10.1186/s13742-015-0047-8.

736   [5]   L. Jiang *et al.*, 'A resource-efficient tool for mixed model association analysis of large-

737         scale data', *Nat. Genet. 2019 5112*, vol. 51, no. 12, pp. 1749–1755, Nov. 2019, doi:

738         10.1038/s41588-019-0530-8.

739   [6]   P. R. Loh *et al.*, 'Efficient Bayesian mixed-model analysis increases association power

740         in large cohorts', *Nat. Genet. 2015 473*, vol. 47, no. 3, pp. 284–290, Feb. 2015, doi:

741         10.1038/ng.3190.

742   [7]   W. Zhou *et al.*, 'Efficiently controlling for case-control imbalance and sample

743         relatedness in large-scale genetic association studies', *Nat. Genet. 2018 509*, vol. 50, no.

744         9, pp. 1335–1341, Aug. 2018, doi: 10.1038/s41588-018-0184-y.

745   [8]   J. Mbatchou *et al.*, 'Computationally efficient whole-genome regression for quantitative

746         and binary traits', *Nat. Genet.*, vol. 53, no. 7, pp. 1097–1103, Jul. 2021, doi:

747         10.1038/s41588-021-00870-7.

748   [9]   N. Homer *et al.*, 'Resolving Individuals Contributing Trace Amounts of DNA to Highly

749         Complex Mixtures Using High-Density SNP Genotyping Microarrays', *PLOS Genet.*,

750         vol. 4, no. 8, p. e1000167, Aug. 2008, doi: 10.1371/JOURNAL.PGEN.1000167.

751   [10]  X. Shi and X. Wu, 'An overview of human genetic privacy', *Ann. N. Y. Acad. Sci.*, vol.

752         1387, no. 1, p. 61, Jan. 2017, doi: 10.1111/NYAS.13211.

753   [11]  E. W. Clayton *et al.*, 'The law of genetic privacy: applications, implications, and

754         limitations', *J. Law Biosci.*, vol. 6, no. 1, pp. 1–36, Oct. 2019, doi:

755         10.1093/JLB/LSZ007.

756   [12]  Z. Wan, J. W. Hazel, E. W. Clayton, Y. Vorobeychik, M. Kantarcioglu, and B. A. Malin,

757         'Sociotechnical safeguards for genomic data privacy', *Nat. Rev. Genet. 2022 237*, vol.

758         23, no. 7, pp. 429–445, Mar. 2022, doi: 10.1038/s41576-022-00455-y.

759   [13]  'UK Biobank'. https://www.ukbiobank.ac.uk/ (accessed Jul. 27, 2023).

760   [14]  'BioBank Japan'. https://biobankjp.org/en/ (accessed Jul. 27, 2023).

761   [15]  'PGC – Psychiatric Genomics Consortium'. https://pgc.unc.edu/ (accessed Jul. 27,

762         2023).

763   [16]  D. Firth, 'Bias Reduction of Maximum Likelihood Estimates', *Source Biom.*, vol. 80, no.

764         1, pp. 27–38, 1993.

765   [17]  P. D. Sasieni, 'From Genotypes to Genes: Doubling the Sample Size', *Biometrics*, vol.

766         53, no. 4, p. 1253, Dec. 1997, doi: 10.2307/2533494.

767   [18]  B. B. Sun *et al.*, 'Genetic associations of protein-coding variants in human disease',

768         *Nature*, vol. 603, no. 7899, pp. 95–102, Mar. 2022, doi: 10.1038/S41586-022-04394-W.

769   [19]  '1000 Genomes | A Deep Catalog of Human Genetic Variation'.

770         https://www.internationalgenome.org/home (accessed Jul. 27, 2023).

771   [20]  'Estonian Biobank', *Tartu Ülikool*, Dec. 13, 2021.

772         https://genomics.ut.ee/en/content/estonian-biobank (accessed Jul. 27, 2023).

773  [21] 'FinnGen research project is an expedition to the frontier of genomics and medicine |

774        FinnGen'. https://www.finngen.fi/en (accessed Jul. 27, 2023).

775  [22] R. A. Fisher, 'On the Interpretation of χ 2 from Contingency Tables, and the Calculation

776        of P', *J. R. Stat. Soc.*, vol. 85, no. 1, p. 87, Jan. 1922, doi: 10.2307/2340521.

777  [23] Q. Wang *et al.*, 'Rare variant contribution to human disease in 281,104 UK Biobank

778        exomes', *Nat. 2021 5977877*, vol. 597, no. 7877, pp. 527–532, Aug. 2021, doi:

779        10.1038/s41586-021-03855-y.

780  [24] K. Ishigaki *et al.*, 'Large-scale genome-wide association study in a Japanese population

781        identifies novel susceptibility loci across different diseases', *Nat. Genet.*, vol. 52, no. 7,

782        pp. 669–679, Jul. 2020, doi: 10.1038/s41588-020-0640-3.

783  [25] S. Ripke *et al.*, 'Biological insights from 108 schizophrenia-associated genetic loci',

784        *Nature*, vol. 511, no. 7510, pp. 421–427, Jul. 2014, doi: 10.1038/nature13595.

785  [26] M. Lam *et al.*, 'Comparative genetic architectures of schizophrenia in East Asian and

786        European populations', *Nat. Genet.*, vol. 51, no. 12, pp. 1670–1678, Dec. 2019, doi:

787        10.1038/s41588-019-0512-x.

788  [27] V. Trubetskoy *et al.*, 'Mapping genomic loci implicates genes and synaptic biology in

789        schizophrenia', *Nature*, vol. 604, no. 7906, pp. 502–508, Apr. 2022, doi:

790        10.1038/s41586-022-04434-5.

791  [28] M. Pirinen, P. Donnelly, and C. C. A. Spencer, 'Including known covariates can reduce

792        power to detect genetic effects in case-control studies', *Nat. Genet. 2012 448*, vol. 44,

793        no. 8, pp. 848–851, Jul. 2012, doi: 10.1038/ng.2346.

794  [29] J. Mefford and J. S. Witte, 'The Covariate's Dilemma', *PLOS Genet.*, vol. 8, no. 11, p.

795        e1003096, Nov. 2012, doi: 10.1371/JOURNAL.PGEN.1003096.

796 [30] L. J. O'Connor, A. P. Schoech, F. Hormozdiari, S. Gazal, N. Patterson, and A. L. Price,

797 'Extreme Polygenicity of Complex Traits Is Explained by Negative Selection', *Am. J.*

798 *Hum. Genet.*, vol. 105, no. 3, p. 456, Sep. 2019, doi: 10.1016/J.AJHG.2019.07.003.

799 [31] N. Chatterjee, B. Wheeler, J. Sampson, P. Hartge, S. J. Chanock, and J. H. Park,

800 'Projecting the performance of risk prediction based on polygenic analyses of genome-

801 wide association studies', *Nat. Genet. 2013 454*, vol. 45, no. 4, pp. 400–405, Mar. 2013,

802 doi: 10.1038/ng.2579.

803 [32] G. Kimmel and R. Shamir, 'A fast method for computing high-significance disease

804 association in large population-based studies', *Am. J. Hum. Genet.*, vol. 79, no. 3, pp.

805 481–492, 2006, doi: 10.1086/507317.

806 [33] P. M. Visscher *et al.*, 'Statistical Power to Detect Genetic (Co)Variance of Complex

807 Traits Using SNP Data in Unrelated Samples', *PLOS Genet.*, vol. 10, no. 4, p. e1004269,

808 2014, doi: 10.1371/JOURNAL.PGEN.1004269.

809 [34] 'Our Future Health', *Our Future Health*. https://ourfuturehealth.org.uk/ (accessed Jul.

810 27, 2023).

811 [35] M. I. Kurki *et al.*, 'FinnGen provides genetic insights from a well-phenotyped isolated

812 population', *Nat. 2023 6137944*, vol. 613, no. 7944, pp. 508–518, Jan. 2023, doi:

813 10.1038/s41586-022-05473-8.

814 [36] J. K. Pritchard and P. Donnelly, 'Case-control studies of association in structured or

815 admixed populations', *Theor. Popul. Biol.*, vol. 60, no. 3, pp. 227–237, Nov. 2001, doi:

816 10.1006/tpbi.2001.1543.

817 [37] R. Puhr, G. Heinze, M. Nold, L. Lusa, and A. Geroldinger, 'Firth's logistic regression

818 with rare events: accurate effect estimates and predictions?', *Stat. Med.*, vol. 36, no. 14,

819 pp. 2302–2317, Jun. 2017, doi: 10.1002/SIM.7273.

820 [38] G. Heinze, M. Ploner, D. Dunkler, H. Southworth, L. Jiricka, and G. Steiner, 'logistf:

821 Firth's Bias-Reduced Logistic Regression'. May 03, 2023. Accessed: Jul. 27, 2023.

822 [Online]. Available: https://cran.r-project.org/web/packages/logistf/index.html

823

# Acknowledgements

# Author contributions

CNF developed the mathematical and statistical methodologies, developed the statistical software and conceptualized the design. ZK developed the statistical software and webtool, designed the methodological analysis pipeline and conducted the real-life analyses. REM contributed to the interpretation of results. HR conceptualized and supervised the study. BBS conceptualized, designed the study, contributed to the method, application, and contextualization of the study. All authors contributed to the writing of the manuscript.

# Competing interests

BBS and HR are employed by Biogen. CNF and ZK are employed by Optima Partners. REM is an advisor to the Epigenetic Clock Development Foundation and Optima Partners.

# Tables

*Table 1. Accuracy, true positive rate, false positive rate and false discovery rate of isGWAS using REGENIE results as gold-standard and a threshold of $p = 5e - 08$ as classification rule. Results are obtained on all 11,079,229 variants used for the analysis of seven diseases in UK Biobank without clumping/finemapping.*

| Disease (ICD code) | Case:control ratio | *TPR* | *FPR* | *Acc* | *FDR* |
|---|---|---|---|---|---|
| **Hypertension (I10)** | 1:2 | 0.625 | 0.000026 | 0.99961 | 0.041 |
| **Asthma (J45)** | 1:6 | 0.982 | 0.000028 | 0.99994 | 0.016 |
| **Atherosclerosis (I25)** | 1:9 | 0.886 | 0.000011 | 0.99995 | 0.043 |
| **Glaucoma (H40)** | 1:26 | 0.891 | 0 | 0.99997 | 0.036 |
| **Stroke (I63)** | 1:56 | NA | 0 | 1 | NA |
| **Colon Cancer (C18)** | 1:94 | 0.944 | 0 | 0.99999 | 0.037 |
| **Thyroid Gland Cancer (C73)** | 1:669 | 1 | 0 | 0.99999 | 0.051 |

# Figures



*Figure 1. Diagram highlighting main differences between isGWAS and other GWAS approaches.*

*Figure 2. Comparative results for Asthma (IC10:J45) from UK Biobank. Subplot (a) is a mirror Manhattan plot comparing* $-\log_{10} P$ *values for isGWAS and REGENIE-Firth and subplot (b) is a locus zoom of the gene FLG2 region +/-250kbp on chromosome 7. Subplot (c) plots* $-\log_{10} P$ *values for isGWAS and REGNIE-Firth with the standard threshold P-value indicated colored by population-level MAF. Subplots (d) and (e) showcase* $\beta - \beta$ *effect size estimates for variants with p-values<0.05 and are coloured by population-level MAF and* $\log_2\left(\frac{SE(isGWAS)}{SE(REGENIE)}\right)$.

*Figure 3. a) Accuracy/TPR/FPR comparing REGENIE-Firth and isGWAS results, where $p = 5e - 08$ threshold has been used as indicator for correct classification accuracy. Results are obtained on all 11,079,229 variants used for the analysis without clumping/finemapping. See Supplementary Table 1 for full results. Manhattan plots b), c), d), e) and f): Comparative results for five diseases from UK Biobank. Mirror Manhattan plots comparing $-\log_{10} P$ values for isGWAS and REGENIE-Firth for six different diseases obtained from UK Biobank. Stroke was excluded from the analysis due to no variants passing significance threshold.*

*Figure 4. Simulation I results. Clockwise from top: a) Mean square error; b) distribution of estimated beta values; c) distribution of associated standard errors and d) distribution of -$\log_{10}$(p-value), for each model - logistic regression, firth regression, 'isGWAS', 'isGWAS_Firth' and, for p-values only, Fisher's Exact Test - and specification of disease prevalence. Panel e) presents the relationship between Firth regression derived -$\log_{10}$(p-value), along the horizontal axis, and the corresponding 'isGWAS_Firth' and Fisher's Exact Test (FET) computed values, on the vertical axis. Results are presented in the range [0,15] as FET regularly failed to converge for very small p-values and coloured according to value of $\pi\delta_{MAF}$, where*

$\pi$ denotes prevalence and $\delta_{MAF} = \frac{(MAF^* - MAF)}{(MAF(1 - MAF))}$. In panels b)-d) a point denotes the median value and error-bars the first and ninth deciles of the range.
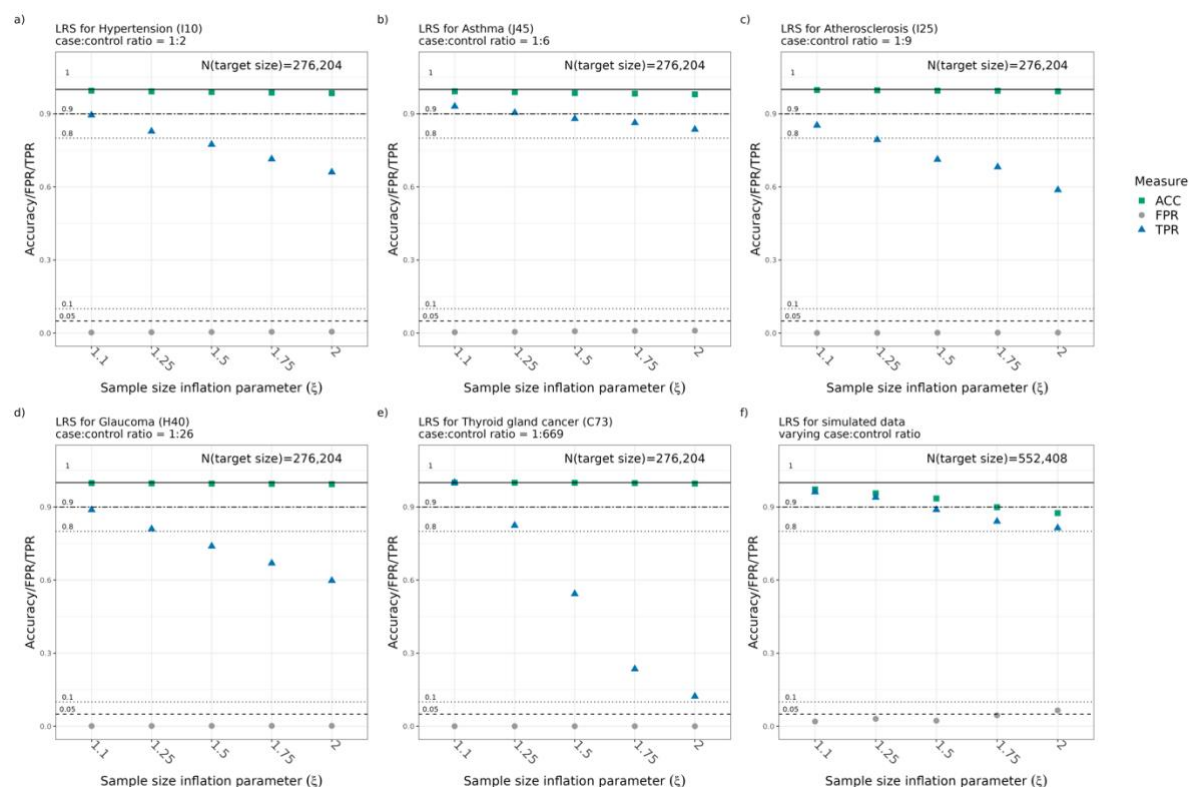


Figure 5. Performance of Leapfrog re-sampler (LRS) benchmarked against results derived from a)-e) the sample of ~276,204 individuals from UK Biobank with five different diseases and b) a simulated sample of 552,408 individuals (i.e., double UK Biobank sample size). For each value of $\xi$, we subset the target sample size down to $\frac{N_{target}}{\xi}$ individuals and deploy the LRS to compute predictions for the target sample $N_{target}$. As the maximum number of UK Biobank samples was 276,204, this was taken as the target. For example, when $\xi = 2$, we subset the full sample to 138,102 individuals and run the LRS to compute predictions of the larger 276,204 sample. We use results from the disease analysis, benchmarking LRS predictions against those computed on the pruned genome sampling uniformly across significance associations, resulting in ~3500 variants per studied disease. Colon cancer and Stroke are excluded from this figure as they don't have significant variants or a very low number of such after pruning. In the right panel we generated 1000 simulated datasets under the null of no genetic association or the alternative (see simulation protocol for details).
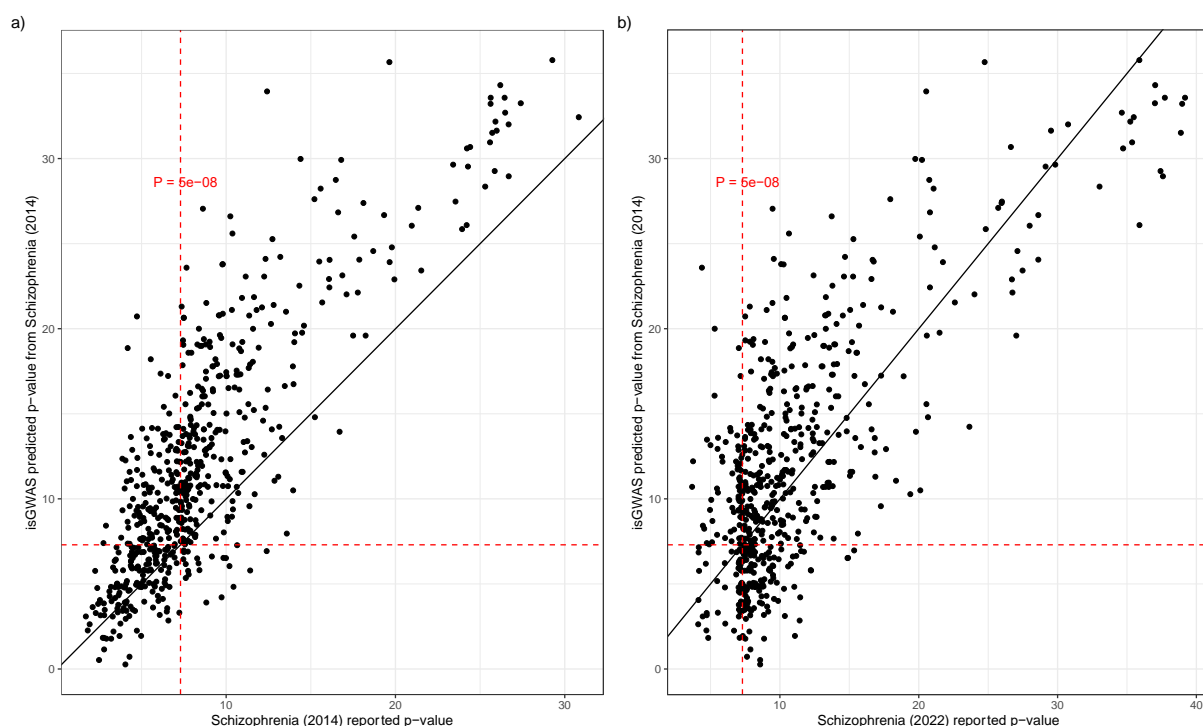
45

*Figure 6. Prediction results for Schizophrenia (2022) using data from Schizophrenia (2014): population-level information for 608 significantly associated loci (P<1e-07) obtained from clumping with parameters ($R^2 = 0.2, p_1 = 1e - 7, p_2 = 1e - 7$) has been used to infer p-values. a) The figure compares reported GWAS Schizophrenia (2014) p-values and isGWAS predicted p-values using population-level information from Schizophrenia (2014) matching for the larger 2022 cohort size. b) The figure compares reported GWAS Schizophrenia (2022) p-values and isGWAS predicted p-values using population-level information from Schizophrenia (2014) matching for the larger 2022 cohort size. The dashed red line represents the threshold P=5e-08.*
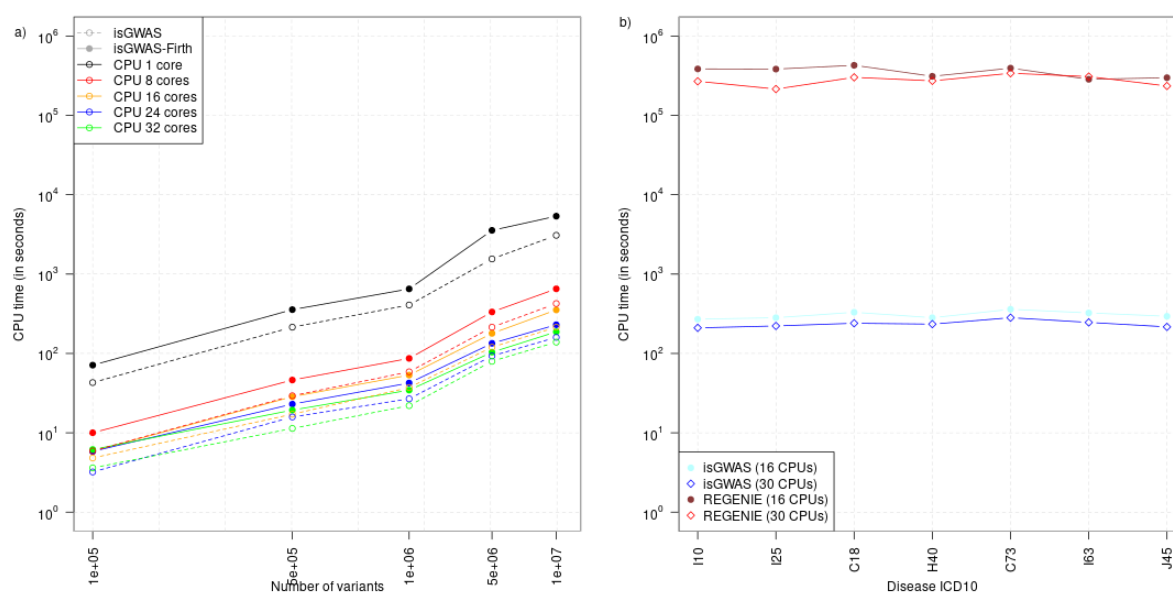
46

*Figure 7. a) Computational CPU time (in seconds) for an increasing number of variants. The results compare performance of isGWAS with and without Firth distributed over different number of CPU cores. The data was obtained from UK Biobank ICD10:C73 disease with low disease prevalence (case-control ratio = 1:669). The x- and y-axis are on $log_{10}$ scale. b) Computational CPU time (in seconds) for seven UK Biobank diseases on 11,079,229 variants for ~335,000 individuals. We compare the performance of isGWAS running on 16 and 30 CPUs vs the performance of REGENIE Step 2 running on 16 CPUs (with 16 threads) and 30 CPUs (with 30 threads). The computation is performed on the same machine. The y-axis is on $log_{10}$ scale.*