

IGREX for quantifying the impact of genetically regulated expression on phenotypes

Mingxuan Cai¹, Lin S. Chen², Jin Liu^{3*}& Can Yang^{1*}

¹Department of Mathematics, The Hong Kong University of Science and Technology

¹Department of Public Health Sciences, The University of Chicago

³Center for Quantitative Medicine, Duke-NUS Medical School

Abstract

1 By leveraging existing GWAS and eQTL resources, transcriptome-wide association studies
2 (TWAS) have achieved many successes in identifying trait-associations of genetically-regulated
3 expression (GREX) levels. TWAS analysis relies on the shared GREX variation across GWAS
4 and the reference eQTL data, which depends on the cellular conditions of the eQTL data.
5 Considering the increasing availability of eQTL data from different conditions and the often
6 unknown trait-relevant cell/tissue-types, we propose a method and tool, IGREX, for precisely
7 quantifying the proportion of phenotypic variation attributed to the GREX component. IGREX
8 takes as input a reference eQTL panel and individual-level or summary-level GWAS data. Using
9 eQTL data of 48 tissue types from the GTEx project as a reference panel, we evaluated the
10 tissue-specific IGREX impact on a wide spectrum of phenotypes. We observed strong GREX
11 effects on immune-related protein biomarkers. By incorporating trans-eQTLs and analyzing
12 genetically-regulated alternative splicing events, we evaluated new potential directions for
13 TWAS analysis.

*Correspondence should be addressed to Jin Liu (jin.liu@duke-nus.edu.sg) and Can Yang (macyang@ust.hk)

14 Introduction

15 Genome-wide association studies (GWAS) have successfully identified tens of thousands of
16 unique associations between single-nucleotide polymorphisms (SNPs) and a wide range of
17 complex traits/diseases (<http://www.ebi.ac.uk/gwas/>). More than 90% of identified risk
18 variants are located in non-coding regions [1], making it challenging to understand their
19 functional mechanisms. Increasing evidence [2, 3, 4, 5, 6, 7, 8, 9] has suggested that many of
20 those risk variants may affect traits/diseases via the modulation of their cis gene expression
21 levels. For example, a study of 18 complex traits revealed an enrichment for expression
22 quantitative trait loci (eQTLs) in 11% of 729 tissue-trait pairs [10]. There is great interest in
23 precisely characterizing the specific role of genetically regulated gene expression (GREX) in
24 human traits and diseases.

25 It is well known that the effects of genetic variation on gene expressions depend on cellular
26 contexts [11]. The rapidly increasing availability of eQTL data from different tissue types,
27 cell types, populations and other conditions provides an unprecedented opportunity to study
28 and evaluate GREX effects in a variety of conditions. For example, the V7 release of the
29 Genotype-Tissue Expression (GTEx) project (<https://gtexportal.org/home/>) has collected
30 gene expression samples from 53 non-diseased tissues across 714 individuals [11]. Multiple
31 blood eQTL resources comprising thousands of individuals are made publicly available [12, 13];
32 and other ongoing projects such as Genetics of DNA Methylation Consortium (GoDMC) and
33 eQTLGen consortium are collecting expression data with sample sizes larger than 10,000
34 [14, 15]. Those data serve as rich eQTL resources for a comprehensive evaluation of GREX
35 effects.

36 The vast amount of publicly available eQTL and GWAS data resources enables an integrative
37 framework, transcriptome-wide association studies (TWAS), for mapping gene-level trait
38 associations and evaluating GREX effects on human traits and diseases. Using a reference
39 eQTL panel (e.g., GTEx), gene-specific expression prediction models can be built based on
40 cis-acting genetic factors. Then the gene expression levels of a GWAS cohort can be predicted
41 based on individual genetic profiles, and the genetically-regulated and predicted expression
42 levels are further associated with the phenotype of interest in the GWAS study to map gene-level
43 trait-associations [16, 17, 18, 19, 20, 21, 22, 23, 24]. Existing methods have been proposed
44 [8, 25], including PrediXcan [16], TWAS [17], FOCUS [19], S-PrediXcan [21], UTMOST [26]

45 and CoMM [22]. Through applications to a wide variety of phenotypes, these methods have
46 successfully identified specific gene-trait associations, whereas a comprehensive and precise
47 evaluation of the impact of GREX variation on various traits and the trait-relevant cellular
48 context is still needed [27].

49 TWAS-types of integrative analysis rely on a key assumption: there exists a steady-state
50 GREX variation shared across reference eQTL data and GWAS data, and the steady-state
51 GREX variation can further induce phenotypic variation. The multi-tissue eQTL data from the
52 GTEx project is commonly used as the reference eQTL panel [16, 21, 26]. The GTEx project
53 has collected data from post-mortem donors and has provided a source of largely non-diseased
54 tissues for general purposes. The GTEx reference may or may not have considerable shared
55 GREX variation with GWAS data of specific phenotypes in specific populations. Given the
56 often unknown disease/trait-relevant tissue types and the increasing availability of eQTL data
57 resources from different conditions, there is a need for new methods and tools that can be
58 used to assess the proportion of the shared GREX variation in the phenotypic variation from a
59 global perspective, and guide the selection of eQTL reference data and tissue-types for specific
60 phenotypes and populations.

61 The heritability measure has been widely used to quantify the impact of genetic variation
62 on phenotypic variation, and has served as a preliminary yet insightful assessment of the
63 potential of genetic studies on various phenotypes [28, 29]. Analogous to the heritability
64 measure, the estimation of proportion of GREX on phenotypic variation can also be used to
65 evaluate the impact of the genetic regulatory effects on phenotypes mediated by expression
66 levels, and inform trait-relevant tissue types or conditions in specific populations. To the
67 best of our knowledge, there are two methods that have been proposed for this purpose
68 [23, 24]. The RhoGE method [23] estimates the proportion of phenotypic variation explained
69 by GREX based on linkage-disequilibrium (LD) score regression (LDSC) [30]. Since it ignores
70 the uncertainty in predicting gene expression levels, the proportion of variance explained by
71 GREX could be substantially under-estimated by RhoGE. The other method, known as the
72 gene expression co-score regression (GECS) [24], requires the analyzed SNPs not being in LD
73 to ensure unbiasedness, which greatly limits its applicability in real data analysis.

74 In this work, we propose a unified framework, named IGREX, for quantifying the impact of
75 genetically regulated expression, while accounting for uncertainty in predicted gene expression

76 levels in the presence of moderate to weak eQTL effects. IGREX requires only summary-level
77 GWAS data as input, greatly enhancing the applicability of the method. We evaluated the
78 performance of IGREX with comprehensive simulation studies, highlighting the importance
79 of accounting for expression estimation uncertainty. Using 48 tissue types from the GTEx
80 project as the reference panel, we applied IGREX to both individual-level and summary-level
81 GWAS data sets, and evaluated the tissue-specific IGREX impact on a wide spectrum of
82 cellular and organismal phenotypes. Our results provide new biological insights into the role
83 of gene expression in the genetic architecture of complex traits. We also demonstrate the
84 reproducibility of results. By incorporating trans-eQTLs and analyzing genetically-regulated
85 alternative splicing events, we evaluated new potential directions for TWAS analysis.

86 Results

87 **Method overview.** IGREX is a two-stage method for quantifying the proportion of phenotypic
88 variation that can be attributed to GREX variation. The method can be applied to both
89 individual-level (IGREX-i) and summary-level (IGREX-s) GWAS data. It first evaluates the
90 posterior distribution of GREX effects based on an eQTL reference panel and then estimates
91 the proportion of variance explained by GREX using the ‘predicted’ gene expression in the
92 GWAS data. Here, we briefly introduce the statistical model of IGREX-i and present additional
93 technical details in the Methods Section.

94 Consider a reference eQTL data set \mathcal{D}_r and an individual-level GWAS data set \mathcal{D}_i . The
95 eQTL data $\mathcal{D}_r = \{\mathbf{Y}, \mathbf{X}_r\}$ is comprised of an $n_r \times G$ gene expression matrix, \mathbf{Y} , and an $n_r \times M$
96 genotype matrix, \mathbf{X}_r , where G is the number of genes, M is the number of SNPs and n_r is the
97 sample size. The GWAS data $\mathcal{D}_i = \{\mathbf{t}, \mathbf{X}\}$ contains a phenotype vector $\mathbf{t} \in \mathbb{R}^n$ and a genotype
98 matrix $\mathbf{X} \in \mathbb{R}^{n \times M}$, where n is the sample size of the GWAS data. Let \mathbf{y}_g and $\mathbf{X}_{r,g}$ be the
99 vector of expression levels of the g -th gene and the genotype matrix corresponding to its local
100 (cis) SNPs from the reference panel, respectively. We first relate \mathbf{y}_g to $\mathbf{X}_{r,g}$ with a linear model:

$$\mathbf{y}_g = \mathbf{X}_{r,g}\boldsymbol{\beta}_g + \mathbf{e}_{r,g}, \quad g = 1, \dots, G, \quad (1)$$

101 where $\boldsymbol{\beta}_g \sim \mathcal{N}(\mathbf{0}, \sigma_{\beta_g}^2 \mathbf{I}_{M_g})$ is the vector of genetic effects of M_g cis SNPs on the expression
102 levels of the g -th gene, and $\mathbf{e}_{r,g} \sim \mathcal{N}(0, \sigma_{r,g}^2 \mathbf{I}_{n_r})$ is the error term. Since we are interested
103 in the steady-state component of gene expression levels regulated by genetic variants, $\boldsymbol{\beta}_g$ is

104 assumed to be the same for individuals in both datasets, \mathcal{D}_r and \mathcal{D}_i . Consequently, the GREX
 105 component of individuals in the GWAS data can be evaluated by $\mathbf{X}_g\boldsymbol{\beta}_g$. Meanwhile, we assume
 106 that the genetic effects on the phenotype of interest \mathbf{t} can be decomposed into two parts, i.e.
 107 the effects mediated via GREX and the effects through alternative pathways not mediated by
 108 gene expression levels:

$$\mathbf{t} = \sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\beta}_g + \mathbf{X} \boldsymbol{\gamma} + \boldsymbol{\epsilon}, \quad (2)$$

109 where $\alpha_g \sim \mathcal{N}(0, \sigma_\alpha^2)$ is the effect size of $\mathbf{X}_g\boldsymbol{\beta}_g$ on \mathbf{t} , $\boldsymbol{\gamma} \sim \mathcal{N}(\mathbf{0}, \sigma_\gamma^2 \mathbf{I}_M)$ is the vector of alternative
 110 genetic effects, and $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_\epsilon^2 \mathbf{I}_n)$ is the error term. In this model, $\sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\beta}_g$ and $\mathbf{X} \boldsymbol{\gamma}$
 111 correspond to the overall impact of the GREX component and the alternative genetic effects on
 112 \mathbf{t} , respectively. Thus, the impact of GREX on the phenotype can be quantified by the proportion
 113 of phenotypic variance explained by the GREX component: $\text{PVE}_{\text{GREX}} = \frac{\text{Var}(\sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\beta}_g)}{\text{Var}(\mathbf{t})}$.

114 To estimate this quantity, we propose a two-stage procedure: In the first stage, we estimate
 115 $\sigma_{\beta_g}^2$ and $\sigma_{r,g}^2$ using an efficient algorithm and evaluate the posterior distribution $\boldsymbol{\beta}_g | \mathbf{y}_g, \mathbf{X}_{r,g} \sim$
 116 $\mathcal{N}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$ for all genes. In the second stage, by treating the posterior obtained in the first
 117 stage as the prior distribution of $\boldsymbol{\beta}_g$ in model (2), we can obtain estimated values of σ_α^2 , σ_γ^2
 118 and σ_ϵ^2 using either method of moments (MoM) or restricted maximum likelihood (REML).
 119 Following this procedure, the resulting estimate of PVE_{GREX} is obtained (with details in the
 120 Methods Section) by

$$\widehat{\text{PVE}}_{\text{GREX}} = \frac{\text{tr}(\sum_{g=1}^G \hat{\sigma}_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T)}{\text{tr}(\sum_{g=1}^G \hat{\sigma}_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T + \hat{\sigma}_\gamma^2 \mathbf{X} \mathbf{X}^T + \hat{\sigma}_\epsilon^2 \mathbf{I}_n)}.$$

121 In the above estimation, the substitution of posterior $\boldsymbol{\beta}_g | \mathbf{y}_g, \mathbf{X}_{r,g}$ accounts for the posterior
 122 variance $\boldsymbol{\Sigma}_g$ and naturally results in the adjustment of estimation uncertainty associated with
 123 $\boldsymbol{\beta}_g$. This is important because in the GWAS data, the gene expression levels are not directly
 124 measured, but rather are predicted or imputed based on genetic variants. It is known that
 125 the prediction accuracy and uncertainty vary substantially among genes. For most of the
 126 genes in the genome, the genetically regulated expression variation accounts for only a small to
 127 moderate proportion of total expression variation. Thus, the prediction may not be accurate
 128 and could be subject to high uncertainty. In contrast, our model accounts for the estimation
 129 uncertainty by $\boldsymbol{\Sigma}_g$ and can yield unbiased estimation for $\widehat{\text{PVE}}_{\text{GREX}}$. In addition, the standard
 130 error of $\widehat{\text{PVE}}_{\text{GREX}}$ can be obtained based on the delta method (see Supplementary Note). The
 131 IGREX framework can also be used to test $H_0 : \text{PVE}_{\text{GREX}} = 0$ for the phenotype of interest in

132 specific populations given an eQTL reference with a specific tissue type or cellular context.

133 In real applications, individual-level GWAS data may not be accessible. We have further
 134 developed IGREX-s which requires only summary-level GWAS data as input (See Methods).
 135 Based on MoM, IGREX-s can approximate IGREX-i while requiring only SNP-level z -scores
 136 from GWAS and a reference genotype matrix $\tilde{\mathbf{X}} \in \mathbb{R}^{m \times M}$ of a similar LD pattern to \mathbf{X} , where
 137 m is the number of samples in the reference panel. Using simulations, we showed that with a
 138 few hundreds of samples in the eQTL reference data, the estimation of IGREX-s with summary
 139 statistics well approximates IGREX-i using individual level data. In practice, $\tilde{\mathbf{X}}$ can be \mathbf{X}_r or
 140 a subset of \mathbf{X} . The estimate of PVE_{GREX} given by IGREX-s is

$$\widehat{\text{PVE}}_{\text{GREX}} = \frac{\hat{\sigma}_\alpha^2}{\hat{\sigma}_t^2} \text{tr} \left(\sum_{g=1}^G (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \hat{\mathbf{R}}_g \right),$$

141 where $\hat{\mathbf{R}}_g = \tilde{\mathbf{X}}_g^T \tilde{\mathbf{X}}_g / (m - 1)$ is the estimated LD matrix associated with the g -th gene and
 142 $\tilde{\mathbf{X}}_g$ is the corresponding columns of $\tilde{\mathbf{X}}$. IGREX also allows for the adjustment of covariates
 143 including sex, age and genotype principal components (See details in Supplementary Note).

144 **Simulation studies.** We conducted extensive simulation studies to evaluate the performance
 145 of IGREX. For all the simulated data, we fixed $n = 4,000$, $G = 200$, $M = 20,000$ (i.e., 100 cis
 146 SNPs for each gene). The total phenotypic heritability was set as $h_t^2 = \frac{\text{Var}(\sum_{g=1}^G \alpha_g \mathbf{x}_g^T \boldsymbol{\beta}_g + \mathbf{x}^T \boldsymbol{\beta}_g)}{\text{Var}(t)} =$
 147 0.5, where $\text{PVE}_{\text{GREX}} = 0.2$ and the proportion explained by the alternative genetic effects,
 148 $\text{PVE}_{\text{Alternative}} = \frac{\text{Var}(\mathbf{x}_g^T \boldsymbol{\gamma})}{\text{Var}(t)} = 0.3$ (results for other scenarios are shown in Supplementary Fig.
 149 1-3). To simulate the genotype data, we first sampled the minor allele frequencies (MAF) from
 150 uniform distribution $\mathcal{U}(0.05, 0.5)$ and data matrices from normal distribution $\mathcal{N}(\mathbf{0}, \Sigma(\rho))$, where
 151 $\Sigma_{jj'} = \rho^{|j-j'|}$ characterizes the LD patterns between SNPs. Then, the genotype matrices \mathbf{X}_r and
 152 \mathbf{X} were obtained by categorizing the entries of generated data matrices into 0, 1, 2 according
 153 to MAF. Given the genotype matrices, $\boldsymbol{\beta}_g$ and α_g , the gene expression \mathbf{y}_g and phenotype \mathbf{t}
 154 were simulated following models (1) and (2). We will discuss the details for generating $\boldsymbol{\beta}_g$ and
 155 α_g later. To assess IGREX-s, we calculated the z -score of each SNP and randomly subsetted
 156 $m = 500$ samples from \mathbf{X} for estimating LD matrix $\hat{\mathbf{R}}_g$ (results for other settings of m are
 157 shown in Supplementary Fig. 4).

158 We first evaluated the estimation performance of IGREX for different settings of eQTL
 159 reference data. Specifically, we varied n_r at $\{800, 1000, 2000\}$, $\text{PVE}_y = \frac{\text{Var}(\mathbf{x}^T \boldsymbol{\beta}_g)}{\text{Var}(\mathbf{y}_g)}$ at $\{0.1, 0.2, 0.3\}$,
 160 where PVE_y quantifies the gene expression heritability explained by its local SNPs. To mimic the

161 scenario in which the expression estimation uncertainty was incorrectly ignored, we obtained the
162 posterior mean of β_g in the first stage, and replaced the true effect size β_g by its posterior mean
163 μ_g while specifying the posterior variance to be $\Sigma_g = \mathbf{0}$ in the second stage, and then conducted
164 REML and MoM as before. We denoted these methods as REML₀ and MoM₀. The simulation
165 results summarized in Fig.1a showed that both PVE_{GREX} and PVE_{Alternative} were accurately
166 estimated using REML-based IGREX-i in all settings. The MoM-based IGREX-i slightly
167 underestimated PVE_{GREX} when both sample size n_r and PVE_y were very small, but steadily
168 achieved similar performance as REML-based estimation when either n_r or PVE_y increased. In
169 all settings, IGREX-s well approximated MoM, producing nearly identical estimation results. In
170 contrast, both REML₀ and MoM₀ did not account for estimation uncertainty in the expression
171 prediction, and they showed poor estimation performance even when sample size was large and
172 PVE_y value was high.

173 Next we conducted simulations to evaluate the situation that the IGREX model was mis-
174 specified. Here we considered the situation where genetic effects β_g and α were sparse while
175 we assumed dense effect sizes in the IGREX model. This was designed to mimic the real data
176 situation that the architecture of eQTL signals is often sparse [31]. Let π_α and π_β be the sparsity
177 of α and β_g , i.e., $\pi_\alpha = (\# \text{ Nonzero entries in } \alpha)/G$ and $\pi_\beta = (\# \text{ Nonzero entries in } \beta_g)/M_g$,
178 respectively. To evaluate the influence of different sparsity patterns on our method, we varied
179 π_α and π_β at $\{0.2, 0.5, 0.8\}$. The nonzero entries in α and β_g were simulated from a normal
180 distribution. As shown in Fig. 1b-c, all three methods of IGREX produced accurate estimates
181 in the presence of sparse genetic effects, implying the robustness of IGREX to model mis-
182 specification. Moreover, the estimation performance was not influenced by the degree of sparsity.
183 Next, we investigated the influence of LD patterns by letting ρ vary at $\{0.1, 0.3, 0.5, 0.8\}$. From
184 Fig.1d, we observed that IGREX produced accurate estimation in the presence of LD. In
185 contrast, REML₀ and MoM₀ consistently underestimated PVE_{GREX} as a result of ignoring
186 estimation uncertainty.

187 We also compared IGREX with an existing method in the literature, RhoGE [23]. RhoGE
188 is an LDSC-based approach for estimating PVE_{GREX}. However, this method does not adjust
189 for estimation uncertainty. The results are shown in Fig. 1e. As expected, IGREX yielded
190 unbiased estimation while RhoGE substantially underestimated PVE_{GREX} in most settings. It
191 achieved similar accuracy as IGREX only when the genetically regulated expression accounted

192 for most of the expression variation, $PVE_y \geq 0.9$. In other words, RhoGE only works well
193 when the genetically-predicted expression levels are very close to the true underlying expression
194 levels for most of the genes, which may not be realistic for real data analysis.

195 **Real data applications with individual-level GWAS data.** With eQTL data of 48
196 human tissues from the GTEx project as reference, we applied IGREX to two individual-level
197 GWAS datasets, the Northern Finland Birth Cohorts program 1966 (NFBC) [32] and the
198 Wellcome Trust Case Control Consortium (WTCCC) [33]. The details of the datasets and the
199 data pre-processing procedures are described in the Methods Section.

200 In analyzing the NFBC data, we focused on six quantitative traits with statistically
201 significant heritability, based on 5,123 individuals and 309,245 genotyped SNPs. Those six
202 traits are Glucose ($h_t^2 = 14.2\% \pm 5.3\%$), high-density lipoprotein cholesterol (HDL, $h_t^2 =$
203 $32.9\% \pm 5.6\%$), low-density lipoprotein cholesterol (LDL, $h_t^2 = 29.0\% \pm 5.5\%$), triglycerides (TG,
204 $h_t^2 = 13.6\% \pm 5.3\%$), total cholesterol (TC, $h_t^2 = 20.1\% \pm 5.4\%$) and systolic blood pressure
205 (SysBP, $h_t^2 = 17.1\% \pm 5.4\%$). Fig. 2a-b shows the tissue-specific \widehat{PVE}_{GREX} estimates of the six
206 traits. The REML and MoM methods yielded similar estimates in most of the tissues.

207 IGREX can also be used to inform trait-relevant tissue types. By testing $H_0 : PVE_{GREX} = 0$
208 in each tissue type, we observed significant GREX components in liver for both LDL and
209 TC. As shown in Fig. 2a, \widehat{PVE}_{GREX} for LDL in liver is as high as 14.3% (with standard
210 error 2.6%), capturing 52.6% of total heritability defined as PVE_{GREX}/h_t^2 ; and TC also has
211 $\widehat{PVE}_{GREX} = 13.7\%$ (with standard error 2.5%) in liver, which captures 79.4% of total heritability
212 (see Supplementary Fig. 6). It is known that LDL synthesized in liver is an important lipoprotein
213 particle for transporting cholesterol in the blood [34, 35]. Our findings suggest that genetic
214 variants affect LDL through regulating their corresponding gene targets and liver is the most
215 relevant tissue involved in gene regulation. Next, we analyzed the impact of ignoring the
216 estimation uncertainty (with the complete results given in the Supplementary Fig. 5). As
217 shown in Fig. 2c-d, the \widehat{PVE}_{GREX} declined substantially as a result of ignoring expression
218 estimation uncertainty. In Fig. 2e, we compared the estimates based on individual level data
219 using IGREX-i versus those based on IGREX-s with summary statistics. For all six of the
220 traits, the IGREX-s estimates well approximated the estimates using the individual level data,
221 which is consistent with our simulation results.

222 Next we investigated the role of GREX in complex human traits and diseases, using the

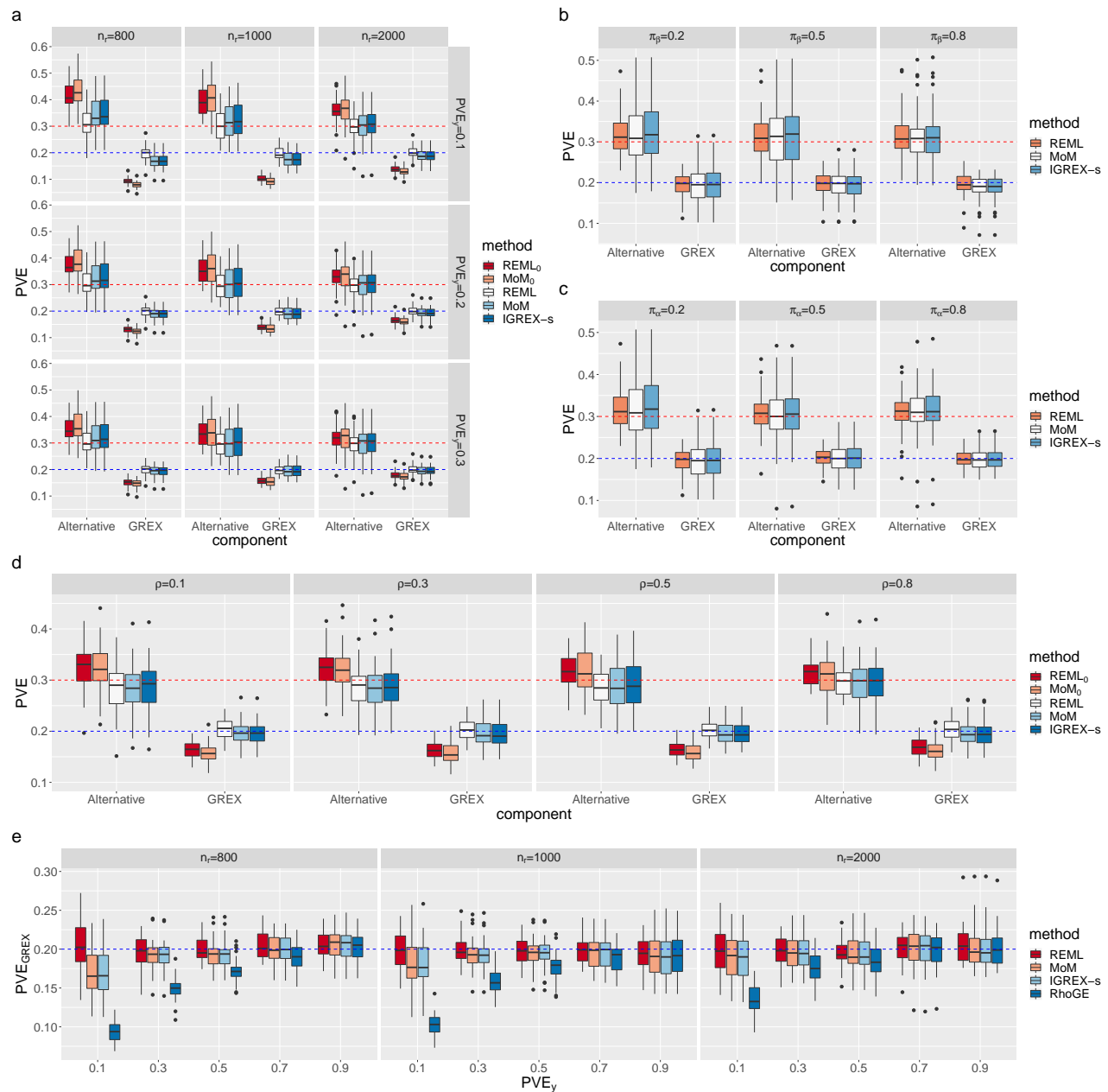


Figure 1: Simulation studies to compare estimation accuracies of IGREX with other methods. REML and MoM in the legend are abbreviations of the IGREX-*i* estimation methods. The blue and red dashed lines represent the true values of PVE_{GREX} and $PVE_{Alternative}$, respectively. We averaged the results over 30 replications and generated box plots for evaluating the estimation performance of: **a** the three models of IGREX, REML₀ and MoM₀ when n_r was varied at $\{800, 1000, 2000\}$ and PVE_y was varied at $\{0.1, 0.2, 0.3\}$; **b** the three models of IGREX when $\pi_\alpha = 0.2$ and π_β was varied at $\{0.2, 0.5, 0.8\}$; **c** the three models of IGREX when $\pi_\beta = 0.2$ and π_α were varied at $\{0.2, 0.5, 0.8\}$; **d** the three models of IGREX, REML₀ and MoM₀ when ρ was varied at $\{0.1, 0.3, 0.5, 0.8\}$; **e** the three models of IGREX and RhoGE when n_r was varied at $\{800, 1000, 2000\}$.

223 WTCCC dataset [33]. We applied IGREX to estimate the tissue-specific PVE_{GREX} of seven
 224 diseases including bipolar disorder (BD), coronary artery disease (CAD), Crohn's disease (CD),

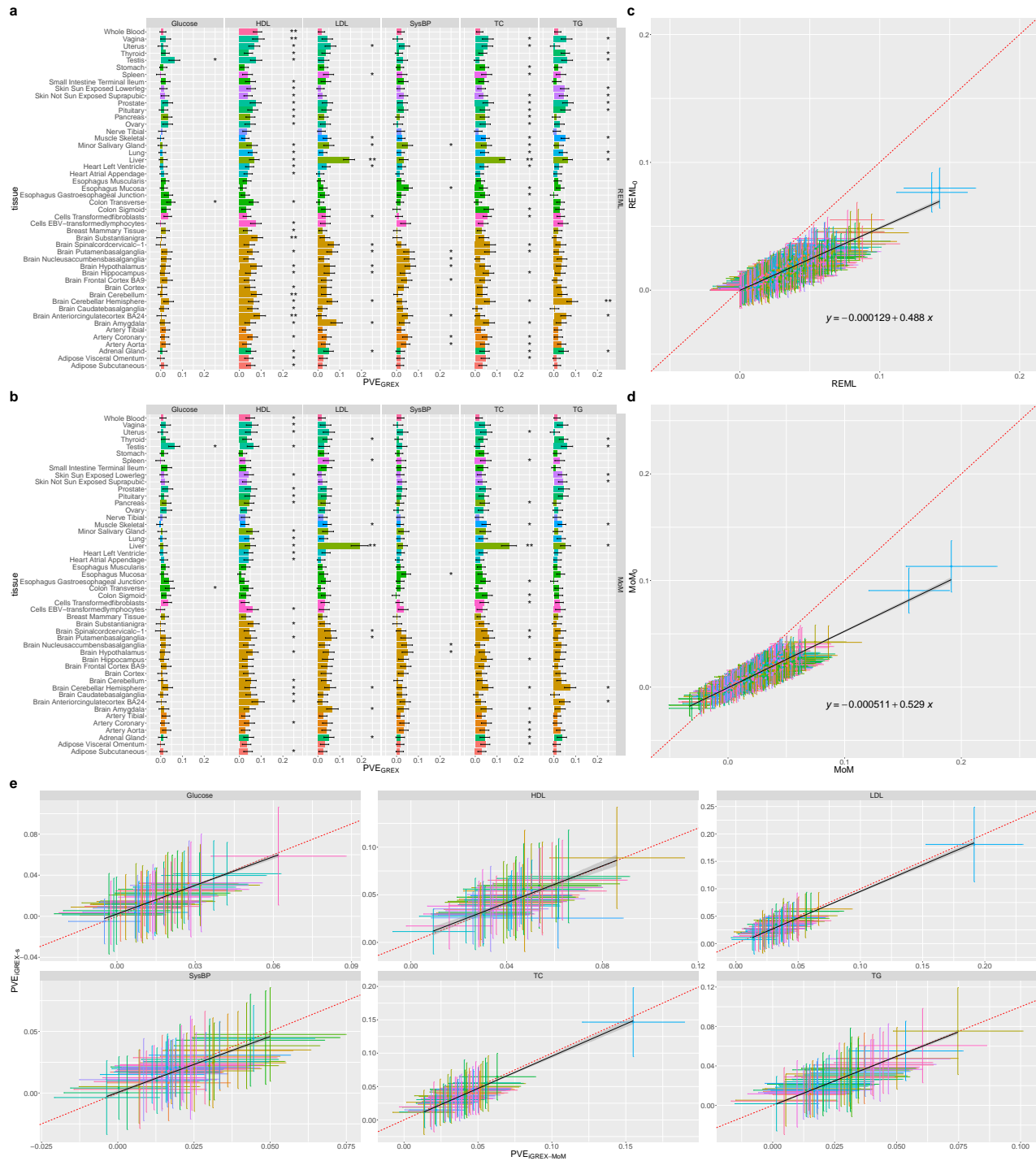


Figure 2: Tissue-specific \widehat{PVE}_{GREX} of the six traits from NFBC data set. (a-b) \widehat{PVE}_{GREX} obtained by REML and MoM. Tissues are colored according to their categories. The number of asterisks represents the significance level: p -value < 0.05 is annotated by *; p -value $< 0.05/48$ is annotated by **. (c-d) All pairs of estimates generated by REML and MoM against their counterparts without accounting for uncertainty. A regression line is fitted and the estimated coefficients are given in the plot. (e) Each panel is a plot of \widehat{PVE}_{GREX} generated by IGREX-s against those generated by MoM for all 48 tissues in one of the six traits.

225 hypertension (HT), rheumatoid arthritis (RA), type 1 diabetes (T1D) and type 2 diabetes
 226 (T2D). The estimates of PVE_{GREX}/h_t^2 obtained by REML are shown in Supplementary Fig. 8.
 227 The top GREX components measured by PVE_{GREX}/h_t^2 are 12.8% for BD in amygdala, 21.2%
 228 for CAD in spinal cord, 18.4% for CD in amygdala, 16.7% for HT in spleen and 17.9% for T2D
 229 in anterior cingulate cortex. The average estimates of PVE_{GREX}/h_t^2 across 48 tissues for RA
 230 and T1D are as high as 34.1% and 71.2%, respectively. Both RA and T1D are autoimmune
 231 diseases, with well-established strong associations in the major histocompatibility complex
 232 (MHC) region [33, 36]. After removing the MHC region, we observed a substantial reduction in
 233 the PVE_{GREX}/h_t^2 estimates: the mean \widehat{PVE}_{GREX} dropped from 34.1% to 7.6% for RA and from
 234 71.2% to 11.7% for T1D, as shown in Fig. 3a. Additionally, the tissue-specific comparisons
 235 presented in Fig. 3b showed an extensive reduction of PVE_{GREX} in all tissue types for T1D
 236 and RA, while such changes were not observed for other traits. This finding suggests the heavy
 237 involvement of GREX variation in the immune functions related to the MHC region for both
 238 RA and T1D. Here we illustrate that IGREX can be used to inform disease/trait-relevant
 239 tissue types or cellular contexts.

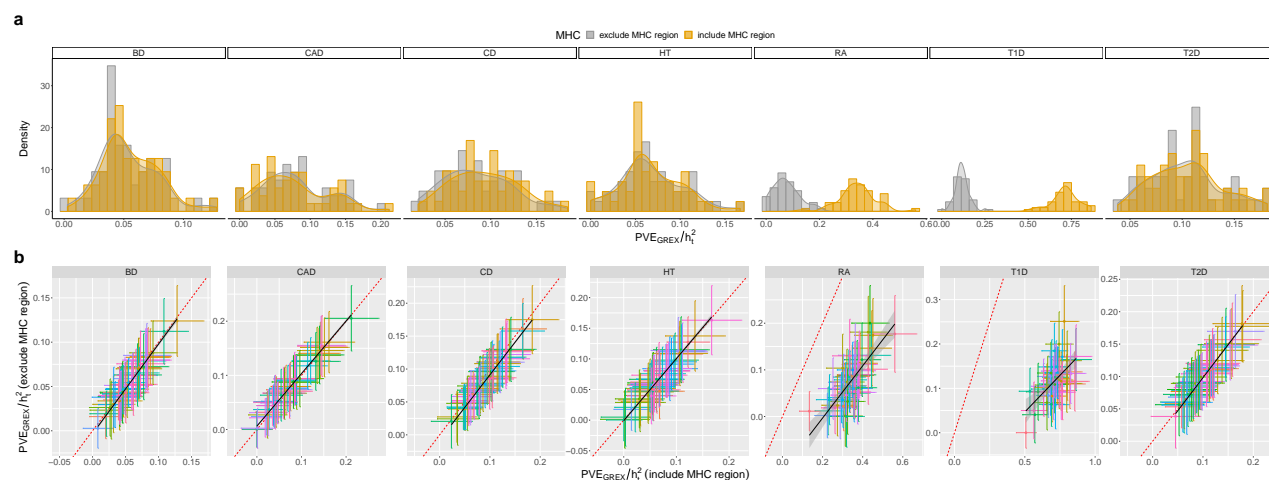


Figure 3: Percentage of heritability explained by GREX (PVE_{GREX}/h_t^2) of the seven traits from WTCCC data. (a) The distributions of estimated PVE_{GREX}/h_t^2 across 48 GTEx tissues. (b) Tissue-specific comparisons of PVE_{GREX}/h_t^2 estimated by whole genome with those estimated by excluding the MHC region.

240 **Analysis of a wide spectrum of phenotypes using IGREX-s with summary-level**
 241 **GWAS data** The vast amount of publicly available summary-level GWAS data and their easy
 242 accessibility allow us to conduct a comprehensive evaluation of the impact of GREX on a wide
 243 spectrum of phenotypes using IGREX-s, from molecular traits such as proteins and metabolites

244 to various complex phenotypes including schizophrenia, height, and body mass index (BMI).
245 In the following analysis, we used the genotypes of the 635 GTEx samples as the LD reference
246 $\tilde{\mathbf{X}}$ in the IGREX-s estimation.

247 First, we estimated PVE_{GREX} in 249 proteins with significantly nonzero heritabilities
248 using summary statistics from a plasma protein quantitative trait loci (pQTL) study [38], as
249 summarized in Fig. 4a. In Supplementary Fig. 10, the heritabilities estimated by IGREX
250 ($\hat{h}_t^2 = \widehat{\text{PVE}}_{\text{GREX}} + \widehat{\text{PVE}}_{\text{Alternative}}$) are shown to be highly consistent with those estimates obtained
251 using MoM [37]. From this perspective, heritability can be attributed to two components:
252 the GREX component and its alternative effects. Then, we grouped 48 tissue types into 16
253 groups by their functions and tested the significance of tissue-specific GREX effects on the 249
254 proteins. We observed a significant GREX contribution in many tissue-protein pairs (Fig. 4b
255 and Supplementary Fig. 11-13). In particular, 9 out of the 249 proteins had significant GREX
256 components in at least one tissue type at 0.05 level after Bonferroni correction. As shown in Fig.
257 4d-e, some proteins, including CD96, DEFB119, MICB and PDE4D, exhibit cross-tissue GREX
258 impacts; meanwhile other proteins, namely CFB, CXCL11, EVI2B, IDUA and LRPAP1, have
259 tissue-specific GREX effect patterns. We found these tissue-specific patterns to be consistent
260 with protein functions. For example, the CFB protein, which is implicated in the growth of
261 preactivated B-lymphocytes, is found to be most associated with GREX in EBV-transformed
262 lymphocytes ($\widehat{\text{PVE}}_{\text{GREX}} = 22.7\%$). As another example, the CXCL11 protein has the highest
263 $\widehat{\text{PVE}}_{\text{GREX}} = 20.0\%$ in pancreas, and the *CXCL11* gene is often over-expressed in pancreas tissue
264 [39]. We also noted that 6 out of the 9 proteins were immune-related, echoing our previous
265 implications of the important role of GREX in immune processes. In addition to the proteins,
266 metabolic traits are also important intermediate traits for complex biological processes. We
267 applied IGREX-s to a summary level data set of circulating metabolites [40], and studied the
268 impact of GREX on metabolic traits. The results are discussed in the Supplementary Materials.

269 Then we applied IGREX-s to the summary data of complex human traits. Here we analyzed
270 three traits: schizophrenia (SCZ), height, and BMI. We considered four datasets of schizophre-
271 nia with increasing and overlapping samples: SCZ subset [41], SCZ1 [42], SCZ1+Sweden
272 (SCZ1Swe)[43] and SCZ2 [44]. We found that the estimated $\text{PVE}_{\text{GREX}}/h_t^2$ in all four SCZ
273 datasets have higher values in the brain tissues than in other tissue types (Fig. 5b). As
274 expected, the statistical power increases with sample size of GWAS (Fig. 5a). Additionally, we

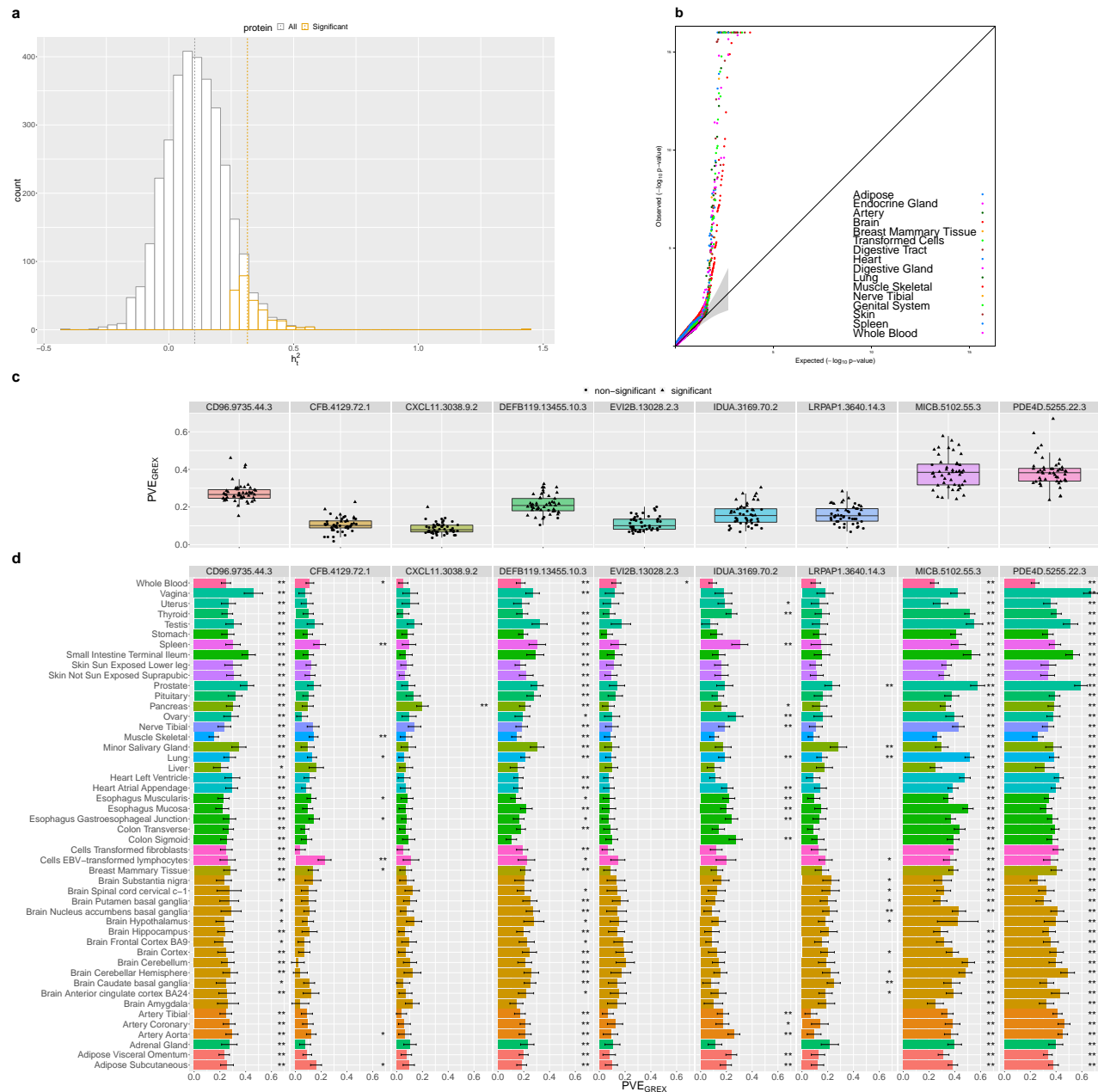


Figure 4: Analysis of plasma pQTL summary statistics. (a) The distribution of estimated heritabilities of 3,283 proteins estimated using [37]. The whole study is colored in grey, while the 249 proteins with significant heritabilities are colored in yellow. Dashed lines represent the means of corresponding distributions. (b) QQ-plot of PVE_{GREX} p -values of tissue-protein pairs. GTEx tissues are categorized into 16 types and colored accordingly. (c) The Manhattan plot of the protein encoding genes in aorta, cerebellum, liver and whole blood. Each point represents a tissue-protein pair. (d) \widehat{PVE}_{GREX} in the 9 proteins whose \widehat{PVE}_{GREX} are significant in at least one tissue at 0.05 level using Bonferroni correction. (e) \widehat{PVE}_{GREX} obtained by IGREX-s. Tissues are colored according to their categories. The number of asterisks represents the significance level: p -value < 0.05/48 is annotated by *; p -value < 0.05/(48 * 9) is annotated by **.

275 also analyzed the human height and BMI phenotypes using pairs of independent GWAS data
 276 for replication purposes. The obtained estimates, \widehat{PVE}_{GREX} , from pairs of independent GWAS

277 data are highly consistent. Although the analysis results are reproducible in several different
 278 data sets, we noted the estimated percentages of heritability explained by GREX for all three
 279 complex traits are less than 10% (8.7% for schizophrenia, 8.7% for height and 3.7% for BMI in
 280 the most expressed tissue types. See Fig. 5c and Supplementary Fig. 15).

281 The relatively low GREX contribution to complex traits other than lipid or molecular
 282 traits can be attributed to multiple reasons. First, it is known that trans-acting genetic effects
 283 can explain a substantial proportion of expression variation [8, 12]. However, trans-eQTL
 284 effects are often tissue-specific and can be harder to detect and replicate across studies [45].
 285 In TWAS-types of analysis, generally the prediction of gene expression is based on only cis
 286 genetic variants of each gene. As such, the \widehat{PVE}_{GREX} values reported here, also based on only
 287 cis genetic variants, may be underestimated. In the next section, we will further explore the
 288 contribution of trans-eQTLs. Second, the genetic effects on gene expression may not be steady
 289 across the reference GTEx data with largely non-diseased tissues for general purposes and the
 290 GWAS data with diseased individuals from specific populations [46]. From this perspective,
 291 before analyzing specific complex traits and diseases via TWAS, it would be helpful to first
 292 estimate the impact of GREX and select the most informative available eQTL reference data.

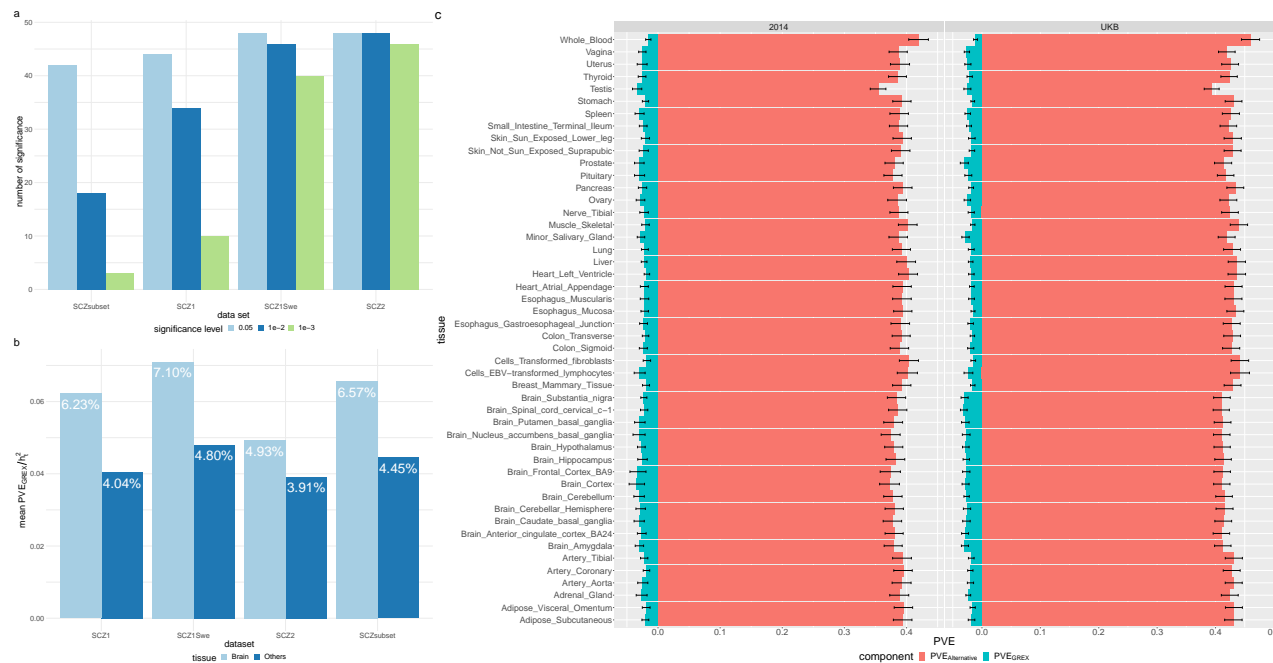


Figure 5: Analyses of complex traits: schizophrenia and height. (a) Number of significant GREX components revealed under different significance levels for the four schizophrenia datasets. (b) Mean estimated percentages of heritability for schizophrenia explained by GREX in brain tissues and in other tissues. (c) \widehat{PVE}_{GREX} and $\widehat{PVE}_{Alternative}$ of height estimated using height2014 and UKB datasets, respectively.

293 **Additional insights on GREX considering trans-eQTLs and genetically-regulated**
294 **alternative splicing events.** The cis-acting genetic effects on local gene expression levels are
295 often shared across tissue types and are often replicable across studies [47]. It is also reported
296 that a substantial proportion (up to 70%) of gene expression heritability can be attributed to
297 trans-acting genetic effects which act predominantly in a tissue-specific manner and have a
298 lower rate of replication across studies [48, 49]. More recently, the eQTLGen consortium [15]
299 has conducted a blood-eQTL meta-analysis and has reported 6,298 (31%) trans-eQTL genes
300 for 10,317 trait-associated SNPs using 31,684 blood samples from 37 datasets. The results
301 suggest that trans-eQTLs are prevalent in the genome, while it is still underpowered to detect
302 them for tissues other than whole blood given the often tissue-specific nature of trans-genetic
303 effects and the limited sample sizes for most tissue types.

304 Although it is still unrealistic to account for all trans-eQTLs in the estimation of PVE_{GREX}
305 due to the limitation of sample sizes, it is possible to explore the potential by incorporating the
306 blood-based trans-eQTLs reported by the eQTLGen consortium and re-estimating PVE_{GREX}/h_t^2 .
307 We first analyzed 13 datasets comprised of 12 phenotypes that have significant PVE_{GREX}/h_t^2
308 estimates in the whole blood, including 7 proteins, 1 lipid trait and 4 complex diseases (with
309 2 SCZ datasets). We observed an increasing trend of PVE_{GREX}/h_t^2 in the blood for all 13
310 datasets (Fig. 6a), by accounting for only $\sim 1,700$ unique trans-eQTLs that are not cis-eQTLs.
311 As a comparison, we applied the same procedure to 13 GTEx brain tissues of the two largest
312 SCZ datasets, and did not observe an increase in PVE_{GREX}/h_t^2 (Fig. 6b). This is not surprising
313 because the trans-eQTLs incorporated above were detected and reported based on whole
314 blood samples and may not be trans-eQTLs in the brain tissues. Our results suggest that
315 the estimation of GREX impacts on traits can be further boosted by incorporating robust
316 trans-eQTLs from the same tissue types.

317 In addition to the gene expression level, we also evaluated the effects of alternative splicing
318 on complex trait heritability. We applied IGREX to quantify the impact of genetically regulated
319 alternative splicing on multiple phenotypes. Alternative splicing is an important gene regulatory
320 process that results in multiple transcripts from a single multi-exon gene. It is commonly
321 observed in humans and plays an essential role in cellular differentiation [50, 51]. Differential
322 variations in splicing may also result in phenotypic variation and contribute to the development
323 of complex diseases including cancer [52, 53, 54]. In a recent work, by extending the TWAS

324 framework to analyze splicing events and associating 40 complex traits with genetically-predicted
 325 splicing quantification, novel putative disease-associated genes were detected [55]. Here, using
 326 multi-tissue splicing quantification data from GTEX as reference, we applied IGREX to study
 327 the impact of genetically-regulated splicing events on four trait-tissue pairs that were found to
 328 have a high PVE_{GREX}/h_t^2 . We estimated the proportion of phenotypic variation explained by
 329 genetically-regulated splicing to be 12.5%, 13.5%, 1.0% and 1.1% for LDL in liver, TC in liver,
 330 SCZ in amygdala and SCZ in cerebellar hemisphere, respectively. Unlike eQTLs that are often
 331 found to be near transcription starting sites, most of the sQTLs were found to be enriched
 332 within gene bodies, in particular within the introns they regulate, and have little to no effects
 333 on cis gene expression levels [55, 56]. In other words, sQTLs are often independent of eQTLs.
 334 Therefore, integrating genetically-regulated splicing quantification may partially explain the
 335 phenotypic variation attributed to alternative genetic factors, $PVE_{Alternative}$. We argue that
 336 with the proper multi-omics reference data, similar analyses can be conducted to quantify
 337 the impact of genetically-regulated methylation, protein, and other multi-omics variation on
 338 phenotype [51].

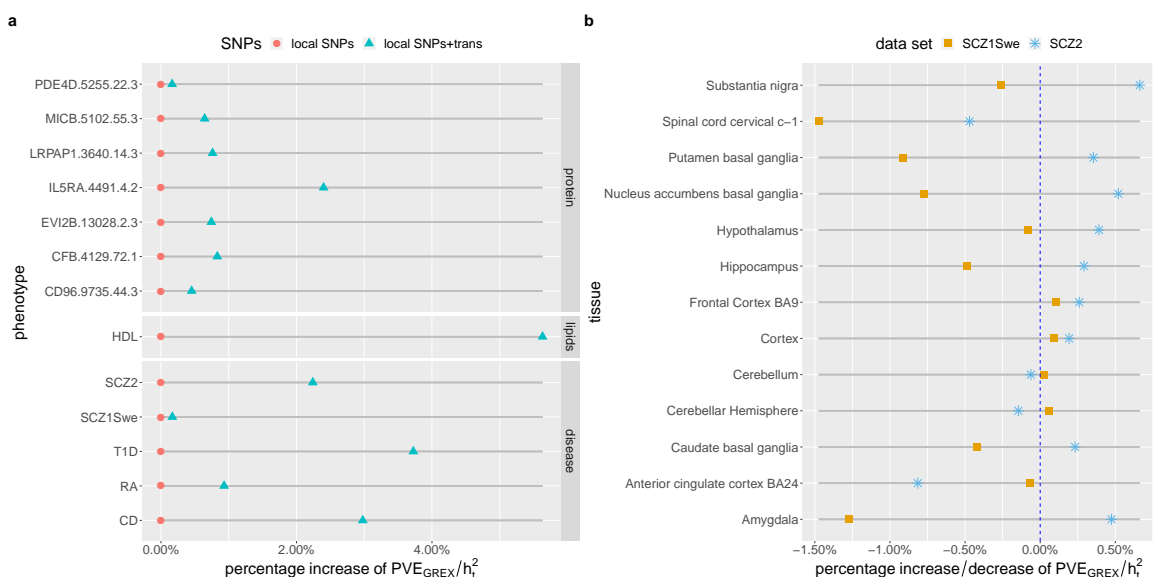


Figure 6: Comparison of PVE_{GREX}/h_t^2 's estimated with only local SNPs and those estimated with additional trans-eQTLs. (a) Estimated PVE_{GREX}/h_t^2 of 13 datasets in blood. All these datasets have significant PVE_{GREX}/h_t^2 in blood at 0.05 nominal level using only local SNPs. (b) Estimated PVE_{GREX}/h_t^2 of two largest SCZ data sets in 13 GTEx brain tissues. All these tissues have significant PVE_{GREX}/h_t^2 at 0.05 nominal level in both datasets using only local SNPs.

339 Discussion

340 In this work, we proposed a method, IGREX, for integrating GWAS and eQTL reference data
341 to quantify the GREX impact on phenotype. IGREX can be applied to both individual-level
342 and summary-level GWAS data, and was shown to achieve estimation accuracy even when the
343 eQTL effects are weak. IGREX can be used in many ways: it can inform the role of GREX
344 variation in various phenotypes and/or the role of GREX in known pathways; it can guide the
345 selection of eQTL reference data and suggest trait-relevant tissues/cell-types/contexts; and it
346 is generally applicable to the integration of GWAS with other omics data types to examine the
347 role of genetically-regulated multi-omics traits.

348 IGREX is closely related to several existing methods and here we briefly discuss the
349 connections and distinctions. By also integrating an eQTL reference and GWAS data, methods
350 including TWAS [17], PrediXcan [16], and the more general MetaXcan [21] aim to identify
351 specific trait-associated genes. In contrast, IGREX estimates the impact of genetically regulated
352 expression from a global perspective by quantifying the phenotypic variation that can be
353 attributed to the GREX component. Since both the TWAS-type of analyses and IGREX rely
354 on the shared GREX variation across eQTL and GWAS data, we argue that with the increasing
355 availability of eQTL resources in different populations, conditions and contexts, the proper
356 selection of eQTL reference panels via IGREX will greatly promote the chances of successes in
357 the subsequent TWAS-type of analyses.

358 There are also existing methods, such as RhoGE, designed for identifying and estimating
359 correlations between gene expression and complex traits. RhoGe provides an LDSC-based
360 approach for estimating PVE_{GREX} . Unlike IGREX, this method does not adjust for estimation
361 uncertainty. Consequently, it significantly underestimates the PVE_{GREX} when the eQTL effects
362 on expression levels are weak or moderate. In fact, RhoGE estimated the PVE_{GREX} for the
363 majority of 1,350 tissue-trait pairs to be almost negligible, with the first quantile, the median,
364 and the third quantile being 0.00125%, 0.162% and 0.616%, respectively [23]. In contrast, as
365 demonstrated via simulation studies, IGREX can accurately estimate PVE_{GREX} in various
366 scenarios by accounting for the estimation uncertainty.

367 Based on estimating PVE_{GREX} for a wide-array of tissue-trait pairs, we observed a stronger
368 impact of GREX on molecular intermediate traits and lipid traits in trait-relevant tissue types.
369 We also observed a relatively low PVE_{GREX} for complex traits in general. The big picture

370 suggests the attenuated impact on downstream phenotypes (e.g, height and SCZ), which is
371 consistent with the result from a pioneer study [57]. However, we noted that the PVE_{GREX}
372 estimates could be improved. A substantial amount of expression heritability is explained
373 by trans-acting genetic factors while current TWAS and IGREX analyses are mainly using
374 only cis-eQTLs. We explored the potential of incorporating trans-eQTLs in TWAS analysis
375 by re-estimating PVE_{GREX} for selected traits in blood tissues with significant trans-eQTLs
376 independently derived from the blood-based eQTLGen Consortium. We observed consistent
377 increases in PVE_{GREX} for blood-related traits. In contrast, such an increase was not observed in
378 the PVE_{GREX} estimates for other tissue types, again illustrating the importance of considering
379 trait-relevant tissue types/conditions in the TWAS-type of analyses. Additionally, we extended
380 the IGREX analysis to quantify the impact of genetically-regulated alternative splicing events
381 on selected traits. Our results suggested the potential for extending TWAS-type of analysis to
382 integrate reference multi-omics QTL data with GWAS in mapping novel disease/trait-associated
383 genes with mechanisms via other omics traits (such as splicing, methylation, protein, etc.).

384 A key assumption in applying IGREX or TWAS methods with a general-purpose eQTL
385 data as reference is the existence of steady-state component in GREX, i.e., the genetic effects
386 on gene expression β_g are shared across the eQTL reference and GWAS data. However, there
387 are many situations in which this assumption is violated. For example, it has been observed
388 that CAD-risk SNPs have a larger overlap with cis-eQTLs isolated from disease-relevant tissues
389 than those from GTEx tissues [46], implying the existence of a dynamic component. In the
390 presence of this dynamic component, the accuracy of \widehat{PVE}_{GREX} based on GTEx is reduced.
391 In those cases, we suggest exploring other trait-relevant or condition-specific eQTL reference
392 panels using IGREX for a better understanding of the role of GREX and before conducting
393 TWAS analysis.

394 Methods

395 **The IGREX-i for individual-level GWAS data.** First, let $\mathcal{D}_r = \{\mathbf{Y}, \mathbf{X}_r\}$ denote the
396 reference data set from an eQTL study, where $\mathbf{Y} \in \mathbb{R}^{n_r \times G}$ is the gene expression matrix,
397 $\mathbf{X}_r \in \mathbb{R}^{n_r \times M}$ is the genotype matrix, n_r is the sample size of the eQTL study, G is the
398 number of genes and M is the number of single-nucleotide polymorphisms (SNPs). Suppose
399 we have individual-level GWAS data $\mathcal{D}_i = \{\mathbf{t}, \mathbf{X}\}$ comprised of phenotype vector $\mathbf{t} \in \mathbb{R}^n$ and

400 genotype matrix $\mathbf{X} \in \mathbb{R}^{n \times M}$, where n is the GWAS sample size. For $g = 1, \dots, G$, we let the
 401 g -th gene expression vector $\mathbf{y}_g \in \mathbb{R}^{n_r}$ denote the corresponding column of \mathbf{Y} , local genotype
 402 matrices $\mathbf{X}_{r,g} \in \mathbb{R}^{n_r \times M_g}$ and $\mathbf{X}_g \in \mathbb{R}^{n \times M_g}$ denote the corresponding M_g columns in \mathbf{X}_r and
 403 \mathbf{X} , respectively, where M_g is the number of local SNPs for g -th gene. To make the notation
 404 uncluttered, we further assume that $\mathbf{X}_{r,g}$ and \mathbf{X}_g have been standardized and both \mathbf{y}_g and \mathbf{t}
 405 have been properly adjusted for covariates. The complete model that accounts for covariates is
 406 described in the Supplementary Materials. Now, we consider linear model (1) that associates
 407 the gene expression vector \mathbf{y}_g to $\mathbf{X}_{r,g}$:

$$\mathbf{y}_g = \mathbf{X}_{r,g}\boldsymbol{\beta}_g + \mathbf{e}_{r,g},$$

408 where $\boldsymbol{\beta}_g$ is an $M_g \times 1$ vector of genetic effects on the gene expression levels, $\mathbf{e}_{r,g} \sim \mathcal{N}(0, \sigma_{r,g}^2 \mathbf{I}_{n_r})$
 409 is a vector of independent noise and \mathbf{I} is the identity matrix with the subscript being its
 410 size. Assuming that there is a steady-state component in gene expression regulated by genetic
 411 variants, individuals in \mathcal{D}_r and \mathcal{D}_i share the same $\boldsymbol{\beta}_g$. Hence, the GREX in \mathcal{D}_i can be evaluated
 412 by $\mathbf{X}_g\boldsymbol{\beta}_g$. Then, we assume that the phenotype \mathbf{t} can be decomposed into two parts, i.e., the
 413 genetic effects via GREX and the genetic effects through alternative pathways, as in model (2):

$$\mathbf{t} = \sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\beta}_g + \mathbf{X}\boldsymbol{\gamma} + \boldsymbol{\epsilon},$$

414 where α_g is the effect of $\mathbf{X}_g\boldsymbol{\beta}_g$ on \mathbf{t} , $\boldsymbol{\gamma}$ is an $n \times 1$ vector of alternative genetic effects and
 415 $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \sigma_{\epsilon}^2 \mathbf{I}_n)$ is a vector of independent errors. The term $\sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\beta}_g$ can be viewed as
 416 the overall impact of GREX on the phenotype and $\mathbf{X}\boldsymbol{\gamma}$ represents the alternative impact of
 417 genotypes on the phenotype. Given a genotype vector $\mathbf{x} \in \mathbb{R}^M$ and a phenotype $t \in \mathbb{R}$, the
 418 impact of GREX can be quantified by the proportion of variance explained by the GREX
 419 component:

$$\text{PVE}_{\text{GREX}} = \frac{\text{Var}(\sum_{g=1}^G \alpha_g \mathbf{x}_g^T \boldsymbol{\beta}_g)}{\text{Var}(t)}, \quad (3)$$

420 where \mathbf{x}_g is the subvector of genotype \mathbf{x} corresponding to the g -th gene.

421 To estimate PVE_{GREX} , we introduce the following probabilistic structure for the effects in
 422 model (1) and (2):

$$\boldsymbol{\beta}_g \sim \mathcal{N}(\mathbf{0}, \sigma_{\beta_g}^2 \mathbf{I}_{M_g}), \quad \alpha_g \sim \mathcal{N}(\mathbf{0}, \sigma_{\alpha}^2), \quad \boldsymbol{\gamma} \sim \mathcal{N}(\mathbf{0}, \sigma_{\gamma}^2 \mathbf{I}_M), \quad (4)$$

423 which is motivated by a recent theoretical justification [58] for heritability estimation on a
 424 mis-specified linear mixed model (LMM). This prior specification in (4) provides a great com-

425 putational advantage as well as a stable performance for IGREX under model mis-specification,
426 as demonstrated in the simulation studies.

427 The proposed method for individual-level GWAS data, IGREX-i, provides a two-stage
428 framework for estimating PVE_{GREX} . In the first stage, we estimate the parameters $\sigma_{\beta_g}^2$ and $\sigma_{r,g}^2$
429 in model (1) via a fast expectation-maximization (EM)-type algorithm, the parameter-expanded
430 EM (PX-EM) algorithm [59]. Based on the estimates, denoted as $\hat{\sigma}_{\beta_g}^2$ and $\hat{\sigma}_{r,g}^2$, the posterior
431 distribution of β_g is given by

$$\beta_g | \mathbf{y}_g, \mathbf{X}_{r,g} \sim \mathcal{N}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g), \text{ where } \boldsymbol{\Sigma}_g = \left(\frac{1}{\hat{\sigma}_{r,g}^2} \mathbf{X}_{r,g}^T \mathbf{X}_{r,g} + \frac{1}{\hat{\sigma}_{\beta_g}^2} \mathbf{I}_{M_g} \right)^{-1}, \boldsymbol{\mu}_g = \boldsymbol{\Sigma}_g \frac{1}{\hat{\sigma}_{r,g}^2} \mathbf{X}_{r,g}^T \mathbf{y}_g. \quad (5)$$

432 In the second stage, we treat the posterior distribution obtained in (5) as the prior distribution
433 of β_g in model (2). This substitution naturally accounts for the uncertainty in estimating
434 β_g which has been captured by $\boldsymbol{\Sigma}_g$. To evaluate the covariance of \mathbf{t} , we first note that
435 $\mathbb{E}(\mathbf{t} | \boldsymbol{\alpha}) = \sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\mu}_g$ and $\text{Cov}(\mathbf{t} | \boldsymbol{\alpha}) = \sum_{g=1}^G \alpha_g^2 \mathbf{X}_g \boldsymbol{\Sigma}_g \mathbf{X}_g^T + \sigma_\gamma^2 \mathbf{X} \mathbf{X}^T + \sigma_\epsilon^2 \mathbf{I}_n$; then, using the
436 law of total expectation and total variance, we obtain $\mathbb{E}(\mathbf{t}) = \mathbb{E}(\mathbb{E}(\mathbf{t} | \boldsymbol{\alpha})) = \mathbf{0}$ and

$$\text{Cov}(\mathbf{t}) = \text{Cov}(\mathbb{E}(\mathbf{t} | \boldsymbol{\alpha})) + \mathbb{E}(\text{Cov}(\mathbf{t} | \boldsymbol{\alpha})) = \sum_{g=1}^G \sigma_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T + \sigma_\gamma^2 \mathbf{X} \mathbf{X}^T + \sigma_\epsilon^2 \mathbf{I}_n, \quad (6)$$

437 respectively. By observing the form of (6), it is clear that the i -th diagonal element of
438 $\sum_{g=1}^G \sigma_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T$ and $\sigma_\gamma^2 \mathbf{X} \mathbf{X}^T$ represents the variance explained by GREX and
439 alternative genetic effects, respectively. Therefore, the PVE_{GREX} defined in (3) can be estimated
440 by

$$\widehat{\text{PVE}}_{\text{GREX}} = \frac{\text{tr}(\sum_{g=1}^G \hat{\sigma}_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T)}{\text{tr}(\sum_{g=1}^G \hat{\sigma}_\alpha^2 \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T + \hat{\sigma}_\gamma^2 \mathbf{X} \mathbf{X}^T + \hat{\sigma}_\epsilon^2 \mathbf{I}_n)}, \quad (7)$$

441 where $\hat{\sigma}_\alpha^2$, $\hat{\sigma}_\gamma^2$ and $\hat{\sigma}_\epsilon^2$ are the estimated values of σ_α^2 , σ_γ^2 and σ_ϵ^2 , respectively.

442 IGREX-i provides two methods for estimating the parameters and $\widehat{\text{PVE}}_{\text{GREX}}$ in the sec-
443 ond stage. Let $\boldsymbol{\psi} = [\sigma_\alpha^2, \sigma_\gamma^2, \sigma_\epsilon^2]^T$ be the vector of parameters to be estimated, $\mathbf{K}_\alpha =$
444 $\sum_{g=1}^G \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T$ and $\mathbf{K}_\gamma = \mathbf{X} \mathbf{X}^T$. The first method is based on MoM, which minimizes
445 the distance between the second moment of \mathbf{t} at the population level and that at the sample
446 level $f(\boldsymbol{\psi}) = \|\mathbf{t} \mathbf{t}^T - (\sigma_\alpha^2 \mathbf{K}_\alpha + \sigma_\gamma^2 \mathbf{K}_\gamma + \sigma_\epsilon^2 \mathbf{I}_n)\|^2$. By setting $\frac{\partial f(\boldsymbol{\psi})}{\partial \sigma_\alpha^2} = \frac{\partial f(\boldsymbol{\psi})}{\partial \sigma_\gamma^2} = \frac{\partial f(\boldsymbol{\psi})}{\partial \sigma_\epsilon^2} = 0$, we obtain
447 the estimating equation

$$\mathbf{S} \boldsymbol{\psi} = \mathbf{q}, \quad (8)$$

$$\text{with } \mathbf{S} = \begin{bmatrix} \text{tr}(\mathbf{K}_\alpha^2) & \text{tr}(\mathbf{K}_\alpha \mathbf{K}_\gamma) & \text{tr}(\mathbf{K}_\alpha) \\ \text{tr}(\mathbf{K}_\alpha \mathbf{K}_\gamma) & \text{tr}(\mathbf{K}_\gamma^2) & \text{tr}(\mathbf{K}_\gamma) \\ \text{tr}(\mathbf{K}_\alpha) & \text{tr}(\mathbf{K}_\gamma) & n \end{bmatrix}, \boldsymbol{\psi} = \begin{bmatrix} \sigma_\alpha^2 \\ \sigma_\gamma^2 \\ \sigma_\epsilon^2 \end{bmatrix}, \mathbf{q} = \begin{bmatrix} \mathbf{t}^T \mathbf{K}_\alpha \mathbf{t} \\ \mathbf{t}^T \mathbf{K}_\gamma \mathbf{t} \\ \mathbf{t}^T \mathbf{t} \end{bmatrix}.$$

449 The solution of Equation (8) is given by $\hat{\boldsymbol{\psi}} = \mathbf{S}^{-1}\mathbf{q}$. And the variance-covariance matrix of $\hat{\boldsymbol{\psi}}$
 450 is given by $\text{Cov}(\hat{\boldsymbol{\psi}}) = \mathbf{S}^{-1}\text{Cov}(\mathbf{q})\mathbf{S}^{-1}$ using the sandwich estimator. Then, the standard error
 451 of $\widehat{\text{PVE}}_{\text{GREX}}$ can be obtained by the delta method (see Supplementary Materials). The second
 452 method applies the restricted maximum likelihood (REML) by further assuming the normal
 453 distribution of \mathbf{t} : $\mathbf{t} \sim \mathcal{N}(\mathbf{0}, \sigma_{\alpha}^2\mathbf{K}_{\alpha} + \sigma_{\gamma}^2\mathbf{K}_{\gamma} + \sigma_{\epsilon}^2\mathbf{I}_n)$. The variance components are estimated by
 454 the Minorization-Maximization (MM) algorithm [60].

455 **The IGREX-s for summary-level GWAS data.** The special formulation of method of
 456 moments allows IGREX to be extended (IGREX-s) to handle summary-level GWAS data (i.e.
 457 z -scores) when the individual-level data \mathcal{D}_i is not accessible. Suppose we only have the z -scores
 458 from summary-level GWAS data $\{z_j\}_{j=1}^M$ generated from \mathcal{D}_i . The definition of the z -score is
 459 $z_j = \frac{(\mathbf{x}_j^T\mathbf{x}_j)^{-1}\mathbf{x}_j^T\mathbf{t}}{\sqrt{\hat{\sigma}_j^2(\mathbf{x}_j^T\mathbf{x}_j)^{-1}}}$, where \mathbf{x}_j is the j -th column of \mathbf{X} and $\hat{\sigma}_j^2$ is the estimate of residual variance
 460 by regressing \mathbf{x}_j on \mathbf{t} . By assuming that z -scores are calculated from a standardized genotype
 461 matrix \mathbf{X} , we have $\mathbf{x}_j^T\mathbf{x}_j = n$. Besides, the polygenicity assumption implies that $\hat{\sigma}_j^2 \approx \hat{\sigma}_t^2$, where
 462 $\hat{\sigma}_t^2$ is the estimate of $\text{Var}(t)$. Hence, we have

$$z_j \approx \frac{\mathbf{x}_j^T\mathbf{t}}{\sqrt{n\hat{\sigma}_t^2}}, \quad (9)$$

463 and $\widehat{\text{PVE}}_{\text{GREX}}$ defined in (3) can be estimated by

$$\widehat{\text{PVE}}_{\text{GREX}} = \frac{\frac{1}{n}\text{tr}(\sum_{g=1}^G \hat{\sigma}_{\alpha}^2\mathbf{X}_g(\boldsymbol{\mu}_g\boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g)\mathbf{X}_g^T)}{\hat{\sigma}_t^2} \approx \frac{\hat{\sigma}_{\alpha}^2}{\hat{\sigma}_t^2}\text{tr}\left(\sum_{g=1}^G(\boldsymbol{\mu}_g\boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g)\hat{\mathbf{R}}_g\right), \quad (10)$$

464 where $\hat{\mathbf{R}}_g = \tilde{\mathbf{X}}_g^T\tilde{\mathbf{X}}_g/(m-1)$ is the estimated LD matrix associated with the g -th gene and $\tilde{\mathbf{X}}_g$
 465 is the corresponding columns of a reference genotype matrix $\tilde{\mathbf{X}}$. In practice, $\tilde{\mathbf{X}} \in \mathbb{R}^{m \times M}$ can be
 466 the genotype matrix either from the GTEx Project or the 1000 Genomes Project. Now, we
 467 consider MoM in the estimating equation (8) to obtain $\frac{\hat{\sigma}_{\alpha}^2}{\hat{\sigma}_t^2}$. By eliminating σ_{ϵ}^2 and dividing
 468 both sides by n^2 , we have

$$\begin{bmatrix} \frac{\text{tr}(\mathbf{K}_{\alpha}^2) - \frac{\text{tr}^2(\mathbf{K}_{\alpha})}{n}}{n^2} & \frac{\text{tr}(\mathbf{K}_{\alpha}\mathbf{K}_{\gamma}) - \frac{\text{tr}(\mathbf{K}_{\alpha})\text{tr}(\mathbf{K}_{\gamma})}{n}}{n^2} \\ \frac{\text{tr}(\mathbf{K}_{\alpha}\mathbf{K}_{\gamma}) - \frac{\text{tr}(\mathbf{K}_{\alpha})\text{tr}(\mathbf{K}_{\gamma})}{n}}{n^2} & \frac{\text{tr}(\mathbf{K}_{\gamma}^2) - \frac{\text{tr}^2(\mathbf{K}_{\gamma})}{n}}{n^2} \end{bmatrix} \begin{bmatrix} \sigma_{\alpha}^2 \\ \sigma_{\gamma}^2 \end{bmatrix} = \begin{bmatrix} \frac{1}{n^2}\mathbf{t}^T\mathbf{K}_{\alpha}\mathbf{t} - \frac{\text{tr}(\mathbf{K}_{\alpha})}{n^3}\mathbf{t}^T\mathbf{t} \\ \frac{1}{n^2}\mathbf{t}^T\mathbf{K}_{\gamma}\mathbf{t} - \frac{\text{tr}(\mathbf{K}_{\gamma})}{n^3}\mathbf{t}^T\mathbf{t} \end{bmatrix}. \quad (11)$$

469 The terms on the left hand side do not involve \mathbf{t} and thus can be approximated using $\tilde{\mathbf{X}}$
 470 [37]. For example, $\frac{\text{tr}(\mathbf{K}_{\alpha}^2) - \frac{\text{tr}^2(\mathbf{K}_{\alpha})}{n}}{n^2}$ can be well approximated by $\frac{\text{tr}(\tilde{\mathbf{K}}_{\alpha}^2) - \frac{\text{tr}^2(\tilde{\mathbf{K}}_{\alpha})}{m}}{m^2}$, where $\tilde{\mathbf{K}}_{\alpha} =$
 471 $\sum_{g=1}^G \tilde{\mathbf{X}}_g(\boldsymbol{\mu}_g\boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g)\tilde{\mathbf{X}}_g^T$. Other terms on the left hand side can be approximated in the same
 472 way. For the right hand side, each term can be approximated using $\hat{\mathbf{R}}_g$ and z -scores from

473 approximation (9): $\mathbf{t}^T \mathbf{K}_\alpha \mathbf{t} \approx n \hat{\sigma}_t^2 \sum_g \mathbf{z}_g^T (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{z}_g$, where $\mathbf{z}_g \in \mathbb{R}^{M_g}$ is the vector of z -scores
 474 corresponding to the g -th gene; $\frac{\text{tr}(\mathbf{K}_\alpha)}{n} \mathbf{t}^T \mathbf{t} \approx n \hat{\sigma}_t^2 \text{tr}(\sum_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \hat{\mathbf{R}}_g)$; $\mathbf{t}^T \mathbf{K}_\gamma \mathbf{t} \approx n \hat{\sigma}_t^2 \sum_{j=1}^M z_j^2$;
 475 and $\frac{\text{tr}(\mathbf{K}_\gamma)}{n} \mathbf{t}^T \mathbf{t} \approx n \hat{\sigma}_t^2$. With these approximations, Equation (11) becomes

$$\begin{bmatrix} \frac{\text{tr}(\tilde{\mathbf{K}}_\alpha^2) - \frac{\text{tr}^2(\tilde{\mathbf{K}}_\alpha)}{m}}{m^2} & \frac{\text{tr}(\tilde{\mathbf{K}}_\alpha \tilde{\mathbf{K}}_\gamma) - \frac{\text{tr}(\tilde{\mathbf{K}}_\alpha) \text{tr}(\tilde{\mathbf{K}}_\gamma)}{m}}{m^2} \\ \frac{\text{tr}(\tilde{\mathbf{K}}_\alpha \tilde{\mathbf{K}}_\gamma) - \frac{\text{tr}(\tilde{\mathbf{K}}_\alpha) \text{tr}(\tilde{\mathbf{K}}_\gamma)}{m}}{m^2} & \frac{\text{tr}(\tilde{\mathbf{K}}_\gamma^2) - \frac{\text{tr}^2(\tilde{\mathbf{K}}_\gamma)}{m}}{m^2} \end{bmatrix} \begin{bmatrix} \frac{\hat{\sigma}_\alpha^2}{\hat{\sigma}_t^2} \\ \frac{\hat{\sigma}_\gamma^2}{\hat{\sigma}_t^2} \end{bmatrix} = \begin{bmatrix} \frac{\sum_g \mathbf{z}_g^T (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{z}_g - \text{tr}(\sum_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \hat{\mathbf{R}}_g)}{\sum_{j=1}^M \frac{z_j^2 - 1}{n}} \end{bmatrix}.$$

476 Then, $\frac{\hat{\sigma}_\alpha^2}{\hat{\sigma}_t^2}$ can be obtained by solving this equation. Plugging this estimate into Equation (10)
 477 gives the $\widehat{\text{PVE}}_{\text{GREX}}$. The standard errors of $\widehat{\text{PVE}}_{\text{GREX}}$ can be estimated by the block jackknife
 478 method [61] (Supplementary Materials).

479 IGREX can incorporate fixed effects to adjust for possible confounding factors, such as
 480 population structure. Details are provided in the Supplementary Note.

481 **GTEEx eQTL dataset.** We used the gene expression data from the V7 release of GTEx
 482 Consortium as our reference dataset. We analyzed the 48 tissues with number of genotyped
 483 samples ≥ 70 , which are collected from 620 donors with total sample size 10,294. The
 484 sample size of each tissue ranges from 80 to 491 (details provided in Supplementary Table
 485 4). We set the mappability cutoff at 0.9 to filter gene expression measures with lower quality,
 486 leaving 16,333 – 27,378 genes to be included in our analysis. Based on the third phase of the
 487 International HapMap project phase 3 (HapMap3), 1,189,556 SNPs were included from the
 488 GTEx genotyped data for analysis. For each gene, we included only the SNPs within 500kb of
 489 the transcription start and end of each protein coding genes. In real data analysis, we used the
 490 covariates provided by the GTEx consortium, including genotype principal components (PCs),
 491 Probabilistic Estimation of Expression Residuals (PEER) factors, genotyping platform and sex
 492 (as described in <https://gtexportal.org/home/documentationPage>).

493 Additionally, the GTEx genotype data was used as an LD reference panel when applying
 494 IGREX-s to GWAS summary statistics. In this application, we used top 5 PCs as covariates.

495 **Individual level GWAS datasets.** The NFBC dataset is comprised of 5,402 individuals
 496 with ten continuous phenotypes related to cardiovascular disease including Glucose, body mass
 497 index (BMI), C-reactive protein (CRP), insulin, high-density lipoprotein cholesterol (HDL),
 498 low-density lipoprotein cholesterol (LDL), triglycerides (TG), total cholesterol (TC), diastolic
 499 blood pressure (DiaBP) and systolic blood pressure (SysBP). There are 364,590 genotyped
 500 SNPs in this dataset. We first excluded the individuals whose reported sex differed from their

501 sex determined from the X chromosome. We then excluded the SNPs with minor allele frequency
502 less than 1%, with missing values in more than 1% of the individuals or with Hardy-Weinberg
503 equilibrium (HWE) p -value below 0.0001. This quality control process yields 5, 123 individuals
504 with 319, 147 SNPs in NFBC dataset for our analysis. We evaluated the genetic relatedness
505 matrix (GRM) using the processed genotype data and selected the top 20 PCs as covariates in
506 the study.

507 The WTCCC dataset contains seven disease phenotypes including bipolar disorder (BD),
508 coronary artery disease (CAD), Crohn's disease (CD), hypertension (HT), rheumatoid arthritis
509 (RA), type 1 diabetes (T1D) and type 2 diabetes (T2D). It includes $\sim 2,000$ cases per phenotype
510 and 3,004 controls with 490,032 genotyped SNPs. We first removed the individuals with
511 genotyping rate less than 5%. Then we excluded the SNPs satisfying at least one of the following:
512 minor allele frequency less than 5%; genotypes missing in more than 1% samples; HWE p -value
513 is below 0.001. We also removed the individuals with estimated genetic correlation larger than
514 2.5%. After quality control, around 4,700 individuals with 300,000 SNPs were retained for our
515 analysis (See Supplementary Table 1). Based on the obtained data, we calculated the GRM
516 and extracted top 20 PCs as covariates to be included in our analysis.

517 **GWAS summary statistics.** We analyzed ten summary level GWAS datasets: human plasma
518 pQTL data [38], circulating metabolite data [40], four schizophrenia datasets [41, 42, 43, 44],
519 two independent height datasets [62] and European ancestry of BMI datasets with sample age
520 ≤ 50 separated by men and women [63]. The SNPs with missing information (i.e. chromosome,
521 minor allele, allele frequency) were first removed. Following the practice of LDSC [30], we
522 checked the χ^2 statistic of each SNP and excluded those with extreme values ($\chi^2 > 80$) to
523 prevent the outliers that may unduly affect the results. The detailed information is provided in
524 Supplementary Table 2. After pre-processing, the remaining SNPs were further matched with
525 reference data, and this step is automatically conducted in our IGREX software.

526 **The eQTLGen summary data.** We used the trans-eQTLs in blood provided by the
527 eQTLGen Consortium [15]. The trans-eQTL analysis were restricted to known complex trait-
528 associated SNPs. The significant trans-eQTLs were identified by controlling the FDR at 0.05.
529 There were 5,4786 gene-SNP pairs composed of 6,298 genes and 3,853 SNPs. The remaining
530 pairs after matching with both reference and GWAS datasets are summarized in Supplementary
531 Table 5.

532 **Data availability.** The GTEx gene expression data was downloaded from GTEx Consor-
533 tium website <https://gtexportal.org/home/datasets>. The GTEx genotype data can be
534 accessed from dbGAP with accession number phs000424.v7.p2. The HapMap3 genotype data is
535 available at <ftp://ftp.ncbi.nlm.nih.gov/hapmap/>. The NFBC study was downloaded from
536 dbGAP using accession number phs000276.v1.p1. The WTCCC data was obtained from its
537 consortium website https://www.wtccc.org.uk/info/access_to_data_samples.html. The
538 GWAS summary statistics can be accessed using the links provided in Supplementary Table 2.
539 The eQTLGen data can be downloaded from <http://www.eqtlgen.org>.

540 **Software.** The R software package IGREX is publicly available on GitHub repository: <https://github.com/mxcai/IGREX>.

542 **Acknowledgements.** We thank Mr. Kevin J. Gleason for proof-reading the work. This work
543 was supported in part by the National Science Funding of China [61501389]; the Hong Kong
544 Research Grant Council [12316116, 12301417 and 16307818]; The Hong Kong University of
545 Science and Technology [startup grant R9405 and IGN17SC02]; Duke-NUS Medical School
546 WBS [R-913-200-098-263]; Ministry of Education, Singapore. AcRF Tier 2 [MOE2016-T2-2-
547 029, MOE2018-T2-1-046 and MOE2018-T2-2-006]. LSC was independently supported by the
548 National Institutes of Health (NIH) grant R01GM108711. The computational work for this
549 article was (fully or partially) performed on resources of the National Supercomputing Centre,
550 Singapore (<https://www.nscg.sg>).

551 References

- 552 [1] Matthew T. Maurano, Richard Humbert, Eric Rynes, Robert E. Thurman, Eric Haugen,
553 Hao Wang, Alex P. Reynolds, Richard Sandstrom, Hongzhu Qu, Jennifer Brody, Anthony
554 Shafer, Fidencio Neri, Kristen Lee, Tanya Kutuyavin, Sandra Stehling-Sun, Audra K.
555 Johnson, Theresa K. Canfield, Erika Giste, Morgan Diegel, Daniel Bates, R. Scott Hansen,
556 Shane Neph, Peter J. Sabo, Shelly Heimfeld, Antony Raubitschek, Steven Ziegler, Chris
557 Cotsapas, Nona Sotoodehnia, Ian Glass, Shamil R. Sunyaev, Rajinder Kaul, and John A.
558 Stamatoyannopoulos. Systematic localization of common disease-associated variation in
559 regulatory DNA. *Science*, 337(6099):1190–1195, 2012.
- 560 [2] William Cookson, Liming Liang, Gonçalo Abecasis, Miriam Moffatt, and Mark Lathrop.
561 Mapping complex disease traits with global gene expression. *Nature Reviews Genetics*,
562 10(3):184, 2009.
- 563 [3] Mark M Pomerantz, Nasim Ahmadiyeh, LI Jia, Paula Herman, Michael P Verzi, Har-
564 shavardhan Doddapaneni, Christine A Beckwith, Jennifer A Chan, Adam Hills, Matt
565 Davis, et al. The 8q24 cancer risk variant rs6983267 shows long-range interaction with
566 myc in colorectal cancer. *Nature genetics*, 41(8):882, 2009.
- 567 [4] Kiran Musunuru, Alanna Strong, Maria Frank-Kamenetsky, Noemi E Lee, Tim Ahfeldt,
568 Katherine V Sachs, Xiaoyu Li, Hui Li, Nicolas Kuperwasser, Vera M Ruda, et al. From non-
569 coding variant to phenotype via sort1 at the 1p13 cholesterol locus. *Nature*, 466(7307):714,
570 2010.
- 571 [5] Olivier Harismendy, Dimple Notani, Xiaoyuan Song, Nazli G Rahim, Bogdan Tanasa,
572 Nathaniel Heintzman, Bing Ren, Xiang-Dong Fu, Eric J Topol, Michael G Rosenfeld, et al.
573 9p21 DNA variants associated with coronary artery disease impair interferon- γ signalling
574 response. *Nature*, 470(7333):264, 2011.
- 575 [6] Dan L Nicolae, Eric Gamazon, Wei Zhang, Shiwei Duan, M Eileen Dolan, and Nancy J
576 Cox. Trait-associated SNPs are more likely to be eqtls: annotation to enhance discovery
577 from gwas. *PLoS genetics*, 6(4):e1000888, 2010.
- 578 [7] Lucia A. Hindorff, Praveen Sethupathy, Heather A. Junkins, Erin M. Ramos, Jayashri P.

- 579 Mehta, Francis S. Collins, and Teri A. Manolio. Potential etiologic and functional implica-
580 tions of genome-wide association loci for human diseases and traits. *Proceedings of the*
581 *National Academy of Sciences*, 106(23):9362–9367, 2009.
- 582 [8] Frank W Albert and Leonid Kruglyak. The role of regulatory variation in complex traits
583 and disease. *Nature Reviews Genetics*, 16(4):197, 2015.
- 584 [9] Kevin J. Gleason, Fan Yang, Brandon L. Pierce, Xin He, and Lin S. Chen. Primo:
585 integration of multiple GWAS and omics QTL summary statistics for elucidation of
586 molecular mechanisms of trait-associated snps and detection of pleiotropy in complex
587 traits. *bioRxiv*, 579581, 2019.
- 588 [10] Eric R Gamazon, Ayellet V Segrè, Martijn van de Bunt, Xiaoquan Wen, Hualin S Xi,
589 Farhad Hormozdiari, Halit Ongen, Anuar Konkashbaev, Eske M Derks, François Aguet,
590 et al. Using an atlas of gene regulation across 44 human tissues to inform complex
591 disease-and trait-associated variation. *Nature genetics*, 50(7):956, 2018.
- 592 [11] GTEx Consortium et al. Genetic effects on gene expression across human tissues. *Nature*,
593 550(7675):204, 2017.
- 594 [12] Luke R Lloyd-Jones, Alexander Holloway, Allan McRae, Jian Yang, Kerrin Small, Jing
595 Zhao, Biao Zeng, Andrew Bakshi, Andres Metspalu, Manolis Dermitzakis, et al. The
596 genetic architecture of gene expression in peripheral blood. *The American Journal of*
597 *Human Genetics*, 100(2):228–237, 2017.
- 598 [13] Harm-Jan Westra, Marjolein J Peters, Tõnu Esko, Hanieh Yaghootkar, Claudia Schurmann,
599 Johannes Kettunen, Mark W Christiansen, Benjamin P Fairfax, Katharina Schramm,
600 Joseph E Powell, et al. Systematic identification of trans eQTLs as putative drivers of
601 known disease associations. *Nature genetics*, 45(10):1238, 2013.
- 602 [14] Ting Qi, Yang Wu, Jian Zeng, Futao Zhang, Angli Xue, Longda Jiang, Zhihong Zhu,
603 Kathryn Kemper, Loic Yengo, Zhili Zheng, et al. Identifying gene targets for brain-related
604 traits using transcriptomic and methylomic data from blood. *Nature communications*,
605 9(1):2282, 2018.
- 606 [15] Urmo Võsa, Annique Claringbould, Harm-Jan Westra, Marc Jan Bonder, Patrick Deelen,
607 Biao Zeng, Holger Kirsten, Ashis Saha, Roman Kreuzhuber, Silva Kasela, et al. Unraveling

- 608 the polygenic architecture of complex traits using blood eQTL meta-analysis. *bioRxiv*,
609 page 447367, 2018.
- 610 [16] Eric R Gamazon, Heather E Wheeler, Kanaan P Shah, Sahar V Mozaffari, Keston Aquino-
611 Michaels, Robert J Carroll, Anne E Eyler, Joshua C Denny, Dan L Nicolae, Nancy J Cox,
612 et al. A gene-based association method for mapping traits using reference transcriptome
613 data. *Nature genetics*, 47(9):1091, 2015.
- 614 [17] Alexander Gusev, Arthur Ko, Huwenbo Shi, Gaurav Bhatia, Wonil Chung, Brenda WJH
615 Penninx, Rick Jansen, Eco JC De Geus, Dorret I Boomsma, Fred A Wright, et al.
616 Integrative approaches for large-scale transcriptome-wide association studies. *Nature*
617 *genetics*, 48(3):245, 2016.
- 618 [18] Zhihong Zhu, Futao Zhang, Han Hu, Andrew Bakshi, Matthew R Robinson, Joseph E
619 Powell, Grant W Montgomery, Michael E Goddard, Naomi R Wray, Peter M Visscher,
620 et al. Integration of summary data from GWAS and eQTL studies predicts complex trait
621 gene targets. *Nature genetics*, 48(5):481, 2016.
- 622 [19] Nicholas Mancuso, Malika K Freund, Ruth Johnson, Huwenbo Shi, Gleb Kichaev, Alexan-
623 der Gusev, and Bogdan Pasaniuc. Probabilistic fine-mapping of transcriptome-wide
624 association studies. *Nature Genetics*, 51(4):675, 2019.
- 625 [20] Kunal Bhutani, Abhishek Sarkar, Yongjin Park, Manolis Kellis, and Nicholas J Schork.
626 Modeling prediction error improves power of transcriptome-wide association studies.
627 *bioRxiv*, page 108316, 2017.
- 628 [21] Alvaro N Barbeira, Scott P Dickinson, Rodrigo Bonazzola, Jiamao Zheng, Heather E
629 Wheeler, Jason M Torres, Eric S Torstenson, Kanaan P Shah, Tzintzuni Garcia, Todd L
630 Edwards, et al. Exploring the phenotypic consequences of tissue specific gene expression
631 variation inferred from GWAS summary statistics. *Nature communications*, 9(1):1825,
632 2018.
- 633 [22] Can Yang, Xiang Wan, Xinyi Lin, Mengjie Chen, Xiang Zhou, and Jin Liu. CoMM:
634 a collaborative mixed model to dissecting genetic contributions to complex traits by
635 leveraging regulatory information. *Bioinformatics*, 35(1644-1652):865, 2018.

- 636 [23] Nicholas Mancuso, Huwenbo Shi, Pagé Goddard, Gleb Kichaev, Alexander Gusev, and
637 Bogdan Pasaniuc. Integrating gene expression with summary association statistics to
638 identify genes associated with 30 complex traits. *The American Journal of Human Genetics*,
639 100(3):473–487, 2017.
- 640 [24] Luke J O’Connor, Alexander Gusev, Xuanyao Liu, Po-Ru Loh, Hilary K Finucane, and
641 Alkes L Price. Estimating the proportion of disease heritability mediated by gene expression
642 levels. *BioRxiv*, page 118018, 2017.
- 643 [25] Daniel J Schaid, Wenan Chen, and Nicholas B Larson. From genome-wide associations to
644 candidate causal variants by statistical fine-mapping. *Nature Reviews Genetics*, 19(8):491,
645 2018.
- 646 [26] Yiming Hu, Mo Li, Qiongshi Lu, Haoyi Weng, Jiawei Wang, Seyedeh M. Zekavat, Zhaolong
647 Yu, Boyang Li, Jianlei Gu, Sydney Muchnik, Yu Shi, Brian W. Kunkle, Shubhabrata
648 Mukherjee, Pradeep Natarajan, Adam Naj, Amanda Kuzma, Yi Zhao, Paul K. Crane, Hui
649 Lu, Hongyu Zhao, and Alzheimer’s Disease Genetics Consortium. A statistical framework
650 for cross-tissue transcriptome-wide association analysis. *Nature Genetics*, 51(3):568–576,
651 2019.
- 652 [27] Michael Wainberg, Nasa Sinnott-Armstrong, Nicholas Mancuso, Alvaro N. Barbeira,
653 David A. Knowles, David Golan, Raili Ermel, Arno Ruusalepp, Thomas Quertermous,
654 Ke Hao, Johan L. M. Björkegren, Hae Kyung Im, Bogdan Pasaniuc, Manuel A. Rivas, and
655 Anshul Kundaje. Opportunities and challenges for transcriptome-wide association studies.
656 *Nature Genetics*, 51(4):592–599, 2019.
- 657 [28] Kanix Wang, Hallie Gaitsch, Hoifung Poon, Nancy J Cox, and Andrey Rzhetsky. Clas-
658 sification of common human diseases derived from shared genetic and environmental
659 determinants. *Nature genetics*, 49(9):1319, 2017.
- 660 [29] Chirag M Lakhani, Braden T Tierney, Arjun K Manrai, Jian Yang, Peter M Visscher, and
661 Chirag J Patel. Repurposing large health insurance claims data to estimate genetic and
662 environmental contributions in 560 phenotypes. *Nature genetics*, 51(2):327, 2019.
- 663 [30] Brendan K Bulik-Sullivan, Po-Ru Loh, Hilary K Finucane, Stephan Ripke, Jian Yang, Nick
664 Patterson, Mark J Daly, Alkes L Price, Benjamin M Neale, Schizophrenia Working Group

- 665 of the Psychiatric Genomics Consortium, et al. LD score regression distinguishes con-
666 founding from polygenicity in genome-wide association studies. *Nature genetics*, 47(3):291,
667 2015.
- 668 [31] Heather E Wheeler, Kaanan P Shah, Jonathon Brenner, Tzintzuni Garcia, Keston Aquino-
669 Michaels, Nancy J Cox, Dan L Nicolae, Hae Kyung Im, GTEx Consortium, et al. Survey
670 of the heritability and sparse architecture of gene expression traits across human tissues.
671 *PLoS genetics*, 12(11):e1006423, 2016.
- 672 [32] Chiara Sabatti, Anna-Liisa Hartikainen, Anneli Pouta, Samuli Ripatti, Jae Brodsky,
673 Chris G Jones, Noah A Zaitlen, Teppo Varilo, Marika Kaakinen, Ulla Sovio, et al. Genome-
674 wide association analysis of metabolic traits in a birth cohort from a founder population.
675 *Nature genetics*, 41(1):35, 2009.
- 676 [33] Wellcome Trust Case Control Consortium et al. Genome-wide association study of 14,000
677 cases of seven common diseases and 3,000 shared controls. *Nature*, 447(7145):661, 2007.
- 678 [34] John M Dietschy, Stephen D Turley, and David K Spady. Role of liver in the maintenance
679 of cholesterol and low density lipoprotein homeostasis in different animal species, including
680 humans. *Journal of lipid research*, 34(10):1637–1659, 1993.
- 681 [35] Petri T. Kovanen, Michael Scott Brown, and Joseph L Goldstein. Increased binding of low
682 density lipoprotein to liver membranes from rats treated with 17 alpha-ethinyl estradiol.
683 *The Journal of biological chemistry*, 254 22:11367–73, 1979.
- 684 [36] Tao Feng and Xiaofeng Zhu. Genome-wide searching of rare genetic variants in WTCCC
685 data. *Human genetics*, 128(3):269–280, 2010.
- 686 [37] Xiang Zhou. A unified framework for variance component estimation with summary
687 statistics in genome-wide association studies. *The annals of applied statistics*, 11(4):2027,
688 2017.
- 689 [38] Benjamin B Sun, Joseph C Maranville, James E Peters, David Stacey, James R Staley,
690 James Blackshaw, Stephen Burgess, Tao Jiang, Ellie Paige, Praveen Surendran, et al.
691 Genomic atlas of the human plasma proteome. *Nature*, 558(7708):73, 2018.

- 692 [39] Katherine E. Cole, Christine A. Strick, Timothy J. Paradis, Kevin T. Ogborne, Marcel
693 Loetscher, Ronald P. Gladue, Wen Lin, James G. Boyd, Bernhard Moser, Douglas E. Wood,
694 Barbara G. Sahagan, and Kuldeep Neote. Interferon-inducible T cell alpha chemoattractant
695 (I-TAC): A novel non-ELR CXC chemokine with potent activity on activated T cells
696 through selective high affinity binding to CXCR3. *Journal of Experimental Medicine*,
697 187(12):2009–2021, 1998.
- 698 [40] Johannes Kettunen, Ayşe Demirkan, Peter Würtz, Harmen HM Draisma, Toomas Haller,
699 Rajesh Rawal, Anika Vaarhorst, Antti J Kangas, Leo-Pekka Lyytikäinen, Matti Pirinen,
700 et al. Genome-wide study for circulating metabolites identifies 62 loci and reveals novel
701 systemic effects of LPA. *Nature communications*, 7:11122, 2016.
- 702 [41] Cross Disorder Group of the Psychiatric Genomics Consortium. Identification of risk loci
703 with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet*,
704 381(9875):1360–1360, 2013.
- 705 [42] S Ripke, AR Sanders, KS Kendler, DF Levinson, P Sklar, PA Holmans, DY Lin, J Duan,
706 RA Ophoff, OA Andreassen, et al. Schizophrenia psychiatric genome-wide association
707 study (gwas) consortium genome-wide association study identifies five new schizophrenia
708 loci. *Nature Genetics*, 43:969–976, 2011.
- 709 [43] Stephan Ripke, Colm O’Dushlaine, Kimberly Chambert, Jennifer L Moran, Anna K Kähler,
710 Susanne Akterin, Sarah E Bergen, Ann L Collins, James J Crowley, Menachem Fromer,
711 et al. Genome-wide association analysis identifies 13 new risk loci for schizophrenia. *Nature*
712 *genetics*, 45(10):1150, 2013.
- 713 [44] Stephan Ripke, Benjamin M Neale, Aiden Corvin, James TR Walters, Kai-How Farh,
714 Peter A Holmans, Phil Lee, Brendan Bulik-Sullivan, David A Collier, Hailiang Huang, et al.
715 Biological insights from 108 schizophrenia-associated genetic loci. *Nature*, 511(7510):421,
716 2014.
- 717 [45] Chen Yao, Roby Joehanes, Andrew D Johnson, Tianxiao Huan, Chunyu Liu, Jane E
718 Freedman, Peter J Munson, David E Hill, Marc Vidal, and Daniel Levy. Dynamic role of
719 trans regulation of gene expression in relation to complex traits. *The American Journal*
720 *of Human Genetics*, 100(4):571–580, 2017.

- 721 [46] Oscar Franzén, Raili Ermel, Ariella Cohain, Nicholas K Akers, Antonio Di Narzo, Husain A
722 Talukdar, Hassan Foroughi-Asl, Claudia Giambartolomei, John F Fullard, Katyayani
723 Sukhavasi, et al. Cardiometabolic risk loci share downstream cis-and trans-gene regulation
724 across tissues and diseases. *Science*, 353(6301):827–830, 2016.
- 725 [47] Center LDACCAnalysis Working Group Coordinating, Common Fund NIH, GTEx Con-
726 sortium, Statistical Methods groupsAnalysis Working Group, et al. Genetic effects on
727 gene expression across human tissues. *Nature*, 550(7675):204, 2017.
- 728 [48] Elin Grundberg, Kerrin S Small, Åsa K Hedman, Alexandra C Nica, Alfonso Buil, Sarah
729 Keildson, Jordana T Bell, Tsun-Po Yang, Eshwar Meduri, Amy Barrett, et al. Mapping cis-
730 and trans-regulatory effects across multiple tissues in twins. *Nature genetics*, 44(10):1084,
731 2012.
- 732 [49] Fred A Wright, Patrick F Sullivan, Andrew I Brooks, Fei Zou, Wei Sun, Kai Xia, Vered
733 Madar, Rick Jansen, Wonil Chung, Yi-Hui Zhou, et al. Heritability and genomics of gene
734 expression in peripheral blood. *Nature genetics*, 46(5):430, 2014.
- 735 [50] Walter Gilbert. Why genes in pieces? *Nature*, 271(5645):501, 1978.
- 736 [51] Arianne J Matlin, Francis Clark, and Christopher WJ Smith. Understanding alternative
737 splicing: towards a cellular code. *Nature reviews Molecular cell biology*, 6(5):386, 2005.
- 738 [52] Yang I Li, Bryce van de Geijn, Anil Raj, David A Knowles, Allegra A Petti, David Golan,
739 Yoav Gilad, and Jonathan K Pritchard. Rna splicing is a primary link between genetic
740 variation and disease. *Science*, 352(6285):600–604, 2016.
- 741 [53] Atsushi Takata, Naomichi Matsumoto, and Tadafumi Kato. Genome-wide identification
742 of splicing qtls in the human brain and their enrichment among schizophrenia-associated
743 loci. *Nature communications*, 8:14519, 2017.
- 744 [54] Rolf I Skotheim and Matthias Nees. Alternative splicing in cancer: noise, functional, or
745 systematic? *The international journal of biochemistry & cell biology*, 39(7-8):1432–1449,
746 2007.
- 747 [55] Yang I Li, David A Knowles, Jack Humphrey, Alvaro N Barbeira, Scott P Dickinson,

- 748 Hae Kyung Im, and Jonathan K Pritchard. Annotation-free quantification of rna splicing
749 using leafcutter. *Nature genetics*, 50(1):151, 2018.
- 750 [56] Maria Gutierrez-Arcelus, Halit Ongen, Tuuli Lappalainen, Stephen B Montgomery, Alfonso
751 Buil, Alisa Yurovsky, Julien Bryois, Ismael Padioleau, Luciana Romano, Alexandra
752 Planchon, et al. Tissue-specific effects of genetic and epigenetic variation on gene regulation
753 and splicing. *PLoS genetics*, 11(1):e1004958, 2015.
- 754 [57] Alexis Battle, Zia Khan, Sidney H Wang, Amy Mitrano, Michael J Ford, Jonathan K
755 Pritchard, and Yoav Gilad. Impact of regulatory variation from RNA to protein. *Science*,
756 347(6222):664–667, 2014.
- 757 [58] Jiming Jiang, Cong Li, Debashis Paul, Can Yang, Hongyu Zhao, et al. On high-dimensional
758 misspecified mixed model analysis in genome-wide association study. *The Annals of*
759 *Statistics*, 44(5):2127–2160, 2016.
- 760 [59] Chuanhai Liu, Donald B Rubin, and Ying Nian Wu. Parameter expansion to accelerate
761 EM: the PX-EM algorithm. *Biometrika*, 85(4):755–770, 1998.
- 762 [60] Hua Zhou, Liuyi Hu, Jin Zhou, and Kenneth Lange. MM algorithms for variance compo-
763 nents models. *Journal of Computational and Graphical Statistics*, pages 1–12, 2019.
- 764 [61] M. H. Quenouille. Notes on bias in estimation. *Biometrika*, 43(3/4):353–360, 1956.
- 765 [62] Andrew R Wood, Tonu Esko, Jian Yang, Sailaja Vedantam, Tune H Pers, Stefan Gustafsson,
766 Audrey Y Chu, Karol Estrada, Jian’an Luan, Zoltán Kutalik, et al. Defining the role
767 of common variation in the genomic and biological architecture of adult human height.
768 *Nature genetics*, 46(11):1173, 2014.
- 769 [63] Thomas W Winkler, Anne E Justice, Mariaelisa Graff, Llilda Barata, Mary F Feitosa,
770 Su Chu, Jacek Czajkowski, Tõnu Esko, Tove Fall, Tuomas O Kilpeläinen, et al. The
771 influence of age and sex on genetic associations with adult body size and shape: a large-scale
772 genome-wide interaction study. *PLoS genetics*, 11(10):e1005378, 2015.