

# Functional Architectures of Local and Distal Regulation of Gene Expression in Multiple Human Tissues

Xuanyao Liu,<sup>1,\*</sup> Hilary K. Finucane,<sup>1,2</sup> Alexander Gusev,<sup>1</sup> Gaurav Bhatia,<sup>1</sup> Steven Gazal,<sup>1</sup> Luke O'Connor,<sup>1,3</sup> Brendan Bulik-Sullivan,<sup>4,5,6</sup> Fred A. Wright,<sup>7</sup> Patrick F. Sullivan,<sup>8,9,10</sup> Benjamin M. Neale,<sup>4,5,6</sup> and Alkes L. Price<sup>1,4,11,\*</sup>

Genetic variants that modulate gene expression levels play an important role in the etiology of human diseases and complex traits. Although large-scale eQTL mapping studies routinely identify many local eQTLs, the molecular mechanisms by which genetic variants regulate expression remain unclear, particularly for distal eQTLs, which these studies are not well powered to detect. Here, we leveraged all variants (not just those that pass stringent significance thresholds) to analyze the functional architecture of local and distal regulation of gene expression in 15 human tissues by employing an extension of stratified LD-score regression that produces robust results in simulations. The top enriched functional categories in local regulation of peripheral-blood gene expression included coding regions (11.41×), conserved regions (4.67×), and four histone marks ( $p < 5 \times 10^{-5}$  for all enrichments); local enrichments were similar across the 15 tissues. We also observed substantial enrichments for distal regulation of peripheral-blood gene expression: coding regions (4.47×), conserved regions (4.51×), and two histone marks ( $p < 3 \times 10^{-7}$  for all enrichments). Analyses of the genetic correlation of gene expression across tissues confirmed that local regulation of gene expression is largely shared across tissues but that distal regulation is highly tissue specific. Our results elucidate the functional components of the genetic architecture of local and distal regulation of gene expression.

## Introduction

Our understanding of the functional elements of the human genome has benefitted greatly from the explosion of functional data generated by the ENCODE project and the Roadmap Epigenomics Consortium.<sup>1,2</sup> In particular, researchers have gained new insights into the functional effects of genetic variants on many complex diseases and traits.<sup>3–12</sup> In parallel, large-scale expression quantitative trait locus (eQTL) mapping studies in multiple human tissues have revealed a large number of genetic variants that affect gene expression<sup>13–19</sup> (reviewed by Albert and Kruglyak<sup>20</sup>). Gene expression serves as an important intermediate cellular phenotype that affects complex diseases and traits,<sup>21–24</sup> and the functional effects of eQTLs provide another lens through which researchers can investigate molecular mechanisms.<sup>9,13–20,25–27</sup>

However, the underlying functional mechanisms of eQTLs are still largely unclear. On one hand, previous studies have produced different functional characterizations of local eQTLs (Table S1), possibly because of differences in the sets of annotations analyzed and/or the sample-size dependence of approaches that assess enrichment by using only top eQTLs. On the other hand, functional characterization of distal eQTLs has been limited,<sup>15,16</sup> given that most studies are under-powered to detect distal eQTLs.

In this study, we extended a recently developed method, stratified linkage disequilibrium (LD)-score regression,<sup>10</sup> to partition the heritability of local and distal regulation of gene expression across different functional categories. Stratified LD-score regression makes use of summary association statistics of all genetic variants (not just the top significant variants) and estimates the heritability explained by each functional category while accounting for LD to other functional categories; this approach is more powerful than other methods for detecting functional enrichment (Figure 7 from Finucane et al.<sup>10</sup>). We extended the method to produce aggregate estimates across all genes for both local and distal regulation of gene expression; our current simulations confirm that this extension to gene expression data produces robust enrichment results. By applying this method to large gene expression datasets in multiple human tissues, we aimed to comprehensively assess the functional enrichments of genetic variants on local and distal regulation of gene expression and shed light on the underlying molecular mechanisms.

## Material and Methods

### Gene Expression Datasets

We analyzed gene expression in 15 human tissues: peripheral blood from Wright et al.,<sup>16</sup> 11 tissues with a sample size larger than 200

<sup>1</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA; <sup>2</sup>Department of Mathematics, Massachusetts Institute of Technology, Cambridge, MA 02142, USA; <sup>3</sup>Program in Bioinformatics and Integrative Genomics, Harvard University, Boston, MA 02115, USA; <sup>4</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA; <sup>5</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA; <sup>6</sup>Analytic and Translational Genetics Unit, Department of Medicine, Massachusetts General Hospital, Boston, MA 02114, USA; <sup>7</sup>Bioinformatics Research Center, Departments of Statistics and Biological Sciences, North Carolina State University, Raleigh, NC 27695, USA; <sup>8</sup>Department of Genetics, University of North Carolina, Chapel Hill, NC 27599, USA; <sup>9</sup>Department of Psychiatry, University of North Carolina, Chapel Hill, NC 27599, USA; <sup>10</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institute, Stockholm 17177, Sweden; <sup>11</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA

\*Correspondence: [xuli@hsph.harvard.edu](mailto:xuli@hsph.harvard.edu) (X.L.), [aprice@hsph.harvard.edu](mailto:aprice@hsph.harvard.edu) (A.L.P.)

<http://dx.doi.org/10.1016/j.ajhg.2017.03.002>

© 2017 American Society of Human Genetics.

**Table 1. Gene Expression Datasets**

Dataset	Tissue	Sample Size	Data Type	No. of SNPs	No. of Probes and/or Genes
Wright et al. <sup>16</sup>	peripheral blood	3,754	expression array	1,142,515	42,044
GTEX <sup>19</sup>	adipose subcutaneous	298	RNA-seq	1,145,366	26,213
GTEX <sup>19</sup>	artery tibial	285	RNA-seq	1,141,287	24,383
GTEX <sup>19</sup>	cells – transformed fibroblasts	272	RNA-seq	1,145,366	22,963
GTEX <sup>19</sup>	esophagus mucosa	241	RNA-seq	1,131,019	25,070
GTEX <sup>19</sup>	esophagus muscularis	218	RNA-seq	1,130,356	24,416
GTEX <sup>19</sup>	lung	278	RNA-seq	1,144,671	27,671
GTEX <sup>19</sup>	muscle skeletal	361	RNA-seq	1,136,801	23,109
GTEX <sup>19</sup>	nerve tibial	256	RNA-seq	1,145,068	26,808
GTEX <sup>19</sup>	skin – sun exposed	302	RNA-seq	1,147,848	26,849
GTEX <sup>19</sup>	thyroid	278	RNA-seq	1,147,844	27,497
GTEX <sup>19</sup>	whole blood	338	RNA-seq	1,114,337	23,164
MuTHER <sup>13</sup>	adipose	776	expression array	878,954	22,058
MuTHER <sup>13</sup>	skin	667	expression array	878,954	22,058
MuTHER <sup>13</sup>	LCL	777	expression array	878,954	22,058

We analyzed gene expression data spanning 15 human tissues from three datasets. For each tissue, we list the sample size, data type, number of SNPs analyzed, and number of probes and/or genes analyzed. We note that stratified LD-score regression restricts to HapMap 3 SNPs from the target dataset as a proxy for SNPs with high-quality imputation.

from the Genotype-Tissue Expression (GTEx) project,<sup>19</sup> and adipose, skin, and lymphoblastoid cell lines (LCLs) from the MuTHER cohort<sup>13</sup> (Table 1). Our analyses required summary association statistics for genome-wide SNPs. For the Wright et al. dataset, we used summary statistics computed from the Netherlands Twin Registry (NTR) and Netherlands Study of Depression and Anxiety (NESDA) cohorts.<sup>16</sup> Genotype and expression quality control and genotype imputation were performed as previously described.<sup>16</sup> Probe sequences were mapped to the human genome (UCSC Genome Browser hg19), and probes with sequences that did not map, mapped to multiple locations, or overlapped a polymorphic SNP (HapMap 3 and 1000 Genomes Project data) were removed. For each gene, multiple probes were included if they passed quality control. For the NTR cohort, t-statistics were computed for each equally split twin set, and combined Z statistics were calculated with empirical correlations among monozygotic and dizygotic twins as previously described.<sup>28</sup> Meta-analyzed Z statistics for the NTR and NESDA cohorts were computed with inverse-variance weighting by sample size. For the GTEx dataset, we used version 6 of the publicly available GTEx summary statistics in local regions<sup>19</sup> (see Web Resources). Genotype and expression quality control was performed as previously described.<sup>19</sup> Only reads that were uniquely mapped, had proper pairs, and were contained 100% within exon boundaries were included in gene-level read count. One transcript per gene was used in our analyses. For the MuTHER dataset, we recomputed local and distal summary statistics as described in Grundberg et al.<sup>13</sup> The use of raw genotypes and expression profiles was approved by the King's College London Department of Twin Registry. Quality control of genotype and expression is described in Grundberg et al.<sup>13</sup> Only uniquely mapping probes with no mismatches and either an Ensembl or RefSeq ID were retained for analysis. Probes encompassing a polymorphic SNP (1000 Genomes Project release June 2010) were excluded. For

each gene, multiple probes were included if they passed quality control. Summary statistics were calculated with a two-step mixed-model-based score test with the GenABEL and ProbABEL packages<sup>29,30</sup> (see Web Resources). The first step fits a mixed model. The fixed effects include age and batch for adipose and LCLs and age, batch, and sample processing for skin. We built the kinship matrix by randomly choosing 10,000 SNPs from the dataset. This step was performed with the “polygenic()” function of the GenABEL software. The second step performs a score test by using the ProbABEL software. This step was performed with the `-mmscore` option of the ProbABEL software.

### Baseline Functional Categories

The 57 functional categories that we analyzed consist of 53 baseline categories<sup>10</sup> (derived from 24 main annotations) that were determined to be important for complex traits and an additional four categories (derived from two additional main annotations). The 26 main annotations were collected from various sources<sup>2,5,6,31–37</sup> and included coding regions, untranslated regions (UTRs), promoters, intronic regions, histone marks, DNase I hypersensitivity sites (DHSs), predicted enhancers, conserved regions, and other annotations (see below). We derived the 57 categories from the 26 main annotations by (1) adding a 500 bp window around each main annotation as an additional category to keep heritability estimates from being inflated by heritability in flanking regions (see Finucane et al.<sup>10</sup>), (2) adding 100 bp windows around chromatin immunoprecipitation sequencing (ChIP-seq) peaks for DHS, H3K4me1, H3K4me3, and H3K9ac annotations, and (3) adding a category containing all SNPs.

All 57 functional categories are publicly available (see Web Resources). The 26 main annotations are described in detail in Table S2. We briefly describe a representative set of 14 main annotations,

ordered by annotation size (percentage of SNPs in the 1000 Genomes European reference genome), that are included in our main figures (based on analyses that include all 57 annotations). 5' UTRs (0.5% of SNPs) and coding regions (1.5%) were derived from RefSeq gene models and post-processed as previously described.<sup>6</sup> Transcription starting sites (TSSs, 1.9%) included combined chromHMM and Segway annotations for six cell lines obtained from Hoffman et al.<sup>34</sup> Conserved regions (2.6%) in mammals were obtained from Lindblad-Toh et al.<sup>32</sup> and post-processed as previously described.<sup>33</sup> Promoters (3.1%) were also derived from RefSeq gene models and post-processed as previously described.<sup>6</sup> Enhancers (6.3%) were combined chromHMM and Segway annotations for six cell lines obtained from Hoffman et al.<sup>34</sup> H3K9ac annotations (12.6%) were a union across cell types, and the H3K9ac marks for each cell type were obtained from Roadmap Epigenomics<sup>1</sup> and post-processed as previously described.<sup>5</sup> Transcription factor binding sites (TFBSs, 13.2%) were obtained from ENCODE<sup>2</sup> and post-processed as previously described.<sup>6</sup> H3K4me3 annotations (13.3%) were a union across cell types, and the H3K4me3 marks for each cell type were obtained from Roadmap Epigenomics<sup>1</sup> and post-processed as previously described.<sup>5</sup> DHSs (16.8%) were a union across cell types, and the DHSs for each cell type were obtained from ENCODE<sup>2</sup> and Roadmap<sup>1</sup> and post-processed as previously described.<sup>5</sup> Super enhancers (Hnisz) (16.8%) were a union across cell types and a subset of closely spaced H3K27ac annotations from Hnisz et al.,<sup>36</sup> given that super enhancers generally refer to sets of enhancers in close genomic proximity.<sup>38</sup> H3K27ac (PGC2) marks (26.9%) were obtained from Roadmap<sup>1</sup> and post-processed as previously described.<sup>37</sup> H3K4me1 annotations (42.7%) were a union across cell types, and the H3K4me1 marks for each cell type were obtained from Roadmap<sup>1</sup> and post-processed as previously described.<sup>5</sup> Repressed annotations (46.1%) were an intersection of chromHMM and Segway annotations from six cell types.<sup>34</sup> We finally note that the two additional main annotations (not included in Finucane et al.<sup>10</sup>) consisted of super enhancers and typical enhancers from Vahedi et al.<sup>39</sup>

### Extension of Stratified LD-Score Regression

In a simple linear model,

$$y_i = \sum_j X_{ij} \beta_j + \varepsilon_i, \quad (\text{Equation 1})$$

where  $y_i$  is a quantitative phenotype in individual  $i$ ,  $X_{ij}$  is the standardized genotype of individual  $i$  at SNP  $j$ ,  $\beta_j$  is the effect size of SNP  $j$ , and  $\varepsilon_i$  is mean-zero noise. The total SNP heritability is defined as

$$h_g^2(\text{total}) = \sum_j \beta_j^2, \quad (\text{Equation 2})$$

and the SNP heritability of category  $C$  is defined as

$$h_g^2(C) = \sum_{j \in C} \beta_j^2. \quad (\text{Equation 3})$$

Stratified LD-score regression<sup>10</sup> (see [Web Resources](#)) relies on the fact that LD to a functional category enriched with heritability will increase the  $\chi^2$  association statistics of a SNP more than LD to other categories. More precisely,

$$E[\chi^2] = N \sum_C \tau_C l(j, C) + Na + 1, \quad (\text{Equation 4})$$

where  $N$  is the sample size,  $l(j, C)$  is the LD score of SNP  $j$  to category  $C$ , defined as  $l(j, C) = \sum_{k \in C} r^2(j, k)$ , and  $a$  measures the contribution

of confounding biases. (In this study, we employed constrained-intercept LD-score regression,<sup>40</sup> in which  $a$  is fixed at 0.) Performing multiple linear regression of  $\chi^2$  on  $l(j, C)$  gives us an estimate  $\widehat{\tau}_C$  of the coefficient  $\tau_C$ , which represents the per-SNP contribution to heritability of each category  $C$ . We estimate  $h_g^2(C)$  via

$$\widehat{h}_g^2(C) = \sum_{j \in C} \widehat{\text{Var}}(\beta_j) = \sum_{j \in C} \sum_{C' \cap C} \widehat{\tau}_{C'}. \quad (\text{Equation 5})$$

We applied stratified LD-score regression for both local and distal regions of each gene. We defined local regions as the regions within 1 Mb of the TSS of each gene and defined distal regions as the rest of the genome (to be consistent with previous studies<sup>16</sup>). To test whether the definition of local regions would affect our results, we also considered a different definition of local regions (within 2 Mb of the TSS) and determined that the estimates were not sensitive to this choice (see [Results](#)). We used the 1000 Genomes (phase 1) Europeans<sup>41</sup> as a reference panel to calculate LD scores. Thus, the LD score  $l(j, C)$  for regression SNP  $j$  is computed with reference SNP  $k$  from 1000 Genomes with minor allele count  $> 5$ . In local and distal analyses, reference SNPs were restricted to SNPs in local and distal regions, respectively, and LD scores were calculated as  $l(j, C) = \sum_{k \in C, k \text{ in local(distal) regions}} r^2(j, k)$ , such that we did not include the effects of SNPs outside the local and distal regions, respectively. Although we had access to individual-level genotype data, we used 1000 Genomes instead of in-sample LD as a reference panel to calculate LD scores because stratified LD-score regression requires LD scores computed with all 1000 Genomes reference SNPs with minor allele count  $> 5$ . Following Finucane et al.,<sup>10</sup> we excluded SNPs with  $\chi^2$  statistics  $> 80$  to reduce variance. We evaluated different  $\chi^2$  thresholds (excluding SNPs with  $\chi^2 > 25, 40, 80, \text{ or } 300$ ). In both local and distal analyses, the estimated enrichments were not sensitive to the choice of threshold (see [Results](#)). Following Finucane et al.,<sup>10</sup> we included in our regression only SNPs that appear in HapMap 3, which we used as a proxy for well-imputed SNPs.

To obtain a genome-wide estimate of the proportion of heritability of a category,  $\text{prop}_-h_g^2(C)$ , for either local or distal regions, we first calculated the average  $\widehat{\tau}_C$  and  $\overline{\tau}_C$  across all genes:

$$\overline{\tau}_C = \sum_{\text{gene } i} \widehat{\tau}_{C,i}. \quad (\text{Equation 6})$$

We included only genes whose total heritability estimate,  $\widehat{h}_g^2(\text{total})$ , was larger than 0. We applied this threshold both because negative heritability is biologically infeasible and because this reduced estimation noise and resulted in more stable estimates (see [Results](#)). We then computed the average category-specific heritability,  $\overline{h}_g^2(C)$ , and divided by the average total heritability  $\overline{h}_g^2(\text{total})$ :

$$\text{prop}_-h_g^2(C) = \frac{\overline{h}_g^2(C)}{\overline{h}_g^2(\text{total})} = \frac{\sum_C \overline{\tau}_C M_{C \wedge C}}{\sum_C \overline{\tau}_C M_C}, \quad (\text{Equation 7})$$

where  $\overline{h}_g^2(C)$  denotes the average estimated heritability of category  $C$ ,  $\overline{h}_g^2(\text{total})$  denotes the average total estimated heritability,  $M_C$  is the number of reference SNPs in category  $C$ , and  $M_{C \wedge C'}$  is the number of overlapping SNPs between categories  $C'$  and  $C$ .

The enrichment of heritability is defined as

$$\text{enrichment}(C) = \frac{\text{prop}_-h_g^2(C)}{\text{prop}_-\text{SNPs}(C)}, \quad (\text{Equation 8})$$

where  $\text{prop}_-\text{SNPs}(C)$  is the proportion of reference SNPs that lie in category  $C$ .

Standard errors (SEs) were computed via block jackknife. In detail, we computed the SE of  $\text{prop}_g \cdot h_g^2(C)$  by partitioning the genes by genomic location into 200 adjacent blocks and jackknifing on genes. This accounts for possible correlations between nearby probes (analogous to the block jackknife on SNPs employed by stratified LD-score regression<sup>10</sup>). We computed the SE of enrichment(C) by dividing the SE of  $\text{prop}_g \cdot h_g^2(C)$  by that of  $\text{prop\_SNPs}(C)$ . We computed the statistical significance of enrichment by using a normal approximation. We used the significance threshold of  $0.05/n_C$ , where  $n_C$  is the number of categories analyzed, to correct for multiple testing.

We also computed an area-under-the-curve (AUC) metric, which quantifies the fact that larger categories (i.e., spanning a larger fraction of the genome) are more informative than smaller categories at a given enrichment level. In detail, for each category, we calculated the area A under the curve  $y = f(x)$ , where  $y$  is  $\text{prop}_g \cdot h_g^2(C)$  and  $x$  is  $\text{prop\_SNPs}(C)$  ( $0 \leq x \leq 1$ ). We defined the AUC as A if  $A \geq 0.5$  or as  $1 - A$  if  $A < 0.5$  (so that the AUC of a category would be equal to the AUC of its complement). The SE of the AUC is calculated as the SE of  $\text{prop}_g \cdot h_g^2(C)$  divided by 2.

## Simulations

We performed null simulations to assess type I error and causal simulations to assess bias in estimates of local enrichment from our extension of stratified LD-score regression. We focused our simulations on analyses of local enrichment because analyses of distal enrichment are very similar to the original version of stratified LD-score regression, which has been shown by previous simulations to produce robust results (Figures 1 and 2 from Finucane et al.<sup>10</sup>). Simulations were performed with genotypes from UK10K.<sup>42</sup> Quality control included removing SNPs with minor allele frequency  $< 0.01$ , missingness  $> 0.01$ , or Hardy-Weinberg equilibrium  $p < 10^{-6}$ . We randomly downsampled to one million SNPs to match the SNP density of the real datasets analyzed (Table 1). We simulated 42,000 gene expression phenotypes (corresponding to 42,000 Wright et al. probes; Table 1) by using genotypes from local regions, defined as within 1 Mb of the TSS of a gene. We assumed a non-infinitesimal, additive model in which 5% of SNPs (in local regions) are causal. In null simulations, local SNPs affect gene expression phenotypes, but no functional categories were enriched ( $\tau_{\text{all\_SNPs}} = 2 \times 10^{-4}$  and  $\tau_C = 0$  for all other categories), and the local heritability of each gene was set to 0.149. In causal simulations, super enhancers from Vahedi et al.<sup>39</sup> (super enhancer [Vahedi], 2.1% of SNPs) and H3K27ac (PGC2) (26.9% of SNPs) were chosen as representative causal enriched categories, whereby  $\tau_{\text{all\_SNPs}} = 5 \times 10^{-5}$  and  $\tau_{\text{super enhancer (Vahedi)}} = \tau_{\text{H3K27ac (Hnisz)}} = 5 \times 10^{-4}$ . We used super enhancer (Vahedi) instead of super enhancer (Hnisz) to represent small functional annotations (i.e., spanning a small fraction of the genome) and thus were able to assess the robustness of our methods for small annotations. The local heritability of each gene was set to 0.221 in causal simulations. In both null and causal simulations, we assessed the accuracy of both enrichment estimates and block-jackknife SEs via ten rounds of simulations.

## Extension of Cross-Trait LD-Score Regression

Cross-trait LD-score regression<sup>40</sup> relies on the fact that SNPs with high LD scores will have a higher product of Z scores (for two

genetically correlated traits) on average than SNPs with low LD scores. More precisely,

$$E[Z_{1j}Z_{2j}] = \frac{\sqrt{N_1 N_2} \rho_g l_j}{M} + \frac{\rho N_s}{\sqrt{N_1 N_2}}, \quad (\text{Equation 9})$$

where  $N_i$  is the sample size for study  $i$ ,  $\rho_g$  is genetic covariance,  $M$  is the number of SNPs,  $l_j$  is the LD score of SNP  $j$ , defined as  $l_j = \sum_k r^2(j, k)$ ,  $N_s$  is the number of overlapping samples in the two studies, and  $\rho$  is the phenotypic correlation among the overlapping samples.

In the simple model defined by Equation 1, let  $\beta_j$  be the effect size of trait 1 at SNP  $j$ , and let  $\gamma_j$  be the effect size of trait 2 at SNP  $j$ . The genetic covariance between trait 1 and trait 2 is defined as

$$\rho_g = \sum_j \beta_j \gamma_j. \quad (\text{Equation 10})$$

Genetic correlation is defined as

$$r_g = \frac{\rho_g}{\sqrt{h_{g1}^2 h_{g2}^2}}. \quad (\text{Equation 11})$$

Regressing the product of Z scores of two traits on  $l_j$  gives  $\widehat{\rho}_g$ , an estimate of  $\rho_g$ . We can also estimate  $h_{g1}^2$  and  $h_{g2}^2$  from standard LD-score regression;<sup>43</sup> the genetic correlation can be estimated from Equation 11.

We extended cross-trait LD-score regression to estimate, for a given pair of tissues, the aggregate genetic correlation of the expression over a large set of common probes between two tissues. We included all probes with positive heritability estimates in each of the two tissues. We estimated genetic correlation separately for local and distal regions. For each pair of tissues and each common probe  $i$ , we estimated the genetic covariance ( $\widehat{\rho}_{g,i}$ ), as well as the total heritability of probe  $i$  in each tissue ( $\widehat{h}_{g1,i}^2$  and  $\widehat{h}_{g2,i}^2$ ). The aggregate genetic correlation across all shared probes is estimated as

$$\widehat{r}_g = \frac{\overline{\widehat{\rho}_g}}{\sqrt{\overline{h_{g1}^2} \overline{h_{g2}^2}}}, \quad (\text{Equation 12})$$

where  $\overline{\widehat{\rho}_g}$ ,  $\overline{h_{g1}^2}$ , and  $\overline{h_{g2}^2}$  are the averages of  $\widehat{\rho}_g$ ,  $\widehat{h}_{g1,i}^2$ , and  $\widehat{h}_{g2,i}^2$ , respectively, taken over probes  $i$ , whose  $\widehat{h}_{g1,i}^2$  and  $\widehat{h}_{g2,i}^2$  are both greater than 0. We estimated the SEs of  $\widehat{r}_g$  by dividing the probes by genomic locations into 200 blocks and performing a block jackknife on the probes as in Bulik-Sullivan et al.<sup>40</sup>

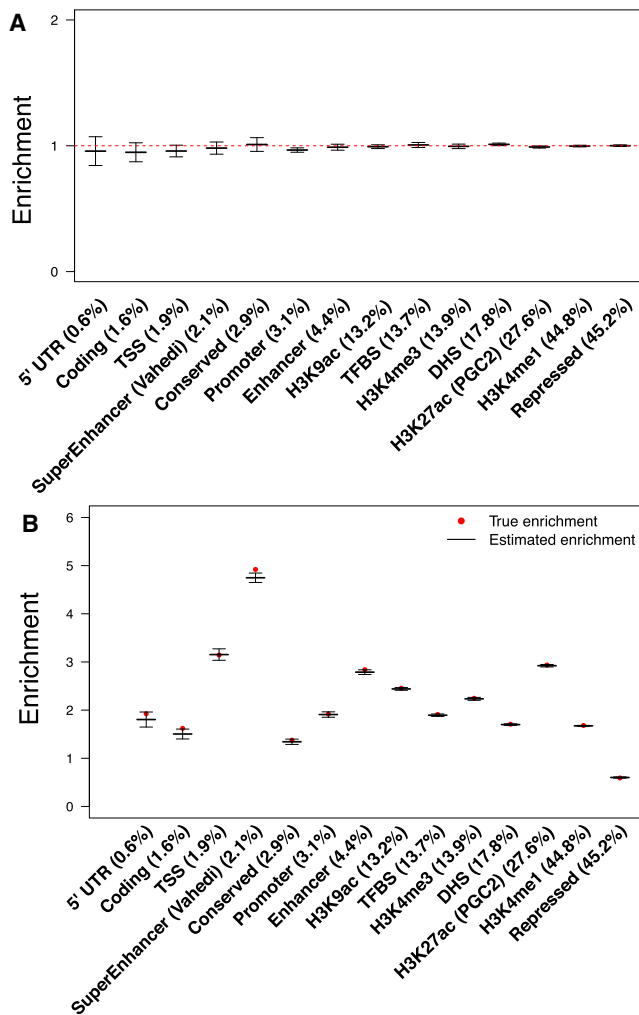
## Software Availability

Open-source software implementing our extensions of stratified LD-score regression and cross-trait LD-score regression is publicly available as part of the LD-score regression software (see [Web Resources](#)).

## Results

### Simulations

We performed null simulations, in which local SNPs affect gene expression phenotypes but no functional categories were enriched, to assess type I error of local enrichment estimates of our extension of LD-score regression (see [Material and Methods](#)). Type I error was well calibrated across ten simulations:  $\tau_{\text{all\_SNPs}}$  was accurately estimated



**Figure 1. Simulations Assessing Type I Error and Bias of Local Enrichment Estimates**

(A) Null simulations demonstrate well-calibrated type I error, given that estimated functional enrichments (average across ten simulations) are not statistically different from 1 after Bonferroni correction. Error bars represent 95% confidence intervals based on empirical SEs of the average enrichment across ten simulations. (B) Causal simulations demonstrate unbiased estimates of functional enrichments. Red dots represent the true expected enrichments. Center bars represent estimated enrichments (average across ten simulations). Error bars represent 95% confidence intervals based on empirical SEs of the average enrichment across ten simulations. Although some estimated enrichments lie just outside the 95% confidence intervals, they are not statistically different from the true enrichments after Bonferroni correction. Results are displayed for a representative set of 14 categories; numerical results for all 57 categories (for null and causal simulations) are reported in Tables S4 and S5.

(Table S3), and the enrichments of all 57 categories were not statistically different from the true enrichment of 1 (Figure 1A and Table S4).

We also performed causal simulations to assess bias in local enrichment estimates by using super enhancer (Vahedi) and H3K27ac (PGC2) as the causal enriched categories (see Material and Methods). We note that true local enrichments of non-causal functional categories are different

from 1 because of overlap with causal categories (Table S5). We observed unbiased estimates of local enrichment across ten simulations for all 57 functional categories, including those occupying less than 1% of the whole genome (Figure 1B and Table S5).

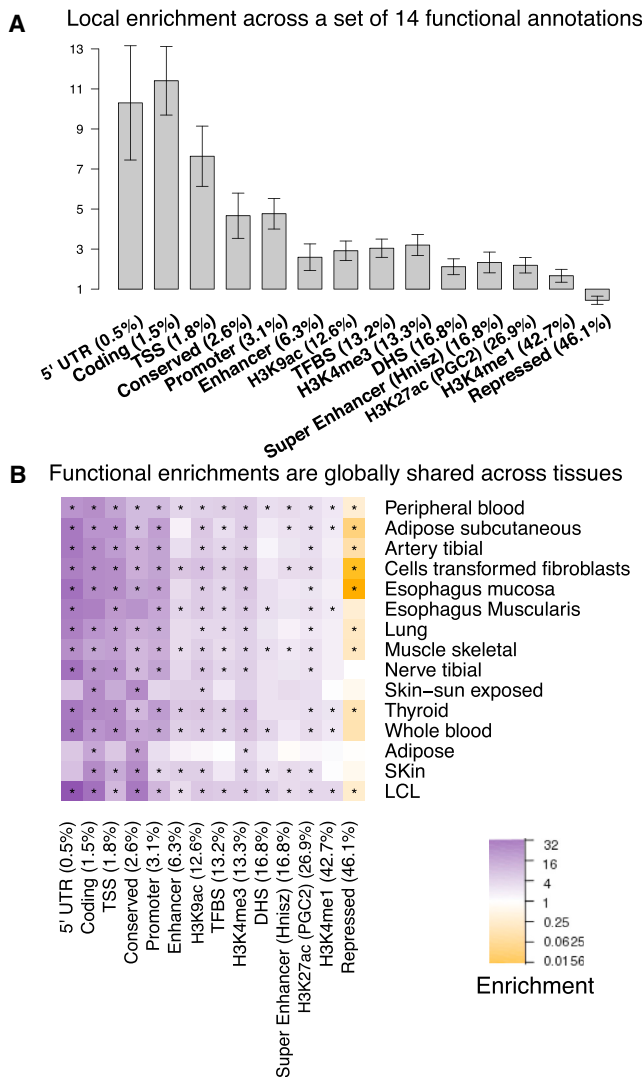
To evaluate whether the block-jackknife SEs were well calibrated, we compared them with empirical standard deviations across null and causal simulations. We determined that block-jackknife SEs were well calibrated: on average across 57 categories, they were 1.044 $\times$  larger than the empirical standard deviations from ten null simulations (Table S4) and 1.004 $\times$  larger than the empirical standard deviations from ten causal simulations (Table S5).

The average local heritability across genes in ten null simulations was estimated to be 0.151 (SE = 0.0001; range = 0.1505–0.1520; actual  $\bar{h}_g^2 = 0.149$ ), and the average local heritability across genes in ten causal simulations was estimated to be 0.203 (SE = 0.0003; range = 0.2019–0.2044; actual  $\bar{h}_g^2 = 0.221$ ), indicating close to unbiased estimates of total local heritability. More than 99% of the simulated genes were estimated to have  $h_g^2(\text{total})$  larger than 0 in the causal simulations. To assess whether restricting the analysis to genes with positive estimated heritability would create bias, we performed additional simulations in which the local heritability of each gene was set to 0.022, causing more genes to have negative estimated heritability. Our results showed that choosing different thresholds on  $\bar{h}_g^2(\text{total})$  did not affect our estimates of local enrichment (Figure S1).

### Functional Architectures of Local Regulation of Gene Expression in 15 Human Tissues

We partitioned local gene expression heritability across functional categories in three datasets spanning 15 human tissues<sup>13,16,19</sup> (Table 1; see Material and Methods). We analyzed 57 functional categories: 53 baseline categories from Finucane et al.<sup>10</sup> and four categories based on super enhancers and typical enhancers from Vahedi et al.<sup>39</sup> (Table S2; see Material and Methods). We estimated the enrichment of each functional category, defined as the proportion of heritability in that category divided by the proportion of SNPs in that category (see Material and Methods).

We first analyzed the Wright et al. gene expression array dataset, which had the largest sample size ( $n = 3,754$ ) and included only a single tissue type, peripheral blood.<sup>16</sup> Many functional categories were significantly enriched (Figure 2A; Table S6); several of these have been implicated in previous studies<sup>9,13–20,25–27</sup> (Table S1), but some have not. We observed that conserved regions were significantly enriched (4.66 $\times$ ; SE = 0.57;  $p = 1.98 \times 10^{-10}$ ). Although the function of conserved regions in gene regulatory programs has previously been reported in yeast,<sup>44</sup> previous evidence of functional enrichments of conserved regions on gene expression in humans is limited.<sup>25,26</sup> We further determined that the enrichment observed in conserved regions is largely attributed to conserved coding regions (15.77 $\times$ ; SE = 1.48;  $p < 10^{-12}$ ; Figure S2A). However, this



**Figure 2. Functional Enrichments for Local Regulation of Gene Expression**

(A) Local enrichment of each category in peripheral blood (Wright et al. dataset,  $n = 3,754$ ). Error bars represent 95% confidence intervals. Results for the AUC metric are displayed in Figure S6.

(B) Local enrichment of each category across 15 tissues. Purple shading indicates enriched categories (enrichment  $> 1$ ), orange shading indicates depleted categories (enrichment  $< 1$ ), and asterisks indicate significant enrichment or depletion after correction for 57 hypotheses tested. Sample sizes of each dataset are reported in Table 1. Results are displayed for a representative set of 14 categories; numerical results for all 57 categories are reported in Table S6, Tables S7 and S8. Estimates of the total local  $h_g^2$  for each tissue (average across all genes) are provided in Table S18. A description of each functional category is provided in Table S2.

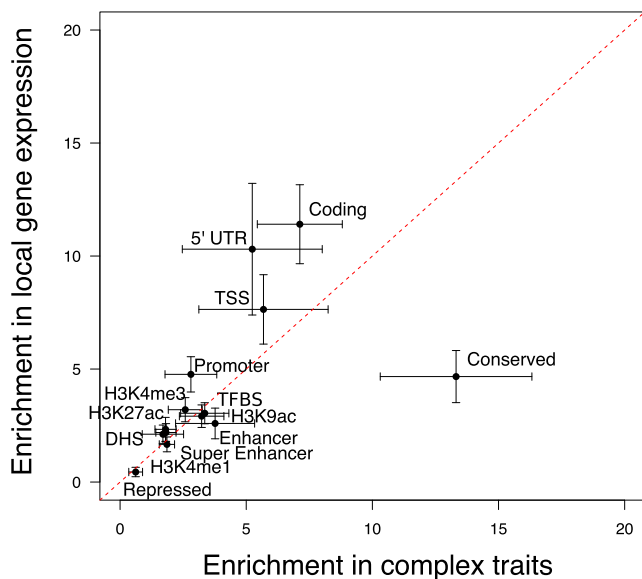
is not the case in complex traits, for which overlapping conserved coding regions and conserved non-coding regions were equally enriched (Figure S2B). Super enhancer (Hnisz) was also significantly enriched ( $2.33\times$ ;  $SE = 0.26$ ;  $p = 4.82 \times 10^{-7}$ ), supporting the role of super enhancers in local regulation of gene expression.

We also confirmed and quantified functional enrichments reported in previous studies of local regulation of gene expression in humans. We observed a large enrichment

in coding regions ( $11.41\times$ ;  $SE = 0.87$ ;  $p < 10^{-12}$ ), which confirmed previous findings<sup>9,14</sup> (Table S1) and is consistent with a recent study reporting that exonic regions are often involved in transcription factor binding<sup>45</sup> or contain splicing signals.<sup>46</sup> This suggests that the impact of coding variants on complex traits could sometimes be due to their effect on expression levels rather than changes in protein sequences. The histone marks H3K4me3, H3K9ac, H3K4me1, and H3K27ac (PGC2) were significantly enriched ( $1.66\times$ – $3.20\times$ ;  $SE = 0.16$ – $0.27$ ;  $p = 4.93 \times 10^{-5}$  to  $1.11 \times 10^{-16}$ ), consistent with previous findings<sup>14,19,26</sup> (Table S1) and confirming the role of histone marks in local regulation of gene expression. 5' UTRs were also significantly enriched ( $10.30\times$ ;  $SE = 1.46$ ;  $p = 1.68 \times 10^{-10}$ ). This enrichment could be driven by the promoter ( $4.77\times$ ;  $SE = 0.39$ ;  $p < 10^{-12}$ ), which overlaps the 5' UTR and directly affects transcription and other regulatory sequences in the 5' UTR, such as upstream open reading frames.<sup>47,48</sup> We also observed significant enrichments at DHSs, enhancers, and TFBSs, consistent with previous studies (Table S1).

We analyzed additional RNA sequencing (RNA-seq) (GTEx) and gene expression array (MuTHER) datasets spanning a total of 15 tissues (Table 1; see Material and Methods). The heritability enrichments were highly consistent across the 15 tissues, despite the widely varying sample sizes and different assays (Figure 2B; Tables S6–S8), which indicates that the functional architecture of local regulation of gene expression is consistent across different tissues. We note that in contrast to stratified LD-score regression, methods for assessing functional enrichment with only top eQTLs could produce enrichment estimates that are highly dependent on sample size (see Discussion).

We compared the functional enrichments that we estimated for local regulation of gene expression in peripheral blood with functional enrichments that we previously reported for a meta-analysis of nine independent complex traits<sup>10</sup> for 53 baseline functional categories. We observed a moderately strong correlation (inverse-variance-weighted Pearson  $r = 0.66$ ; Figure 3). The enrichments for local regulation of gene expression were comparable to the enrichments for complex traits for most functional categories: enrichments for only 3 and 2 out of 53 categories were significantly smaller and larger, respectively, for local regulation of gene expression after Bonferroni correction (Table S9). In particular, conserved regions exhibited a significantly lower enrichment in local regulation of gene expression, suggesting that variants in conserved regions could affect complex traits through mechanisms other than local regulation of gene expression. Notably, because of the large number of genes in each gene expression dataset, analyzing gene expression as an intermediate phenotype generally resulted in smaller SEs than analyses of complex traits in very large sample sizes, leading to enrichments that were more statistically significant (Figure 3; Table S9). Thus, gene expression data can be a particularly valuable means of assessing functionally important genomic regions.



**Figure 3. Comparison of Functional Enrichments for Local Regulation of Gene Expression in Peripheral Blood and Nine Complex Traits**

Enrichments for complex traits are meta-analyzed enrichments of nine complex traits and diseases from Finucane et al.<sup>10</sup> Error bars represent 95% confidence intervals. The red dashed line represents  $y = x$ . Results are displayed for a representative set of 14 categories; numerical results for all 53 categories are reported in Table S9. H3K27ac (PGC2) is denoted as H3K27ac in the figure. A description of each functional category is provided in Table S2.

We performed several secondary analyses that did not substantially change our results. First, we included distance from the TSS ( $\pm 10$ , 20, and 30 kb from the TSS of the corresponding gene) as additional functional categories to assess whether distance from the TSS might explain some of the observed enrichments. The annotation defined by  $\pm 10$  kb to the TSS exhibited the largest increase in per-SNP heritability when conditioned on other annotations ( $\tau$ ) among all annotations in the model, consistent with previous work emphasizing the importance of distance to TSS.<sup>25</sup> However, no significant differences were observed in  $\tau$  estimates for other functional categories after the distance-to-TSS annotations were included in the model (Figure S3). This indicates that the enrichments we observed for other functional categories are independent of the effect of distance to TSS. Second, we evaluated whether estimates of local enrichment are sensitive to the definition of local regions. We extended local regions from 1 to 2 Mb around the TSS and obtained comparable results (Figure S4). Third, we evaluated whether estimates of local enrichment are sensitive to the choice of threshold on  $\chi^2$  statistics (see Material and Methods). We applied different  $\chi^2$  thresholds (excluding SNPs with  $\chi^2 > 40$ , 80 [default], and 300) and obtained comparable results (Table S10). Fourth, we evaluated the impact of the choice of  $\widehat{h}_g^2(\text{total})$  threshold on estimates of local enrichment (see Material and Methods). We determined that applying a threshold on  $\widehat{h}_g^2(\text{total})$  reduces the

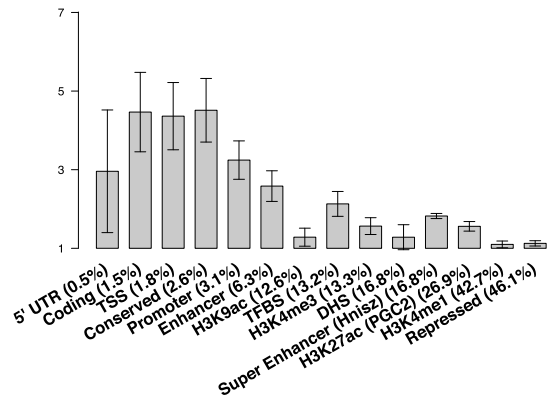
estimation noise and that enrichment estimates are stable as long as the probes with extremely negative estimates are removed (Figure S5 and Table S11). We also modified the analysis by including only probes whose heritability estimates were significantly positive ( $p < 0.05$  before Bonferroni correction). The enrichment was also consistent with the estimates obtained after application of  $\widehat{h}_g^2(\text{total})$  thresholds (Figure S5).

### Functional Architectures of Distal Regulation of Gene Expression in Four Human Tissues

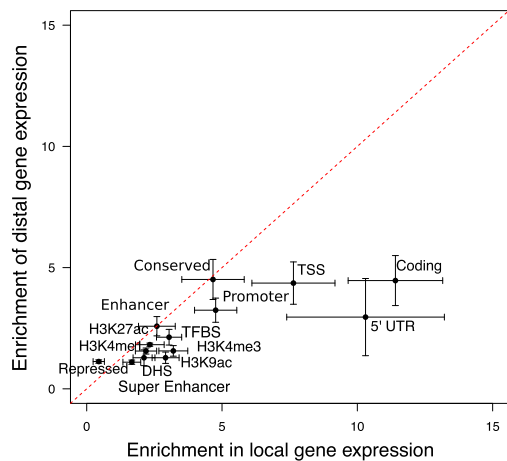
Functional characterization of distal regulation of gene expression has been limited because of the low statistical power to identify distal eQTLs. We partitioned distal gene expression heritability across functional categories in four human tissues. We first analyzed the Wright et al. gene expression array dataset.<sup>16</sup> Many functional categories were significantly enriched in the distal analysis (Figure 4A; Figure S6 and Table S12). In particular, we again observed significant enrichments at conserved regions ( $4.51\times$ ; SE = 0.41;  $p < 10^{-12}$ ), coding regions ( $4.47\times$ ; SE = 0.52;  $p = 1.79 \times 10^{-11}$ ), and super enhancer (Hnisz) regions ( $1.82\times$ ; SE = 0.03;  $p < 10^{-12}$ ). To test the hypothesis that only SNPs in expressed genes should affect expression levels of other genes distally, we added an additional coding annotation by considering coding regions of highly expressed genes in whole blood (RPKM  $> 5$  in GTEx whole blood; 8% of genes). We observed a substantially larger enrichment of  $17.23\times$  (SE = 5.76;  $p = 0.005$ ; see Table S13), which is consistent with an important contribution of expressed genes (such as transcription factors) in the distal regulatory control of gene expression. To our knowledge, the enrichment in distal coding regions of expressed genes has not been reported in previous studies in humans or model organisms.<sup>20</sup> In addition, two histone marks were significantly enriched: H3K27ac (PGC2) ( $1.56\times$ ; SE = 0.06;  $p < 10^{-12}$ ) and H3K4me3 ( $1.56\times$ ; SE = 0.11;  $p = 2.29 \times 10^{-7}$ ). H3K4me1 and H3K9ac were not significant after correction for 57 hypotheses tested, but broadly defined H3K4me1 regions (H3K4me1 extended by 500 bp; 60.9% of SNPs) explained 98.0% of distal heritability ( $1.61\times$ ; SE = 0.02;  $p < 10^{-12}$ ). These results suggest that most SNPs that affect distal gene regulation lie near regions marked by H3K4me1. We note that previous studies of distal eQTLs in blood reported distal enrichments only in 5' UTRs<sup>16</sup> (whose enrichment in our analyses was not statistically significant after correction for multiple testing:  $2.96\times$ ; SE = 0.80;  $p = 0.013$ ) and in enhancer regions of myeloid and lymphoid cell lines<sup>15</sup> (we similarly detected distal enrichment in enhancers as defined by Hoffman et al.<sup>34</sup>:  $2.58\times$ ; SE = 0.20;  $p < 10^{-12}$ ). We are not aware of any other previous results on distal enrichment.

We compared the enrichments in distal regulation of gene expression with the local enrichments estimated above across the 57 categories and observed a strong correlation (inverse-variance-weighted Pearson  $r = 0.90$ ; Figure 4B; Table S12). The enrichments in distal and local

**A** Significant enrichment observed in distal gene expression regulation



**B**



**Figure 4. Functional Enrichments for Distal Regulation of Gene Expression**

(A) The distal enrichment of each category in peripheral blood (Wright et al.<sup>16</sup> dataset,  $n = 3,754$ ). Error bars represent 95% confidence intervals. Results for the area under curve (AUC) metric are displayed in Figure S6.

(B) Comparison of functional enrichments for distal gene expression regulation versus local gene expression regulation in peripheral blood. Error bars represent 95% confidence intervals. The red dashed line represents  $y = x$ . H3K27ac (PGC2) is denoted as H3K27ac in the figure. A description of each functional category is provided in Table S2.

regulation of gene expression were comparable for most functional categories: enrichments for only 11 and 4 out of 57 categories were significantly smaller and larger, respectively, for distal regulation of gene expression after Bonferroni correction (Table S12). This suggests that the dearth of previously reported functional enrichments for distal regulation of gene expression is due to the low power of approaches based on top distal eQTLs (which most studies are underpowered to detect) and not due to the absence of functional enrichments.

We performed two secondary analyses. First, we evaluated the impact of applying a  $\widehat{h}_g^2(\text{total})$  threshold on estimates of distal enrichment (see Material and Methods). Similar to local analyses, we determined that applying a threshold produces quantitatively smaller but more precise estimates (Figure S5C and Table S11C). Second, we estimated functional enrichments separately for intra- and

inter-chromosomal distal regions. We observed that the enrichment of inter-chromosomal distal regions was comparable to the total distal enrichment, indicating that functional elements on chromosomes different from those of the expressed gene are actively involved in regulation of gene expression (Figure S7 and Table S14). In particular, this implies that the definition of distal regions (which includes intra-chromosomal regions  $> 1$  Mb from the TSS) has little effect on enrichment estimates. In some cases, larger enrichments were observed for inter-chromosomal distal regions than for intra-chromosomal distal regions, but many of these differences were not statistically significant and could be due to estimation noise.

We also analyzed distal enrichment in the MuTHER gene expression array dataset (Table 1) and observed many significant enrichments (Figure S8 and Table S15). We did not include the GTEx dataset in the distal analysis because of its smaller sample size.

**Genetic Correlation of Gene Expression between Different Tissues**

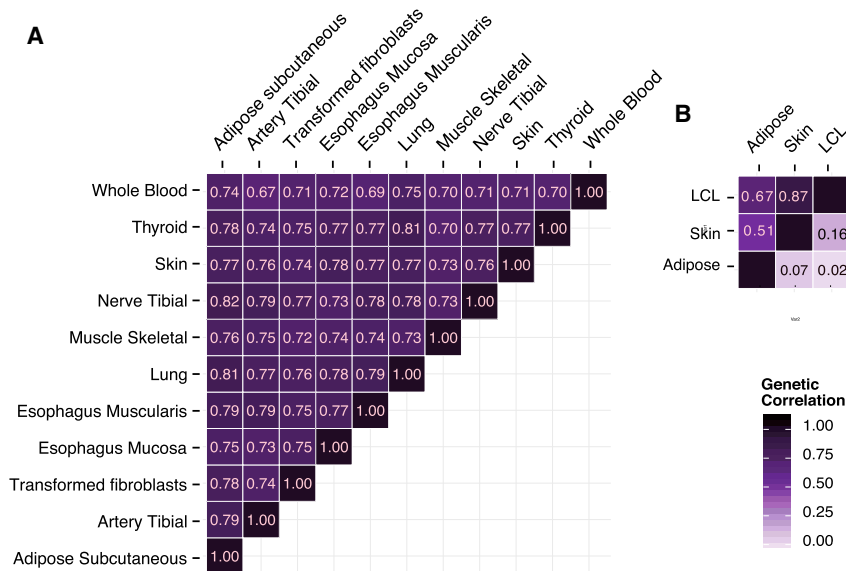
We extended cross-trait LD-score regression<sup>40</sup> to estimate the pairwise genetic correlations of local gene expression between different tissues (see Material and Methods). Pairwise genetic correlations were estimated separately in 11 GTEx tissues and in three MuTHER tissues (Figure 5; Tables S16 and S17). The average pairwise genetic correlation was 0.75 (average SE = 0.02). The lowest genetic correlation was observed between adipose tissue and skin in MuTHER data ( $r = 0.51$ ; SE = 0.17), although it was not statistically significant in comparison with correlations from other MuTHER data. The remaining 57 pairwise correlations were all larger than 0.675, indicating that local regulation of gene expression is highly correlated across tissues, consistent with previous studies.<sup>13,19,49–51</sup>

We also estimated the pairwise genetic correlations of distal gene expression between the three MuTHER tissues (Figure 5; Table S17). Interestingly, the average pairwise genetic correlation was much smaller at 0.08 (average SE = 0.01), indicating that distal regulation of gene expression is highly tissue specific. This is consistent with previous work,<sup>49</sup> although relatively few previous studies have investigated the sharing of distal regulation of gene expression across tissues because of the low power to detect distal eQTLs. We also note that this does not contradict our finding of consistent distal enrichments across tissues (Figure S8), because it is possible that different tissues have different distal eQTLs that nonetheless reside in the same functional categories.

**Discussion**

In this study, we comprehensively investigated functional enrichments for both local and distal regulation of gene expression in multiple human tissues by applying an extension of stratified LD-score regression<sup>10</sup> to large gene





**Figure 5. Pairwise Local and Distal Genetic Correlation across Tissues**

(A) Pairwise local genetic correlations across 11 GTEx tissues.

(B) Local (upper left) and distal (lower right) genetic correlations across three MuTHER tissues. Numerical results are reported in Tables S16 and S17.

expression datasets. We detected widespread functional enrichments for both local and distal gene regulation, including enrichments at coding regions, conserved regions, super enhancers, and several histone marks; some of the local enrichments and most of the distal enrichments were not identified in previous studies (Table S1). We also confirmed that local regulation of gene expression is highly genetically correlated across tissues, whereas distal regulation is highly tissue specific.<sup>49</sup>

The functional enrichments that we detected for local regulation of gene expression were generally more statistically significant than enrichments that we previously reported for analyses of complex traits in very large sample sizes.<sup>10</sup> This emphasizes the value of studying gene expression as an intermediate phenotype for studying complex diseases and traits, particularly in analyses of functional enrichment. Our systematic investigation of enrichment of local regulation of gene expression across 15 tissues identified highly consistent enrichments across tissues, despite the widely varying samples sizes and different assays. This conclusion was possible because the heritability approach employed by stratified LD-score regression produces enrichment estimates that are independent of sample size,<sup>10</sup> in the sense that small sample size does not bias point estimates (although small sample size could limit power to detect significant enrichments). On the other hand, methods for assessing functional enrichment by using only top eQTLs could be highly dependent on sample size because the enrichment of associated variants in regulatory annotations could vary with effect size (see Table 1 from Sveinbjornsson et al.<sup>11</sup>). In addition, our results on enrichment of distal regulation of gene expression represent a substantial advance over previous results on functional enrichment of distal eQTLs, which were limited by the small number of individually significant distal eQTLs detected by previous studies. Our results highlight the advantages of leveraging genome-wide polygenic signals

over restricting to top eQTLs in efforts to identify functional enrichments.

Our work has several limitations. First, stratified LD-score regression models only additive effects and cannot capture non-additive effects or epistasis, which could play an important role in regulating gene expression.<sup>52–56</sup> Second, stratified LD-score regression analyzes summary-level data and thus does not take advantage

of the additional information available in individual-level data. Although functional-enrichment analyses of individual-level data can be performed with restricted maximum likelihood (REML) and its extensions,<sup>57–59</sup> those methods are applicable only to a small number of non-overlapping functional annotations; to our knowledge, all current methods that are applicable to a large number of overlapping functional annotations are based on summary statistics,<sup>60</sup> whereas analyzing one annotation at a time can produce severely biased results (see Figure 2b from Finucane et al.<sup>10</sup>). Third, stratified LD-score regression is designed for highly polygenic traits and does not take full advantage of non-infinitesimal genetic architectures, which are a particularly likely characteristic of local regulation of gene expression.<sup>61</sup> Our highly consistent local enrichments across 15 tissues indicate that the method does produce robust results for analyses of local gene expression, but methods that account for non-infinitesimal genetic architectures might produce even more precise estimates. However, to our knowledge, existing methods for heritability analysis that model non-infinitesimal genetic architectures<sup>62,63</sup> are not applicable to enrichment analyses involving a large number of overlapping functional annotations. Fourth, stratified LD-score regression is designed to partition the heritability explained by common variants, but rare variants could also play an important role in regulating gene expression.<sup>64</sup> Fifth, the functional enrichments that we inferred are relative local and distal  $h_g^2$  values that are small in absolute terms (Tables S18 and S19); however, other studies have also inferred low values of gene expression heritability.<sup>13,16,24</sup> The low average estimates of heritability can be attributed to environmental noise, including noise in measurements of gene expression. (The fact that individual estimates are sometimes negative can be attributed to estimation noise; we did not constrain our estimates to the plausible 0–1 range, which could lead to bias in the average of the estimates.) However, the low inferred heritability of

gene expression has not precluded important biological discoveries. Sixth, our results on functional enrichment were based on eQTLs and did not consider splicing QTLs (sQTLs), a rich area for future investigation.<sup>18,46,65,66</sup> Seventh, we detected no significant cell-type-specific local enrichments and only limited cell-type-specific distal enrichments (see Tables S20–S22), although similar analyses have detected strong cell-type-specific enrichments for complex traits.<sup>10</sup> The absence of local cell-type-specific enrichments is consistent with our observation that local functional enrichments are highly consistent across different tissues, and future analyses might need to restrict to appropriate gene sets (and/or consider sQTLs) to detect cell-type-specific signals. Despite these limitations, our findings shed light on the genetic architecture and molecular mechanisms underlying the regulation of gene expression and demonstrate that gene expression is an appropriate intermediate phenotype for analyzing functional enrichments of complex diseases and traits.

### Supplemental Data

Supplemental Data include 8 figures and 23 tables and can be found with this article online at <http://dx.doi.org/10.1016/j.ajhg.2017.03.002>.

### Conflicts of Interest

P.F.S. is a member of the scientific advisory board of Pfizer and Lundbeck.

### Acknowledgments

We are grateful to Soumya Raychaudhuri for helpful discussions. This research was funded by NIH grant R01 MH107649. P.F.S. acknowledges funding from the US National Institute of Mental Health (U01 MH109528) for the Psychiatric Genomics Consortium.

Received: November 29, 2016

Accepted: February 24, 2017

Published: March 23, 2017

### Web Resources

Annotations for 57 functional categories, <http://data.broadinstitute.org/alkesgroup/LDSCORE/>

GenABEL and ProbABEL packages, <http://www.genable.org>

GTEx Portal, <http://gtexportal.org/home/datasets>

LDSC (LD Score), [https://github.com/Nealelab/ldsc/tree/gene\\_expression\\_ldsc](https://github.com/Nealelab/ldsc/tree/gene_expression_ldsc)

UCSC Genome Browser, <https://genome.ucsc.edu/>

### References

- Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A., Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., Ziller, M.J., et al.; Roadmap Epigenomics Consortium (2015). Integrative analysis of 111 reference human epigenomes. *Nature* 518, 317–330.
- ENCODE Project Consortium (2012). An integrated encyclopedia of DNA elements in the human genome. *Nature* 489, 57–74.
- Ernst, J., Kheradpour, P., Mikkelsen, T.S., Shores, N., Ward, L.D., Epstein, C.B., Zhang, X., Wang, L., Issner, R., Coyne, M., et al. (2011). Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* 473, 43–49.
- 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.
- Trynka, G., Sandor, C., Han, B., Xu, H., Stranger, B.E., Liu, X.S., and Raychaudhuri, S. (2013). Chromatin marks identify critical cell types for fine mapping complex trait variants. *Nat. Genet.* 45, 124–130.
- Gusev, A., Lee, S.H., Trynka, G., Finucane, H., Vilhjálmsson, B.J., Xu, H., Zang, C., Ripke, S., Bulik-Sullivan, B., Stahl, E., et al.; Schizophrenia Working Group of the Psychiatric Genomics Consortium; SWE-SCZ Consortium; Schizophrenia Working Group of the Psychiatric Genomics Consortium; and SWE-SCZ Consortium (2014). Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *Am. J. Hum. Genet.* 95, 535–552.
- Pickrell, J.K. (2014). Joint analysis of functional genomic data and genome-wide association studies of 18 human traits. *Am. J. Hum. Genet.* 94, 559–573.
- Kichaev, G., Yang, W.-Y., Lindstrom, S., Hormozdiari, F., Eskin, E., Price, A.L., Kraft, P., and Pasaniuc, B. (2014). Integrating functional data to prioritize causal variants in statistical fine-mapping studies. *PLoS Genet.* 10, e1004722.
- Farh, K.K.-H., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W.J., Beik, S., Shores, N., Whitton, H., Ryan, R.J.H., Shishkin, A.A., et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* 518, 337–343.
- Finucane, H.K., Bulik-Sullivan, B., Gusev, A., Trynka, G., Reshef, Y., Loh, P.-R., Anttila, V., Xu, H., Zang, C., Farh, K., et al.; ReproGen Consortium; Schizophrenia Working Group of the Psychiatric Genomics Consortium; and RACI Consortium (2015). Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.* 47, 1228–1235.
- Sveinbjornsson, G., Albrechtsen, A., Zink, F., Gudjonsson, S.A., Oddson, A., Måsson, G., Holm, H., Kong, A., Thorsteinsdottir, U., Sulem, P., et al. (2016). Weighting sequence variants based on their annotation increases power of whole-genome association studies. *Nat. Genet.* 48, 314–317.
- Schork, A.J., Thompson, W.K., Pham, P., Torkamani, A., Roddey, J.C., Sullivan, P.F., Kelsoe, J.R., O'Donovan, M.C., Furberg, H., Schork, N.J., et al.; Tobacco and Genetics Consortium; Bipolar Disorder Psychiatric Genomics Consortium; and Schizophrenia Psychiatric Genomics Consortium (2013). All SNPs are not created equal: genome-wide association studies reveal a consistent pattern of enrichment among functionally annotated SNPs. *PLoS Genet.* 9, e1003449.
- 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.; Multiple Tissue Human Expression Resource (MuTHER) Consortium (2012). Mapping cis- and trans-regulatory effects across multiple tissues in twins. *Nat. Genet.* 44, 1084–1089.
- Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P.A., Monlong, J., Rivas, M.A., González-Porta, M., Kurbatova, N.,

- Griebel, T., Ferreira, P.G., et al.; Geuvadis Consortium (2013). Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* 501, 506–511.
15. 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.
  16. 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.
  17. Battle, A., Mostafavi, S., Zhu, X., Potash, J.B., Weissman, M.M., McCormick, C., Haudenschild, C.D., Beckman, K.B., Shi, J., Mei, R., et al. (2014). Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome Res.* 24, 14–24.
  18. Zhang, X., Joehanes, R., Chen, B.H., Huan, T., Ying, S., Munson, P.J., Johnson, A.D., Levy, D., and O'Donnell, C.J. (2015). Identification of common genetic variants controlling transcript isoform variation in human whole blood. *Nat. Genet.* 47, 345–352.
  19. Deluca, D.S., Segre, A.V., Sullivan, T.J., Young, T.R., Gelfand, E.T., Trowbridge, C.A., Maller, J.B., Tukiainen, T., et al.; GTEx Consortium (2015). Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* 348, 648–660.
  20. Albert, F.W., and Kruglyak, L. (2015). The role of regulatory variation in complex traits and disease. *Nat. Rev. Genet.* 16, 197–212.
  21. Davis, L.K., Yu, D., Keenan, C.L., Gamazon, E.R., Konkashbaev, A.I., Derks, E.M., Neale, B.M., Yang, J., Lee, S.H., Evans, P., et al. (2013). Partitioning the heritability of Tourette syndrome and obsessive compulsive disorder reveals differences in genetic architecture. *PLoS Genet.* 9, e1003864.
  22. Torres, J.M., Gamazon, E.R., Parra, E.J., Below, J.E., Valladares-Salgado, A., Wacher, N., Cruz, M., Hanis, C.L., and Cox, N.J. (2014). Cross-tissue and tissue-specific eQTLs: partitioning the heritability of a complex trait. *Am. J. Hum. Genet.* 95, 521–534.
  23. 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., Im, H.K.; and GTEx Consortium (2015). A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* 47, 1091–1098.
  24. Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B.W.J.H., Jansen, R., de Geus, E.J.C., Boomsma, D.I., Wright, F.A., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* 48, 245–252.
  25. Veyrieras, J.-B., Kudaravalli, S., Kim, S.Y., Dermitzakis, E.T., Gilad, Y., Stephens, M., and Pritchard, J.K. (2008). High-resolution mapping of expression-QTLs yields insight into human gene regulation. *PLoS Genet.* 4, e1000214–e1000215.
  26. Gaffney, D.J., Veyrieras, J.-B., Degner, J.F., Pique-Regi, R., Pai, A.A., Crawford, G.E., Stephens, M., Gilad, Y., and Pritchard, J.K. (2012). Dissecting the regulatory architecture of gene expression QTLs. *Genome Biol.* 13, R7.
  27. Battle, A., Khan, Z., Wang, S.H., Mitrano, A., Ford, M.J., Pritchard, J.K., and Gilad, Y. (2015). Genomic variation. Impact of regulatory variation from RNA to protein. *Science* 347, 664–667.
  28. Yin, Z., Xia, K., Chung, W., Sullivan, P.F., and Zou, F. (2015). Fast eQTL Analysis for Twin Studies. *Genet. Epidemiol.* 39, 357–365.
  29. Aulchenko, Y.S., de Koning, D.-J., and Haley, C. (2007). Genomewide rapid association using mixed model and regression: a fast and simple method for genomewide pedigree-based quantitative trait loci association analysis. *Genetics* 177, 577–585.
  30. Chen, W.-M., and Abecasis, G.R. (2007). Family-based association tests for genomewide association scans. *Am. J. Hum. Genet.* 81, 913–926.
  31. Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Res.* 12, 996–1006.
  32. Lindblad-Toh, K., Garber, M., Zuk, O., Lin, M.F., Parker, B.J., Washietl, S., Kheradpour, P., Ernst, J., Jordan, G., Mauceli, E., et al.; Broad Institute Sequencing Platform and Whole Genome Assembly Team; Baylor College of Medicine Human Genome Sequencing Center Sequencing Team; and Genome Institute at Washington University (2011). A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* 478, 476–482.
  33. Ward, L.D., and Kellis, M. (2012). Evidence of abundant purifying selection in humans for recently acquired regulatory functions. *Science* 337, 1675–1678.
  34. Hoffman, M.M., Ernst, J., Wilder, S.P., Kundaje, A., Harris, R.S., Libbrecht, M., Giardine, B., Ellenbogen, P.M., Bilmes, J.A., Birney, E., et al. (2013). Integrative annotation of chromatin elements from ENCODE data. *Nucleic Acids Res.* 41, 827–841.
  35. Andersson, R., Gebhard, C., Miguel-Escalada, I., Hoof, I., Bornholdt, J., Boyd, M., Chen, Y., Zhao, X., Schmidl, C., Suzuki, T., et al.; FANTOM Consortium (2014). An atlas of active enhancers across human cell types and tissues. *Nature* 507, 455–461.
  36. Hnisz, D., Abraham, B.J., Lee, T.I., Lau, A., Saint-André, V., Sigova, A.A., Hoke, H.A., and Young, R.A. (2013). Super-enhancers in the control of cell identity and disease. *Cell* 155, 934–947.
  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., et al.; Schizophrenia Working Group of the Psychiatric Genomics Consortium (2014). Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511, 421–427.
  38. Pott, S., and Lieb, J.D. (2015). What are super-enhancers? *Nat. Genet.* 47, 8–12.
  39. Vahedi, G., Kanno, Y., Furumoto, Y., Jiang, K., Parker, S.C.J., Erdos, M.R., Davis, S.R., Roychoudhuri, R., Restifo, N.P., Gaidina, M., et al. (2015). Super-enhancers delineate disease-associated regulatory nodes in T cells. *Nature* 520, 558–562.
  40. Bulik-Sullivan, B., Finucane, H.K., Anttila, V., Gusev, A., Day, F.R., Loh, P.-R., Duncan, L., Perry, J.R., Patterson, N., Robinson, E.B., et al.; ReproGen Consortium; Psychiatric Genomics Consortium; and Genetic Consortium for Anorexia Nervosa of the Wellcome Trust Case Control Consortium 3 (2015). An atlas of genetic correlations across human diseases and traits. *Nat. Genet.* 47, 1236–1241.
  41. Abecasis, G.R., Auton, A., Brooks, L.D., DePristo, M.A., Durbin, R.M., Handsaker, R.E., Kang, H.M., Marth, G.T., McVean, G.A.; and 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature* 491, 56–65.
  42. Walter, K., Min, J.L., Huang, J., Crooks, L., Memari, Y., McCarthy, S., Perry, J.R.B., Xu, C., Futema, M., Lawson, D., et al.; UK10K Consortium (2015). The UK10K project identifies rare variants in health and disease. *Nature* 526, 82–90.

43. 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.; and Schizophrenia Working Group of the Psychiatric Genomics Consortium (2015). LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* *47*, 291–295.
44. Lee, S.-I., Dudley, A.M., Drubin, D., Silver, P.A., Krogan, N.J., Pe'er, D., and Koller, D. (2009). Learning a prior on regulatory potential from eQTL data. *PLoS Genet.* *5*, e1000358.
45. Stergachis, A.B., Haugen, E., Shafer, A., Fu, W., Vernot, B., Reynolds, A., Raubitschek, A., Ziegler, S., LeProust, E.M., Akey, J.M., and Stamatoyannopoulos, J.A. (2013). Exonic transcription factor binding directs codon choice and affects protein evolution. *Science* *342*, 1367–1372.
46. Li, Y.I., van de Geijn, B., Raj, A., Knowles, D.A., Petti, A.A., Golani, D., Gilad, Y., and Pritchard, J.K. (2016). RNA splicing is a primary link between genetic variation and disease. *Science* *352*, 600–604.
47. Kervestin, S., and Jacobson, A. (2012). NMD: a multifaceted response to premature translational termination. *Nat. Rev. Mol. Cell Biol.* *13*, 700–712.
48. Cenik, C., Cenik, E.S., Byeon, G.W., Grubert, F., Candille, S.I., Spacek, D., Alsallakh, B., Tilgner, H., Araya, C.L., Tang, H., et al. (2015). Integrative analysis of RNA, translation, and protein levels reveals distinct regulatory variation across humans. *Genome Res.* *25*, 1610–1621.
49. Price, A.L., Helgason, A., Thorleifsson, G., McCarroll, S.A., Kong, A., and Stefansson, K. (2011). Single-tissue and cross-tissue heritability of gene expression via identity-by-descent in related or unrelated individuals. *PLoS Genet.* *7*, e1001317–e1001319.
50. Nica, A.C., Parts, L., Glass, D., Nisbet, J., Barrett, A., Sekowska, M., Travers, M., Potter, S., Grundberg, E., Small, K., et al.; MuTHER Consortium (2011). The architecture of gene regulatory variation across multiple human tissues: the MuTHER study. *PLoS Genet.* *7*, e1002003.
51. Flutre, T., Wen, X., Pritchard, J., and Stephens, M. (2013). A statistical framework for joint eQTL analysis in multiple tissues. *PLoS Genet.* *9*, e1003486.
52. Lappalainen, T., Montgomery, S.B., Nica, A.C., and Dermitzakis, E.T. (2011). Epistatic selection between coding and regulatory variation in human evolution and disease. *Am. J. Hum. Genet.* *89*, 459–463.
53. Hemani, G., Shakhbazov, K., Westra, H.-J., Esko, T., Henders, A.K., McRae, A.F., Yang, J., Gibson, G., Martin, N.G., Metspalu, A., et al. (2014). Detection and replication of epistasis influencing transcription in humans. *Nature* *508*, 249–253.
54. Wood, A.R., Tuke, M.A., Nalls, M.A., Hernandez, D.G., Bandinelli, S., Singleton, A.B., Melzer, D., Ferrucci, L., Frayling, T.M., and Weedon, M.N. (2014). Another explanation for apparent epistasis. *Nature* *514*, E3–E5.
55. Buil, A., Brown, A.A., Lappalainen, T., Viñuela, A., Davies, M.N., Zheng, H.-F., Richards, J.B., Glass, D., Small, K.S., Durbin, R., et al. (2015). Gene-gene and gene-environment interactions detected by transcriptome sequence analysis in twins. *Nat. Genet.* *47*, 88–91.
56. Fish, A.E., Capra, J.A., and Bush, W.S. (2016). Are Interactions between cis-Regulatory Variants Evidence for Biological Epistasis or Statistical Artifacts? *Am. J. Hum. Genet.* *99*, 817–830.
57. Yang, J., Benyamin, B., McEvoy, B.P., Gordon, S., Henders, A.K., Nyholt, D.R., Madden, P.A., Heath, A.C., Martin, N.G., Montgomery, G.W., et al. (2010). Common SNPs explain a large proportion of the heritability for human height. *Nat. Genet.* *42*, 565–569.
58. Yang, J., Manolio, T.A., Pasquale, L.R., Boerwinkle, E., Caporaso, N., Cunningham, J.M., de Andrade, M., Feenstra, B., Feingold, E., Hayes, M.G., et al. (2011). Genome partitioning of genetic variation for complex traits using common SNPs. *Nat. Genet.* *43*, 519–525.
59. Loh, P.-R., Bhatia, G., Gusev, A., Finucane, H.K., Bulik-Sullivan, B.K., Pollack, S.J., de Candia, T.R., Lee, S.H., Wray, N.R., Kendler, K.S., et al.; Schizophrenia Working Group of Psychiatric Genomics Consortium (2015). Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nat. Genet.* *47*, 1385–1392.
60. Pasaniuc, B., and Price, A.L. (2017). Dissecting the genetics of complex traits using summary association statistics. *Nat. Rev. Genet.* *18*, 117–127.
61. Wheeler, H.E., Shah, K.P., Brenner, J., Garcia, T., Aquino-Michaels, K., Cox, N.J., Nicolae, D.L., Im, H.K.; and GTEx Consortium (2016). Survey of the Heritability and Sparse Architecture of Gene Expression Traits across Human Tissues. *PLoS Genet.* *12*, e1006423.
62. Zhou, X., Carbonetto, P., and Stephens, M. (2013). Polygenic modeling with bayesian sparse linear mixed models. *PLoS Genet.* *9*, e1003264.
63. Moser, G., Lee, S.H., Hayes, B.J., Goddard, M.E., Wray, N.R., and Visscher, P.M. (2015). Simultaneous discovery, estimation and prediction analysis of complex traits using a bayesian mixture model. *PLoS Genet.* *11*, e1004969.
64. Zhao, J., Akinsanmi, I., Arafat, D., Cradick, T.J., Lee, C.M., Banskota, S., Marigorta, U.M., Bao, G., and Gibson, G. (2016). A Burden of Rare Variants Associated with Extremes of Gene Expression in Human Peripheral Blood. *Am. J. Hum. Genet.* *98*, 299–309.
65. Ongen, H., and Dermitzakis, E.T. (2015). Alternative Splicing QTLs in European and African Populations. *Am. J. Hum. Genet.* *97*, 567–575.
66. 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.