

Controlling for Background Genetic Effects Using Polygenic Scores Improves the Power of Genome-wide Association Studies

Declan Bennett

National University of Ireland, Galway

Dónal O'Shea

National University of Ireland, Galway

John Ferguson

National University of Ireland, Galway

Derek Morris

National University of Ireland, Galway

Cathal Seoighe (✉ Cathal.Seoighe@nuigalway.ie)

National University of Ireland, Galway

Research Article

Keywords: complex diseases, genetic effects, polygenic score, GWAS results

Posted Date: September 17th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-873301/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Controlling for background genetic effects using polygenic scores improves the power of genome-wide association studies

Declan Bennett¹, Donal O'Shea^{1,2}, John Ferguson^{1,3}, Derek Morris⁴, and Cathal Seoighe^{1,*}

¹School of mathematics, Statistics and Applied Mathematics, National University of Ireland Galway, H91TK33, Ireland

²SFI Centre for Research Training in Genomics Data Science, National University of Ireland Galway, H91TK33, Ireland

³Biostatistics Unit, Clinical Research Facility, National University of Ireland Galway, H91TK33, Ireland

⁴Centre for Neuroimaging, Cognition and Genomics, Discipline of Biochemistry, National University of Ireland Galway, H91CF50, Ireland.

*cathal.seoighe@nuigalway.ie

ABSTRACT

Ongoing increases in the size of human genotype and phenotype collections offer the promise of improved understanding of the genetics of complex diseases. In addition to the biological insights that can be gained from the nature of the variants that contribute to the genetic component of complex trait variability, these data bring forward the prospect of predicting complex traits and the risk of complex genetic diseases from genotype data. Here we show that advances in phenotype prediction can be applied to improve the power of genome-wide association studies. We demonstrate a simple and efficient method to model genetic background effects using polygenic scores derived from SNPs that are not on the same chromosome as the target SNP. Using simulated and real data we found that this can result in a substantial increase in the number of variants passing genome-wide significance thresholds. This increase in power to detect trait-associated variants also translates into an increase in the accuracy with which the resulting polygenic score predicts the phenotype from genotype data. Our results suggest that advances in methods for phenotype prediction can be exploited to improve the control of background genetic effects, leading to more accurate GWAS results and further improvements in phenotype prediction.

1 Introduction

2 Linear mixed effects models (LMMs) are routinely applied to detect associations between SNPs and phenotypes in genome-wide
3 association studies (GWAS) and many methods have been developed that enable these models to be applied efficiently to the
4 large scale datasets that are typically now encountered in studies of complex traits¹⁻⁹. Compared to fixed effects models for
5 GWAS¹⁰, LMMs can be designed that have the advantage of being applicable to samples that include related individuals^{1,11,12}.
6 LMMs for this purpose typically include a random effect with covariance proportional to the kinship matrix that indicates
7 the degree of relatedness between pairs of individuals in the sample¹². The relatedness of individuals in the sample may be
8 known *a priori* or may be derived from the genotype data by constructing a genetic relationship matrix (GRM), with entries
9 corresponding to the genotypic covariance between pairs of individuals. When the entries of the GRM below a specified
10 threshold are set to zero, the GRM is approximately equivalent to a family kinship matrix, with the degree of relatedness that
11 the matrix captures controlled by this threshold. Thresholding the matrix to capture close family relationships (or cryptic
12 relatedness¹³) allows specialized computational methods for sparse matrices to be applied so that model fitting remains tractable

13 for studies that include large numbers of individuals⁹. This is the approach taken by fastGWA⁹, a recently developed tool that
14 has been shown to generate correctly calibrated statistical results efficiently for biobank-scale GWAS.

15 In addition to enabling application to samples containing related individuals, LMMs can also account for genetic background
16 effects^{11,14}. When a statistical model is used to test for a relationship between a given SNP (the test SNP) and a phenotype,
17 genetic variants in the genome that are not in linkage disequilibrium with the test SNP may also make a substantial contribution
18 to the phenotypic variation. If this contribution to phenotypic variation is not accounted for it contributes to the error term in
19 the model. If the trait of interest is both highly polygenic and highly heritable this noise may be substantial. Failure to account
20 for sources of variance in the response in a statistical model can reduce the power to detect a relationship of interest^{15,16}. A
21 LMM with a full GRM (i.e. derived from all SNPs in the data and with no threshold applied on the level of genetic correlation
22 between individuals) is equivalent to a model in which all variants are assumed to have a causal effect on the phenotype,
23 with effect sizes consisting of independent samples from a Gaussian distribution⁸. This is typically not a good fit to the true
24 effect size distribution, and instead, the software package BOLT-LMM⁸ uses a spike-and-slab Gaussian mixture for the effect
25 size distribution, with a component (the spike) close to zero corresponding to weak genome-wide effects and accounting for
26 family relationships, and component with larger variance (the slab) corresponding to variants with large effects⁸. Fitting this
27 more sophisticated model requires specialist numeric methods, that are relatively computationally intensive. Consequently
28 BOLT-LMM is much more computationally intensive than fastGWA⁹.

29 The full GRM is an $N \times N$ matrix, where N is the number of individuals in the study. The memory requirement of
30 BOLT-LMM is kept tractable by not explicitly evaluating the GRM but rather BOLT-LMM solves the mixed model equations
31 by computing the product of the inverse GRM and the phenotype vectors. Nonetheless, the overall compute time and memory
32 requirements of BOLT-LMM are a function of both N and the number of model SNPs, M , that contribute to the GRM (with
33 $O((NM)^{1.5})$ compute time and $\frac{NM}{4}$ bytes of memory required. Various options have been explored for which SNPs to include in
34 the calculation of the GRM¹¹. Including SNPs in LD with the target SNP results in loss of power, as the effect of the target SNP
35 is partially accounted for by the random effect through the GRM. This has been referred to as proximal contamination¹⁴. On
36 the other hand, including all (or most) SNPs that are not in LD with the target SNP, e.g. using a Leave One Chromosome Out
37 (LOCO) approach, can result in dilution of the extent to which the relevant part of the genetic background is captured by the
38 GRM. In the latter case, SNPs that are not relevant, in that they do not capture direct genetic effects or tag relevant population
39 structure effects, effectively add noise to the GRM¹⁴. Alternatively, the GRM can be built from only the SNPs that are found
40 using a linear model to be associated with the phenotype. Although this results in an increase in statistical power^{11,17,18}, it
41 does not fully control for population structure and is not recommended if population structure is of substantial concern^{8,11}.
42 Methods have been developed that incorporate principal components into the GRM calculation built from significant SNPs;
43 however, most of these methods are not suited to large biobank-scale data, without access to cloud computing or large compute
44 farms¹⁹⁻²¹. Background genetic effects can also be included in the statistical model as fixed effects and this is the recommended
45 approach when there are SNPs with large effect sizes¹¹. A model fitting approach to determine the SNPs to include as fixed

46 effects has been developed, and this also results in increased power in GWAS¹⁴.

47 As the genomic architecture of complex diseases is uncovered with the help of large biobanks, there is an advancing
48 prospect of predicting quantitative phenotypes and the risk of complex diseases from genotype data. Recent years have seen
49 substantial success and emerging clinical utility in phenotype prediction from polygenic scores (PGS)^{22,23}. PGS are constructed
50 from weighted sums of allele dosages, with the weights corresponding to the effects size of the variants. Risk variants (variants
51 associated with the phenotype) are typically inferred from the largest available GWAS, generally a meta-analysis. The clinical
52 potential of PGS has already been shown in complex diseases such as coronary artery disease (CAD), diabetes and cancer^{23–25}.
53 In CAD, the identification of individuals with similar risk to those with rare high-risk monogenic variants has been reported²⁴.
54 Similarly, in breast cancer, pathogenic variants in BRCA1/2 account for 25% of familial risk of the disease with genome wide
55 variants accounting for a further 18% of the risk^{26,27}. It is likely that in the future specialist machine learning methods will
56 be developed to predict phenotype from genotype²³, potentially achieving higher accuracy by incorporating the possibility of
57 non-additive effects.

58 Here, we set out an approach to GWAS that seeks to separate the model fitting at the test locus and estimation of the
59 genetic background effect. After carrying out an initial round of GWAS using an existing method, we derive a PGS for
60 each chromosome, using the summary statistics for SNPs on the remaining chromosomes. We refer to this as the Leave
61 One Chromosome Out Poly Genic Score (LOCO PGS). We then perform a second round of GWAS, including the relevant
62 LOCO PGS as a fixed effect to account for the contribution to the variation in the phenotype of SNPs that are not on the same
63 chromosome as the test locus. We tested this approach in two ways. Firstly, using simulated data we tested for an improvement
64 in power on the task of recovering known causal variants as a function of study size, number of causal variants and trait
65 heritability. In addition, we applied the method to standing height data from the UK Biobank and determined the number and
66 characteristics of additional variants that were detected. For an objective assessment of performance on real data, where the
67 true associations are unknown, we divided the data into test and training sets and predicted the phenotype in the test set. The
68 improvement in performance on the critical task of complex phenotype prediction illustrates the utility of the PGS as a means of
69 accounting for off target genetic effects. This straightforward, modular approach to accounting for genetic background effects
70 in GWAS has the advantage of leveraging advances in phenotype prediction as they become available. It also offers significant
71 improvements in speed relative to existing methods that correct for genetic background.

72 **Results**

73 We incorporated the LOCO PGS as a fixed effect in a linear mixed model using the existing tools, GCTA fastGWA, BOLT-LMM
74 and REGENIE^{8,9,28}. We refer to the methods that result from including the LOCO PGS fixed effect by appending PGS and the
75 name of the method used to calculate the PGS to the name of the original tool. For example, fastGWA with a LOCO PGS
76 fixed effect, calculated using the pruning and thresholding (P&T²⁹) or LDpred2³⁰ methods are referred to as fastGWA-PGS-PT
77 and fastGWA-PGS-LDpred2, respectively. We simulated data to evaluate the impact of including the LOCO PGS as a fixed

78 effect in GWAS. The simulations consisted initially of a normally-distributed continuous trait in 100,000 individuals. The
79 trait had a narrow-sense heritability (h^2) of 0.5 and there were 1,000 causal SNPs with normally-distributed effects on the
80 trait (see Methods for details). To check the validity of our approach we performed simulations under the null model of no
81 association between genotype and phenotype and found that the method was well calibrated (Fig. S1). This was the case both
82 for the P&T method of calculating the LOCO PGS (with a fixed P value threshold of 5×10^{-5}) and for the LDpred2 method
83 and was in-line with our expectations, as the LOCO PGS is approximately uncorrelated with the genotype of the tested SNP
84 (see Supplementary Material for a mathematical justification). The median false positive rate rose slightly when we used high
85 P-value thresholds ($P < 0.05$ and $P < 0.5$) with the P&T method to calculate the LOCO PGS (Fig. S1). In this case the majority
86 of the variants contributing to the LOCO-PGS are likely to be false positives and this may cause the method to become unstable,
87 as some individual simulations had false positive rates as high as 10% (in the case of the threshold of $P < 0.5$). We therefore
88 recommend against the use of these high thresholds if applying the P&T method to calculate the LOCO PGS. This issue did not
89 arise for LDpred2, which does not require a P-value threshold to be specified (Fig. S1).

90 In 100 simulations we found that including a LOCO PGS resulted in a substantial improvement in power to detect the known
91 causal SNPs (Fig. 1). When we included the PGS obtained using P&T as a fixed effect with fastGWA (i.e. fastGWA-PGS-PT)
92 we recovered 82 additional causal variants, on average, below the conventional P-value threshold of 5×10^{-8} compared to
93 fastGWA (corresponding to a relative increase in power of 18.4%; $p = 3.0 \times 10^{-32}$ from a paired T-test; Table S1-S3). The
94 performance was further improved when we used LDpred2 to calculate the LOCO PGS (fastGWA-PGS-LDpred2). This resulted
95 in the recovery of, on average, 115 more causal variants than fastGWA alone (relative increase of 25.9%; $p = 2.3 \times 10^{-36}$). We
96 also simulated case control data for a binary traits with h^2 of 0.5 and 1,000 causal loci, with disease prevalence, k , of 0.1 and
97 0.3. As with the quantitative trait simulations, inclusion of a LOCO PGS fixed effect always resulted in an increase in the
98 average number of casual loci recovered, with an average of 28 more causal loci recovered for a disease prevalence of $k=0.1$ (p
99 $= 0.19$) while, $k = 0.3$ recovered on average 48 more causal loci ($p= 0.03$) (Fig. S2 and Table S4 & S5).

100 The contribution to phenotype variance of background SNPs can also be modelled as a random effect in a linear mixed
101 model. This approach is applied by BOLT-LMM, which uses a normal mixture random effect, with a component corresponding
102 to SNPs with large effects. The running time of BOLT-LMM is proportional to $MN^{1.5}$ and the memory requirement is
103 approximately $MN/4$ bytes, where N is the number of individuals in the dataset and M is the number of SNPs included in the
104 GRM⁸. When we ran BOLT-LMM with a subset of 165,683 SNPs (see Methods for how these were selected) we found that
105 including the LOCO PGS as a fixed effect resulted in a substantial gain in power (Fig. 1), likely resulting from inability of the
106 reduced GRM to account fully for genetic background. No further improvement was obtained by adding the LOCO PGS fixed
107 effect to BOLT-LMM with a GRM consisting of all of the 664,393 directly genotyped SNPs (Fig. S3); however, the power
108 obtained with the smaller GRM with the PGS fixed effect was close to the power obtained with the larger GRM, but with a
109 much lower memory requirement (Table 1). The highest power of all methods was achieved with fastGWA-PGS-LDpred2,
110 which slightly exceeded the power of BOLT-LMM, even when all variants contributed to the GRM (i.e. BOLT-LMM-665 in

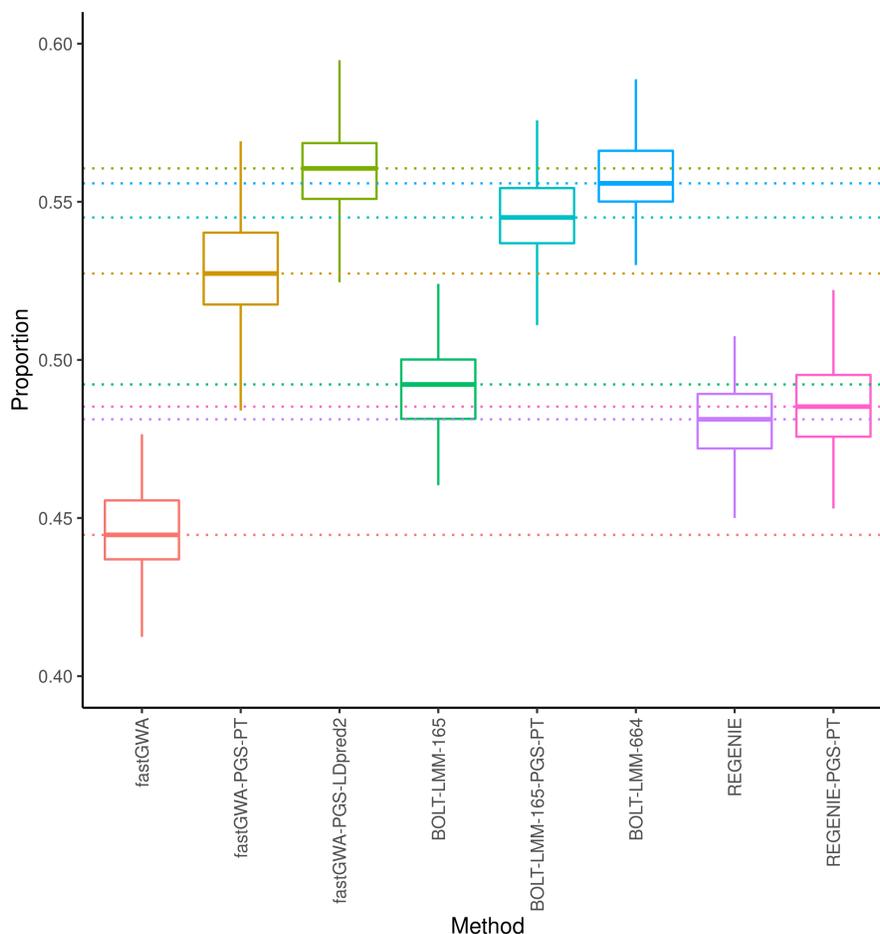


Figure 1. The proportion of causal variants recovered in 100 simulations. The boxplot shows the median (center line), upper and lower quartiles (hinges) and the maximum and minimum values not more than 1.5 times the interquartile range from the corresponding hinge (whiskers). The simulations consisted of 100,000 individuals and a continuous trait, with narrow-sense heritability of 0.5 and 1,000 causal variants. BOLT-LMM-165 denotes BOLT-LMM with a GRM derived from 165,684 variants resulting from strict LD-pruning. BOLT-LMM-664 refers to the use of BOLT-LMM with a GRM derived from all 664,393 variants in the simulations. Methods that include PGS in the name involved the use of a LOCO PGS fixed effect, derived either from pruning and thresholding (methods ending in PT) or using LDpred2.

111 Fig. 1). Recently, a new fast method, REGENIE²⁸, has been released that also includes control of the genetic background effect
 112 based on prediction of the phenotype from SNPs that are not on the same chromosome as the test SNP. In our simulations
 113 the performance of REGENIE was higher than fastGWA but well behind fastGWA-PGS-LDpred2 and BOLT-LMM-664.
 114 REGENIE showed no improvement when the LOCO PGS was added as a fixed effect, suggesting that it accounts adequately
 115 for the genetic background effect.

116 We calculated receiver operator characteristic (ROC) curves to investigate whether the increased number of causal variants
 117 recovered when we included the LOCO PGS as a fixed effect reflected a reduction in P-values across the board for the
 118 phenotype-associated variants or also an improvement in the ordering of the variants, when the variants are ordered by the
 119 evidence of an association with the phenotype. Over 100 simulations we found that the area under the ROC curve (AUC) was
 120 always higher for fastGWA-PGS-LDpred2 than for fastGWA without the LOCO PGS fixed effect (Fig. 2). This was also the

Table 1. Pipeline computation time and memory for simulations consisting of 100,000 individuals and 664,393 variants. Analyses were performed on a single compute node with 32 Xeon(R) CPU D-1541 CPUs and 128GB of RAM. Note that REGENIE was omitted from the table, as the simulation is based on a single phenotype and would unfairly disadvantage REGENIE, which is optimized for the task of performing association analyses on multiple phenotypes simultaneously.

Method	CPU Time (s)			Total (CPU Time)	Max Memory (GB)
	GWAS	LOCO PGS	GWAS(22 chr)		
fastGWA	501.2	0.0	0.0	501.2	0.5
fastGWA-PGS-PT	501.2	245.8	2,953.3	3,700.2	0.7
fastGWA-PGS-LDpred2	501.2	58,880.0	2,953.3	62,334.5	6.3
BOLT-LMM-165	92,108.0	0.0	0.0	92,108.0	3.9
BOLT-LMM-165-PGS-PT	92,108.0	245.8	614,514.4	706,868.2	3.9
BOLT-LMM-664	119,202.0	0.0	0.0	119,202.0	15.5

121 case for 99 of the 100 simulations when we added the PGS fixed effect to BOLT-LMM-165. The difference in sensitivity as a
 122 function of specificity (Table S6) showed that the sensitivity was consistently higher at a given specificity when the LOCO
 123 PGS-LDpred2 was included as a fixed effect, indicating an improvement in the ordering of the SNPs. The increase in mean
 124 sensitivity was up to 0.073 in the case of fastGWA-PGS-LDpred2 vs fastGWA, corresponding to a relative increase of 11.6% (at
 125 a specificity of 0.9988) over fastGWA. The addition of the LOCO PGS fixed effect led to a smaller but still consistent increase
 126 in sensitivity for BOLT-LMM-165. In this case, the greatest increase in the mean sensitivity was 0.028, corresponding to a
 127 4.2% relative increase in sensitivity (at a specificity of 0.9991)

128 In addition to increasing the statistical power to detect causal variants, including the PGS fixed effect also resulted
 129 in an improvement in effect size estimates (Fig. S4). We found that when a fixed effect PGS was incorporated into the
 130 association study the median squared error (MEDSE) of the effect size estimate was substantially reduced (Fig. S4, Table S7-9).
 131 Interestingly, the MEDSE of the effect size estimate was largest across all methods for BOLT-LMM with the reduced GRM
 132 (Fig. S4).

133 Effects of trait heritability, number of causal variants and sample size

134 We simulated data over a range of values of sample size, h^2 and of the number of causal SNPs to investigate how these
 135 parameters affect the impact of including the LOCO PGS as a fixed effect on GWAS power. For this analysis we used the
 136 P&T method to calculate the LOCO PGS, due to its lower computational cost (Table 1). For the larger sample size, a small
 137 improvement in power was obtained even for the lowest values of h^2 (0.1) simulated, with a statistically significant improvement
 138 for $h^2 \geq 0.2$ (Fig. 3). The improvement was not statistically significant at this value of h^2 when only 100,000 samples were
 139 used in the simulation, but even in this case the number of causal variants recovered was always at least as large and typically
 140 larger when the PGS fixed effect was included in the model (Tables S10, S11). This was somewhat surprising, given that it is
 141 assumed that large sample sizes are required for accurate phenotype prediction from PGS³¹.

142 The improvement in power resulting from the inclusion of the PGS fixed effect increased consistently with increasing
 143 numbers of causal variants in the case of the larger sample size. This was not the case for the smaller sample size, for which the
 144 improvement decreased or was lost altogether when the number of causal variants was large (Fig. 3). This is likely due to the

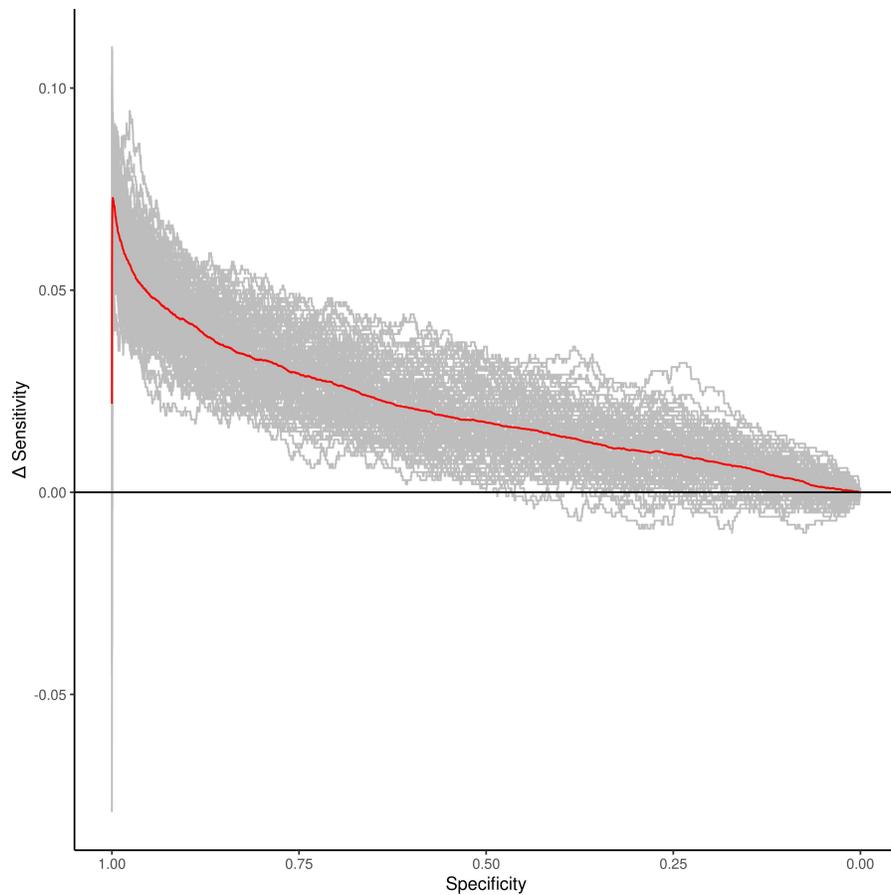


Figure 2. Difference in sensitivity (between fastGWA-PGS-LDpred2 and fastGWA) as a function of specificity for 100 simulations of a continuous trait with narrow-sense heritability of 0.5 and 1,000 causal variants in 100,000 individuals. The specificity (x-axis) is discretized in bins of size 0.0001. Each grey line shows the results of one simulation. The red line shows the mean difference over all simulations.

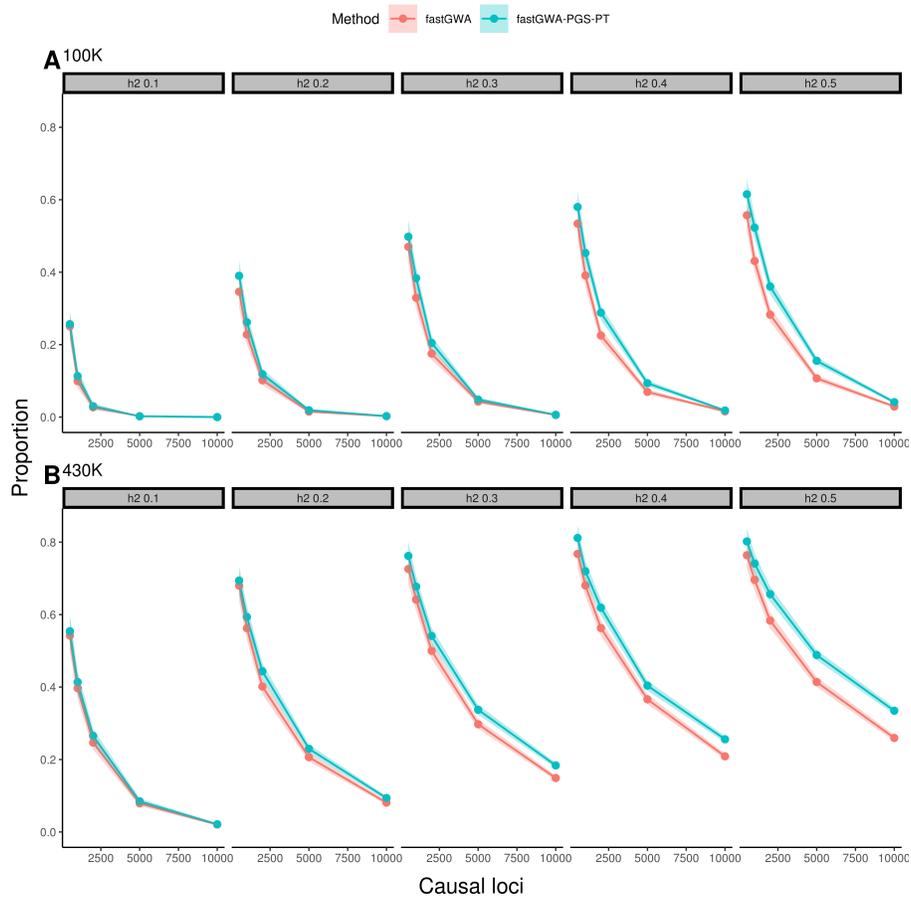


Figure 3. Proportion of causal variants recovered in simulations of a quantitative trait over a range of values of h^2 and the number of causal loci. Simulations on the top (A) and bottom (B) panels were based on 100,000 and 430,000 randomly sampled individuals from the UK Biobank, respectively.

145 loss of power to detect true causal variants and to estimate their effect sizes accurately when the genetic effect is distributed
146 over too large a number of causal variants, resulting in the inability to correct for the genetic background using the PGS.
147 This suggests that larger sample sizes would be required for highly polygenic traits in order to obtain a benefit from using
148 the LOCO PGS fixed effect. However, the larger sample size simulated is comparable in scale to the UK Biobank and with
149 a sample of this size our simulations suggest that a significant improvement in power can be obtained, even for a trait with
150 10,000 independent causal loci. For the case control simulation (N=100,000), a more modest increase in power was observed as
151 heritability increased, whereas the power to recover smaller effect loci decreased dramatically compared to the quantitative
152 simulation. However, we found that for all except three simulations the inclusion of a fixed effect LOCO PGS improved the
153 power to detect associated loci (Fig. S5, Table S12).

154 Application to UK Biobank phenotypes

155 We assessed the impact of including the LOCO PGS fixed effect on the performance of fastGWA on real data using standing
156 height, BMI, and heel bone mineral density (HBMD) in individuals of British ancestry ($N_{height}=395,133$, $N_{BMI}=395,149$ &
157 $N_{HBMD}=229,191$) from the UK Biobank. The distribution of P-values obtained from fastGWA with the LOCO PGS fixed effect
158 was lower than that obtained using fastGWA, regardless of the method used to calculate the PGS (Fig. S6-8). At a genome-wide
159 significance level of 5×10^{-8} inclusion of a LOCO PGS always increased the number of independent loci recovered, compared
160 to fastGWA (Table 2). We also applied BOLT-LMM to the real data. In this case we used all 556,516 lightly pruned HAPMAP3
161 variants for the GRM (see Methods for details). Across height, HBMD, and BMI, BOLT-LMM identified the largest number
162 of independent associated loci. Including the PGS fixed effect resulted in substantial increases in the number of independent
163 associated loci, compared to fastGWA alone for all phenotypes (Table 2, Table S13).

Table 2. Number of independent significant loci identified and resulting phenotype prediction model fit. R^2 Full is the coefficient of determination of a model that includes the PGS, sex, age & 10 PCs as covariates while R^2 PGS is the coefficient for a model that includes only the PGS. BOLT-LMM was applied with a GRM consisting of 556,516 variants.

Method	Significant loci	R^2 full	95% CI	R^2 PGS	95% CI	Spearman's ρ	Phenotype
fastGWA	1,381	0.696	0.689, 0.702	0.165	0.158, 0.170	0.382	Height
fastGWA-PGS-PT	1,583	0.701	0.694, 0.707	0.173	0.166, 0.179	0.391	
fastGWA-PGS-LDpred2	1,717	0.703	0.696, 0.709	0.176	0.170, 0.182	0.395	
BOLT-LMM	1,804	0.703	0.697, 0.709	0.170	0.164, 0.176	0.388	
fastGWA	450	0.151	0.146, 0.158	0.130	0.124, 0.135	0.351	BMI
fastGWA-PGS-PT	493	0.153	0.147, 0.159	0.130	0.125, 0.136	0.351	
fastGWA-PGS-LDpred2	500	0.151	0.146, 0.157	0.127	0.121, 0.133	0.346	
BOLT-LMM	583	0.155	0.150, 0.162	0.134	0.128, 0.139	0.356	
fastGWA	324	0.216	0.204, 0.232	0.158	0.144, 0.171	0.427	HBMD
fastGWA-PGS-PT	365	0.221	0.208, 0.238	0.164	0.152, 0.178	0.439	
fastGWA-PGS-LDpred2	385	0.225	0.210, 0.241	0.167	0.154, 0.182	0.444	
BOLT-LMM	393	0.223	0.209, 0.238	0.165	0.152, 0.178	0.437	

164 One way to determine objectively whether fastGWA with a LOCO PGS fixed effect outperforms fastGWA on real data
165 is to apply the methods on the key task of phenotype prediction. We used summary statistics from the 3 analyses above to
166 calculate PGS scores using LDpred2 and P&T (see Methods for details of how the independent training and test datasets were

167 determined). For two of the three phenotypes (height and HBMD), the PGS fixed effect resulted in an increase in the correlation
168 between the PGS and the phenotype in the test data (Table 2). In both cases the highest correlation with the phenotype was
169 obtained using fastGWA-PGS-LDpred2, which out-performed BOLT-LMM on this task. For the remaining phenotype (BMI),
170 the addition of the PGS fixed effect resulted in no change or a slightly worse correlation with the phenotype in the test data.
171 In this case the highest performance was obtained by BOLT-LMM (but at a substantial computational cost; Table 1). However,
172 even in this case, we found that including only the SNPs with low P-values in the polygenic score (as implemented by the P&T
173 method) resulted in an improvement over fastGWA (Fig. S9).

174 Discussion

175 Omitting covariates that are associated with a response and independent of an effect of interest can result in a reduction in the
176 efficiency of the estimation of the effect of interest^{15,16}. Complex traits are associated with the genotype of many loci across the
177 genome, but the effects of genetic variants other than the variant being tested are often not fully modelled by GWAS methods.
178 We evaluated a simple two-stage approach to accounting for this genetic background effect that consists of performing an
179 initial GWAS and using the summary statistics to calculate a polygenic score and then including the polygenic score, derived
180 from SNPs not on the same chromosome as the target SNP, as a fixed effect in a second round of association testing. Using
181 simulated data, we found that this led to a substantial improvement in power of fastGWA, an efficient tool for biobank scale
182 GWAS that does not fully control for genetic background effects. When we included the LOCO polygenic score as a fixed
183 effect with fastGWA (which we refer to as fastGWA-PGS), the power exceeded that of REGENIE²⁸, a recent, computationally
184 efficient tool for GWAS that uses ridge regression to control for genetic background effects. When BOLT-LMM⁸ was used
185 with a GRM derived from all of the simulated variants, the LOCO PGS fixed effect did not provide any boost in power (Fig.
186 S3); however, the equivalent (or slightly improved) performance of fastGWA-PGS-LDpred2 (Fig. 1) was achieved at a much
187 lower computational cost (Table 1). Furthermore, we note, that our simulations were favourable to BOLT-LMM because the
188 LOCO PGS was calculated from the same set of variants that were used in the GRM of BOLT-LMM. In practice, in the case of
189 P&T millions of variants can be included in the LOCO PGS calculations, but the number of model variants, M , that can be
190 included in the GRM of BOLT-LMM is constrained by memory and compute time, both of which scale at least linearly with M .
191 A further key advantage of the approach that we propose is that it is modular. Any phenotype prediction method can be used to
192 predict the combined effect of the LOCO genetic variants on the phenotype. As methods for phenotype prediction improve, we
193 anticipate that the performance of this approach will increase.

194 The increase in power using the PGS fixed effect was largest for simulated phenotypes with high heritability and a large
195 number of causal variants (Fig. 3). In these cases the many background SNPs collectively explain a substantial proportion
196 of the phenotypic variance and summarizing the contribution of these background SNPs to the phenotype via the LOCO
197 PGS is likely to result in a better estimate of the effect of the target SNP and its standard error. The boost in performance
198 derived from including the LOCO PGS as a fixed effect also depended on study sizes. For example, when the number of

199 causal variants became large (10,000) there was no substantial boost in performance in the simulation that included 100,000
200 individuals, presumably because in this case the study size was not sufficient to identify and accurately estimate the effects
201 of the causal variants. Even with this large number of causal variants the larger simulation (with 430,000 individuals) still
202 showed a significant improvement arising from the LOCO PGS fixed effect (Fig. 3). Across all the simulation parameters
203 we investigated, the performance of fastGWA with a LOCO PGS fixed effect (calculated using the P&T method) was never
204 worse than fastGWA without the fixed effect included. We also note that we calculated the P&T LOCO PGS using SNPs that
205 were selected based on a fixed P-value threshold. Further increases in power may be possible by optimizing the SNPs that are
206 used to calculate the PGS separately for each omitted chromosome, though care should be taken to avoid very high values
207 of the threshold, which may result in an elevation in the false positive rate (SFig. 1). Thresholding based on a P-value was
208 not required for LDpred2, which may help to explain why we achieved significantly better power when the LOCO PGS was
209 calculated using this method rather than pruning and thresholding (Fig. 1).

210 We also applied the method to real data (standing height, heel bone mineral density (HBMD) and body mass index (BMI)
211 in individuals of British ancestry in the UK Biobank). Consistent with the simulation results, we found more independent
212 trait-associated loci using fastGWA-PGS-LDpred2 than with fastGWA alone for all three traits (30%, 19%, & 11% more
213 for height, HBMD, and BMI, respectively; Table 2). Although, BOLT-LMM recovered the largest number of independent
214 significant loci across all UK Biobank traits, this did not always translate into better correlation between a PGS calculated
215 from the resulting summary statistics and the phenotype in the test dataset. In fact, the highest correlation was obtained by
216 fastGWA-PGS-LDpred2 for two of the three traits. This could be explained by a higher proportion of true positives among the
217 loci detected using the PGS-based methods or a more accurate estimate of the effects sizes by these methods, as suggested by
218 Fig. S4. For BMI, the correlation was in fact lower between the PGS and the phenotype in the test dataset when the LOCO
219 PGS fixed effect was used (Table2). However, even in this case a larger number of significant variants were recovered than with
220 fastGWA.

221 The use of polygenic scores for phenotype prediction from genotype is an increasingly important application of the
222 results of GWAS³². High polygenic scores can capture a substantial component of the risk of complex diseases^{24,33} and guide
223 interventions that can confer health benefits to individuals and reduce the stress on health systems³⁴. Performing GWAS on a
224 subset of samples and predicting on the remainder, we observed an increase in the correlation of the PGS with the phenotype
225 when we included the LOCO PGS as a fixed effect in two out of three traits considered, consistent with improved effect size
226 estimates (Fig. S4). Our results suggest that a modular approach that integrates advances in phenotype prediction with efficient
227 GWAS methods can have a significant impact on the power of GWAS and that this can, in turn, lead to more accurate phenotype
228 prediction. A recent study showed that models that allow unequal a priori contribution of SNPs to trait heritability can lead to
229 substantial improvements in the accuracy of trait³⁵. Although not explored in this work, the incorporation of external PGS
230 instruments from large meta-analyses in the first round of GWAS may also provide an additional gain in performance, similar
231 to what is proposed in Bulik-Sullivan³⁶. Indeed, our results show that GWAS summary statistics can be used to account for

232 genetic background effects, with results matching the performance of methods such as BOLT-LMM that require individual-level
233 data for this purpose. As new efficient methods emerge from these and further insights, they can be easily substituted for the
234 calculation of the LOCO PGS fixed effect. The current fast pace of methodological innovation in phenotype prediction supports
235 the use, at least for the time being, of the simple modular approach to modeling genetic background effects evaluated here.

236 **Conclusion**

237 The tasks of detecting trait-associated variants and predicting the trait in a new sample from the summary statistics of these
238 variants are closely intertwined. Improved performance on the trait-association task can result in more associated variants and
239 better estimates of their effect sizes, resulting in improvement on the prediction task. On the other hand, improved methods for
240 phenotype prediction can help to control for background genetic effects in methods that identify the trait-associated variants
241 and their effects. The method that we have explored here consists of incorporating a LOCO PGS as a fixed-effect covariate
242 to control for these background genetic effects; however, any method for phenotype prediction could play this role, once
243 its application is restricted to variants that are not linked to the target SNP. We show here that incorporating the PGS as a
244 fixed-effect covariate results in increased power to detect trait-associated variants in GWAS. The resulting trait-associated
245 variants and effect size estimates can lead to an improvement in the PGS, as illustrated by improved performance in the task of
246 predicting the phenotype in a test dataset.

247 **Methods**

248 **Simulations**

249 ***Genotype QC***

250 The use of the UK Biobank Materials falls within UK Biobank's generic Research Tissue Bank (RTB) approval from the NHS
251 North West Research Ethics Committee, UK. The simulated genotype data was based on autosomal genotyped data from
252 the UK Biobank. To limit the effects of population stratification only individuals reporting white British ancestry (data field
253 21000; code 1001; N=443,076) were included in these analyses. The genotype data for the simulation analysis was based
254 on directly genotyped variants with minor allele frequency (MAF) greater than 0.05%. Variants with genotype missingness
255 greater than 2% or that failed a test for Hardy-Weinberg equilibrium (HWE) at $\alpha=0.0001$ were excluded, resulting in a total of
256 664,393 genetic variants. There were 429,359 samples remaining following filtering. The sparse GRM required by fastGWA
257 was created by setting entries corresponding to sample pairs with an estimated relatedness of less than 0.05 to 0. To account for
258 population structure in the association studies, principal component analysis (PCA) was performed on a set of 165,684 variants
259 LD-pruned with an R^2 greater than 0.1 in a sliding window of size 500bp, sliding by 200bp. This set was also used as the basis
260 of the BOLT-LMM analyses with the reduced GRM size (referred to as BOLT-LMM-165 in Results). All genotype QC was
261 implemented in plink2³⁷.

262 Based on the above genotype data, we simulated a continuous phenotype using the GCTA software suite³⁸. The initial

simulation (Fig. 1) consisted of 100,000 individuals, 1,000 randomly sampled causal variants and $h^2 = 0.5$. This simulation was repeated 100 times with the 664,393 variants remaining after variant filtering for the GRM calculation. Power was calculated as the proportion of the causal variants recovered. Further simulations were carried out to investigate the effects of varying the number of causal SNPs (500, 1000, 2000, 5000 & 10,000), h^2 (0.1, 0.2, 0.3 0.4, 0.5) and the sample size (100,000 & 430,000) on method performance. In each case all parameters other than the ones being varied were the same as the initial simulation, and one simulation was performed per set of parameter values. The pROC R package was used to generate receiver operating characteristic (ROC) curves, variants within 1 Mb of the causal variants were removed.³⁹ We applied the same simulation strategy to binary traits with two levels of disease prevalence, 0.1 & 0.3, using 1,000 causal loci with $h^2 = 0.5$, and 100,000 samples. To calculate the false positive rate we performed 100 simulations with 100,000 samples, $h^2 = 0.5$ and 1000 causal variants restricted to the even chromosomes.

Simulation association tests

Association testing was performed using fastGWA, REGENIE and BOLT-LMM. To account for known sources of covariation (technical batch effects, population structure, biological effects) 10 PCs, sex, age, genotype batch and assessment centre were included as fixed-effect covariates in statistical models. For the PGS method we first performed GWAS (using fastGWA, REGENIE or BOLT-LMM) and calculated PGS scores on a Leave One Chromosome Out (LOCO) basis. This resulted in 22 sets of PGS values (one for each autosomal chromosome, calculated from the summary statistics of variants on all other autosomal chromosomes). Two PGS strategies were used in this study, pruning and thresholding (P+T), denoted with the suffix PGS-PT and LDpred2, denoted by the suffix PGS-LDpred2. The LOCO PGS-PT were calculated using PRSice2 (version 2.2.12 (2020-02-20))²⁹. To decrease computation time and reduce the likelihood of over-fitting a P-value threshold of 5×10^{-5} was chosen, a priori, for the LOCO PGS-PT calculation. Association testing was then performed using fastGWA in a chromosome-wise manner, with the corresponding LOCO PGS included as a fixed effect. The bigsnpr R package (bigsnpr v1.6.1 & R v3.6.1) was used to calculate the LOCO PGS-LDpred2 fixed effects³⁰. To reduce computation time, 22 LOCO genotype objects containing the SNP correlations were precomputed.

Application to the UK Biobank

UK Biobank association tests

The genotype selection, quality control and genetic relationship matrix were performed following the QC procedure in Jiang *et al.*⁹. The genetic relationship matrix used with fastGWA and BOLT-LMM was calculated for all European individuals (N=458,686), using a set of 556,516 lightly pruned HAPMAP3 variants (R^2 greater than 0.9 in a 100 variant sliding window of size 1,000 & MAF > 0.01)⁹. Association summary statistics were generated from a set of 1.1 million HAPMAP3 variants (MAF > 0.01, HWE $\alpha = 1 \times 10^{-6}$ and missingness < 0.05)⁹. Principal components were calculated using a set of 34,775 variants (LD-pruned with $R^2 = 0.05$ in a sliding window of size 1,000bp, sliding by 50bp)⁴⁰. To identify white British samples with similar genetic backgrounds we clustered samples based on the first 6 principal components⁴⁰, resulting in a subset of 406,319 white-British samples. Sample pairs that had a KING kinship coefficient above 0.05, with one member of the pair within the

296 white-British group and the other in the group self-reporting as white European were removed. This left 399,135 white British
297 and 46,406 other European samples^{40,41}. To account for known sources of phenotype and genotype variation, 10 PCs, age, sex,
298 genotype batch and assessment centre were included as fixed-effect covariates for the BOLT-LMM and fastGWA analyses.
299 PRSice2 and LDpred2 were used to calculate the LOCO PGS. Independent loci were identified using the clumping algorithm in
300 plink2 (P-value threshold = 5×10^{-9} , window size = 5Mb, and LD R^2 threshold = 0.01).

301 **UK Biobank phenotype prediction**

302 To test the performance of fastGWA with a LOCO PGS fixed effect on the task of predicting standing height, BMI and HBMD,
303 the UK Biobank data was partitioned into training and test datasets. The test data consisted of white British individuals with
304 similar genetic background described above and the polygenic score predictions were tested on the remaining independent
305 European samples. Summary statistics were generated using fastGWA, fastGWA-PGS-PT, fastGWA-PGS-LDpred2 and
306 BOLT-LMM. We used LDpred2 and PRSice2 to predict the phenotypic values in the test set. LDpred2 requires LD correlation
307 data and we used a pre-computed set built on the 1.1 million HAPMAP3 variants for this purpose. The model fit was assessed
308 for each method by fitting a linear model to the values of the phenotype in the test set as a function of their predicted values,
309 accounting for known sources of phenotypic variation, i.e sex, age, PC's. We report both the proportion of variation explained
310 collectively by the PGS, sex, age, the first 4 principal components and assessment centre as well as the R^2 using only the PGS
311 in the regression model.

312 **Data Availability**

313 All genotype and phenotype data analyzed are available, subject to application, from the UK Biobank (application 23739).
314 Code to implement fastGWA-PGS-PT and fastGWA-PGS-LDpred2 as described in this work is available under MIT license
315 from <https://github.com/declan93/PGS-LMM/>.

316 **References**

- 317 **1.** Chen, W. M. & Abecasis, G. R. Family-based association tests for genomewide association scans. *Am J Hum Genet* **81**,
318 913–926 (2007).
- 319 **2.** Svishcheva, G. R., Axenovich, T. I., Belonogova, N. M., van Duijn, C. M. & Aulchenko, Y. S. Rapid variance components-
320 based method for whole-genome association analysis. *Nat Genet* **44**, 1166–1170 (2012).
- 321 **3.** Jakobsdottir, J. & McPeck, M. S. MASTOR: mixed-model association mapping of quantitative traits in samples with
322 related individuals. *Am J Hum Genet* **92**, 652–666 (2013).
- 323 **4.** Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nature*
324 *genetics* **42**, 348 (2010).
- 325 **5.** Zhang, Z. *et al.* Mixed linear model approach adapted for genome-wide association studies. *Nat Genet* **42**, 355–360 (2010).

- 326 **6.** Lippert, C. *et al.* FaST linear mixed models for genome-wide association studies. *Nat Methods* **8**, 833–835 (2011).
- 327 **7.** Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association studies. *Nature genetics* **44**,
328 821–824 (2012).
- 329 **8.** Loh, P.-R. *et al.* Efficient bayesian mixed-model analysis increases association power in large cohorts. *Nature genetics* **47**,
330 284 (2015).
- 331 **9.** Jiang, L. *et al.* A resource-efficient tool for mixed model association analysis of large-scale data. Tech. Rep., Nature
332 Publishing Group (2019).
- 333 **10.** Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*
334 **81**, 559–575 (2007).
- 335 **11.** Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M. & Price, A. L. Advantages and pitfalls in the application of
336 mixed-model association methods. *Nature genetics* **46**, 100 (2014).
- 337 **12.** Eu-Ahsunthornwattana, J. *et al.* Comparison of methods to account for relatedness in genome-wide association studies
338 with family-based data. *PLoS Genet* **10**, e1004445 (2014).
- 339 **13.** Devlin, B. & Roeder, K. Genomic control for association studies. *Biometrics* **55**, 997–1004 (1999).
- 340 **14.** Listgarten, J. *et al.* Improved linear mixed models for genome-wide association studies. *Nature methods* **9**, 525 (2012).
- 341 **15.** Fisher, R. A. *The Design of Experiments* (Oliver and Boyd, 1935).
- 342 **16.** Neuhaus, J. M. Estimation efficiency with omitted covariates in generalized linear models. *Journal of the American*
343 *Statistical Association* **93**, 1124–1129 (1998).
- 344 **17.** Listgarten, J., Lippert, C. & Heckerman, D. Fast-lmm-select for addressing confounding from spatial structure and rare
345 variants. *Nature genetics* **45**, 470–471 (2013).
- 346 **18.** Lippert, C. *et al.* The benefits of selecting phenotype-specific variants for applications of mixed models in genomics.
347 *Scientific reports* **3**, 1815 (2013).
- 348 **19.** Tucker, G., Price, A. L. & Berger, B. Improving the power of gwas and avoiding confounding from population stratification
349 with pc-select. *Genetics* **197**, 1045–1049 (2014).
- 350 **20.** Canela-Xandri, O., Rawlik, K. & Tenesa, A. An atlas of genetic associations in uk biobank. *Nature genetics* **50**, 1593–1599
351 (2018).
- 352 **21.** Kadie, C. & Heckerman, D. Ludicrous speed linear mixed models for genome-wide association studies. *bioRxiv* (2019).
353 URL <https://www.biorxiv.org/content/early/2019/12/07/154682>. <https://www.biorxiv.org/content/early/2019/12/07/154682.full.pdf>.
354
- 355 **22.** Tam, V. *et al.* Benefits and limitations of genome-wide association studies. *Nature Reviews Genetics* **20**, 467–484 (2019).

- 356 **23.** Torkamani, A., Wineinger, N. E. & Topol, E. J. The personal and clinical utility of polygenic risk scores. *Nature Reviews*
357 *Genetics* **19**, 581 (2018).
- 358 **24.** Khera, A. V. *et al.* Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to
359 monogenic mutations. *Nature genetics* **50**, 1219–1224 (2018).
- 360 **25.** Yanes, T., Young, M.-A., Meiser, B. & James, P. A. Clinical applications of polygenic breast cancer risk: a critical review
361 and perspectives of an emerging field. *Breast Cancer Research* **22**, 1–10 (2020).
- 362 **26.** Michailidou, K. *et al.* Association analysis identifies 65 new breast cancer risk loci. *Nature* **551**, 92 (2017).
- 363 **27.** Bahcall, O. Common variation and heritability estimates for breast, ovarian and prostate cancers. *Nat Genet* **10** (2013).
- 364 **28.** Mbatchou, J. *et al.* Computationally efficient whole genome regression for quantitative and binary traits. *bioRxiv* (2020).
365 URL <https://www.biorxiv.org/content/early/2020/06/22/2020.06.19.162354>. [https://](https://www.biorxiv.org/content/early/2020/06/22/2020.06.19.162354.full.pdf)
366 www.biorxiv.org/content/early/2020/06/22/2020.06.19.162354.full.pdf.
- 367 **29.** Choi, S. W. & O'Reilly, P. F. Prsice-2: Polygenic risk score software for biobank-scale data. *Gigascience* **8**, giz082 (2019).
- 368 **30.** Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics* **36**, 5424–5431
369 (2020). URL <https://doi.org/10.1093/bioinformatics/btaa1029>. [https://academic.oup.](https://academic.oup.com/bioinformatics/article-pdf/36/22-23/5424/36855825/btaa1029.pdf)
370 [com/bioinformatics/article-pdf/36/22-23/5424/36855825/btaa1029.pdf](https://academic.oup.com/bioinformatics/article-pdf/36/22-23/5424/36855825/btaa1029.pdf).
- 371 **31.** Dudbridge, F. Power and predictive accuracy of polygenic risk scores. *PLoS genetics* **9** (2013).
- 372 **32.** Martin, A. R., Daly, M. J., Robinson, E. B., Hyman, S. E. & Neale, B. M. Predicting polygenic risk of psychiatric disorders.
373 *Biological psychiatry* **86**, 97–109 (2019).
- 374 **33.** Mars, N. *et al.* Polygenic and clinical risk scores and their impact on age at onset and prediction of cardiometabolic
375 diseases and common cancers. *Nature Medicine* 1–9 (2020).
- 376 **34.** Gibson, G. On the utilization of polygenic risk scores for therapeutic targeting. *PLoS genetics* **15** (2019).
- 377 **35.** Zhang, Q., Prive, F., Vilhjalmsón, B. J. & Speed, D. Improved genetic prediction of complex traits from individual-level
378 data or summary statistics. *bioRxiv* (2020).
- 379 **36.** Bulik-Sullivan, B. Mixed models for meta-analysis and sequencing. *bioRxiv* (2015). URL [https://www.biorxiv.](https://www.biorxiv.org/content/early/2015/05/29/020115)
380 [org/content/early/2015/05/29/020115](https://www.biorxiv.org/content/early/2015/05/29/020115). [https://www.biorxiv.org/content/early/2015/](https://www.biorxiv.org/content/early/2015/05/29/020115.full.pdf)
381 [05/29/020115.full.pdf](https://www.biorxiv.org/content/early/2015/05/29/020115.full.pdf).
- 382 **37.** Chang, C. C. *et al.* Second-generation plink: rising to the challenge of larger and richer datasets. *Gigascience* **4**, s13742–015
383 (2015).
- 384 **38.** Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. Gcta: a tool for genome-wide complex trait analysis. *The American*
385 *Journal of Human Genetics* **88**, 76–82 (2011).

386 **39.** Robin, X. *et al.* proc: an open-source package for r and s+ to analyze and compare roc curves. *BMC bioinformatics* **12**, 77
387 (2011).

388 **40.** Bycroft, C. *et al.* The uk biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).

389 **41.** Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies. *Bioinformatics* **26**, 2867–2873
390 (2010).

391 **Acknowledgements**

392 This research has been conducted using the UK Biobank Resource under application number 23739. This publication has
393 emanated from research conducted with the financial support of Science Foundation Ireland under grant number 16/IA/4612.
394 DOS was funded through Science Foundation Ireland grant number 18/CRT/6214.

395 **Author contributions statement**

396 CS conceived and supervised the project and performed analyses. DB implemented the pipeline and performed analyses. DB
397 and CS wrote the manuscript, with input from DM. DM advised on application of the method to human phenotypes. DOS
398 performed analysis. JF provided the mathematical justification of the method.

399 **Additional information**

400 **Competing interests**

401 The authors declare that they have no competing interests.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementarySciRep.pdf](#)