

Supplementary Materials

1. Algorithm details

1.1 Step 0. Constructing the sparse GRM

In the sparse GRM, denoted by ψ_s , GRM elements below a user-specified relative coefficient cutoff are zeroed out with close family structures preserved. To improve the test accuracy for rare variants, SAIGE-GENE approximates the variance of score statistics calculated with the full GRM ψ_f using the variance calculated with the sparse GRM ψ_s and the ratios of these two variance estimates estimated using a subset of genetic markers.

To construct the sparse GRM ψ_s , a small subset of randomly select markers were used to identify related sample pairs whose relative coefficient pass the use-specified cutoff, which is to find out the indices of non-zero elements in ψ_s . Next, the values of the nonzero elements in ψ_s are then estimated using the full set of genetic markers that are used in Step 1 for ψ_f . This step is only needed for once for each data set or biobank and parallel computation is allowed. Once the sparse GRM is constructed for a data set, it can be re-used in SAIGE-GENE for all phenotypes.

1.2 Step 1. Fitting the null generalized linear mixed model

The same model fitting framework and cutting-edge computation approaches used in the original SAIGE¹ are used in SAIGE-GENE to fit the null GLMM for large sample sizes. These include estimating model parameters using the AI-REML approach², solving linear systems by the preconditioned conjugate gradient method³, using Hutchinson's randomized trace estimator^{4,5} to obtain traces of matrices, and allowing for parallel computation for the vector multiplication. For details of the likelihood, parameter estimates and information matrices, please refer to the Supplementary Note in the SAIGE paper¹. In addition, SAIGE-GENE can estimate variance component parameters, thus heritability, by fitting a null GLMM using the sparse GRM. The estimated variance component parameters can then be used as initial values for the model fitting using the full GRM. By plotting the heritability estimates using the sparse GRM versus using the full GRM for 24 quantitative traits with sample size larger than or equal to 10,000 from the UK Biobank (**Supplementary Figure 3**), we have shown that for most phenotypes, the sparse GRM can provide similar variance component estimates, thus heritability, as the full GRM does. Therefore, comparing to arbitrary initial values, values estimated using the sparse GRM tend to be closer to the true parameter values. For real-data analysis, robust performance of convergence has also been observed for some phenotypes, such as waist hip ratio (N = 408,144) in the UKB. Using initial values 0.5 for heritability, step 1 did not even converge after 6,300 CPU hours, while using initial values estimated with the sparse GRM, it took 1836 CPU hours to finish the step 1.

1.3 Step 2. Gene-based association tests

Test statistics of the Burden, SKAT and SKAT-O tests for a gene can be constructed based on the score test statistics from the marginal model for individual variants in the gene. Suppose there are q variants in the region or gene to test. The score test statistics for variant j ($j=1, \dots, q$) under $H_0: \beta_j = 0$ is $T_j = g_j^T(Y - \hat{\mu})$

where g_j and Y are $N \times 1$ genotype and phenotype vectors, respectively, and $\hat{\mu}$ is the estimated mean of Y under the null hypothesis.

Let u_j denote a threshold indicator or weight for variant j and U be a diagonal matrix with u_j as the j th element. The Burden test statistics can be written as $Q_{Burden} = \left(\sum_{j=1}^q u_j T_j \right)^2$. Suppose $\tilde{G} = G - X(X^T \hat{W} X)^{-1} X^T \hat{W} G$, where $G = (g_1, \dots, g_q)$ is the $N \times q$ genotype matrix of the q genetic variants, and $\hat{P} = \hat{\Sigma}^{-1} - \hat{\Sigma}^{-1} X (X^T \hat{\Sigma}^{-1} X)^{-1} X^T \hat{\Sigma}^{-1}$ with $\hat{\Sigma} = \hat{W}^{-1} + \hat{\tau} \psi$. Under the null hypothesis of no genetic effects, Q_{Burden} followed $\lambda_B \chi_1^2$, where $\lambda_B = J^T U \tilde{G}^T \hat{P} \tilde{G} U J$ and J is a $q \times 1$ vector with all elements being unity and χ_1^2 is a chi-squared distribution with 1 degree of freedom⁶. The SKAT test⁷ can be written as $Q_{SKAT} = \sum_{j=1}^q u_j^2 T_j^2$, which follows a mixture of chi-square distribution $\sum_{j=1}^q \lambda_{Sj} \chi_1^2$, where λ_{Sj} are the eigenvalues of $U \tilde{G}^T \hat{P} \tilde{G} U$. The SKAT-O test developed by Lee et al in 2012⁸ uses a linear combination of the Burden and SKAT tests statistics $Q_{SKATO} = (1 - \rho) Q_{SKAT} + \rho Q_{Burden}$, $0 \leq \rho \leq 1$. To conduct the test, the minimum p-value from grid of ρ is calculated and the p-value of the minimum p-value is estimated through numerical integration. Following the suggestion in Lee et al⁹, we use a grid of eight values of $\rho = (0, 0.1^2, 0.2^2, 0.3^2, 0.4^2, 0.5^2, 0.5, 1)$ to find the minimum p-value.

1.3.1 Estimating $\tilde{G}^T \hat{P} \tilde{G}$

For each gene, given \hat{P} , calculation of $\tilde{G}^T \hat{P} \tilde{G}$ can be computationally expensive. Suppose $\tilde{g} = g - X(X^T \hat{W} X)^{-1} X^T \hat{W} g$, which represents a covariate adjusted single variant genotype $N \times 1$ vector. To reduce computation cost, an approximation approach has been used in SAIGE¹, BOLT-LMM¹⁰ and GRAMMAR-GAMMAR¹¹, in which the ratio between $\tilde{g}^T \hat{P} \tilde{g}$ and $\tilde{g}^T \tilde{g}$ is estimated by a small subset of randomly selected genetic markers that has been shown to be approximately constant for all variants¹. Given the ratio $\hat{r} = \tilde{g}^T \hat{P} \tilde{g} / \tilde{g}^T \tilde{g}$, $\tilde{g}^T \hat{P} \tilde{g}$ for all other variants can be easily obtained as $\hat{r} \tilde{g}^T \tilde{g}$. However, the variations of estimated \hat{r} for extremely rare variants are large and including some closely related samples in the denominator helps reduce the variation of \hat{r} as shown in **Supplementary Figure 2**. It can also be observed from the plots in **Supplementary Figure 2** that the variance ratio for those extremely rare variants could be quite different from the ratio for more frequent variants, SAIGE-GENE estimates variance ratios for different MAC categories. By default, MAC categories are set to be MAC equals to 1, 2, 3, 4, 5, 6 to 10, 11 to 20, and is greater than 20. For each MAC category, a ratio \hat{r}_s is estimated as the average of the ratios computed from 30 randomly selected markers, among which every marker has a ratio $\tilde{g}^T \hat{P} \tilde{g} / \tilde{g}^T \hat{P}_s \tilde{g}$, where $\hat{P}_s = \hat{\Sigma}_s^{-1} - \hat{\Sigma}_s^{-1} X (X^T \hat{\Sigma}_s^{-1} X)^{-1} X^T \hat{\Sigma}_s^{-1}$ and $\hat{\Sigma}_s = \hat{W}^{-1} + \tau \psi_s$. ψ_s is a sparse GRM that preserves closely related samples. The coefficient of variance (CV) of \hat{r} is used to evaluate the numerical stability of the \hat{r} estimation. As in SAIGE, the default value of CV threshold is 0.001. If CV of \hat{r} is larger than the threshold, SAIGE-GENE will increase the number of markers by 10 to estimate \hat{r} until the estimation is stable with CV below or equal to the threshold. Once the variance ratios have been estimated for different MAC categories. For each genetic marker in genes or regions that are to be tested in Step 2, a \hat{r}_s can be obtained according to its MAC. Let \hat{R}_s be a $q \times 1$ vector whose j th element is the ratio \hat{r}_s for the j th marker in the tested gene. For the tested gene with q markers, $\tilde{G}^T \hat{P} \tilde{G}$ can be approximated as $\hat{R}_s^{\frac{1}{2}} \tilde{G}^T \hat{P}_s \tilde{G} \hat{R}_s^{\frac{1}{2}}$.

Products of $\hat{\Sigma}_s^{-1}$ and other vectors or matrices are obtained using the sparse LU decomposition through the solve function in R¹². Note that the $N \times N$ matrix \hat{P}_s is not a sparse matrix because $\hat{\Sigma}_s^{-1}$ is not sparse

and \tilde{G} is also a dense matrix, which is converted from the sparse matrix G , as $\tilde{G} = G - X(X^T \hat{W} X)^{-1} X^T \hat{W} G$. It can be shown that

$$\begin{aligned} \tilde{G}^T \hat{P}_s \tilde{G} &= (G^T - G^T \hat{W} X (X^T \hat{W} X)^{-1} X^T) \left(\hat{\Sigma}_s^{-1} - \hat{\Sigma}_s^{-1} X (X^T \hat{\Sigma}_s^{-1} X)^{-1} X^T \hat{\Sigma}_s^{-1} \right) (G \\ &\quad - X (X^T \hat{W} X)^{-1} X^T \hat{W} G) = G^T \hat{P}_s G \end{aligned}$$

Computation of $G^T \hat{P}_s G$ is more computational efficient than that of $\tilde{G}^T \hat{P}_s \tilde{G}$ as G can be stored as a sparse matrix.

1.3.2 Conditional analysis

To test whether the association signals from a tested gene or region are independent from a given marker or multiple markers, the conditional analysis based on summary statistics from unconditional association tests with the linkage disequilibrium r^2 among testing and conditioning markers¹³ have been implemented in SAIGE-GENE.

Let G be the genotypes for a gene to be tested for association, which contains q markers, and G_2 be the genotypes for the conditioning markers, which contains q_2 markers. Let β denote a $q \times 1$ coefficient vector of the genetic effect for the gene to be tested and β_2 be a $q_2 \times 1$ coefficient vector of the genetic effect for the conditioning markers. The genotype matrix with the non-genetic covariates projected out $\tilde{G} = G - X(X^T \hat{W} X)^{-1} X^T \hat{W} G$ and $\tilde{G}_2 = G_2 - X(X^T \hat{W} X)^{-1} X^T \hat{W} G_2$. In the unconditioning association tests, the test statistics $T = \tilde{G}^T (Y - \hat{\mu})$ and $T_2 = \tilde{G}_2^T (Y - \hat{\mu})$. In conditional analysis, under the null hypothesis, $E(T) = E(\tilde{G}^T P(\tilde{G}_2 \beta_2)) = \tilde{G}^T \hat{P} \tilde{G}_2 \beta_2$ and $E(T_2) = E(\tilde{G}_2^T P(\tilde{G}_2 \beta_2)) = \tilde{G}_2^T \hat{P}_s \tilde{G}_2 \beta_2$. T and T_2 jointly follow the multivariate normal with mean $(E(T), E(T_2))$ and variance $S = \begin{bmatrix} \tilde{G}^T \hat{P} \tilde{G} & \tilde{G}^T \hat{P} \tilde{G}_2 \\ \tilde{G}_2^T \hat{P} \tilde{G} & \tilde{G}_2^T \hat{P}_s \tilde{G}_2 \end{bmatrix}$.

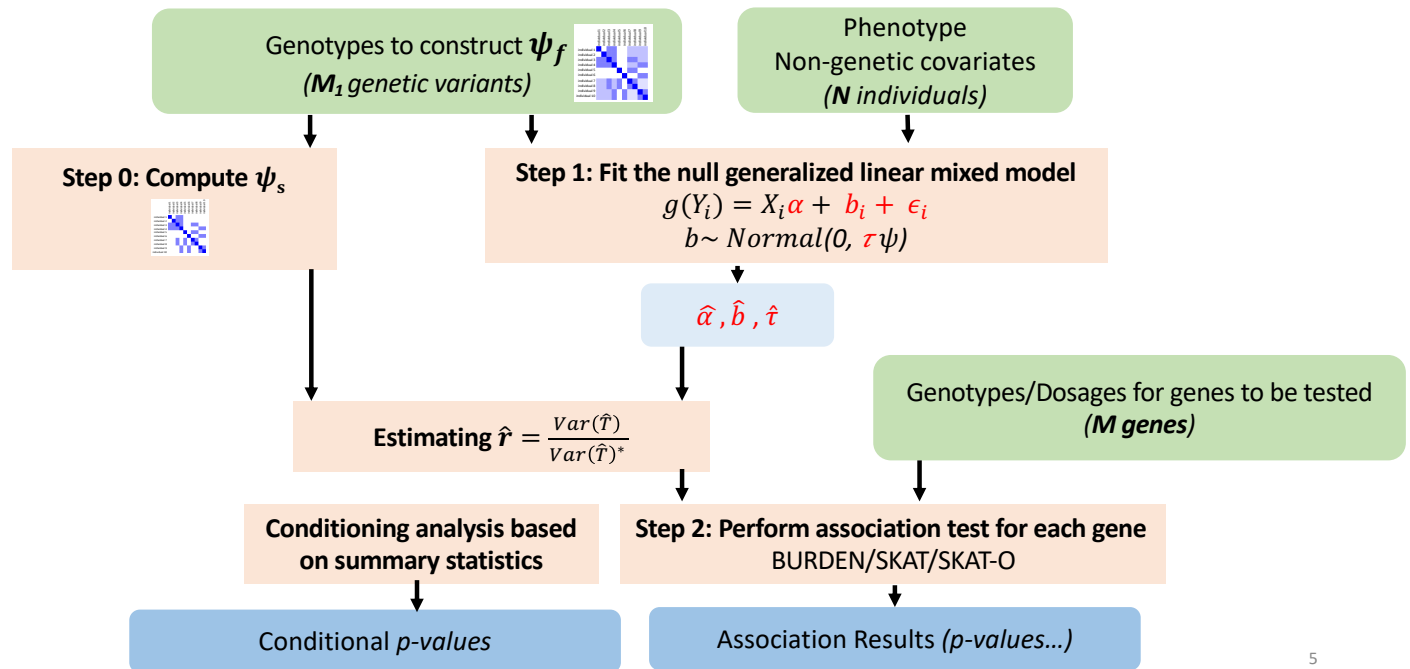
Thus under the null hypothesis of $\beta = 0$, the $T|T_2$ follows the conditional distribution $E(T|T_2) = \tilde{G}^T \hat{P} \tilde{G}_2 (\tilde{G}_2^T \hat{P} \tilde{G}_2)^{-1} T_2$ and $\text{var}(T|T_2) = \tilde{G}^T \hat{P} \tilde{G} - \tilde{G}^T \hat{P} \tilde{G}_2 (\tilde{G}_2^T \hat{P} \tilde{G}_2)^{-1} \tilde{G}_2^T \hat{P} \tilde{G} = G^T \hat{P} G - G^T \hat{P} G_2 (G_2^T \hat{P} G_2)^{-1} G_2^T \hat{P} G$. The test statistic of the conditional analysis can be written as $(T - E(T|T_2))^2 / \text{var}(T|T_2)$, which follows the χ^2 distribution with one degree of freedom. Similar to the unconditional analysis, we approximate these terms by $\tilde{G}^T \hat{P}_s \tilde{G}$, $\tilde{G}^T \hat{P}_s \tilde{G}_2$, $\tilde{G}_2^T \hat{P}_s \tilde{G}_2$, and $\tilde{G}_2^T \hat{P}_s \tilde{G}$ and the corresponding variance ratio matrices.

1.3.3 Adjust for test statistics variance using the genome control inflation factor from single-variant association tests.

To further control for type I error rates, SAIGE-GENE allows for using the genome control inflation factor from single-variant association tests to adjust for gene- or region- based tests. Let λ_{GC} be the genome control inflation factor from the single-variant association tests, the $\tilde{G}^T \hat{P} \tilde{G}$ for each tested gene is multiplied by λ_{GC} . As **Supplementary Table 3** shows, this approach has successfully attenuated the type I error inflation

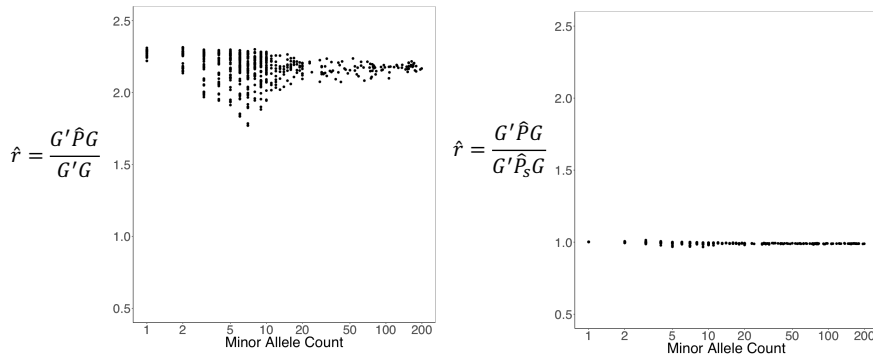
2. Supplementary figures

Supplementary Figure 1. Workflow of SAIGE-GENE.

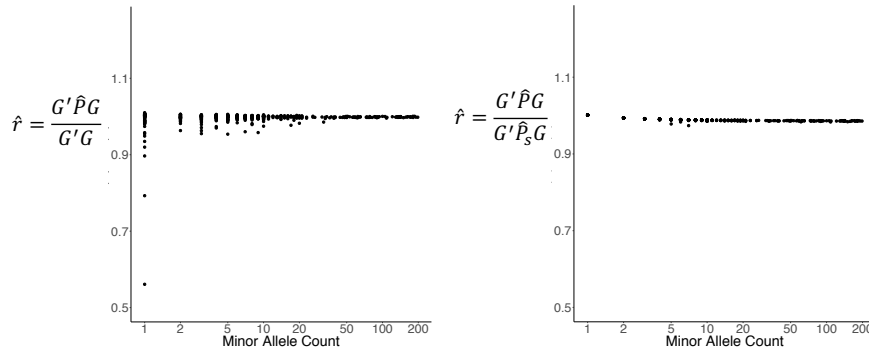


Supplementary Figure 2. Plots of the variance ratio of the score statistics by MAC for rare variants with and without the full GRM for sample relatedness (left) and with the full GRM and a sparse GRM for closely related samples(right). A. 500 families and 5,000 independent individuals were simulated with $h^2= 0.2$ based on the pedigree structure shown in **Supplementary Figure 9**. The sparse GRM was constructed using a coefficient of relatedness cutoff 0.2. B. 20,000 samples with White British ancestry were randomly selected from the UK Biobank and the null model was fitted for the automated read pulse rate. The sparse GRM was constructed using a coefficient of relatedness cutoff 0.125. C. 20,000 samples were randomly selected form the HUNT study and the null model was fitted for HDL. The sparse GRM was constructed using a coefficient of relatedness cutoff 0.125.

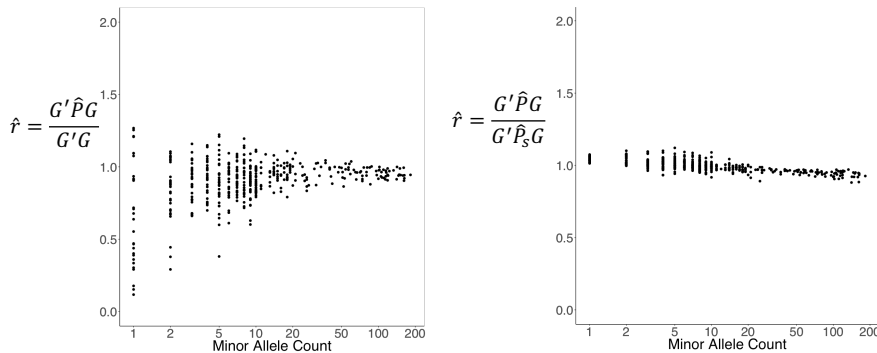
A. Simulation: 500 families and 5,000 independent individuals



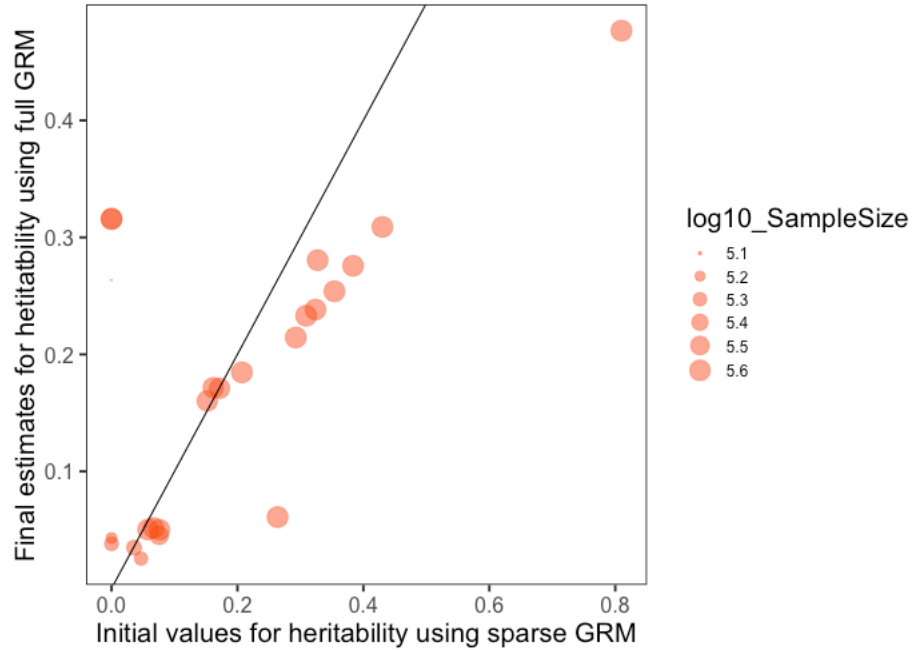
B. UK Biobank: Pulse rate automated read (mean), randomly selected 20,000 individuals



C. HUNT: HDL, randomly selected 20,000 individuals

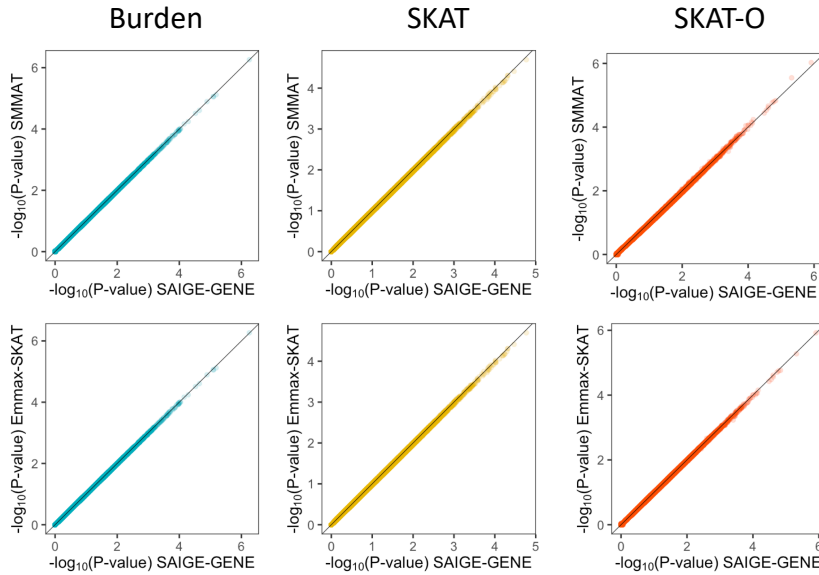


Supplementary Figure 3. Comparing heritability estimates using the sparse GRM to heritability estimates using the full GRM for 24 quantitative traits in the UK Biobank with sample size (N) $\geq 100,000$. The sparse GRM was constructed using a coefficient of relatedness cutoff 0.125, corresponding to up to 3rd degree relative pairs.

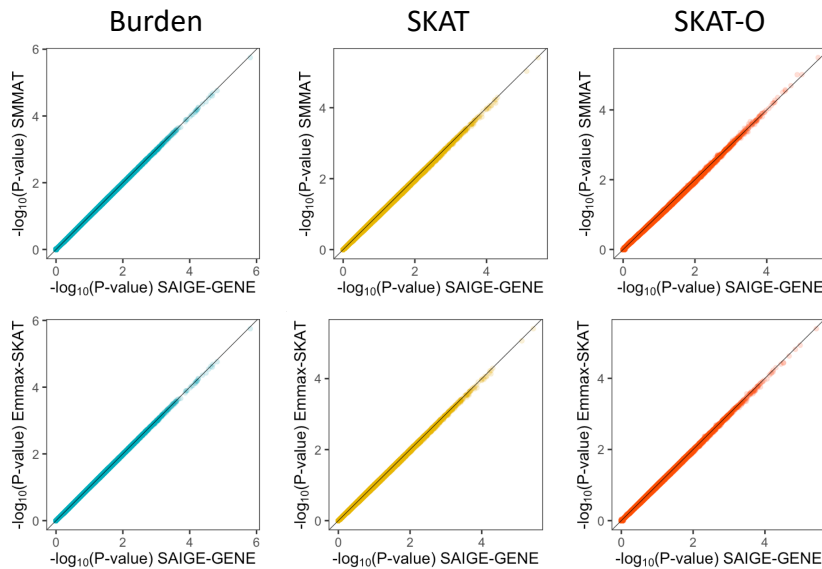


Supplementary Figure 4. Scatter plots of association p-values from SAIGE-GENE versus SMMAT¹⁴ and EmmaX-SKAT for the Burden, SKAT, and SKAT-O tests based on simulation data on the $-\log_{10}$ scale. 1,000,000 genes were tested with 1000 families, each having 10 members, as shown in the **Supplementary Figure 9**. The Pearson's correlation coefficients $r^2 > 0.99$ for $-\log_{10}(\text{P-values})$ between SAIGE and SMMAT and between SAIGE and EmmaX-SKAT. A. $h^2 = 0.2$, B. $h^2 = 0.4$

A. $h^2 = 0.2$

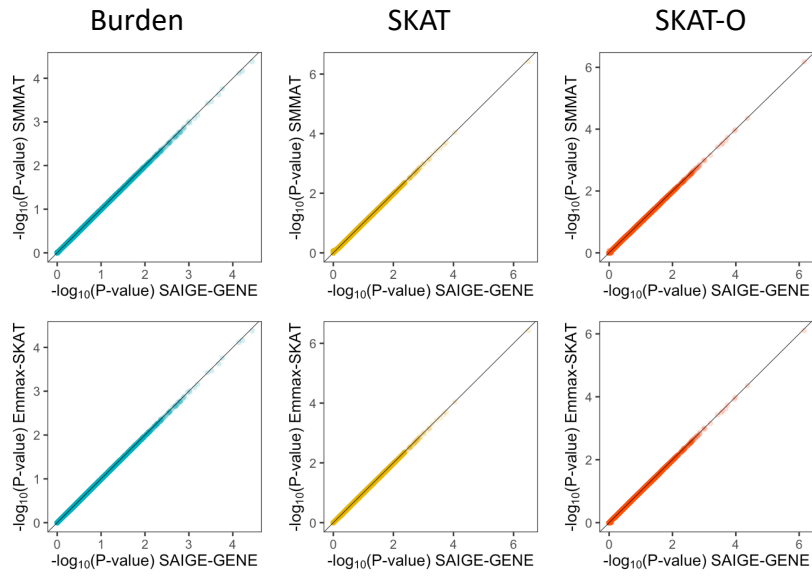


B. $h^2 = 0.4$

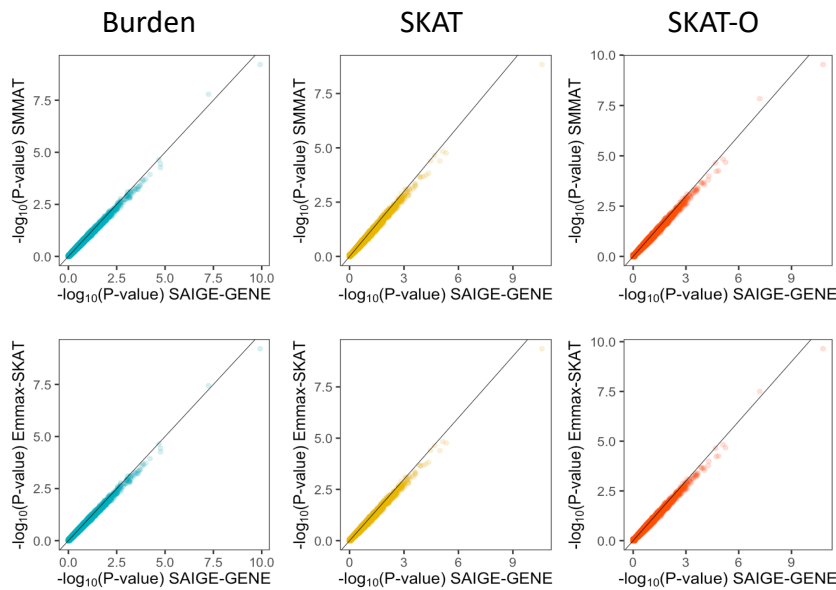


Supplementary Figure 5. Scatter plots of association p-values from SAIGE-GENE versus SMMAT and EmmaX-SKAT for the Burden, SKAT, and SKAT-O tests based on real data analysis on the $-\log_{10}$ scale. 12,000 genes were tested for A. automated read pulse rate using 20,000 randomly selected white British samples in the HRC-imputed UK Biobank; B. HDL using 20,000 randomly selected samples in HUNT. Missense and stop-gain variants with $MAF \leq 1\%$ were included. The Pearson's correlation coefficients $r^2 > 0.99$ for $-\log_{10}(P\text{-values})$ between SAIGE and SMMAT and between SAIGE and EmmaX-SKAT.

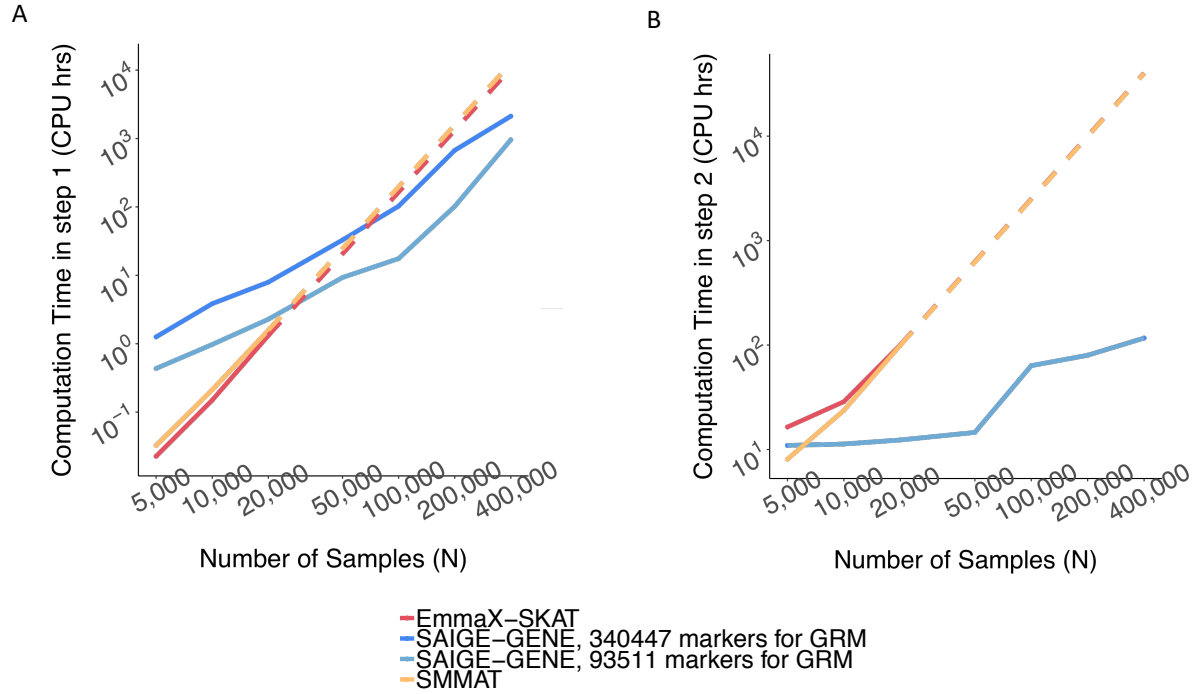
A. automated read pulse rate in the UK Biobank



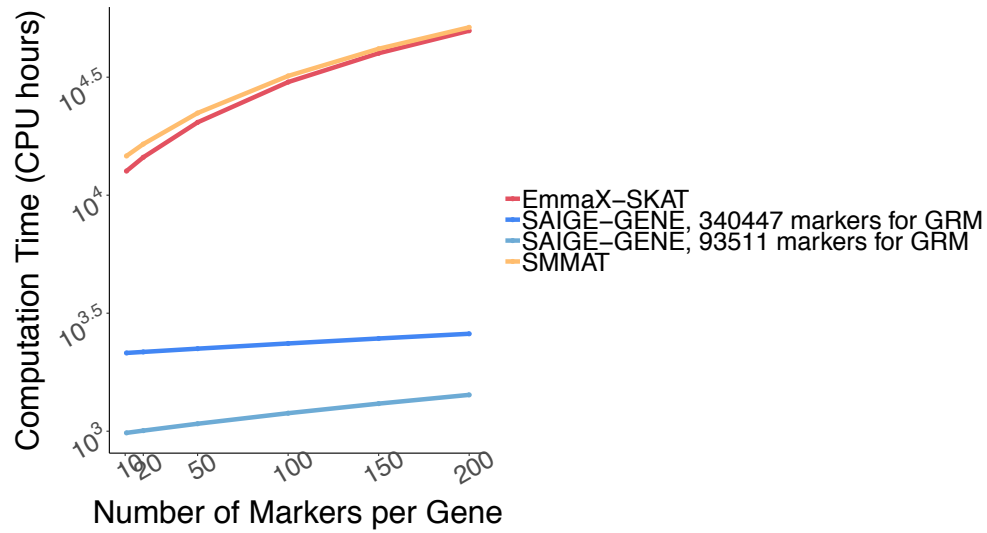
B. HDL in the HUNT study



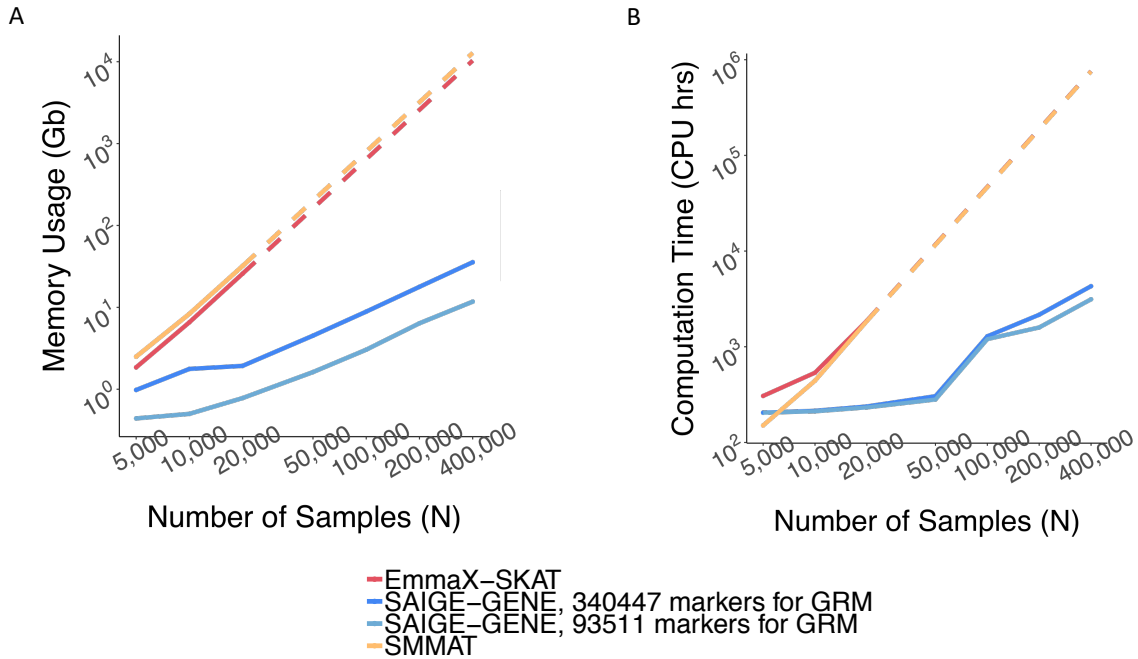
Supplementary Figure 6. Empirical computation time for A. step 1 for fitting a null mixed model and B. step 2 for association tests, respectively by sample sizes (N) for gene-based tests for 15,342 genes, each containing 50 rare variants. Benchmarking was performed on randomly sub-sampled UK Biobank data with 408,144 White British participants for waist-to-hip ratio. The reported run times and memory are medians of five runs with samples randomly selected from the full sample set using different sampling seeds. The reported computation time and memory for EmmaX-SKAT and SMMAT is the projected computation time when $N > 20,000$. As the number of tested markers varies by sample sizes, the computation time is projected for 50 markers per gene for plotting. Numerical data are provided in **Supplementary Table 1.**



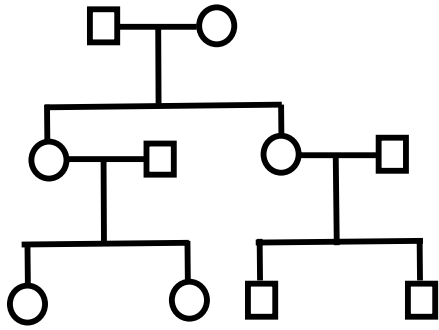
Supplementary Figure 7. Log-log plot of the estimated run time as a function of number of markers per gene. Benchmarking was performed on randomly sub-sampled 400,000 UK Biobank data with 408,144 white British participants for waist-to-hip ratio on 15,342 genes. Run times are medians of five runs with samples randomly selected from the full sample set using different sampling seeds. The computation time for other different number of markers per gene is projected based on the benchmarked time.



Supplementary Figure 8. Log-log plots of the estimated A. run time and B. memory usage as a function of sample size (N) for genome-wide tests for 286,000 chunks, each containing 50 variants on average, given that there are 14.3 million markers in the HRC-imputed UKB with $MAF \leq 1\%$ and imputation info score ≥ 0.8 . Numerical data are provided in **Supplementary Table 1**. Benchmarking was performed on randomly sub-sampled UK Biobank data with 408,144 white British participants for waist-to-hip ratio. run times are medians of five runs with samples randomly selected from the full sample set using different sampling seeds.



Supplementary Figure 9. Pedigree of families, each with 10 members, in the simulation study.



3. Supplementary tables

Supplementary Table 1. The estimated run time (A) and memory use (B) across different sample sizes. Benchmarking was performed on randomly sub-sampled UK Biobank data with 408,144 white British participants and 15,342 genes were tested for waist hip ratio. For simplicity, the number of markers in the gene was 50 regardless of sample sizes. The run times and memory are medians of five runs with samples randomly selected from the full sample set using different sampling seeds. The computation cost for genome-wide region-based tests were projected from the exome-wide gene-based tests results given that there are 14.3 million markers in the HRC-imputed UKB with $MAF \leq 1\%$ and imputation info score ≥ 0.8 . Total 286,000 chunks are tested with 50 markers in each chunk.

	sampleSize	step1(CPU hrs)	step2(CPU hrs)	Total Time(CPU hrs)	Memory (Gb)	Program	
Exome-wide gene-based tests	5,000	0.43	10.94	11.37	0.44	SAIGE-GENE; 93,511 markers for GRM	
	10,000	0.98	11.29	12.26	0.50	SAIGE-GENE; 93,511 markers for GRM	
	20,000	2.28	12.34	14.61	0.78	SAIGE-GENE; 93,511 markers for GRM	
	50,000	9.29	14.57	23.87	1.62	SAIGE-GENE; 93,511 markers for GRM	
	100,000	17.58	63.51	81.09	3.04	SAIGE-GENE; 93,511 markers for GRM	
	200,000	101.85	79.87	181.71	6.39	SAIGE-GENE; 93,511 markers for GRM	
	400,000	958.98	116.65	1075.63	11.74	SAIGE-GENE; 93,511 markers for GRM	
	5,000	1.26	10.94	12.20	0.98	SAIGE-GENE; 340,447 markers for GRM	
	10,000	3.85	11.29	15.14	1.77	SAIGE-GENE; 340,447 markers for GRM	
	20,000	7.95	12.34	20.28	1.93	SAIGE-GENE; 340,447 markers for GRM	
	50,000	32.88	14.57	47.45	4.51	SAIGE-GENE; 340,447 markers for GRM	
	100,000	102.71	63.51	166.22	8.87	SAIGE-GENE; 340,447 markers for GRM	
	200,000	672.29	79.87	752.16	17.82	SAIGE-GENE; 340,447 markers for GRM	
	400,000	2120.89	116.65	2237.54	35.59	SAIGE-GENE; 340,447 markers for GRM	
	5,000	0.02	16.42	16.45	1.85	EmmaX-SKAT	
	10,000	0.15	28.72	28.87	6.57	EmmaX-SKAT	
	20,000	1.32	100.32	101.64	25.82	EmmaX-SKAT	
	50,000	20.62	626.98	647.60	161.37	EmmaX-SKAT	
	100,000	164.93	2507.93	2672.86	645.50	EmmaX-SKAT	
	200,000	1319.44	10031.73	11351.17	2581.99	EmmaX-SKAT	
	400,000	10555.55	40126.91	50682.46	10327.96	EmmaX-SKAT	
	5,000	0.03	8.07	8.11	2.50	SMMAT	
	10,000	0.21	23.68	23.90	8.34	SMMAT	
	20,000	1.57	99.44	101.01	32.03	SMMAT	
	50,000	24.57	621.47	646.04	200.22	SMMAT	
	100,000	196.56	2485.89	2682.46	800.87	SMMAT	
	200,000	1572.50	9943.57	11516.07	3203.49	SMMAT	
	400,000	12580.00	39774.28	52354.28	12813.95	SMMAT	
	Genome-wide region-based tests	5,000	0.43	203.89	204.33	0.44	SAIGE-GENE; 93,511 markers for GRM
		10,000	0.98	210.42	211.40	0.50	SAIGE-GENE; 93,511 markers for GRM
20,000		2.28	229.96	232.24	0.78	SAIGE-GENE; 93,511 markers for GRM	
50,000		9.29	271.67	280.97	1.62	SAIGE-GENE; 93,511 markers for GRM	
100,000		17.58	1183.98	1201.56	3.04	SAIGE-GENE; 93,511 markers for GRM	
200,000		101.85	1488.86	1590.70	6.39	SAIGE-GENE; 93,511 markers for GRM	
400,000		958.98	2174.55	3133.53	11.74	SAIGE-GENE; 93,511 markers for GRM	
5,000		1.26	203.89	205.15	0.98	SAIGE-GENE; 340,447 markers for GRM	
10,000		3.85	210.42	214.27	1.77	SAIGE-GENE; 340,447 markers for GRM	
20,000		7.95	229.96	237.90	1.93	SAIGE-GENE; 340,447 markers for GRM	
50,000		32.88	271.67	304.55	4.51	SAIGE-GENE; 340,447 markers for GRM	
100,000		102.71	1183.98	1286.69	8.87	SAIGE-GENE; 340,447 markers for GRM	
200,000		672.29	1488.86	2161.15	17.82	SAIGE-GENE; 340,447 markers for GRM	

	400,000	2120.89	2174.55	4295.44	35.59	SAIGE-GENE; 340,447 markers for GRM
	5,000	0.02	306.15	306.17	1.85	EmmaX-SKAT
	10,000	0.15	535.40	535.55	6.57	EmmaX-SKAT
	20,000	1.32	1870.08	1871.40	25.82	EmmaX-SKAT
	50,000	20.62	11687.99	11708.61	161.37	EmmaX-SKAT
	100,000	164.93	46751.96	46916.89	645.50	EmmaX-SKAT
	200,000	1319.44	187007.85	188327.29	2581.99	EmmaX-SKAT
	400,000	10555.55	748031.39	758586.94	10327.96	EmmaX-SKAT
	5,000	0.03	150.51	150.55	2.50	SMMAT
	10,000	0.21	441.50	441.71	8.34	SMMAT
	20,000	1.57	1853.64	1855.22	32.03	SMMAT
	50,000	24.57	11585.28	11609.85	200.22	SMMAT
	100,000	196.56	46341.12	46537.68	800.87	SMMAT
	200,000	1572.50	185364.46	186936.96	3203.49	SMMAT
	400,000	12580.00	741457.85	754037.85	12813.95	SMMAT

Supplementary Table 2. Exome-wide significant genes with p-values $\leq 2.5 \times 10^{-6}$ identified by SAIGE-GENE in the UK Biobank for the 53 quantitative traits.

Phentoype	Gene	P-value	Number of variants	Sample Size
Waist circumference	<i>GPR151</i>	2.15E-10	6	408227
Waist circumference	<i>C16orf70</i>	1.94E-06	2	408227
Hip circumference	<i>SYPL2</i>	2.91E-08	4	408182
Hip circumference	<i>TRAPPC4</i>	5.83E-07	2	408182
Hip circumference	<i>ANO1</i>	5.98E-07	4	408182
Hip circumference	<i>GPR151</i>	6.08E-07	6	408182
Hip circumference	<i>C16orf70</i>	1.14E-06	2	408182
Waist hip ratio	<i>TAS2R46</i>	1.36E-08	2	408144
Waist hip ratio	<i>GPR151</i>	3.00E-08	6	408144
Waist hip ratio	<i>SLC5A3</i>	1.33E-07	4	408144
Height standing	<i>SCMH1</i>	3.41E-39	8	408034
Height standing	<i>ACAN</i>	1.39E-36	33	408034
Height standing	<i>FBN2</i>	3.15E-32	16	408034
Height standing	<i>ZFAT</i>	3.61E-32	13	408034
Height standing	<i>HTRA1</i>	1.66E-29	4	408034
Height standing	<i>NPR3</i>	3.60E-25	3	408034
Height standing	<i>STC2</i>	1.10E-24	4	408034
Height standing	<i>SPSB3</i>	1.10E-22	4	408034
Height standing	<i>NUBP2</i>	2.01E-22	9	408034
Height standing	<i>ATAD2</i>	8.07E-21	6	408034
Height standing	<i>ADAMTS6</i>	1.83E-17	6	408034
Height standing	<i>GRAMD2A</i>	3.41E-16	4	408034
Height standing	<i>GRM4</i>	7.23E-16	2	408034
Height standing	<i>MTMR11</i>	1.94E-15	3	408034
Height standing	<i>CRISPLD2</i>	2.04E-14	11	408034
Height standing	<i>CERCAM</i>	7.05E-14	10	408034
Height standing	<i>PDE3B</i>	1.88E-13	8	408034
Height standing	<i>FNDC3B</i>	1.39E-12	8	408034
Height standing	<i>PKD1</i>	6.40E-12	43	408034
Height standing	<i>FER</i>	6.68E-12	5	408034
Height standing	<i>C16orf70</i>	7.45E-12	2	408034
Height standing	<i>HAPLN3</i>	1.25E-11	6	408034
Height standing	<i>ST3GAL4</i>	5.37E-11	5	408034
Height standing	<i>SHANK1</i>	8.09E-11	5	408034
Height standing	<i>TRAPPC13</i>	1.04E-10	2	408034
Height standing	<i>S1PR5</i>	1.44E-10	5	408034
Height standing	<i>PTCH1</i>	2.65E-10	15	408034
Height standing	<i>COL8A1</i>	3.25E-10	2	408034
Height standing	<i>EXD1</i>	3.84E-10	3	408034
Height standing	<i>ATAD5</i>	4.74E-10	8	408034
Height standing	<i>ESR1</i>	7.28E-10	9	408034
Height standing	<i>CLEC3A</i>	9.08E-10	3	408034
Height standing	<i>PTH1R</i>	1.40E-09	3	408034
Height standing	<i>FGFR3</i>	3.50E-09	4	408034
Height standing	<i>NOX4</i>	4.09E-09	4	408034
Height standing	<i>CYR61</i>	4.40E-09	2	408034
Height standing	<i>TBX3</i>	4.59E-09	2	408034
Height standing	<i>SAMD4A</i>	4.77E-09	7	408034
Height standing	<i>ZCCHC6</i>	1.42E-08	7	408034
Height standing	<i>CTU2</i>	2.09E-08	12	408034
Height standing	<i>LPP</i>	2.11E-08	8	408034
Height standing	<i>LRRC8A</i>	2.87E-08	2	408034

Height standing	<i>DGKH</i>	3.08E-08	9	408034
Height standing	<i>ABCB6</i>	4.74E-08	11	408034
Height standing	<i>PIEZO1</i>	6.63E-08	55	408034
Height standing	<i>ELN</i>	6.67E-08	8	408034
Height standing	<i>ATG10</i>	1.35E-07	4	408034
Height standing	<i>CADM1</i>	1.35E-07	2	408034
Height standing	<i>CPPED1</i>	1.40E-07	7	408034
Height standing	<i>PFDN2</i>	1.41E-07	2	408034
Height standing	<i>PROP1</i>	1.63E-07	4	408034
Height standing	<i>PRSS56</i>	1.68E-07	3	408034
Height standing	<i>GYG1</i>	1.78E-07	5	408034
Height standing	<i>CHGA</i>	2.20E-07	9	408034
Height standing	<i>SERPINE2</i>	2.80E-07	3	408034
Height standing	<i>SPATA5</i>	2.83E-07	5	408034
Height standing	<i>AMOTL1</i>	2.85E-07	6	408034
Height standing	<i>TMEM150B</i>	2.88E-07	3	408034
Height standing	<i>COL11A1</i>	2.91E-07	8	408034
Height standing	<i>PTK7</i>	2.95E-07	6	408034
Height standing	<i>PHC3</i>	3.24E-07	7	408034
Height standing	<i>TARS2</i>	4.22E-07	3	408034
Height standing	<i>VASN</i>	5.25E-07	7	408034
Height standing	<i>TXLNA</i>	5.66E-07	3	408034
Height standing	<i>FAM76A</i>	5.70E-07	3	408034
Height standing	<i>C11orf57</i>	5.71E-07	2	408034
Height standing	<i>FBLN2</i>	6.61E-07	13	408034
Height standing	<i>MC3R</i>	6.92E-07	2	408034
Height standing	<i>GTF2E2</i>	7.05E-07	3	408034
Height standing	<i>FIBIN</i>	1.09E-06	2	408034
Height standing	<i>PREB</i>	1.12E-06	5	408034
Height standing	<i>SIX6</i>	1.16E-06	3	408034
Height standing	<i>DOT1L</i>	1.20E-06	5	408034
Height standing	<i>SETD2</i>	1.30E-06	19	408034
Height standing	<i>USP37</i>	1.40E-06	8	408034
Height standing	<i>BCKDHA</i>	1.91E-06	2	408034
Height standing	<i>HR</i>	1.92E-06	21	408034
Height standing	<i>PAM</i>	1.98E-06	5	408034
Height standing	<i>PLG</i>	2.04E-06	11	408034
Height standing	<i>LLGL1</i>	2.30E-06	9	408034
Height standing	<i>ZNF205</i>	2.42E-06	4	408034
Body mass index	<i>GPR151</i>	1.74E-09	6	407605
Body mass index	<i>SYPL2</i>	6.43E-09	4	407605
Body mass index	<i>TRAPPC4</i>	2.87E-07	2	407605
Body mass index	<i>IQSEC1</i>	7.02E-07	12	407605
Body mass index	<i>HCRTR2</i>	1.50E-06	4	407605
Body mass index	<i>FRMD5</i>	1.92E-06	4	407605
Weight	<i>ZFAT</i>	1.80E-11	13	401786
Weight	<i>C16orf70</i>	7.28E-11	2	401786
Weight	<i>STC2</i>	1.69E-09	4	401786
Weight	<i>GPR151</i>	1.73E-09	6	401786
Weight	<i>ANO1</i>	3.10E-08	4	401786
Weight	<i>SYPL2</i>	6.22E-08	4	401786
Weight	<i>FRMD5</i>	6.24E-08	4	401786
Weight	<i>NUBP2</i>	1.70E-07	9	401786
Weight	<i>SCMH1</i>	1.87E-07	8	401786
Weight	<i>TRAPPC4</i>	1.16E-06	2	401786
Weight	<i>HTRA1</i>	2.00E-06	4	401786
Weight	<i>SPSB3</i>	2.32E-06	4	401786
Whole body water mass	<i>STC2</i>	8.55E-24	4	401782

Whole body water mass	<i>NUBP2</i>	1.52E-19	9	401782
Whole body water mass	<i>ZFAT</i>	1.83E-18	13	401782
Whole body water mass	<i>SCMH1</i>	8.44E-18	8	401782
Whole body water mass	<i>SPSB3</i>	2.40E-17	4	401782
Whole body water mass	<i>HTRA1</i>	2.11E-12	4	401782
Whole body water mass	<i>ANO1</i>	4.31E-12	4	401782
Whole body water mass	<i>PHC3</i>	1.32E-11	7	401782
Whole body water mass	<i>C16orf70</i>	3.33E-11	2	401782
Whole body water mass	<i>ATAD2</i>	2.82E-09	6	401782
Whole body water mass	<i>GRM4</i>	6.47E-09	2	401782
Whole body water mass	<i>ESR1</i>	1.15E-07	9	401782
Whole body water mass	<i>ZMYM6</i>	1.75E-07	11	401782
Whole body water mass	<i>TMEM150B</i>	2.96E-07	3	401782
Whole body water mass	<i>FRMD5</i>	2.98E-07	4	401782
Whole body water mass	<i>LCOR</i>	3.88E-07	11	401782
Whole body water mass	<i>PLEKHJ1</i>	9.18E-07	5	401782
Whole body water mass	<i>GLI3</i>	1.09E-06	8	401782
Whole body water mass	<i>ACAN</i>	2.12E-06	33	401782
Basal metabolic rate	<i>STC2</i>	1.50E-20	4	401771
Basal metabolic rate	<i>ZFAT</i>	2.39E-17	13	401771
Basal metabolic rate	<i>NUBP2</i>	4.62E-17	9	401771
Basal metabolic rate	<i>SCMH1</i>	3.15E-15	8	401771
Basal metabolic rate	<i>SPSB3</i>	5.37E-15	4	401771
Basal metabolic rate	<i>C16orf70</i>	9.08E-12	2	401771
Basal metabolic rate	<i>HTRA1</i>	4.87E-11	4	401771
Basal metabolic rate	<i>ANO1</i>	5.53E-11	4	401771
Basal metabolic rate	<i>PHC3</i>	2.11E-10	7	401771
Basal metabolic rate	<i>GRM4</i>	9.53E-10	2	401771
Basal metabolic rate	<i>FRMD5</i>	1.35E-07	4	401771
Basal metabolic rate	<i>ATAD2</i>	2.27E-07	6	401771
Basal metabolic rate	<i>LCOR</i>	7.41E-07	11	401771
Basal metabolic rate	<i>ZMYM6</i>	8.04E-07	11	401771
Basal metabolic rate	<i>TMEM150B</i>	1.10E-06	3	401771
Basal metabolic rate	<i>ESR1</i>	1.14E-06	9	401771
Basal metabolic rate	<i>GPR151</i>	1.97E-06	6	401771
Whole body fat free mass	<i>STC2</i>	5.52E-24	4	401747
Whole body fat free mass	<i>NUBP2</i>	1.70E-19	9	401747
Whole body fat free mass	<i>ZFAT</i>	5.30E-19	13	401747
Whole body fat free mass	<i>SCMH1</i>	2.90E-18	8	401747
Whole body fat free mass	<i>SPSB3</i>	2.51E-17	4	401747
Whole body fat free mass	<i>HTRA1</i>	1.32E-12	4	401747
Whole body fat free mass	<i>ANO1</i>	6.36E-12	4	401747
Whole body fat free mass	<i>PHC3</i>	7.42E-12	7	401747
Whole body fat free mass	<i>C16orf70</i>	1.18E-10	2	401747
Whole body fat free mass	<i>GRM4</i>	6.51E-09	2	401747
Whole body fat free mass	<i>ATAD2</i>	6.80E-09	6	401747
Whole body fat free mass	<i>ESR1</i>	7.69E-08	9	401747
Whole body fat free mass	<i>ZMYM6</i>	2.27E-07	11	401747
Whole body fat free mass	<i>TMEM150B</i>	5.78E-07	3	401747
Whole body fat free mass	<i>FRMD5</i>	5.88E-07	4	401747
Whole body fat free mass	<i>LCOR</i>	1.06E-06	11	401747
Whole body fat free mass	<i>ACAN</i>	1.31E-06	33	401747
Whole body fat free mass	<i>PLEKHJ1</i>	1.72E-06	5	401747
Whole body fat free mass	<i>GLI3</i>	1.94E-06	8	401747
Whole body fat free mass	<i>PAM</i>	2.01E-06	5	401747
Impedance of whole body	<i>CYR61</i>	3.81E-18	2	401746
Impedance of whole body	<i>POR</i>	2.90E-12	7	401746
Impedance of whole body	<i>ADAMTS3</i>	7.23E-12	8	401746

Impedance of whole body	<i>ANO1</i>	4.91E-11	4	401746
Impedance of whole body	<i>STC2</i>	1.61E-09	4	401746
Impedance of whole body	<i>ECM2</i>	1.95E-08	10	401746
Impedance of whole body	<i>ZNF469</i>	3.93E-08	53	401746
Impedance of whole body	<i>NUBP2</i>	3.02E-07	9	401746
Impedance of whole body	<i>FBN2</i>	3.59E-07	16	401746
Impedance of whole body	<i>FAM198A</i>	1.08E-06	7	401746
Impedance of whole body	<i>ZMYM6</i>	1.18E-06	11	401746
Impedance of whole body	<i>SNED1</i>	1.99E-06	9	401746
Body fat percentage	<i>CYR61</i>	4.99E-12	2	401556
Body fat percentage	<i>GPR151</i>	2.85E-10	6	401556
Body fat percentage	<i>SYPL2</i>	2.81E-07	4	401556
Body fat percentage	<i>C10orf35</i>	3.64E-07	2	401556
Days per week moderate phys activity 10min	<i>RAD51AP1</i>	3.87E-07	7	389204
Blood pressure diastolic automated mean	<i>DBH</i>	3.08E-14	12	385365
Blood pressure diastolic automated mean	<i>SLC9A3R2</i>	2.90E-11	5	385365
Blood pressure diastolic automated mean	<i>ARID1B</i>	1.08E-06	7	385365
Pulse rate automated mean	<i>TBX5</i>	9.69E-35	4	385365
Pulse rate automated mean	<i>MYH6</i>	3.61E-15	14	385365
Pulse rate automated mean	<i>TTN</i>	3.18E-10	368	385365
Pulse rate automated mean	<i>KIF1C</i>	4.78E-10	12	385365
Pulse rate automated mean	<i>ARHGGEF40</i>	7.02E-08	7	385365
Pulse rate automated mean	<i>FNIP1</i>	3.58E-07	8	385365
Pulse rate automated mean	<i>DBH</i>	1.74E-06	12	385365
Blood pressure systolic automated mean	<i>SLC9A3R2</i>	2.36E-12	5	385362
Blood pressure systolic automated mean	<i>ZFAT</i>	7.06E-12	13	385362
Blood pressure systolic automated mean	<i>DBH</i>	2.85E-10	12	385362
Blood pressure systolic automated mean	<i>RRAS</i>	1.33E-07	2	385362
Blood pressure systolic automated mean	<i>NOX4</i>	1.91E-07	4	385362
Blood pressure systolic automated mean	<i>TBX5</i>	2.87E-07	4	385362
Blood pressure systolic automated mean	<i>COL21A1</i>	8.05E-07	14	385362

Supplementary Table 3. Empirical type I error rates for SAIGE-GENE, SAIGE-GENE-GCadj (Gene-based tests in SAIGE-GENE adjusted using the GC lambda of single-variants association results), EmmaX-SKAT^{7,15} and SMMAT¹⁴. h^2 : heritability.

	alpha	500 families and 5000 independent samples						1000 families and no independent samples					
		$h^2=0.2$			$h^2=0.4$			$h^2=0.2$			$h^2=0.4$		
		burden	skat	skato	burden	skat	skato	burden	skat	skato	burden	skat	skato
SAIGE-GENE	0.05	0.05217	0.05496	0.05691	0.05324	0.05853	0.05979	0.05266	0.05619	0.05796	0.05505	0.06249	0.06312
	0.01	0.01067	0.01132	0.01226	0.01112	0.01234	0.01315	0.01089	0.01163	0.01256	0.01174	0.01345	0.01420
	0.001	0.00110	0.00118	0.00133	0.00118	0.00135	0.00147	0.00116	0.00121	0.00135	0.00129	0.00147	0.00159
	0.0001	0.00011	0.00012	0.00013	0.00013	0.00014	0.00016	0.00012	0.00013	0.00014	0.00015	0.00016	0.00017
	2.50E-06	2.72E-06	4.12E-06	2.51E-06	3.10E-06	4.30E-06	3.10E-06	4.00E-06	5.10E-06	5.50E-06	3.50E-06	6.30E-06	5.70E-06
SAIGE-GENE-Gcadj	0.05	0.05055	0.05105	0.05366	0.05185	0.05502	0.05685	0.05058	0.05123	0.05375	0.05104	0.05260	0.05488
	0.01	0.01014	0.01030	0.01129	0.01064	0.01138	0.01227	0.01020	0.01028	0.01133	0.01039	0.01072	0.01172
	0.001	0.00102	0.00104	0.00119	0.00110	0.00121	0.00134	0.00105	0.00103	0.00117	0.00105	0.00109	0.00123
	0.0001	0.00010	0.00010	0.00011	0.00011	0.00013	0.00014	0.00010	0.00011	0.00012	0.00011	0.00012	0.00013
	2.50E-06	2.50E-06	3.50E-06	2.20E-06	2.60E-06	3.80E-06	2.50E-06	3.10E-06	4.00E-06	4.90E-06	2.90E-06	4.00E-06	3.80E-06
EmmaX-SKAT	0.05	0.05145	0.05314	0.05541	0.05283	0.05641	0.05820	0.05152	0.05329	0.05552	0.05271	0.05637	0.05807
	0.01	0.01043	0.01082	0.01181	0.01089	0.01168	0.01261	0.01051	0.01082	0.01184	0.01094	0.01170	0.01266
	0.001	0.00107	0.00111	0.00126	0.00114	0.00124	0.00138	0.00110	0.00110	0.00124	0.00115	0.00122	0.00136
	0.0001	0.00011	0.00011	0.00012	0.00012	0.00013	0.00014	0.00011	0.00012	0.00013	0.00013	0.00013	0.00015
	2.50E-06	2.50E-06	3.90E-06	2.30E-06	2.60E-06	4.10E-06	2.60E-06	3.40E-06	4.50E-06	5.30E-06	3.10E-06	4.90E-06	4.60E-06
SMMAT	0.05	0.05146	0.05312	0.05376	0.05282	0.05640	0.05645	0.05152	0.05329	0.05387	0.05271	0.05637	0.05637
	0.01	0.01043	0.01082	0.01159	0.01088	0.01168	0.01239	0.01051	0.01083	0.01164	0.01094	0.01170	0.01244
	0.001	0.00107	0.00110	0.00131	0.00113	0.00124	0.00143	0.00110	0.00110	0.00129	0.00115	0.00122	0.00141
	0.0001	0.00011	0.00011	0.00014	0.00012	0.00013	0.00016	0.00011	0.00012	0.00014	0.00013	0.00013	0.00016
	2.50E-06	2.54E-06	3.75E-06	2.94E-06	2.63E-06	3.85E-06	3.75E-06	3.40E-06	4.21E-06	6.01E-06	3.10E-06	4.40E-06	5.11E-06

Supplementary Table 4. Empirical type 1 error rates for SAIGE-GENE for binary traits with four different prevalence.

	Alpha	Prev=0.01	Prev=0.1	Prev=0.2	Prev=0.5
SKAT-O	0.05	0.0724177	0.05472287	0.0535236	0.05355896
	0.01	0.0279551	0.01234915	0.0113321	0.01111481
	0.001	0.0083735	0.00162367	0.0012498	0.00110831
	0.0001	0.0027596	0.00024293	0.0001461	0.00010691
	2.50E-06	0.0005338	1.59E-05	5.70E-06	1.70E-06
SKAT	0.05	0.0864831	0.05326219	0.0510109	0.05075666
	0.01	0.0327375	0.0116663	0.0103784	0.01003213
	0.001	0.0094332	0.00154323	0.0011195	0.00096757
	0.0001	0.0030869	0.00024373	0.0001332	9.12E-05
	2.50E-06	0.0006283	1.71E-05	6.60E-06	1.70E-06
BURDEN	0.05	0.0456423	0.05036851	0.0505411	0.05084995
	0.01	0.0111034	0.01022357	0.0101198	0.010203
	0.001	0.0023232	0.00111053	0.001029	0.00102092
	0.0001	0.0005645	0.00013059	0.00011	0.000102
	2.50E-06	7.12E-05	5.32E-06	3.70E-06	2.90E-06

Supplementary Table 5. Empirical power for SAIGE-GENE and EmmaX-SKAT with two different percentages of causal variants (top vs bottom panels) and two different ratios of positive and negative effect directions (left vs right). β : effect size. h^2 : heritability.

$h^2=0.4$	Proportion of causal variants= 0.4					
	$\beta -/+ = 0.8/0.2$			$\beta -/+ = 1/0$		
	Burden	SKAT	SKAT-O	Burden	SKAT	SKAT-O
EmmaX-SKAT	64.70%	90.60%	91.90%	96.80%	91%	97.40%
SAIGE-GENE	64.60%	90.30%	91.80%	96.70%	90.80%	97.20%
	Proportion of causal variants= 0.1					
	$\beta -/+ = 0.8/0.2$			$\beta -/+ = 1/0$		
	Burden	SKAT	SKAT-O	Burden	SKAT	SKAT-O
EmmaX-SKAT	55%	87.30%	86.20%	74.50%	87%	85.90%
SAIGE-GENE	55.30%	86.80%	85.90%	75.10%	87%	85.90%

References

- 1 Zhou, W. *et al.* Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat Genet* **50**, 1335-1341, doi:10.1038/s41588-018-0184-y (2018).
- 2 Gilmour, A. R., Thompson, R. & Cullis, B. R. Average Information REML: An Efficient Algorithm for Variance Parameter Estimation in Linear Mixed Models. *Biometrics* **51**, 1440-1450, doi:10.2307/2533274 (1995).
- 3 Kaasschieter, E. F. Preconditioned conjugate gradients for solving singular systems. *Journal of Computational and Applied Mathematics* **24**, 265-275, doi:[https://doi.org/10.1016/0377-0427\(88\)90358-5](https://doi.org/10.1016/0377-0427(88)90358-5) (1988).
- 4 Hutchinson, M. F. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics - Simulation and Computation* **19**, 433-450, doi:10.1080/03610919008812866 (1990).
- 5 Avron, H. & Toledo, S. Randomized algorithms for estimating the trace of an implicit symmetric positive semi-definite matrix %J J. ACM. **58**, 1-34, doi:10.1145/1944345.1944349 (2011).
- 6 Lee, S., Abecasis, G. R., Boehnke, M. & Lin, X. Rare-variant association analysis: study designs and statistical tests. *Am J Hum Genet* **95**, 5-23, doi:10.1016/j.ajhg.2014.06.009 (2014).
- 7 Wu, M. C. *et al.* Rare-variant association testing for sequencing data with the sequence kernel association test. *Am J Hum Genet* **89**, 82-93, doi:10.1016/j.ajhg.2011.05.029 (2011).
- 8 Lee, S., Wu, M. C. & Lin, X. Optimal tests for rare variant effects in sequencing association studies. *Biostatistics* **13**, 762-775, doi:10.1093/biostatistics/kxs014 (2012).
- 9 Lee, S., Teslovich, T. M., Boehnke, M. & Lin, X. General framework for meta-analysis of rare variants in sequencing association studies. *Am J Hum Genet* **93**, 42-53, doi:10.1016/j.ajhg.2013.05.010 (2013).
- 10 Loh, P. R. *et al.* Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet* **47**, 284-290, doi:10.1038/ng.3190 (2015).
- 11 Svishcheva, G. R., Axenovich, T. I., Belonogova, N. M., van Duijn, C. M. & Aulchenko, Y. S. Rapid variance components-based method for whole-genome association analysis. *Nat Genet* **44**, 1166-1170, doi:10.1038/ng.2410 (2012).
- 12 Davis, T. A. *Direct Methods for Sparse Linear Systems (Fundamentals of Algorithms 2)*. (Society for Industrial and Applied Mathematics, 2006).
- 13 Liu, D. J. *et al.* Meta-analysis of gene-level tests for rare variant association. *Nat Genet* **46**, 200-204, doi:10.1038/ng.2852 (2014).
- 14 Chen, H. *et al.* Efficient variant set mixed model association tests for continuous and binary traits in large-scale whole genome sequencing studies. *bioRxiv* (2018).
- 15 Kang, H. M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nat Genet* **42**, 348-354, doi:10.1038/ng.548 (2010).