

IGREX for quantifying the impact of genetically regulated expression on phenotypes

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

¹Department of Mathematics, The Hong Kong University of Science and Technology, Clear Water Bay, Hong Kong, ²Department of Public Health Sciences, The University of Chicago, IL 60637, USA and ³Center for Quantitative Medicine, Duke-NUS Medical School, 169856, Singapore

Received August 14, 2019; Revised January 08, 2020; Editorial Decision February 01, 2020; Accepted February 05, 2020

ABSTRACT

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

INTRODUCTION

Genome-wide association studies (GWAS) have successfully identified tens of thousands of unique associations between single-nucleotide polymorphisms (SNPs) and a wide range of complex traits/diseases (<http://www.ebi.ac.uk/gwas/>). More than 90% of identified risk variants are located in non-coding regions (1), making it challenging to understand their functional mechanisms. Increasing evidence (2–8) has suggested that many of those risk variants may affect traits/diseases via the modulation of their *cis* gene expression levels. For example, a study of 18 complex traits revealed an enrichment for expression quantitative trait loci (eQTLs) in 11% of 729 tissue-trait pairs (9).

There is great interest in precisely characterizing the specific role of genetically regulated gene expression (GREX) in human traits and diseases.

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

The vast amount of publicly available eQTL and GWAS data resources enables an integrative framework, transcriptome-wide association studies (TWAS), for mapping gene-level trait associations and evaluating GREX effects on human traits and diseases. Using a reference eQTL panel (e.g. GTEx), gene-specific expression prediction models can be built based on *cis*-acting genetic factors. Then the gene expression levels of a GWAS cohort can be predicted based on individual genetic profiles, and the genetically regulated and predicted expression levels are further associated with the phenotype of interest in the GWAS study to map gene-level trait-associations (14–20). Existing methods have been proposed (8,21), including PrediXcan (14), TWAS (15), FOCUS (17), S-PrediXcan (18), UTMOST (22) and CoMM (19). Through applications to a wide variety of phenotypes, these methods have successfully identified specific gene-trait associations, whereas a comprehensive and precise evaluation of the impact of GREX variation on various traits and the trait-relevant cellular context is still needed (23).

*To whom correspondence should be addressed. Tel: +852 2358 7462; Fax: +852 23581643; Email: macyang@ust.hk
Correspondence may also be addressed to Jin Liu. Tel: +65 6576 7376; Fax: +65 6225 1244; Email: jin.liu@duke-nus.edu.sg

TWAS-types of integrative analysis rely on a key assumption: there exists a steady-state GREX variation shared across reference eQTL data and GWAS data, and the steady-state GREX variation can further induce phenotypic variation. The multi-tissue eQTL data from the GTEx project is commonly used as the reference eQTL panel (14,18,22). The GTEx project has collected data from post-mortem donors and has provided a source of largely non-diseased tissues for general purposes. The GTEx reference may or may not have considerable shared GREX variation with GWAS data of specific phenotypes in specific populations. Given the often unknown disease/trait-relevant tissue types and the increasing availability of eQTL data resources from different conditions, there is a need for new methods and tools that can be used to assess the proportion of the shared GREX variation in the phenotypic variation from a global perspective, and guide the selection of eQTL reference data and tissue-types for specific phenotypes and populations.

The heritability measure has been widely used to quantify the impact of genetic variation on phenotypic variation, and has served as a preliminary yet insightful assessment of the potential of genetic studies on various phenotypes (24,25). Analogous to the heritability measure, the estimation of proportion of GREX on phenotypic variation can also be used to evaluate the impact of the genetic regulatory effects on phenotypes mediated by expression levels, and inform trait-relevant tissue types or conditions in specific populations. To the best of our knowledge, there are two methods that have been proposed for this purpose (20). The RhoGE method (20) estimates the proportion of phenotypic variation explained by GREX based on linkage-disequilibrium (LD) score regression (LDSC) (26). Since it ignores the uncertainty in predicting gene expression levels, the proportion of variance explained by GREX could be substantially under-estimated by RhoGE. The other method, known as the gene expression co-score regression (GECS), requires the analyzed SNPs not being in LD to ensure unbiasedness, which greatly limits its applicability in real data analysis.

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

MATERIALS AND METHODS

Datasets and data preprocessing

GTEx eQTL dataset. We used the gene expression data from the V7 release of GTEx Consortium (<https://gtexportal.org/home/datasets>) as our reference dataset. We analyzed the 48 tissues with number of genotyped samples ≥ 70 , which are collected from 620 donors with total sample size 10 294. The sample size of each tissue ranges from 80 to 491 (details provided in Supplementary Table S4). We set the mappability cutoff at 0.9 to filter gene expression measures with lower quality, leaving 16 333–27 378 genes to be included in our analysis. Based on the third phase of the International HapMap project phase 3 (HapMap3), 1 189 556 SNPs were included from the GTEx genotyped data for analysis. For each gene, we included only the SNPs within 500 kb of the transcription start and end of each protein coding genes. In real data analysis, we used the covariates provided by the GTEx consortium, including genotype principal components (PCs), Probabilistic Estimation of Expression Residuals (PEER) factors, genotyping platform and sex (as described in <https://gtexportal.org/home/documentationPage>).

Individual level GWAS datasets. We obtained the individual-level genotype and phenotype data of the Northern Finland Birth Cohorts program 1966 (NFBC) (27) from the database of Genotypes and Phenotypes (db-GaP) with accession number phs000424.v7.p2. This dataset is comprised of 5402 individuals with 10 continuous phenotypes related to cardiovascular disease including Glucose, body mass index (BMI), C-reactive protein (CRP), insulin, high-density lipoprotein cholesterol (HDL), low-density lipoprotein cholesterol (LDL), triglycerides (TG), total cholesterol (TC), diastolic blood pressure (DiaBP) and systolic blood pressure (SysBP). There are 364 590 genotyped SNPs in this dataset. We first excluded the individuals whose reported sex differed from their sex determined from the X chromosome. We then excluded the SNPs with minor allele frequency less than 1%, with missing values in more than 1% of the individuals or with Hardy-Weinberg equilibrium (HWE) P -value below 0.0001. This quality control process following (28) yields 5123 individuals with 319 147 SNPs for our analysis. We evaluated the genetic relatedness matrix (GRM) using the processed genotype data and selected the top 20 PCs as covariates in the study.

Another individual-level GWAS dataset is from the Wellcome Trust Case Control Consortium (WTCCC, <https://www.wtccc.org.uk/>) (29). The WTCCC dataset contains seven disease phenotypes including bipolar disorder (BD), coronary artery disease (CAD), Crohn's disease (CD), hypertension (HT), rheumatoid arthritis (RA), type 1 diabetes (T1D) and type 2 diabetes (T2D). It includes ~2000 cases per phenotype and 3004 controls with 490 032 genotyped SNPs. Following the QC process of (30,31), we first removed the individuals with genotyping rate $< 5\%$. Then we excluded the SNPs satisfying at least one of the following: minor allele frequency $< 5\%$; genotypes missing in more than 1% samples; HWE P -value is below 0.001. We also removed the individuals with estimated genetic correlation

larger than 2.5%. After quality control, around 4700 individuals with 300 000 SNPs were retained for our analysis (See Supplementary Table S1). Based on the obtained data, we calculated the GRM and extracted top 20 PCs as covariates to be included in our analysis.

GWAS summary statistics. Besides the individual-level GWAS data, we analyzed 10 summary level GWAS datasets: human plasma protein quantitative trait loci (pQTL) dataset (32), circulating metabolite data (33), four schizophrenia datasets (34–37), two independent height datasets (38) and BMI datasets from European ancestry with age under 50 separated by men and women (39). The SNPs with missing information (i.e. chromosome, minor allele, allele frequency) were first removed. Following the practice of LDSC (26), we checked the χ^2 statistic of each SNP and excluded those with extreme values ($\chi^2 > 80$) to prevent the outliers that may unduly affect the results. The detailed information and download links are provided in Supplementary Table S2. After pre-processing, the remaining SNPs were further matched with reference data and this step is automatically conducted in our IGREX software.

The trans-eQTLGen summary data. In the analysis involving the trans-eQTLs, we used the SNPs identified in blood tissue provided by the eQTLGen Consortium (<http://www.eqtlgen.org>). The trans-eQTL analysis were restricted to known complex trait-associated SNPs. The significant trans-eQTLs were identified by controlling the FDR at 0.05. There were 59 786 gene-SNP pairs composed of 6298 genes and 3853 SNPs. The remaining pairs after matching with both reference and GWAS datasets are summarized in Supplementary Table S5.

IGREX framework for quantifying the GREX component

IGREX is a two-stage method for quantifying the proportion of phenotypic variation that can be attributed to GREX variation. The method can be applied to both individual-level (IGREX-i) and summary-level (IGREX-s) GWAS data. It first evaluates the posterior distribution of GREX effects based on an eQTL reference panel and then estimates the proportion of variance explained by GREX using the ‘predicted’ gene expression in the GWAS data. The details of both IGREX-i and IGREX-s are described as follows.

The IGREX-i for individual-level GWAS data. Consider a reference eQTL dataset \mathcal{D}_r and an individual-level GWAS dataset \mathcal{D}_i . The 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 $n_r \times M$ genotype matrix, \mathbf{X}_r , where G is the number of genes, M is the number of SNPs and n_r is the 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 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 vector of expression levels of the g -th gene and the genotype matrix corresponding to its local (*cis*) SNPs from the reference panel, respectively. We first relate \mathbf{y}_g to $\mathbf{X}_{r,g}$ with a linear model:

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

where $\boldsymbol{\beta}_g \in \mathbb{R}^{M_g}$ is the vector of genetic effects of M_g *cis* SNPs on the expression 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 in the steady-state component of gene expression levels regulated by genetic variants, $\boldsymbol{\beta}_g$ is assumed to be the same for individuals in both datasets, \mathcal{D}_r and \mathcal{D}_i . Consequently, the GREX component of individuals in the GWAS data can be evaluated by $\mathbf{X}_g \boldsymbol{\beta}_g$, where \mathbf{X}_g denotes the columns of \mathbf{X} corresponding to the *cis*-SNPs of the g -th gene. Meanwhile, we assume that the genetic effects on the phenotype of interest \mathbf{t} can be decomposed into two parts, i.e. the effects mediated via GREX and the effects through alternative pathways not mediated by gene expression levels:

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

where scalar α_g is the effect size of $\mathbf{X}_g \boldsymbol{\beta}_g$ on \mathbf{t} , $\boldsymbol{\gamma} \in \mathbb{R}^M$ is the vector of alternative 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 $\mathbf{X} \boldsymbol{\gamma}$ correspond to the overall impact of the GREX component and the alternative genetic effects on \mathbf{t} , respectively. Thus, the impact of GREX on the phenotype can be quantified by the proportion of phenotypic variance explained by the GREX component:

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

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

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

which is motivated by a recent theoretical justification (40) for heritability estimation on a mis-specified linear mixed model (LMM). This prior specification in (4) provides a great computational advantage as well as a stable performance for IGREX under model mis-specification, as demonstrated in the simulation studies.

The proposed method for individual-level GWAS data, IGREX-i, provides a two-stage framework for estimating PVE_{GREX} . In the first stage, we estimate the parameters $\sigma_{\beta_g}^2$ and $\sigma_{r,g}^2$ in model (1) via a fast expectation-maximization (EM)-type algorithm, the parameter-expanded EM (PX-EM) algorithm (41). Based on the estimates, denoted as $\hat{\sigma}_{\beta_g}^2$ and $\hat{\sigma}_{r,g}^2$, the posterior distribution of $\boldsymbol{\beta}_g$ is given by

$$\boldsymbol{\beta}_g | \mathbf{y}_g, \mathbf{X}_{r,g} \sim \mathcal{N}(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g), \quad (5)$$

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

In the second stage, we estimate PVE_{GREX} by treating the posterior distribution (5) as the prior distribution of $\boldsymbol{\beta}_g$ in model (2). To evaluate the covariance of \mathbf{t} , we first note that $\mathbb{E}(\mathbf{t} | \boldsymbol{\alpha}) = \sum_{g=1}^G \alpha_g \mathbf{X}_g \boldsymbol{\mu}_g$ 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 law of total expectation and total variance, we obtain

$\mathbb{E}(\mathbf{t}) = \mathbb{E}(\mathbb{E}(\mathbf{t}|\boldsymbol{\alpha})) = \mathbf{0}$ and

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

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

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

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. In this estimation, the substitution of posterior $\boldsymbol{\beta}_g | \mathbf{y}_g, \mathbf{X}_{r,g}$ accounts for the posterior variance $\boldsymbol{\Sigma}_g$ and naturally results in the adjustment of estimation uncertainty associated with $\boldsymbol{\beta}_g$. This is important because in the GWAS data, the gene expression levels are not directly measured, but rather are predicted or imputed based on genetic variants. It is known that the prediction accuracy and uncertainty vary substantially among genes. For most of the genes in the genome, the genetically regulated expression variation accounts for only a small to moderate proportion of total expression variation. Thus, the prediction may not be accurate and could be subject to high uncertainty. In contrast, our model accounts for the estimation uncertainty by $\boldsymbol{\Sigma}_g$ and can yield unbiased estimation for $\widehat{\text{PVE}}_{\text{GREX}}$.

IGREX-i provides two methods for estimating the parameters and $\widehat{\text{PVE}}_{\text{GREX}}$ in the second stage. Let $\boldsymbol{\psi} = [\sigma_{\alpha}^2, \sigma_{\gamma}^2, \sigma_{\epsilon}^2]^T$ be the vector of parameters to be estimated, $\mathbf{K}_{\alpha} = \sum_{g=1}^G \mathbf{X}_g (\boldsymbol{\mu}_g \boldsymbol{\mu}_g^T + \boldsymbol{\Sigma}_g) \mathbf{X}_g^T$ and $\mathbf{K}_{\gamma} = \mathbf{X} \mathbf{X}^T$. The first method is based on MoM, which minimizes the distance between the second moment of \mathbf{t} at the population level and that at the sample 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 the estimating equation

$$\begin{aligned} \mathbf{S} \boldsymbol{\psi} &= \mathbf{q}, \text{ with } \mathbf{S} = \begin{bmatrix} \text{tr}(\mathbf{K}_{\alpha}^2) & \text{tr}(\mathbf{K}_{\alpha} \mathbf{K}_{\gamma}) & \text{tr}(\mathbf{K}_{\alpha}) \\ \text{tr}(\mathbf{K}_{\alpha} \mathbf{K}_{\gamma}) & \text{tr}(\mathbf{K}_{\gamma}^2) & \text{tr}(\mathbf{K}_{\gamma}) \\ \text{tr}(\mathbf{K}_{\alpha}) & \text{tr}(\mathbf{K}_{\gamma}) & n \end{bmatrix}, \quad (8) \\ \boldsymbol{\psi} &= \begin{bmatrix} \sigma_{\alpha}^2 \\ \sigma_{\gamma}^2 \\ \sigma_{\epsilon}^2 \end{bmatrix}, \quad \mathbf{q} = \begin{bmatrix} \mathbf{t}^T \mathbf{K}_{\alpha} \mathbf{t} \\ \mathbf{t}^T \mathbf{K}_{\gamma} \mathbf{t} \\ \mathbf{t}^T \mathbf{t} \end{bmatrix}. \end{aligned}$$

The solution of normal equation (8) is given by $\hat{\boldsymbol{\psi}} = \mathbf{S}^{-1} \mathbf{q}$. And the variance-covariance matrix of $\hat{\boldsymbol{\psi}}$ is given by $\text{Cov}(\hat{\boldsymbol{\psi}}) = \mathbf{S}^{-1} \text{Cov}(\mathbf{q}) \mathbf{S}^{-1}$ using the sandwich estimator. The second method applies the restricted maximum likelihood (REML) by further assuming the normal 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 the Minorization-Maximization (MM) algorithm (42), and their standard errors are estimated by the

inverse of Fisher information matrix. In both MoM and REML, the standard error of $\widehat{\text{PVE}}_{\text{GREX}}$ can be obtained by the delta method (see Supplementary Note).

In addition to the point estimate $\widehat{\text{PVE}}_{\text{GREX}}$, the IGREX framework can be also used to test $H_0: \text{PVE}_{\text{GREX}} = 0$ for the phenotype of interest in specific populations given an eQTL reference with a specific tissue type or cellular context. Because this is a test of the boundary point, the test statistic follows a mixture of χ^2 distribution. Usually, the Davies method can be used to produce well-calibrated P -values at a cost of computational inefficiency when the sample size n is large (43). As an approximation of the Davies method, IGREX adopts a normal test using the point estimate of $\widehat{\text{PVE}}_{\text{GREX}}$ and its standard error. Using simulated data, we showed that the normal test provides satisfactory approximation to the exact test (See Supplementary Section 2.4 and Figure S7).

The IGREX-s for summary-level GWAS data. In real applications, individual-level GWAS data \mathcal{D}_i may not be accessible. We have further developed IGREX-s which requires only summary-level GWAS data as input. Based on MoM, IGREX-s can approximate IGREX-i while requiring only SNP-level z -scores 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 m is the number of samples in the reference panel.

Suppose we only have the z -scores from summary-level GWAS data $\{z_j\}_{j=1}^M$ generated from \mathcal{D}_i . The definition of the z -score is $z_j = \frac{(\mathbf{x}_j^T \mathbf{x}_j)^{-1} \mathbf{x}_j^T \mathbf{t}}{\sqrt{\hat{\sigma}_j^2 (\mathbf{x}_j^T \mathbf{x}_j)^{-1}}}$, where \mathbf{x}_j is the j -th column of \mathbf{X} and $\hat{\sigma}_j^2$ is the estimate of residual variance by regressing \mathbf{x}_j on \mathbf{t} . Because z -scores are invariant with respect to the scales of both \mathbf{X} and \mathbf{t} , we can assume that the z -scores are calculated from a standardized genotype matrix $\tilde{\mathbf{X}}$ without loss of generality. Hence, 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 $\hat{\sigma}_t^2$ is the estimate of $\text{Var}(\mathbf{t})$. Hence, we have

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

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

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

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 $\tilde{\mathbf{X}}_g$ is the corresponding columns of a reference genotype matrix $\tilde{\mathbf{X}}$. In practice, $\tilde{\mathbf{X}}$ can be the eQTL reference genotype $\tilde{\mathbf{X}}$, (e.g. the genotype matrix from the GTEx Project), a subset of \mathbf{X} or from the 1000 Genomes Project. Using simulations, we showed that with a few hundreds of samples in the eQTL reference data, the estimation of IGREX-s with summary statistics well approximates IGREX-i using individual level data. Now, we 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 both sides by n^2 , we have

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

The terms on the left hand side do not involve \mathbf{t} and thus can be approximated using $\tilde{\mathbf{X}}$ (44). 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 $\tilde{\mathbf{K}}_\alpha = \sum_{g=1}^G \tilde{\mathbf{X}}_g (\mu_g \mu_g^T + \Sigma_g) \tilde{\mathbf{X}}_g^T$. Other terms on the left hand side can be approximated in the same way. For the right hand side, each term can be approximated using $\hat{\mathbf{R}}_g$ and z -scores from approximation (9): $\mathbf{t}^T \mathbf{K}_\alpha \mathbf{t} \approx n \hat{\sigma}_i^2 \sum_g \mathbf{z}_g^T (\mu_g \mu_g^T + \Sigma_g) \mathbf{z}_g$, where $\mathbf{z}_g \in \mathbb{R}^{M_g}$ is the vector of z -scores corresponding to the g -th gene; $\frac{\text{tr}(\mathbf{K}_\alpha) \mathbf{t}^T \mathbf{t}}{n^3} \approx n \hat{\sigma}_i^2 \text{tr}(\sum_g (\mu_g \mu_g^T + \Sigma_g) \hat{\mathbf{R}}_g)$; $\mathbf{t}^T \mathbf{K}_\gamma \mathbf{t} \approx n \hat{\sigma}_i^2 \sum_{j=1}^M z_j^2$; and $\frac{\text{tr}(\mathbf{K}_\gamma) \mathbf{t}^T \mathbf{t}}{n^3} \approx n \hat{\sigma}_i^2$. With these approximations, Equation (11) becomes

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

Then, $\frac{\hat{\sigma}_\alpha^2}{\hat{\sigma}_i^2}$ can be obtained by solving this equation. Plugging this estimate into Equation (10) gives the $\widehat{\text{PVE}}_{\text{GREX}}$. The standard errors of $\widehat{\text{PVE}}_{\text{GREX}}$ can be estimated by the block jackknife method (45). Given the $\widehat{\text{PVE}}_{\text{GREX}}$ and its standard error, we can test the tissue-wise hypothesis H_0 : $\text{PVE}_{\text{GREX}} = 0$ (details described in the Supplementary Note Section 2.3).

IGREX also allows for the adjustment of covariates including sex, age and genotype principal components (See details in Supplementary Note Section 2.1–2.3).

RESULTS

Simulation studies

We conducted extensive simulation studies to evaluate the performance of IGREX using the genotypes from the NFBC dataset. First, we extracted genotypes \mathbf{X} from the first chromosome of the NFBC dataset, which is comprised of $M = 23\,718$ SNPs. Among the 5123 samples, we randomly subsampled $n = 3000$ as the GWAS individuals and treated a subset of the rest samples as eQTL reference. The total phenotypic heritability was set as $h_i^2 = \frac{\text{Var}(\sum_{g=1}^G \alpha_g \mathbf{x}_g^T \beta_g + \mathbf{x}^T \beta_g)}{\text{Var}(t)} = 0.5$, where $\text{PVE}_{\text{GREX}} = 0.2$ and the proportion explained by the alternative genetic effects, $\text{PVE}_{\text{Alternative}} = \frac{\text{Var}(\mathbf{x}_g^T \beta_g)}{\text{Var}(t)} = 0.3$. Given the genotype matrices, β_g and α_g , the gene expression \mathbf{y}_g and phenotype \mathbf{t} were sim-

ulated following models (1) and (2). We will discuss the details for generating β_g and α_g later. To assess IGREX-s, we calculated the z -score of each SNP and randomly subsampled $m = 500$ samples from \mathbf{X} for estimating LD matrix $\hat{\mathbf{R}}_g$ (results for other settings of m are shown in Supplementary Figure S5).

We first evaluated the estimation performance of IGREX for different settings of eQTL reference data. Specifically, we considered n_r varying at $\{200, 400, 800\}$. Note that the setting $n_r = 200$ mimics the situation in GTEx study, whose average sample size is 214. We further considered larger n_r 's as the sample size of eQTL studies would increase in the future.

We also varied $\text{PVE}_y = \frac{\text{Var}(\mathbf{x}^T \beta_g)}{\text{Var}(y_g)}$ at $\{0.1, 0.2, 0.3\}$, where PVE_y quantifies the gene expression heritability explained by its local SNPs. To mimic the scenario in which the expression estimation uncertainty was incorrectly ignored, we obtained the posterior mean of β_g in the first stage, and replaced the true effect size β_g by its posterior mean μ_g while specifying the posterior variance to be $\Sigma_g = \mathbf{0}$ in the second stage, and then conducted REML and MoM as before. We denote these methods as REML_0 and MoM_0 . The simulation results summarized in Figure 1A show that both PVE_{GREX} and $\text{PVE}_{\text{Alternative}}$ are accurately estimated using IGREX-i in most settings. IGREX slightly underestimates PVE_{GREX} when both n_r and PVE_y are small (i.e. $n_r = 200$, $\text{PVE}_y = 0.1$), while the accuracy steadily increases as PVE_y increases. Additionally, IGREX-s well approximated MoM, producing nearly identical estimation results. In contrast, both REML_0 and MoM_0 does not account for estimation uncertainty in the expression prediction, and they show poor estimation performance even when sample size is large and PVE_y value is high. Additionally, we varied the number of cis-SNPs p_g at $\{20, 50, 100\}$ to investigate its influence. As shown in Figure 1B, the REML-based IGREX-i produces accurate estimates under all scenarios. The MoM-based IGREX-i and IGREX-s slightly underestimate PVE_{GREX} when p_g is large and n_r is small, but they achieves identical performance as REML as n_r increases.

Next we conducted simulations to evaluate the situation that the IGREX model was mis-specified based on the NFBC genotypes. Here we considered the situation where genetic effects β_g and α were sparse while we assumed dense effect sizes in the IGREX model. This was designed to mimic the real data situation that the architecture of eQTL signals is often sparse (46). Let π_α and π_β be the sparsity of α and β_g , i.e. $\pi_\alpha = (\# \text{ Nonzero entries in } \alpha) / G$ and $\pi_\beta = (\# \text{ Nonzero entries in } \beta_g) / M_g$, respectively. To evaluate the influence of different sparsity patterns on our method, we varied π_α and π_β at $\{0.2, 0.5, 0.8\}$. The nonzero entries in α and β_g were simulated form a normal distribution. As shown in Figure 1C and D, all three methods of IGREX produced accurate estimates in the presence of sparse genetic effects, implying the robustness of IGREX to model mis-specification. Moreover, the estimation performance was not influenced by the degree of sparsity. Next, we analyzed the influence of LD patterns by varying the autocorrelation between SNPs $\rho \in \{0.1, 0.3, 0.5, 0.8\}$ and generating \mathbf{X} based on the ρ 's. The simulation details are described in Supplementary Note Section 2.5. From Figure 1E, we observed that IGREX produced accurate estimation

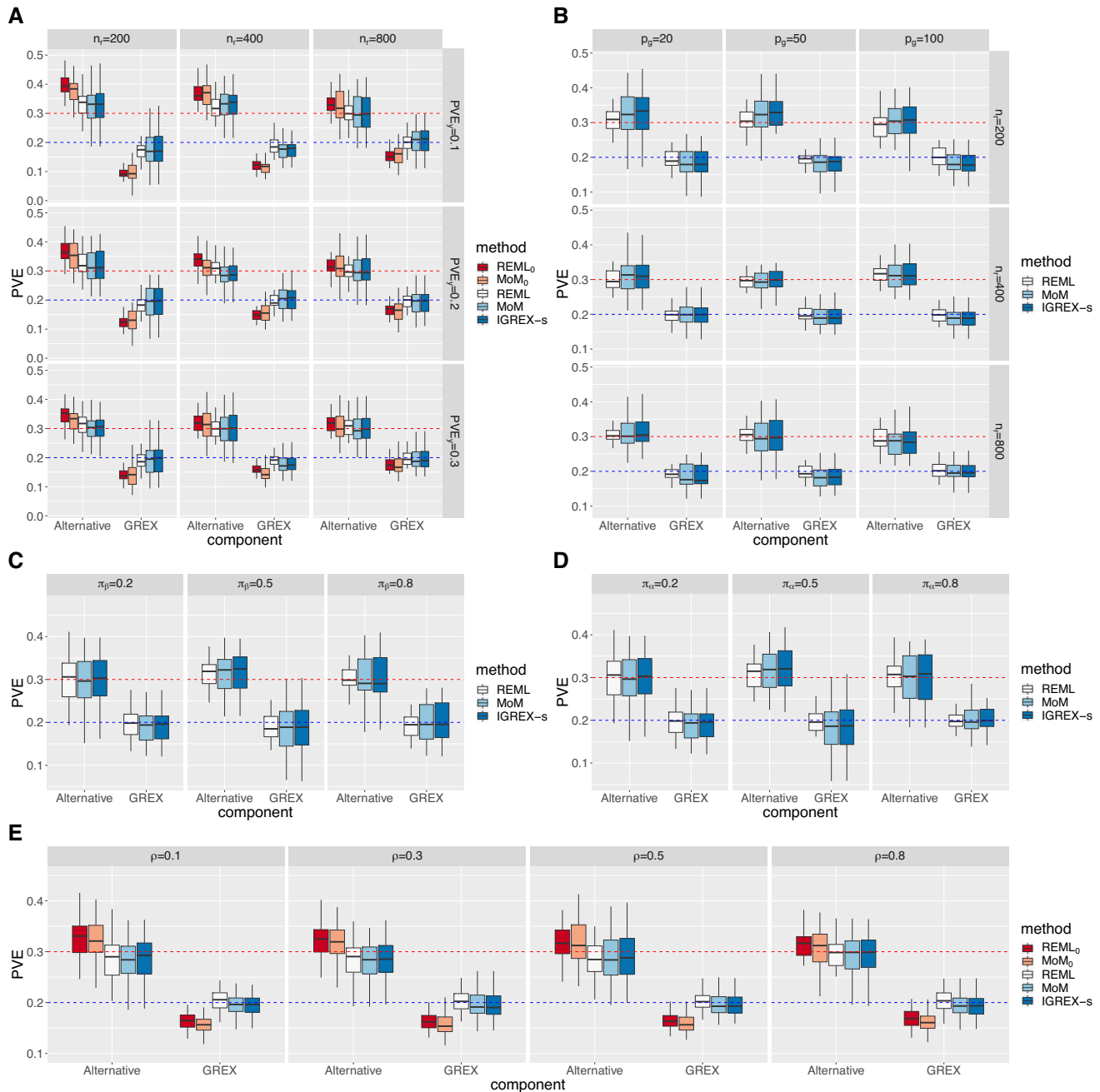


Figure 1. Simulation studies to compare estimation accuracies of IGREX with other methods. REML and MoM in the legend are abbreviations of the IGREX-i estimation methods. The blue and red dashed lines represent the true values of PVE_{GREX} and $PVE_{Alternative}$, respectively. We averaged the results 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 {200, 400, 800} and PVE_y was varied at {0.1, 0.2, 0.3}; (B) the three models of IGREX when n_r was varied at {200, 400, 800} and p_g were varied at {20, 50, 100}; (C) the three models of IGREX when $\pi_\alpha = 0.2$ and π_β was varied at {0.2, 0.5, 0.8}; (D) the three models of IGREX when $\pi_\beta = 0.2$ and π_α were varied at {0.2, 0.5, 0.8}; (E) the three models of IGREX, REML₀ and MoM₀ when ρ was varied at {0.1, 0.3, 0.5, 0.8}.

in various setting of LD. In contrast, REML₀ and MoM₀ consistently underestimated PVE_{GREX} as a result of ignoring estimation uncertainty.

We also compared IGREX with an existing method in the literature, RhoGE (20). RhoGE is an LDSC-based approach for estimating PVE_{GREX} . However, this method does not adjust for estimation uncertainty. The results are shown in Supplementary Figure S6. As expected, IGREX

yielded unbiased estimation while RhoGE substantially underestimated PVE_{GREX} in most settings. It achieved similar accuracy as IGREX only when the genetically regulated expression accounted for most of the expression variation, $PVE_y \geq 0.9$. In other words, RhoGE only works well when the genetically predicted expression levels are very close to the true underlying expression levels for most of the genes, which may not be realistic for real data analysis.

Real data applications with individual-level GWAS data

With eQTL data of 48 human tissues from the GTEx project as reference, we applied IGREX to two individual-level GWAS datasets, the Northern Finland Birth Cohorts program 1966 (NFBC) (27) and the Wellcome Trust Case Control Consortium (WTCCC) (29).

In analyzing the NFBC data, we focused on six quantitative traits with statistically significant heritability, based on 5123 individuals and 309 245 genotyped SNPs. Those six traits are Glucose ($h_t^2 = 14.2\% \pm 5.3\%$), high-density lipoprotein cholesterol (HDL, $h_t^2 = 32.9\% \pm 5.6\%$), low-density lipoprotein cholesterol (LDL, $h_t^2 = 29.0\% \pm 5.5\%$), triglycerides (TG, $h_t^2 = 13.6\% \pm 5.3\%$), total cholesterol (TC, $h_t^2 = 20.1\% \pm 5.4\%$) and systolic blood pressure (SysBP, $h_t^2 = 17.1\% \pm 5.4\%$). Figure 2A and B shows the tissue-specific \widehat{PVE}_{GREX} estimates of the six traits. The REML and MoM methods yielded similar estimates in most of the tissues. Besides, we visualized the relationship between the point estimates of PVE_{GREX} and the eQTL effect sizes for the NFBC dataset in Supplementary Figures S10 and 12. On the one hand, we can observe that the \widehat{PVE}_{GREX} does not increase as eQTL sample size increases. On the other hand, as shown in Supplementary Figures S11 and 13, there is a decreasing trend of the standard error as the eQTL sample size becomes larger. These results suggest that the eQTL sample size only influence the standard errors of \widehat{PVE}_{GREX} . This conclusion is confirmed by the later analysis in pQTL dataset (Supplementary Figures S18 and 19).

IGREX can also be used to inform trait-relevant tissue types. By testing $H_0: PVE_{GREX} = 0$ in each tissue type, we observed significant GREX components in liver for both LDL and TC. As shown in Figure 2A, \widehat{PVE}_{GREX} for LDL in liver is as high as 14.3% (with standard error 2.6%), capturing 52.6% of total heritability defined as PVE_{GREX}/h_t^2 ; and TC also has $\widehat{PVE}_{GREX} = 13.7\%$ (with standard error 2.5%) in liver, which captures 79.4% of total heritability (see Supplementary Figure S9). It is known that LDL synthesized in liver is an important lipoprotein particle for transporting cholesterol in the blood (47,48). Our findings suggest that genetic variants affect LDL through regulating their corresponding gene targets and liver is the most relevant tissue involved in gene regulation. Next, we analyzed the impact of ignoring the estimation uncertainty (with the complete results given in the Supplementary Figure S8). As shown in Figure 2C and D, the \widehat{PVE}_{GREX} declined substantially as a result of ignoring expression estimation uncertainty. In Figure 2E, we compared the estimates based on individual level data using IGREX-i versus those based on IGREX-s with summary statistics, where 500 samples from the NFBC dataset were used as \tilde{X} in IGREX-s. For all six of the traits, the IGREX-s estimates well approximated the estimates using the individual level data, which is consistent with our simulation results. We additionally compared the IGREX-s results obtained by using the GTEx reference panel with those obtained by using subsamples from the NFBC dataset. As shown in Supplementary Figure S15, the GTEx reference panel can also produce satisfactory IGREX-s approximation in practice.

Next we investigated the role of GREX in complex human traits and diseases, using the WTCCC dataset (29). We applied IGREX to estimate the tissue-specific PVE_{GREX} of seven diseases including bipolar disorder (BD), coronary artery disease (CAD), Crohn's disease (CD), hypertension (HT), rheumatoid arthritis (RA), type 1 diabetes (T1D) and type 2 diabetes (T2D). The estimates of PVE_{GREX}/h_t^2 obtained by REML and MoM are shown in Supplementary Figures S16 and 17, respectively. The top GREX components measured by PVE_{GREX}/h_t^2 are 12.8% for BD in amygdala, 21.2% for CAD in spinal cord, 18.4% for CD in amygdala, 16.7% for HT in spleen and 17.9% for T2D in anterior cingulate cortex. The average estimates of PVE_{GREX}/h_t^2 across 48 tissues for RA and T1D are as high as 34.1% and 71.2%, respectively. Both RA and T1D are autoimmune diseases, with well-established strong associations in the major histocompatibility complex (MHC) region (29,49). After removing the MHC region, we observed a substantial reduction in the PVE_{GREX}/h_t^2 estimates: the mean \widehat{PVE}_{GREX} dropped from 34.1% to 7.6% for RA and from 71.2% to 11.7% for T1D, as shown in Figure 3A. Additionally, the tissue-specific comparisons presented in Figure 3B showed an extensive reduction of PVE_{GREX} in all tissue types for T1D and RA, while such changes were not observed for other traits. This finding suggests the heavy involvement of GREX variation in the immune functions related to the MHC region for both RA and T1D. Here we illustrate that IGREX can be used to inform disease/trait-relevant tissue types or cellular contexts.

Analysis of a wide spectrum of phenotypes using IGREX-s with summary-level GWAS data

The vast amount of publicly available summary-level GWAS data and their easy accessibility allow us to conduct a comprehensive evaluation of the impact of GREX on a wide spectrum of phenotypes using IGREX-s, from molecular traits such as proteins and metabolites to various complex phenotypes including schizophrenia, height, and body mass index (BMI). In the following analysis, we used the genotypes of the 635 GTEx samples as the LD reference \tilde{X} in the IGREX-s estimation.

First, we estimated PVE_{GREX} in 249 proteins (244 of which were identified with unique UniProt IDs, see Supplementary Table S3) with significantly non-zero heritabilities using summary statistics from a plasma protein quantitative trait loci (pQTL) study (32), as summarized in Figure 4A. In Supplementary Figure S20, the heritabilities estimated by IGREX-s ($\hat{h}_t^2 = \widehat{PVE}_{GREX} + \widehat{PVE}_{Alternative}$) are shown to be highly consistent with those estimates obtained using MoM (44). From this perspective, heritability can be attributed to two components: the GREX component and its alternative effects. Then, we grouped 48 tissue types into 16 groups by their functions and tested the significance of tissue-specific GREX effects on the 249 proteins. We observed a significant GREX contribution in many tissue-protein pairs (Figure 4B and Supplementary Figures S21–23). In particular, 9 out of the 249 proteins had significant GREX components in at least one tissue type at 0.05 level after Bonferroni correction. As shown in Figure 4C and D, some proteins, including CD96, DEFB119, MICB and PDE4D,

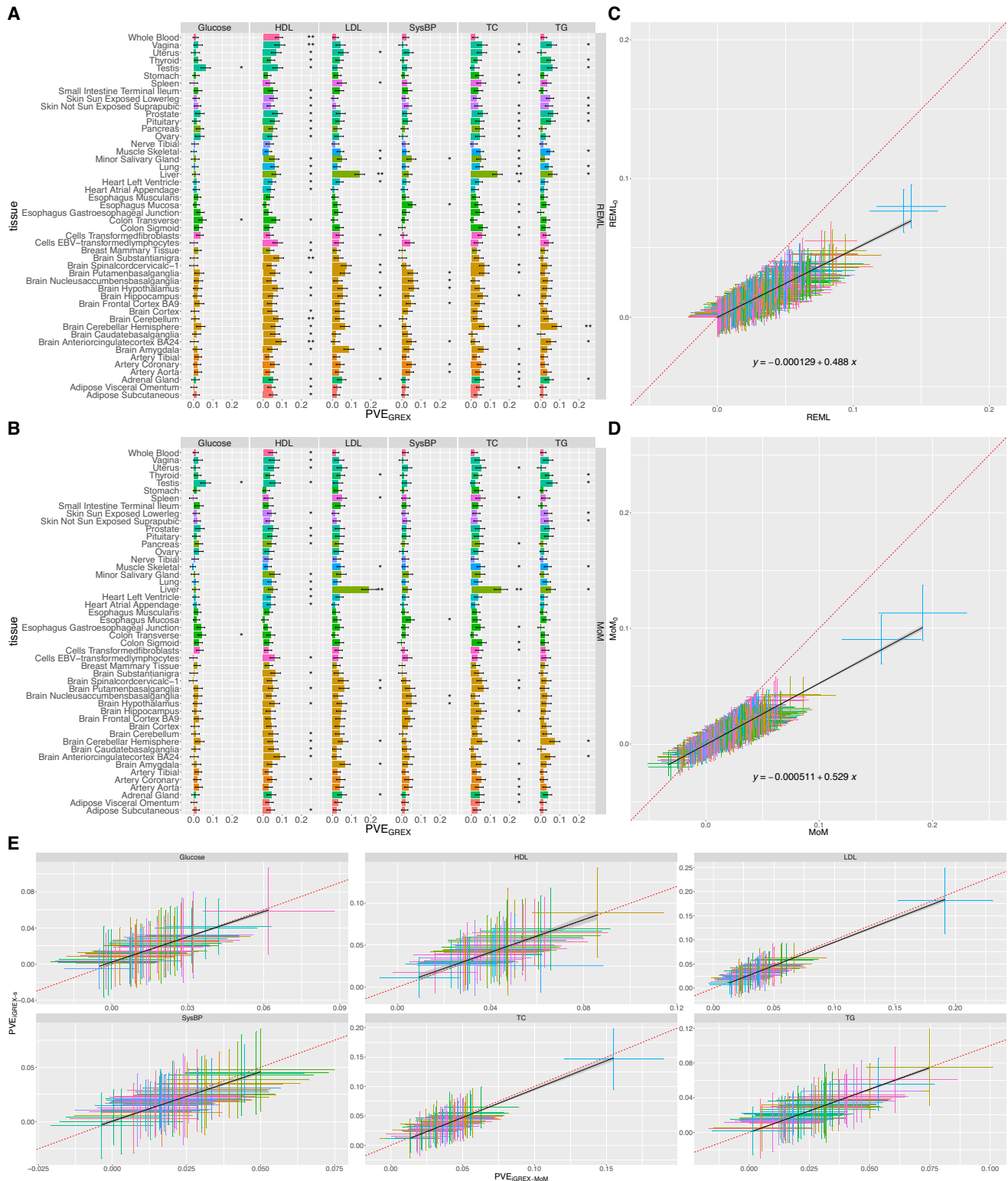


Figure 2. Tissue-specific PVE_{GREX} of the six traits from NFBC dataset. (A and B) PVE_{GREX} obtained by REML and MoM. Tissues are colored according to their categories. The number of asterisks represents the significance level: P -value < 0.05 is annotated by *; P -value $< 0.05/48$ is annotated by **. (C and 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 PVE_{GREX} generated by IGREX-s against those generated by MoM for all 48 tissues in one of the six traits.

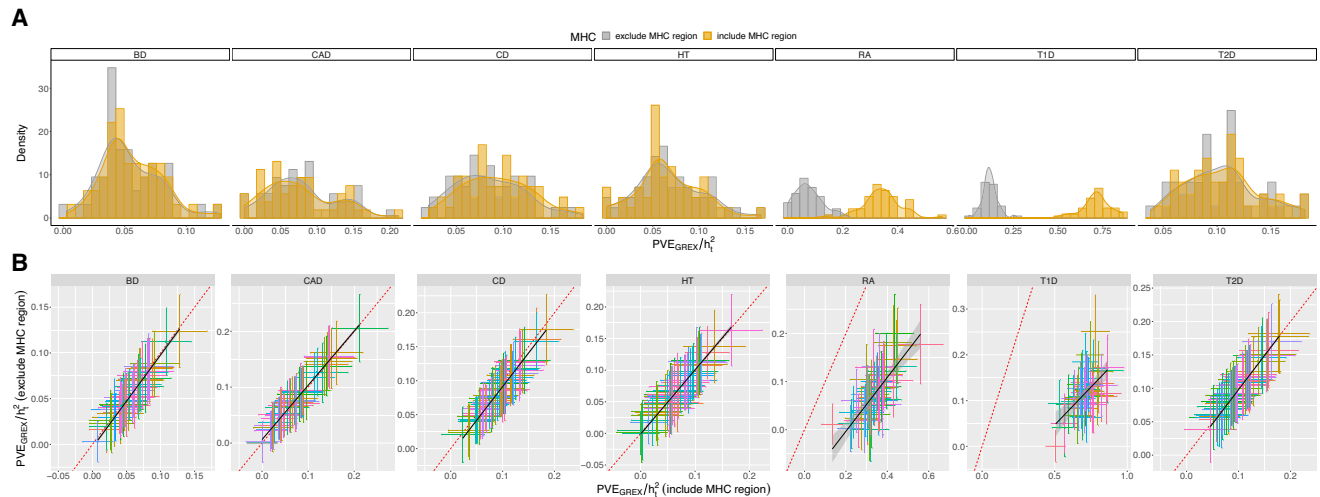


Figure 3. Percentage of heritability explained by GREX (PVE_{GREX}/h_t^2) of the seven traits from WTCCC data. (A) The distributions of estimated PVE_{GREX}/h_t^2 across 48 GTEx tissues. (B) Tissue-specific comparisons of PVE_{GREX}/h_t^2 estimated by whole genome with those estimated by excluding the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

exhibit cross-tissue GREX impacts; meanwhile other proteins, namely CFB, CXCL11, EVI2B, IDUA and LRPAP1, have tissue-specific GREX effect patterns. We found these tissue-specific patterns to be consistent with protein functions. For example, the CFB protein, which is implicated in the growth of preactivated B-lymphocytes, is found to be most associated with GREX in EBV-transformed lymphocytes ($PVE_{GREX} = 22.7\%$). As another example, the CXCL11 protein has the highest $PVE_{GREX} = 20.0\%$ in pancreas, and the *CXCL11* gene is often over-expressed in pancreas tissue (50). We also noted that six out of the nine proteins were immune-related, echoing our previous implications of the important role of GREX in immune processes. Using the pQTL dataset, we have conducted sensitivity analyses for IGREX (with details given in the Supplementary Note Section 2.6), illustrating the robustness of IGREX in various practical scenarios. In addition to the proteins, metabolic traits are also important intermediate traits for complex biological processes. We applied IGREX-s to a summary level dataset of circulating metabolites (33), and studied the impact of GREX on metabolic traits. The results are presented in Supplementary Figure S30 and discussed in the Supplementary Note Section 2.7.

Then we applied IGREX-s to the summary data of complex human traits. Here we analyzed three traits: schizophrenia (SCZ), height, and BMI. We considered four datasets of schizophrenia with increasing and overlapping samples: SCZ subset (34), SCZ1 (35), SCZ1+Sweden (SCZ1Swe) (36) and SCZ2 (37). We found that the estimated PVE_{GREX}/h_t^2 in all four SCZ datasets have higher values in the brain tissues than in other tissue types (Figure 5B and Supplementary Figure S28). As expected, the statistical power increases with sample size of GWAS (Figure 5A). Additionally, we also analyzed the human height and BMI phenotypes using pairs of independent GWAS data for replication purposes. The obtained estimates, PVE_{GREX} , from pairs of independent GWAS data are highly consistent. Although the analysis results are reproducible in sev-

eral different datasets, we noted the estimated percentages of heritability explained by GREX for all three complex traits are $<10\%$ (8.7% for schizophrenia, 8.7% for height and 3.7% for BMI in the most expressed tissue types. See Figure 5C and Supplementary Figure S29).

The relatively low GREX contribution to complex traits other than lipid or molecular traits can be attributed to multiple reasons. First, it is known that trans-acting genetic effects can explain a substantial proportion of expression variation (8,11). However, trans-eQTL effects are often tissue-specific and can be harder to detect and replicate across studies (51). In TWAS-types of analysis, generally the prediction of gene expression is based on only *cis* genetic variants of each gene. As such, the PVE_{GREX} values reported here, also based on only *cis* genetic variants, may be underestimated. In the next section, we will further explore the contribution of trans-eQTLs. Second, the genetic effects on gene expression may not be steady across the reference GTEx data with largely non-diseased tissues for general purposes and the GWAS data with diseased individuals from specific populations (52). From this perspective, before analyzing specific complex traits and diseases via TWAS, it would be helpful to first estimate the impact of GREX and select the most informative available eQTL reference data.

Additional insights on GREX considering trans-eQTLs and genetically regulated alternative splicing events

The *cis*-acting genetic effects on local gene expression levels are often shared across tissue types and are often replicable across studies (10). It is also reported that a substantial proportion (up to 70%) of gene expression heritability can be attributed to trans-acting genetic effects which act predominantly in a tissue-specific manner and have a lower rate of replication across studies (53,54). More recently, the eQTLGen consortium has conducted a blood-eQTL meta-analysis and has reported 6,298 (31%) trans-eQTL genes for 10 317 trait-associated SNPs using 31 684 blood samples

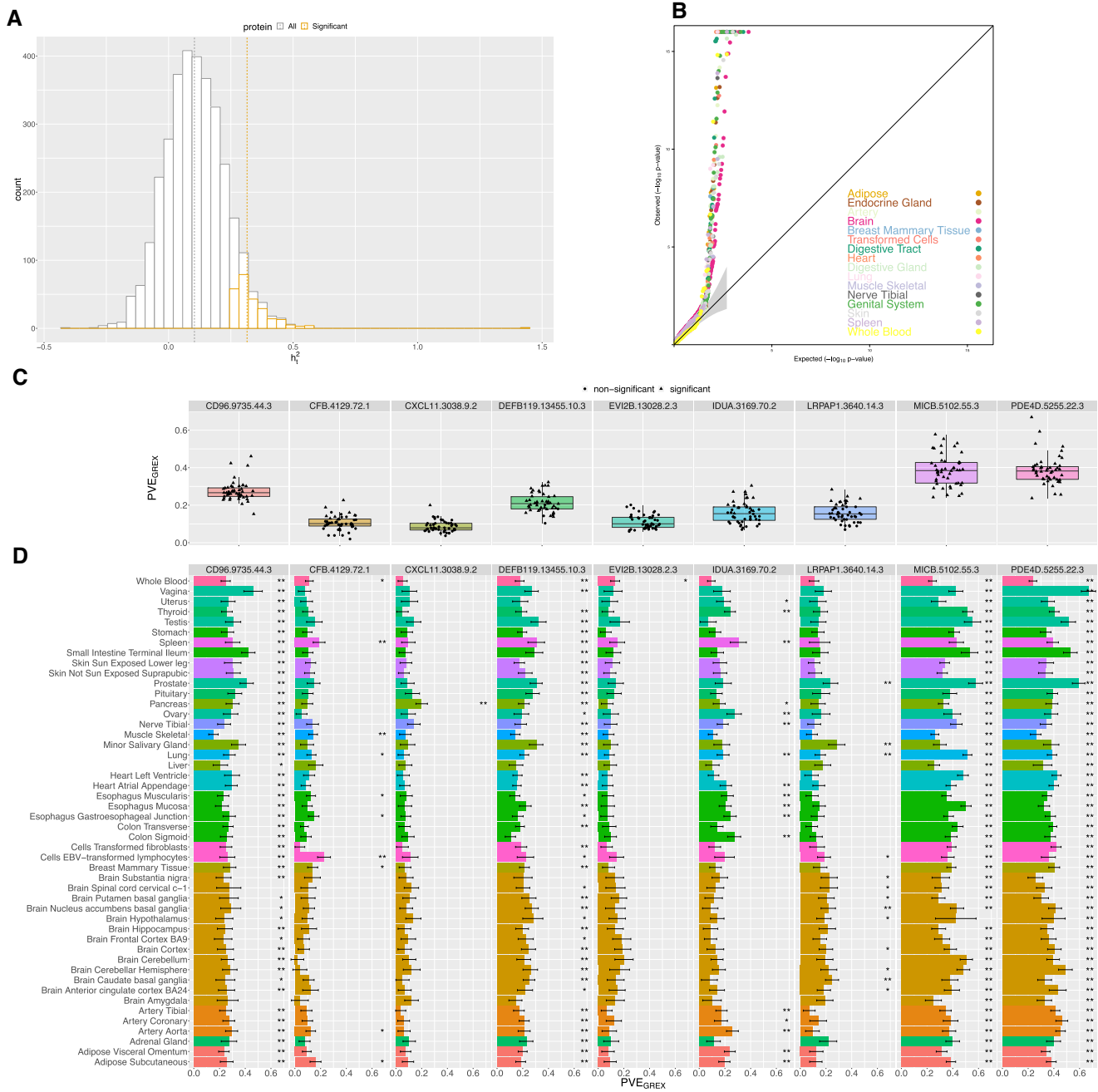


Figure 4. Analysis of plasma pQTL summary statistics. (A) The distribution of estimated heritabilities of 3283 proteins estimated using (44). The whole study is colored in gray, while the 249 proteins with significant heritabilities are colored in yellow. Dashed lines represent the means of corresponding distributions. (B) QQ-plot of PVE_{GREX} P -values of tissue-protein pairs. GTEx tissues are categorized into 16 types and colored accordingly. (C) PVE_{GREX} in the nine proteins whose PVE_{GREX} are significant in at least one tissue at 0.05 level using Bonferroni correction. (D) 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 **.

from 37 datasets. The results suggest that trans-eQTLs are prevalent in the genome, while it is still underpowered to detect them for tissues other than whole blood given the often tissue-specific nature of trans-genetic effects and the limited sample sizes for most tissue types.

Although it is still unrealistic to account for all trans-eQTLs in the estimation of PVE_{GREX} due to the limitation of sample sizes, it is possible to explore the potential by in-

corporating the blood-based trans-eQTLs reported by the eQTLGen consortium and re-estimating PVE_{GREX}/h^2 . We first analyzed 13 datasets comprised of 12 phenotypes that have significant PVE_{GREX}/h^2 estimates in the whole blood, including seven proteins, one lipid trait and four complex diseases (with two SCZ datasets). We observed an increasing trend of PVE_{GREX}/h^2 in the blood for all 13 datasets (Figure 6A), by accounting for only ~1700 unique trans-

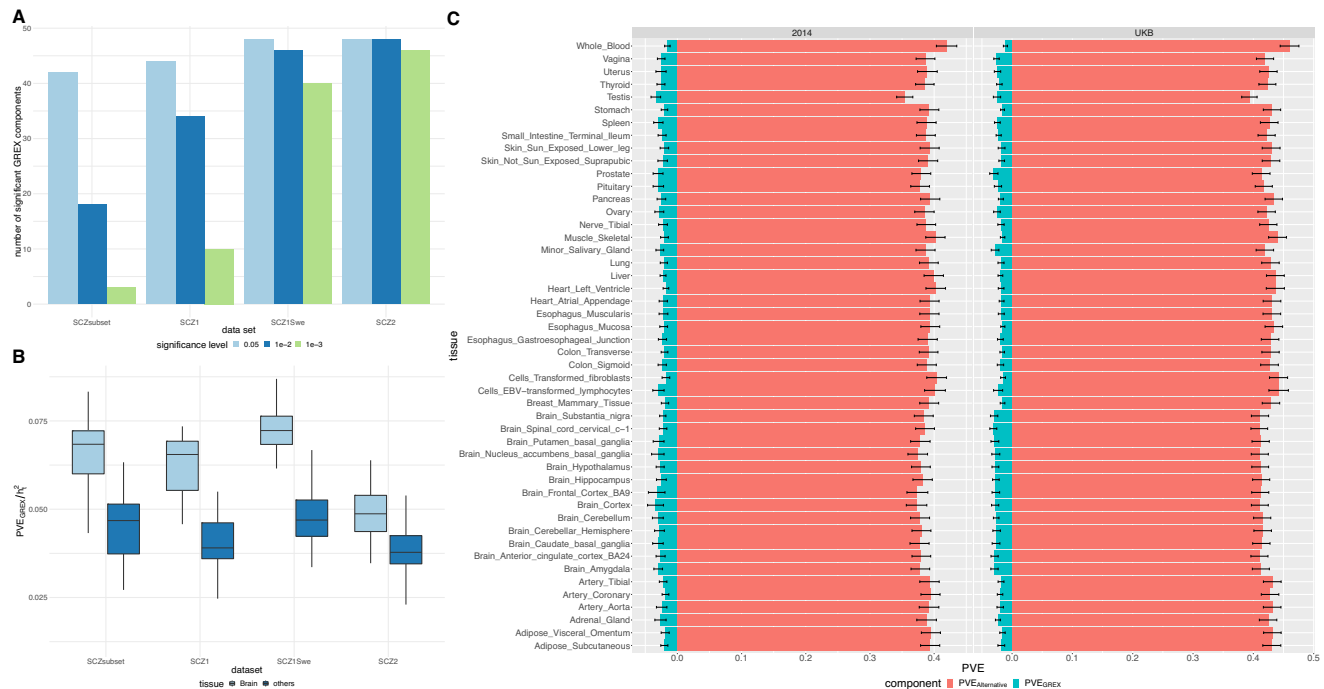


Figure 5. Analyses of complex traits: schizophrenia and height. (A) Number of significant GREX components revealed under different significance levels for the four schizophrenia datasets. (B) Estimated percentages of heritability for schizophrenia explained by GREX in brain tissues and in other tissues. (C) PVE_{GREX} and PVE_{Alternative} of height estimated using height2014 and UKB datasets, respectively.

eQTLs that are not *cis*-eQTLs. As a comparison, we applied the same procedure to 13 GTEx brain tissues of the two largest SCZ datasets, and did not observe an increase in PVE_{GREX}/*h*_T² (Figure 6B). This is not surprising because the trans-eQTLs incorporated above were detected and reported based on whole blood samples and may not be trans-eQTLs in the brain tissues. Our results suggest that the estimation of GREX impacts on traits can be further boosted by incorporating robust trans-eQTLs from the same tissue types.

In addition to the gene expression level, we also evaluated the effects of alternative splicing on complex trait heritability. We applied IGREX to quantify the impact of genetically regulated alternative splicing on multiple phenotypes. Alternative splicing is an important gene regulatory process that results in multiple transcripts from a single multi-exon gene. It is commonly observed in humans and plays an essential role in cellular differentiation (55,56). Differential variations in splicing may also result in phenotypic variation and contribute to the development of complex diseases including cancer (57–59). In a recent work, by extending the TWAS framework to analyze splicing events and associating 40 complex traits with genetically predicted splicing quantification, novel putative disease-associated genes were detected (60). Here, using multi-tissue splicing quantification data from GTEx as reference, we applied IGREX to study the impact of genetically regulated splicing events on four trait-tissue pairs that were found to have a high PVE_{GREX}/*h*_T². We estimated the proportion of phenotypic variation explained by genetically regulated splicing to be 12.5%, 13.5%, 1.0% and 1.1% for LDL in liver, TC in liver, SCZ in amygdala and SCZ in cerebellar hemisphere, respec-

tively. Unlike eQTLs that are often found to be near transcription starting sites, most of the sQTLs were found to be enriched within gene bodies, in particular within the introns they regulate, and have little to no effects on *cis* gene expression levels (60,61). In other words, sQTLs are often independent of eQTLs. Therefore, integrating genetically regulated splicing quantification may partially explain the phenotypic variation attributed to alternative genetic factors, PVE_{Alternative}. We argue that with the proper multi-omics reference data, similar analyses can be conducted to quantify the impact of genetically regulated methylation, protein, and other multi-omics variation on phenotype (56).

DISCUSSION

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

IGREX is closely related to several existing methods and here we briefly discuss the connections and distinctions. By also integrating an eQTL reference and GWAS data, methods including TWAS (15), PrediXcan (14) and the

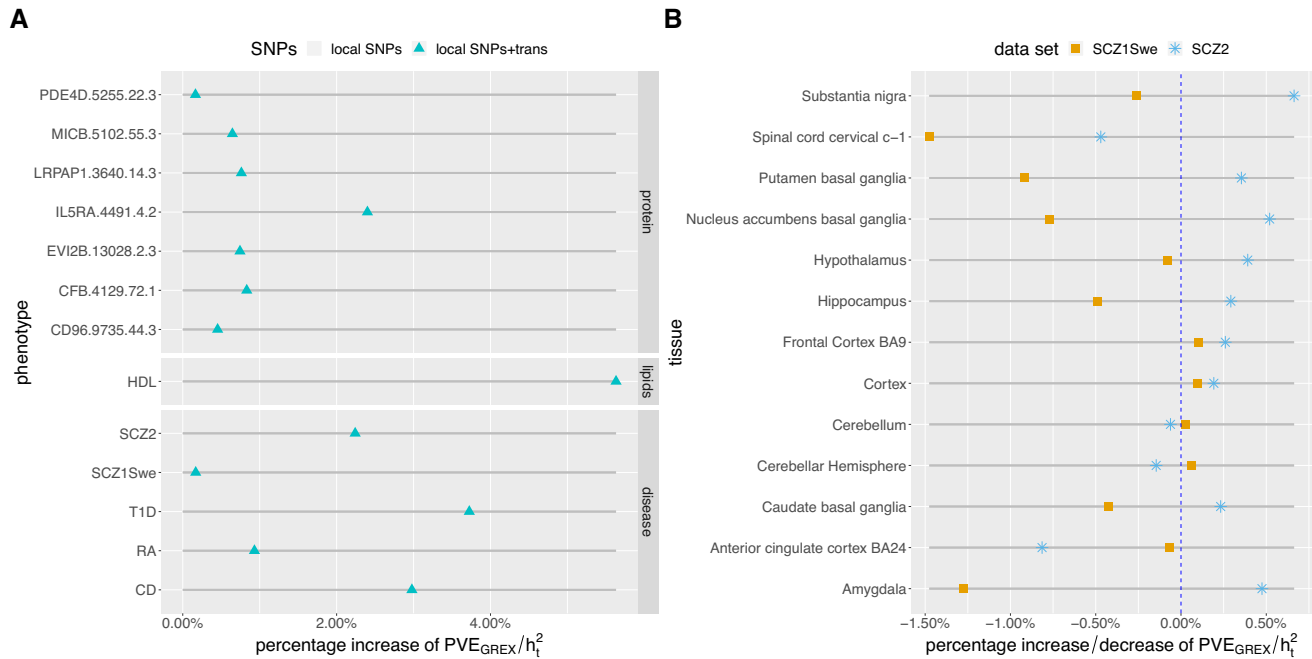


Figure 6. Percentage increase/decrease of PVE_{GREX}/h^2 estimates with *cis+trans* SNPs against those estimated with only *cis*-SNPs. (A) Percentage increase of 13 datasets in blood. All these datasets have significant PVE_{GREX}/h^2 in blood at 0.05 nominal level using only local SNPs. (B) Percentage increase/decrease of two largest SCZ datasets in 13 GTEx brain tissues. All these tissues have significant PVE_{GREX}/h^2 at 0.05 nominal level in both datasets using only local SNPs.

more general MetaXcan (18) aim to identify specific trait-associated genes. In contrast, IGREX estimates the impact of genetically regulated expression from a global perspective by quantifying the phenotypic variation that can be attributed to the GREX component. Since both the TWAS-type of analyses and IGREX rely on the shared GREX variation across eQTL and GWAS data, we argue that with the increasing availability of eQTL resources in different populations, conditions and contexts, the proper selection of eQTL reference panels via IGREX will greatly promote the chances of successes in the subsequent TWAS-type of analyses.

There are also existing methods, such as RhoGE, designed for identifying and estimating correlations between gene expression and complex traits. RhoGe provides an LDSC-based approach for estimating PVE_{GREX} . Unlike IGREX, this method does not adjust for estimation uncertainty. Consequently, it significantly underestimates the PVE_{GREX} when the eQTL effects on expression levels are weak or moderate. In fact, RhoGE estimated the PVE_{GREX} for the majority of 1350 tissue-trait pairs to be almost negligible, with the first quantile, the median, and the third quantile being 0.00125%, 0.162% and 0.616%, respectively (20). In contrast, as demonstrated via simulation studies, IGREX can accurately estimate PVE_{GREX} in various scenarios by accounting for the estimation uncertainty.

Based on estimating PVE_{GREX} for a wide-array of tissue-trait pairs, we observed a stronger impact of GREX on molecular intermediate traits and lipid traits in trait-relevant tissue types. We also observed a relatively low PVE_{GREX} for complex traits in general. The big picture suggests the attenuated impact on downstream phenotypes (e.g.

height and SCZ), which is consistent with the result from a pioneer study (62). However, we noted that the PVE_{GREX} estimates could be improved. A substantial amount of expression heritability is explained by trans-acting genetic factors while current TWAS and IGREX analyses are mainly using only *cis*-eQTLs. We explored the potential of incorporating trans-eQTLs in TWAS analysis by re-estimating PVE_{GREX} for selected traits in blood tissues with significant trans-eQTLs independently derived from the blood-based eQTLGen Consortium. We observed consistent increases in PVE_{GREX} for blood-related traits. In contrast, such an increase was not observed in the PVE_{GREX} estimates for other tissue types, again illustrating the importance of considering trait-relevant tissue types/conditions in the TWAS-type of analyses. Additionally, we extended the IGREX analysis to quantify the impact of genetically regulated alternative splicing events on selected traits. Our results suggested the potential for extending TWAS-type of analysis to integrate reference multi-omics QTL data with GWAS in mapping novel disease/trait-associated genes with mechanisms via other omics traits (such as splicing, methylation, protein, etc.).

IGREX is widely applicable for various GWAS phenotypes because it can handle the GWAS summary statistics. Given the increasing resources in eQTL study, it may be also desirable for developing methods that can handle the eQTL summary statistics, which can potentially boost the power of identifying the GREX component when the weaker distal eQTL effects are considered. One of the possible solution is to first derive the posterior distributions of β_g from the summary statistics of eQTL study using existing methods, and then plug the obtained distribution to the second

stage of IGREX. The RSS method (63) turns out to be a possible candidate for retrieving the posterior distribution of β_g from eQTL summary statistics. While this is a possible extension to IGREX, there will be some computational issues for practical applications and we leave it for further investigation in the future.

A key assumption in applying IGREX or TWAS methods with a general-purpose eQTL data as reference is the existence of steady-state component in GREX, i.e. the genetic effects on gene expression β_g are shared across the eQTL reference and GWAS data. However, there are many situations in which this assumption is violated. For example, it has been observed that CAD-risk SNPs have a larger overlap with *cis*-eQTLs isolated from disease-relevant tissues than those from GTEx tissues (52), implying the existence of a dynamic component. In the presence of this dynamic component, the accuracy of \widehat{PVE}_{GREX} based on GTEx is reduced. In those cases, we suggest exploring other trait-relevant or condition-specific eQTL reference panels using IGREX for a better understanding of the role of GREX and before conducting TWAS analysis.

DATA AVAILABILITY

The GTEx gene expression data were downloaded from GTEx Consortium website <https://gtexportal.org/home/datasets>. The GTEx genotype data can be accessed from dbGAP with accession number phs000424.v7.p2. The HapMap3 genotype data is available at <ftp://ftp.ncbi.nlm.nih.gov/hapmap/>. The NFBC study was downloaded from dbGAP using accession number phs000276.v1.p1. The WTCCC data was obtained from its consortium website https://www.wtccc.org.uk/info/access_to_data_samples.html. The GWAS summary statistics can be accessed using the links provided in Supplementary Table S2. The eQTL-Gen data can be downloaded from <http://www.eqtlgen.org>. The R software package IGREX is publicly available on GitHub repository: <https://github.com/mxcai/iGREX>.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

We thank Prof. Hongyu Zhao for helpful comments and discussions, and Mr Kevin J. Gleason for proof-reading the work. The computational work for this article was (fully or partially) performed on resources of the National Supercomputing Centre, Singapore (<https://www.nsc.sg>).

FUNDING

National Science Funding of China [61501389 to M.C., C.Y.]; Hong Kong Research Grant Council [12316116, 12301417, 16307818, 16301419 to M.C., C.Y.]; Hong Kong University of Science and Technology [R9405, IGN17SC02, Z0428 to M.C., C.Y.]; Duke-NUS Medical School [R-913-200-098-263 to J.L.]; Ministry of Education, Singapore, AcRF Tier 2 [MOE2016-T2-2-029, MOE2018-T2-1-046, MOE2018-T2-2-006 to J.L.]; National Institutes of Health [R01GM108711, U24CA210993-SUB to L.S.C.].

Conflict of interest statement. None declared.

REFERENCES

- Maurano,M.T., Humbert,R., Rynes,E., Thurman,R.E., Haugen,E., Wang,H., Reynolds,A.P., Sandstrom,R., Qu,H., Brody,J. *et al.* (2012) Systematic localization of common disease-associated variation in regulatory DNA. *Science*, **337**, 1190–1195.
- Cookson,W., Liang,L., Abecasis,G., Moffatt,M. and Lathrop,M. (2009) Mapping complex disease traits with global gene expression. *Nat. Rev. Genet.*, **10**, 184–194.
- Pomerantz,M.M., Ahmadiyeh,N., Jia,L., Herman,P., Verzi,M.P., Doddapaneni,H., Beckwith,C.A., Chan,J.A., Hills,A., Davis,M. *et al.* (2009) The 8q24 cancer risk variant rs6983267 shows long-range interaction with MYC in colorectal cancer. *Nat. Genet.*, **41**, 882–884.
- Musunuru,K., Strong,A., Frank-Kamenetsky,M., Lee,N.E., Ahfeldt,T., Sachs,K.V., Li,X., Li,H., Kuperwasser,N., Ruda,V.M. *et al.* (2010) From noncoding variant to phenotype via SORT1 at the 1p13 cholesterol locus. *Nature*, **466**, 714–719.
- Harismendy,O., Notani,D., Song,X., Rahim,N.G., Tanasa,B., Heintzman,N., Ren,B., Fu,X.-D., Topol,E.J., Rosenfeld,M.G. *et al.* (2011) 9p21 DNA variants associated with coronary artery disease impair interferon- γ signalling response. *Nature*, **470**, 264–268.
- Nicolae,D.L., Gamazon,E., Zhang,W., Duan,S., Dolan,M.E. and Cox,N.J. (2010) Trait-associated SNPs are more likely to be eQTLs: annotation to enhance discovery from GWAS. *PLoS Genet.*, **6**, e1000888.
- Hindorf,L.A., Sethupathy,P., Junkins,H.A., Ramos,E.M., Mehta,J.P., Collins,F.S. and Manolio,T.A. (2009) Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl. Acad. Sci. U.S.A.*, **106**, 9362–9367.
- Albert,F.W. and Kruglyak,L. (2015) The role of regulatory variation in complex traits and disease. *Nat. Rev. Genet.*, **16**, 197–212.
- Gamazon,E.R., Segrè,A.V., van de Bunt,M., Wen,X., Xi,H.S., Hormozdiari,F., Ongen,H., Konkashbaev,A., Derks,E.M., Aguet,F. *et al.* (2018) Using an atlas of gene regulation across 44 human tissues to inform complex disease-and trait-associated variation. *Nat. Genet.*, **50**, 956–967.
- GTEx Consortium (2017) Genetic effects on gene expression across human tissues. *Nature*, **550**, 204–213.
- Lloyd-Jones,L.R., Holloway,A., McRae,A., Yang,J., Small,K., Zhao,J., Zeng,B., Bakshi,A., Metspalu,A., Dermizakis,M. *et al.* (2017) The genetic architecture of gene expression in peripheral blood. *Am. J. Hum. Genet.*, **100**, 228–237.
- Westra,H.-J., Peters,M.J., Esko,T., Yaghootkar,H., Schurmann,C., Kettunen,J., Christiansen,M.W., Fairfax,B.P., Schramm,K., Powell,J.E. *et al.* (2013) Systematic identification of trans eQTLs as putative drivers of known disease associations. *Nat. Genet.*, **45**, 1238–1243.
- Qi,T., Wu,Y., Zeng,J., Zhang,F., Xue,A., Jiang,L., Zhu,Z., Kemper,K., Yengo,L., Zheng,Z. *et al.* (2018) Identifying gene targets for brain-related traits using transcriptomic and methylomic data from blood. *Nat. Commun.*, **9**, 2282–2282.
- Gamazon,E.R., Wheeler,H.E., Shah,K.P., Mozaffari,S.V., Aquino-Michaels,K., Carroll,R.J., Eyler,A.E., Denny,J.C., Nicolae,D.L., Cox,N.J. *et al.* (2015) A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.*, **47**, 1091–1098.
- Gusev,A., Ko,A., Shi,H., Bhatia,G., Chung,W., Penninx,B.W., Jansen,R., De Geus,E.J., Boomsma,D.I., Wright,F.A. *et al.* (2016) Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.*, **48**, 245–252.
- Zhu,Z., Zhang,F., Hu,H., Bakshi,A., Robinson,M.R., Powell,J.E., Montgomery,G.W., Goddard,M.E., Wray,N.R., Visscher,P.M. *et al.* (2016) Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.*, **48**, 481–487.
- Mancuso,N., Freund,M.K., Johnson,R., Shi,H., Kichaev,G., Gusev,A. and Pasaniuc,B. (2019) Probabilistic fine-mapping of transcriptome-wide association studies. *Nat. Genet.*, **51**, 675–682.
- Barbeira,A.N., Dickinson,S.P., Bonazzola,R., Zheng,J., Wheeler,H.E., Torres,J.M., Torstenson,E.S., Shah,K.P., Garcia,T., Edwards,T.L. *et al.* (2018) Exploring the phenotypic consequences of

- tissue specific gene expression variation inferred from GWAS summary statistics. *Nat. Commun.*, **9**, 1825–1825.
19. Yang, C., Wan, X., Lin, X., Chen, M., Zhou, X. and Liu, J. (2018) CoMM: a collaborative mixed model to dissecting genetic contributions to complex traits by leveraging regulatory information. *Bioinformatics*, **35**, 1644–1652.
 20. Mancuso, N., Shi, H., Goddard, P., Kichaev, G., Gusev, A. and Pasaniuc, B. (2017) Integrating gene expression with summary association statistics to identify genes associated with 30 complex traits. *Am. J. Hum. Genet.*, **100**, 473–487.
 21. Schaid, D.J., Chen, W. and Larson, N.B. (2018) From genome-wide associations to candidate causal variants by statistical fine-mapping. *Nat. Rev. Genet.*, **19**, 491–504.
 22. Hu, Y., Li, M., Lu, Q., Weng, H., Wang, J., Zekavat, S.M., Yu, Z., Li, B., Gu, J., Muchnik, S. *et al.* (2019) A statistical framework for cross-tissue transcriptome-wide association analysis. *Nat. Genet.*, **51**, 568–576.
 23. Wainberg, M., Sinnott-Armstrong, N., Mancuso, N., Barbeira, A.N., Knowles, D.A., Golan, D., Ermel, R., Ruusalepp, A., Quertermous, T., Hao, K. *et al.* (2019) Opportunities and challenges for transcriptome-wide association studies. *Nat. Genet.*, **51**, 592–599.
 24. Wang, K., Gaitsch, H., Poon, H., Cox, N.J. and Rzhetsky, A. (2017) Classification of common human diseases derived from shared genetic and environmental determinants. *Nat. Genet.*, **49**, 1319–1325.
 25. Lakhani, C.M., Tierney, B.T., Manrai, A.K., Yang, J., Visscher, P.M. and Patel, C.J. (2019) Repurposing large health insurance claims data to estimate genetic and environmental contributions in 560 phenotypes. *Nat. Genet.*, **51**, 327–334.
 26. Bulik-Sullivan, B.K., Loh, P.-R., Finucane, H.K., Ripke, S., Yang, J., Patterson, N., Daly, M.J., Price, A.L., Neale, B.M., Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* (2015) LD score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.*, **47**, 291–295.
 27. Sabatti, C., Hartikainen, A.-L., Pouta, A., Ripatti, S., Brodsky, J., Jones, C.G., Zaitlen, N.A., Varilo, T., Kaakinen, M., Sovio, U. *et al.* (2009) Genome-wide association analysis of metabolic traits in a birth cohort from a founder population. *Nat. Genet.*, **41**, 35–46.
 28. Zhou, X. and Stephens, M. (2014) Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nat. Methods*, **11**, 407–409.
 29. Wellcome, Trust Case Control Consortium (2007) Genome-wide association study of 14,000 cases of seven common diseases and 3,000 shared controls. *Nature*, **447**, 661–678.
 30. Dai, M., Ming, J., Cai, M., Liu, J., Yang, C., Wan, X. and Xu, Z. (2017) IGESS: a statistical approach to integrating individual-level genotype data and summary statistics in genome-wide association studies. *Bioinformatics*, **33**, 2882–2889.
 31. Cai, M., Dai, M., Ming, J., Peng, H., Liu, J. and Yang, C. (2019) BIVAS: a scalable Bayesian method for bi-level variable selection with applications. *J. Comput. Graph. Statist.*, 1–38.
 32. Sun, B.B., Maranville, J.C., Peters, J.E., Stacey, D., Staley, J.R., Blackshaw, J., Burgess, S., Jiang, T., Paige, E., Surendran, P. *et al.* (2018) Genomic atlas of the human plasma proteome. *Nature*, **558**, 73–79.
 33. Kettunen, J., Demirkan, A., Würtz, P., Draisma, H.H., Haller, T., Rawal, R., Vaarhorst, A., Kangas, A.J., Lytikäinen, L.-P., Pirinen, M. *et al.* (2016) Genome-wide study for circulating metabolites identifies 62 loci and reveals novel systemic effects of LPA. *Nat. Commun.*, **7**, 11122–11122.
 34. Cross Disorder Group of the Psychiatric Genomics Consortium, (2013) Identification of risk loci with shared effects on five major psychiatric disorders: a genome-wide analysis. *Lancet*, **381**, 1360–1360.
 35. Ripke, S., Sanders, A., Kendler, K., Levinson, D., Sklar, P., Holmans, P., Lin, D., Duan, J., Ophoff, R., Andreassen, O. *et al.* (2011) Schizophrenia Psychiatric Genome-Wide Association Study (GWAS) Consortium Genome-wide association study identifies five new schizophrenia loci. *Nat. Genet.*, **43**, 969–976.
 36. Ripke, S., O’Dushlaine, C., Chambert, K., Moran, J.L., Kähler, A.K., Akterin, S., Bergen, S.E., Collins, A.L., Crowley, J.J., Fromer, M. *et al.* (2013) Genome-wide association analysis identifies 13 new risk loci for schizophrenia. *Nat. Genet.*, **45**, 1150–1159.
 37. Ripke, S., Neale, B.M., Corvin, A., Walters, J.T., Farh, K.-H., Holmans, P.A., Lee, P., Bulik-Sullivan, B., Collier, D.A., Huang, H. *et al.* (2014) Biological insights from 108 schizophrenia-associated genetic loci. *Nature*, **511**, 421–427.
 38. Wood, A.R., Esko, T., Yang, J., Vedantam, S., Pers, T.H., Gustafsson, S., Chu, A.Y., Estrada, K., Luan, J., Kutalik, Z. *et al.* (2014) Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat. Genet.*, **46**, 1173–1186.
 39. Winkler, T.W., Justice, A.E., Graff, M., Barata, L., Feitosa, M.F., Chu, S., Czajkowski, J., Esko, T., Fall, T., Kilpeläinen, T.O. *et al.* (2015) The influence of age and sex on genetic associations with adult body size and shape: a large-scale genome-wide interaction study. *PLoS Genet.*, **11**, e1005378.
 40. Jiang, J., Li, C., Paul, D., Yang, C. and Zhao, H. (2016) On high-dimensional misspecified mixed model analysis in genome-wide association study. *Ann. Statist.*, **44**, 2127–2160.
 41. Liu, C., Rubin, D.B. and Wu, Y.N. (1998) Parameter expansion to accelerate EM: the PX-EM algorithm. *Biometrika*, **85**, 755–770.
 42. Zhou, H., Hu, L., Zhou, J. and Lange, K. (2019) MM algorithms for variance components models. *J. Comput. Graph. Statist.*, 1–12.
 43. Davies, R.B. (1980) The distribution of a linear combination of χ^2 random variables. *J. R. Stat. Soc. C (Applied Statistics)*, **29**, 323–333.
 44. Zhou, X. (2017) A unified framework for variance component estimation with summary statistics in genome-wide association studies. *Ann. Appl. Statist.*, **11**, 2027–2051.
 45. Quenouille, M.H. (1956) Notes on bias in estimation. *Biometrika*, **43**, 353–360.
 46. Wheeler, H.E., Shah, K.P., Brenner, J., Garcia, T., Aquino-Michaels, K., Cox, N.J., Nicolae, D.L., Im, H.K. and Consortium, G. (2016) Survey of the heritability and sparse architecture of gene expression traits across human tissues. *PLoS Genet.*, **12**, e1006423.
 47. Dietschy, J.M., Turley, S.D. and Spady, D.K. (1993) Role of liver in the maintenance of cholesterol and low density lipoprotein homeostasis in different animal species, including humans. *J. Lipid Res.*, **34**, 1637–1659.
 48. Kovanen, P.T., Brown, M.S. and Goldstein, J.L. (1979) Increased binding of low density lipoprotein to liver membranes from rats treated with 17 alpha-ethinyl estradiol. *J. Biol. Chem.*, **254**, 11367–11373.
 49. Feng, T. and Zhu, X. (2010) Genome-wide searching of rare genetic variants in WTCCC data. *Hum. Genet.*, **128**, 269–280.
 50. Cole, K.E., Strick, C.A., Paradis, T.J., Ogborne, K.T., Loetscher, M., Gladue, R.P., Lin, W., Boyd, J.G., Moser, B., Wood, D.E. *et al.* (1998) Interferon-inducible T Cell Alpha Chemoattractant (I-TAC): A Novel Non-ELR CXC Chemokine with Potent Activity on Activated T Cells through Selective High Affinity Binding to CXCR3. *J. Exp. Med.*, **187**, 2009–2021.
 51. Yao, C., Joehanes, R., Johnson, A.D., Huan, T., Liu, C., Freedman, J.E., Munson, P.J., Hill, D.E., Vidal, M. and Levy, D. (2017) Dynamic role of trans regulation of gene expression in relation to complex traits. *Am. J. Hum. Genet.*, **100**, 571–580.
 52. Franzén, O., Ermel, R., Cohain, A., Akers, N.K., Di Narzo, A., Talukdar, H.A., Foroughi-Asl, H., Giambartolomei, C., Fullard, J.F., Sukhvasi, K. *et al.* (2016) Cardiometabolic risk loci share downstream cis- and trans-gene regulation across tissues and diseases. *Science*, **353**, 827–830.
 53. Grundberg, E., Small, K.S., Hedman, Å.K., Nica, A.C., Buil, A., Keildson, S., Bell, J.T., Yang, T.-P., Meduri, E., Barrett, A. *et al.* (2012) Mapping cis- and trans-regulatory effects across multiple tissues in twins. *Nat. Genet.*, **44**, 1084–1089.
 54. Wright, F.A., Sullivan, P.F., Brooks, A.I., Zou, F., Sun, W., Xia, K., Madar, V., Jansen, R., Chung, W., Zhou, Y.-H. *et al.* (2014) Heritability and genomics of gene expression in peripheral blood. *Nat. Genet.*, **46**, 430–437.
 55. Gilbert, W. (1978) Why genes in pieces? *Nature*, **271**, 501–501.
 56. Matlin, A.J., Clark, F. and Smith, C.W. (2005) Understanding alternative splicing: towards a cellular code. *Nat. Rev. Mol. Cell Biol.*, **6**, 386–398.
 57. Li, Y.I., van de Geijn, B., Raj, A., Knowles, D.A., Petti, A.A., Golan, D., Gilad, Y. and Pritchard, J.K. (2016) RNA splicing is a primary link between genetic variation and disease. *Science*, **352**, 600–604.
 58. Takata, A., Matsumoto, N. and Kato, T. (2017) Genome-wide identification of splicing QTLs in the human brain and their enrichment among schizophrenia-associated loci. *Nat. Commun.*, **8**, 14519–14519.
 59. Skotheim, R.I. and Nees, M. (2007) Alternative splicing in cancer: noise, functional, or systematic? *Int. J. Biochem. Cell Biol.*, **39**, 1432–1449.

60. Li, Y.I., Knowles, D.A., Humphrey, J., Barbeira, A.N., Dickinson, S.P., Im, H.K. and Pritchard, J.K. (2018) Annotation-free quantification of RNA splicing using LeafCutter. *Nat. Genet.*, **50**, 151–158.
61. Gutierrez-Arcelus, M., Ongen, H., Lappalainen, T., Montgomery, S.B., Buil, A., Yurovsky, A., Bryois, J., Padioleau, I., Romano, L., Planchon, A. *et al.* (2015) Tissue-specific effects of genetic and epigenetic variation on gene regulation and splicing. *PLoS Genet.*, **11**, e1004958.
62. Battle, A., Khan, Z., Wang, S.H., Mitrano, A., Ford, M.J., Pritchard, J.K. and Gilad, Y. (2014) Impact of regulatory variation from RNA to protein. *Science*, **347**, 664–667.
63. Zhu, X. and Stephens, M. (2017) Bayesian large-scale multiple regression with summary statistics from genome-wide association studies. *Ann. Appl. Statist.*, **11**, 1561–1592.