

# Trans-ethnic eQTL meta-analysis of human brain reveals regulatory architecture and candidate causal variants for brain-related traits

Biao Zeng<sup>1,2,3</sup>, Jaroslav Bendl<sup>1,2,3</sup>, Roman Kosoy<sup>1,2,3</sup>, John F. Fullard<sup>1,2,3</sup>, Gabriel E. Hoffman<sup>1,2,3</sup> #, Panos Roussos<sup>1,2,3,4,5</sup> #

<sup>1</sup>Pamela Sklar Division of Psychiatric Genomics, <sup>2</sup>Department of Genetics and Genomic Sciences, <sup>3</sup>Icahn Institute for Data Science and Genomic Technology, <sup>4</sup>Department of Psychiatry, Icahn School of Medicine at Mount Sinai, New York, New York, USA.

<sup>5</sup>Mental Illness Research, Education and Clinical Centers, James J. Peters VA Medical Center, Bronx, New York, USA.

#Correspondence:

Gabriel E. Hoffman, PhD; Icahn School of Medicine at Mount Sinai; 1425 Madison Avenue, Icahn Building L3-70B; New York, NY, 10029, USA; [gabriel.hoffman@mssm.edu](mailto:gabriel.hoffman@mssm.edu)

Panos Roussos, MD, PhD; Icahn School of Medicine at Mount Sinai; 1470 Madison Ave, Fl. 9, Rm 107; New York, NY, 10029, USA; [panagiotis.roussos@mssm.edu](mailto:panagiotis.roussos@mssm.edu)

## Abstract

While large-scale genome-wide association studies (GWAS) have identified hundreds of loci associated with neuropsychiatric and neurodegenerative traits, identifying the variants, genes and molecular mechanisms underlying these traits remains challenging. Integrating GWAS results with expression quantitative trait loci (eQTLs) and identifying shared genetic architecture has been widely adopted to nominate genes and candidate causal variants. However, this integrative approach is often limited by the sample size, the statistical power of the eQTL dataset, and the strong linkage disequilibrium between variants. Here we developed the multivariate multiple QTL (mmQTL) approach and applied it to perform a large-scale trans-ethnic eQTL meta-analysis to increase power and fine-mapping resolution. Importantly, this method also increases power to identify conditional eQTL's that are enriched for cell type specific regulatory effects. Analysis of 3,188 RNA-seq samples from 2,029 donors, including 444 non-European individuals, yields an effective sample size of 2,974, which is substantially larger than previous brain eQTL efforts. Joint statistical fine-mapping of eQTL and GWAS identified 301 variant-trait pairs for 23 brain-related traits driven by 189 unique candidate causal variants for 179 unique genes. This integrative analysis identifies novel disease genes and elucidates potential regulatory mechanisms for genes underlying schizophrenia, bipolar disorder and Alzheimer's disease.

## Introduction

Genome-wide association studies (GWAS) have associated hundreds of loci with neuropsychiatric and neurodegenerative traits (Jansen et al., 2019; Nalls et al., 2019; Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2014; Visscher et al., 2017; Wray et al., 2018). Yet elucidating the molecular mechanisms underlying these traits remains challenging since most risk variants are non-coding and highly correlated due to linkage disequilibrium (Schaid et al., 2018; Visscher et al., 2017). Integration of risk loci with expression quantitative trait loci (eQTL) has been widely adopted to identify genes and candidate causal variants (Gallagher and Chen-Plotkin, 2018; GTEx Consortium, 2020; Hormozdiari et al., 2016). Recent work by the Genotype-Tissue Expression (GTEx) consortium across 838 individuals and 49 tissues, detected eQTLs for 95% of protein-coding and >60% of long non-coding RNA genes (GTEx Consortium, 2020). While the power to detect primary (i.e. the most significant association) eQTLs is very high, advances in identifying tissue- and cell-type-specific effects, conditionally independent effects, and candidate causal variants in trait-relevant tissues and cell types promises to further inform the molecular etiology of disease (Dobbyn et al., 2018; GTEx Consortium, 2020; Hormozdiari et al., 2016, 2018; Kim-Hellmuth et al., 2020).

Large-scale efforts have been undertaken to catalogue human brain eQTLs (Fromer et al., 2016; GTEx Consortium, 2020; Jaffe et al., 2018; Ng et al., 2017; Wang et al., 2018a). All these efforts focus on homogenate brain tissue, which is composed of multiple cell types (Cao et al., 2020; Darmanis et al., 2015; Habib et al., 2017; Lake et al., 2018), and, therefore, cell type-specific eQTLs are not fully captured (Fairfax et al., 2014; Raj et al., 2014; van der Wijst et al., 2018). This is an important limitation given that disease variants act through cell-type-specific biological effects (Farh et al., 2015; Finucane et al., 2018; Raj et al., 2014). Initial efforts have performed cell type-specific eQTL analysis in human brain by experimentally purifying specific cell types (Jaffe et al., 2020; de Paiva Lopes et al., 2020; Young et al., 2019), but the sample size of such studies are necessarily limited by the increased experimental costs, and data quality can be affected by the additional experimental steps. An alternative strategy to capture cell type-specific effects is to statistically define conditional- or context-dependent eQTL (Dobbyn et al., 2018; Kim-Hellmuth et al., 2020). While existing studies have sufficient power to detect primary eQTLs, identifying conditionally independent eQTLs that capture more subtle cell type-specific effects requires large sample sizes (Jansen et al., 2017; Zhernakova et al., 2017).

Following eQTL detection, statistical fine-mapping can identify candidate causal variants likely to drive variation in expression (Benner et al., 2016; Hormozdiari et al., 2014, 2016; Schaid et al., 2018). Going one step further, joint statistical fine-mapping integrating GWAS and gene expression traits can define the candidate causal variants that increase disease risk through alterations of gene expression (Hormozdiari et al., 2016). Interpreting and validating such variants can pinpoint genes such as *FURIN* (Schrode et al., 2019), *BIN1* (Nott et al., 2019) and *C4* (Sekar et al., 2016) along with molecular mechanisms that can be further studied in experimental systems. Yet the resolution of statistical fine-mapping for eQTL and GWAS is incomplete due to limited sample sizes and lack of trans-ancestry analysis (Schaid et al., 2018). Sample size of more than 2,000 donors is needed to detect eQTLs and perform GWAS colocalization for identification of causal variants explaining 1% of heritability (Hormozdiari et al., 2016). The largest current human brain eQTL mega-analysis by PsychENCODE included 1,387 unique donors from multiple cohorts (Wang et al., 2018a). Moreover, most eQTL analyses have been limited to European populations, despite the fact that much shorter linkage disequilibrium in individuals of African or African-American ancestry can substantially increase the resolution of statistical fine-mapping (Asimit et al., 2016; Morris, 2011; Schaid et al., 2018; Zaitlen et al., 2010).

Given the limited availability of human brain samples, it is critical to maximize power and fine-mapping resolution by combining existing datasets. Yet differences in study designs have, thus far, hindered such efforts. Trans-ethnic studies have long been challenging in genetics, but linear mixed models can control the false positive rate in the presence of complex population structure (Sul et al., 2018; Yang et al., 2014; Zhou and Stephens, 2012). Moreover, expression measurements from multiple brain regions in GTEx are not statistically independent, so combining these data entails explicit modelling of these correlated measurements from the same set of individuals (Han et al., 2016).

In order to realize the potential of trans-ethnic eQTL fine-mapping and integration with brain-related GWAS results, we developed the multivariate multiple QTL (mmQTL) pipeline and applied it to a combined analysis of brain tissues from PsychENCODE, Religious Orders Study and Memory and Aging Project (ROSMAP) and GTEx. Our pipeline performs eQTL detection with a linear mixed model, identifies conditionally independent eQTL and combines results across datasets with a random effects meta-analysis that models the correlation between multiple brain regions from a shared set of individuals. Joint fine-mapping then identifies candidate causal variants shared between gene expression and GWAS traits. This integrative

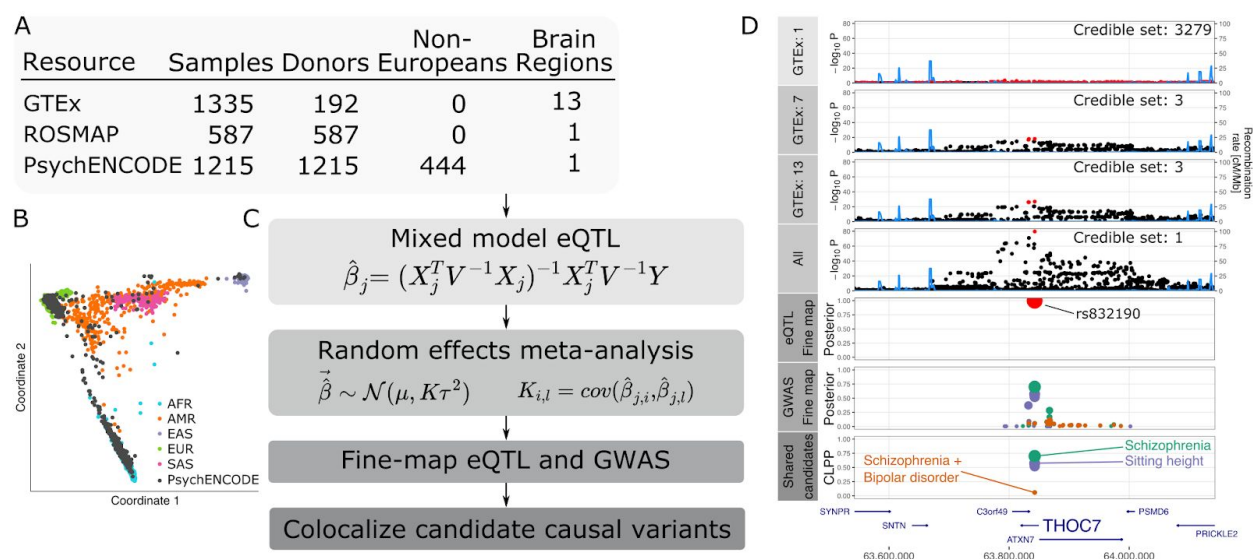
analysis identifies novel disease genes and elucidates potential regulatory mechanisms for genes underlying schizophrenia (SZ), bipolar disorder (BD) and Alzheimer's disease (AD).

## Results

### Analysis overview

We performed a trans-ethnic eQTL meta-analysis on RNA-seq gene expression data from non-overlapping samples from the dorsolateral prefrontal cortex (DLPFC) from PsychENCODE (Wang et al., 2018a) and ROSMAP (Bennett et al., 2018), and 13 brain regions from GTEx (GTEx Consortium et al., 2017) (**Figure 1A**). We accounted for diverse ancestry (**Figure 1B**), repeated measures from a shared set of donors in GTEx and effect size heterogeneity using a linear mixed model for analysis of each dataset, followed by combining these 15 eQTL analyses using a random effects meta-analysis (**Figure 1C**). This statistical framework is implemented in our mmQTL software (**see Methods**). Statistical fine-mapping of the eQTL meta-analysis was integrated with GWAS fine-mapping from CAUSALdb (Wang et al., 2020) to identify candidate causal variants shared between gene expression and neuropsychiatric traits.

For example, results for *THOC7* illustrate that increasing the number of GTEx tissue from 1 to 7 to 13 enhances power and decreases the size of the 95% credible sets, while integration with PsychENCODE and ROSMAP nominates a single candidate causal variant (**Figure 1D**). Integrating GWAS and eQTL results produces colocalization posterior probabilities (CLPP) > 0.05 for SZ, BD and sitting height, and identifies rs832190 and *THOC7* as the candidate causal variant and gene, respectively, for this locus.



**Figure 1: Workflow for trans-ethnic eQTL meta-analysis.** **A)** RNA-seq datasets with details about ancestry and repeated measures. **B)** Multidimensional scaling illustrating diverse ancestry of donors from PsychENCODE resource. **C)** mmQTL workflow is composed of eQTL analysis within each brain region for each resource using a linear mixed model to account for population stratification. Each analysis is then combined using a random effects meta-analysis that accounts for repeated measures from GTEx sample and effect size heterogeneity across brain regions and resources. Statistical fine-mapping is performed on GWAS and combined eQTL results separately. Finally, fine-mapping posterior probabilities from the eQTL analysis and each GWAS are combined to produce colocalization posterior probabilities (CLPP). **D)** Analysis of data for *THOC7* from 1, 7 and 13 GTEx brain tissues, and addition of PsychENCODE and ROSMAP, reduces the size of the 95% credible sets indicated by red points. Statistical fine-mapping for this gene and integration with GWAS nominates a single candidate causal variant, rs832190, affecting SZ, a combined risk for SZ and BD, and sitting height in this region.

## Biologically motivated simulations

Simulations motivated by the scenarios considered here (i.e. diverse ancestry and repeated measures design of the human brain datasets) were used to evaluate mmQTL performance in terms of: 1) controlling the false positive rate, 2) leveraging eQTL effects shared across multiple tissues and 3) reducing the size of the credible set from statistical fine-mapping (**Figure 2**). For the eQTL analysis we considered a linear regression model including 5 genotype PC's and a linear mixed model that counts for the genetic similarity between all pairs of samples (Sul et al.,

2018; Yang et al., 2014; Zhou and Stephens, 2012). The summary statistics for each SNP-gene pair were aggregated across tissues using a fixed- or random-effects meta-analysis, or simply the minimum p-value with a Sidak correction to account for the number of tissues. The first two explicitly account for the repeated measures design by modeling the correlation between summary statistics under the null, while the Sidak-corrected minimum p-values assume independence.

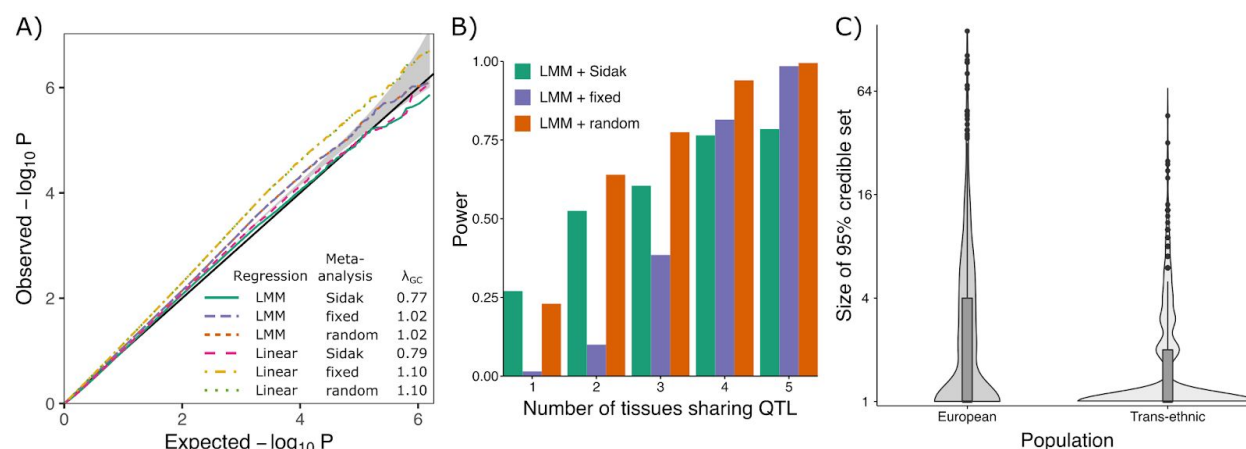
We simulated genotypes for 500 individuals in each of three distinct populations: European, African, and Asian. A single causal eQTL explaining 1-2% of expression variation in up to 5 tissues for these 1,500 individuals was simulated for 800 randomly chosen genes where the number of tissues with a shared effect varied from 1 to 5. Correlation between the same gene expression trait measured in two tissues was simulated to be low ( $r=0.12$ ) or high ( $r=0.45$ ) (**see Methods**).

In a null simulation with all genetic effects set to zero in both the low and high correlation scenarios, the linear mixed model accurately controlled the false positive rate when summary statistics from multiple tissues were aggregated using the Sidak method as well as fixed or random effects meta-analysis (**Figure 2A, Supplementary Figure 1**). As expected, the linear model did not adequately account for the complex population structure and showed an inflated false positive rate. Therefore, it was not included in subsequent simulations.

Power analyses were performed on the same set of samples of diverse ancestry where the number of tissues with a shared eQTL effect varied between 1 and 5 (**Figure 2B**). Using a p-value cutoff of  $10^{-6}$ , the random effects meta-analysis following a linear mixed model eQTL analysis had the highest power under most levels of eQTL sharing across tissues because it models heterogeneity in effect sizes across tissues. The fixed-effect meta-analysis was less powerful because it assumes a shared effect size across tissues. The Sidak corrected minimum p-value only performed best when the eQTL was tissue-specific (i.e. no cross-tissue sharing) since it assumes statistical independence of the results from each tissue.

The mmQTL workflow with linear mixed model followed by a random-effects meta-analysis demonstrated accurate control of the false positive rate while retaining high power under biologically motivated simulations. With the goal of identifying candidate causal variants shared with brain-related traits, we evaluated the benefit of using a dataset of diverse ancestry. A

dataset of 1,500 European individuals was simulated in addition to the trans-ethnic cohort above. One causal variant with effect size  $1\% \pm 0.09\%$  was used to simulate gene expression traits. Statistical fine-mapping of eQTL results from the trans-ethnic cohort produced 95% credible sets containing a mean of 2.0 SNP's compared to a mean of 4.8 for the European only cohort (**Figure 2C**). In the trans-ethnic cohort, 73.0% of genes have a single candidate causal variant compared to 51.6% in the European cohort. Moreover, random-effects meta-analysis reduces the credible set by 10.0% compared to fixed effects meta-analysis.



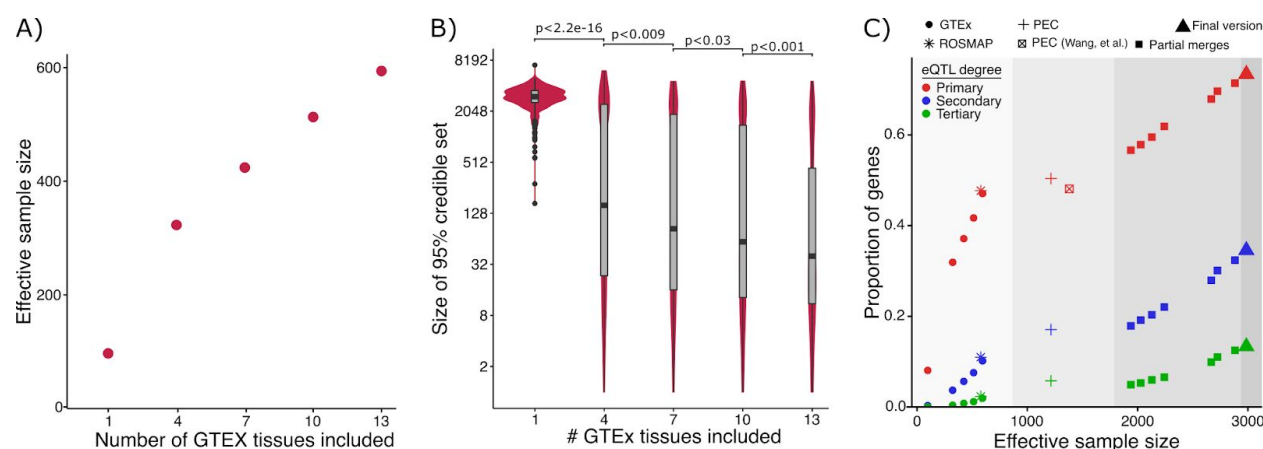
**Figure 2: Biologically motivated simulations demonstrate performance of mmQTL workflow: high correlation scenario. A)** QQ plot of results from null simulation shows that the linear mixed model (LMM) with fixed or random effect meta-analysis accurately controls the false positive rate, while linear regression with 5 genotype principal components did not. The Sidak method was very conservative in both cases.  $\lambda_{GC}$  indicates the genomic control inflation factor. **B)** Power from LMM followed by 3 types of meta-analysis versus the number of tissues sharing an eQTL. **C)** Size of the 95% credible sets from statistical fine-mapping for a dataset of European samples versus a trans-ethnic dataset of the same size.

## Evaluating mmQTL workflow on real data

Here we evaluate the empirical performance of our mmQTL workflow on real data by analysing an increasing number of brain regions ( $k=1,4,7,13$ ) from GTEx (**Figure 3**). As expected, mmQTL is able to borrow information across multiple brain regions using a random-effects meta-analysis so that increasing  $k$  substantially increases the empirical effective sample size ( $N_{eff}$ ) (**Figure 3A**). With  $k=13$ , there are 1,335 RNA-seq samples from 192 individuals producing

empirical  $N_{\text{eff}} = 524$ . Moreover, increasing  $k$  decreases the median size of the 95% credible sets from statistical fine-mapping (**Figure 3B**).

The value of adding each successive study to the meta-analysis was evaluated for primary eQTLs as well as secondary and tertiary conditional eQTLs using a conservative p-value cutoff of  $10^{-6}$  (**Figure 3C, see Methods**). The PsychENCODE study included the largest cohort and yielded 50.4% of genes having genome-wide significant primary eQTLs. Adding data from GTEx and ROSMAP produced a combined eQTL analysis comprising 3,188 RNA-seq samples from 2,029 donors to give  $N_{\text{eff}} = 2,974$ . Powered by this substantial increase in  $N_{\text{eff}}$ , eQTLs were detected for 73% of genes analysed in the final meta-analysis.



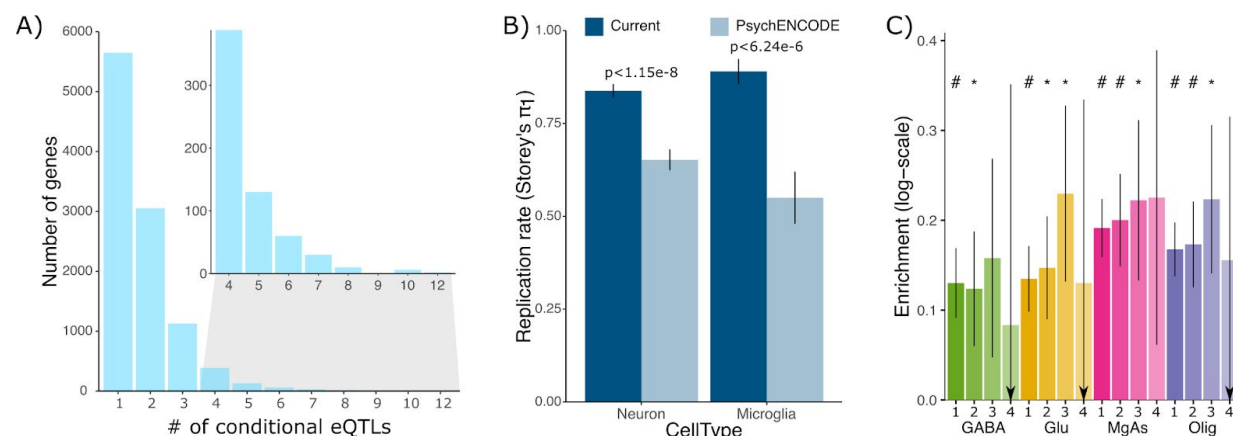
**Figure 3: Evaluation of mmQTL workflow on real data.** **A)** Increasing the number of brain regions from GTEx increases the effective sample size. **B)** Increasing the number of brain regions from GTEx decreases the median 95% credible set size. P-values are shown from one-sided Kolmogorov–Smirnov test between adjacent categories. **C)** Including additional datasets increases the proportion of genes with a detectable primary or conditional eQTL. Colors indicate degree of eQTL. Panel is divided into regions showing 1) GTEx and ROSMAP results; 2) PsychENCODE (PEC) data analyzed here, and published PEC summary statistics (Wang et al., 2018a); 3) adding an increasing number of GTEx brain tissues to the PEC+ROSMAP results; 4) final version merging PEC+ROSMAP+GTEx.

## Properties of brain eQTL meta-analysis

Our brain eQTL meta-analysis identifies 10,456 genes with a genome-wide significant eQTL, including 4,808 with at least one conditional eQTL using a conservative p-value threshold of  $10^{-6}$  (**Figure 4A**). These eQTL results are highly reproducible with estimated replicated rate

$\pi_1=75.6\%$  when evaluated in an independent dataset of bulk brain tissue (Wang et al., 2018b) using Storey's  $\pi_1$  statistic (Storey and Tibshirani, 2003). The increased power from our meta-analysis enables detection of cell-type-specific eQTL not detectable in smaller studies of bulk brain tissue. eQTLs detected in the granule cell layer of the dentate gyrus enriched for excitatory neurons (Jaffe et al., 2020), are replicated in our analysis at  $\pi_1=83.8\%$  compared to  $\pi_1=65.2\%$  in the PsychENCODE analysis (one-sided z-test  $p < 1.15e^{-8}$ ), and eQTLs detected in purified microglia (Kosoy, et al, in preparation) are replicated in our analysis at  $\pi_1=89.0\%$  compared to  $\pi_1=55.0\%$ , from the PsychENCODE analysis (one-sided z-test  $p < 6.24e^{-6}$ ) (**Figure 4B**). Overlaying variants in 95% credible sets with with ATAC-seq regions identified by fluorescence activated nuclei sorting for 4 cell populations (GABAergic neurons, glutamatergic neurons, oligodendrocytes, and a mixture of microglia and astrocytes) (Hauberg et al., 2020) identifies significant enrichment within open chromatin regions for each cell population (**Figure 4C**).

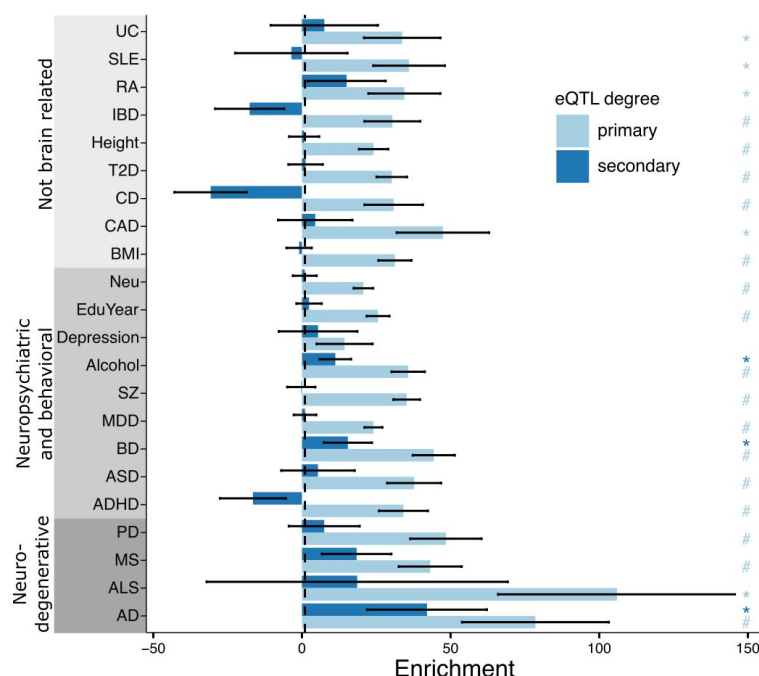
Conditional eQTLs have different properties than primary eQTLs. While primary eQTLs are a median of 24 kb from the transcription start site, conditional eQTL are more distal with median distances of 39 kb for secondary, 53 kb for tertiary and 75 kb for quaternary eQTL ( $p < 0.001$  for all comparisons of adjacent categories using one sided Kolmogorov–Smirnov test) (**Supplementary Figure 2A**). This is consistent with primary eQTLs often affecting promoters and conditional eQTL more often affecting enhancers. In addition, genes with more independent eQTLs have higher cell type specificity in human (Darmanis et al., 2015) (Spearman  $\rho = 0.0408$ ,  $p = 1.85 \times 10^{-6}$ ) and mouse (Zeisel et al., 2015) (Spearman  $\rho = 0.04435$ ,  $p = 2.18 \times 10^{-7}$ ) brain (**Supplementary Figure 2B**). Finally, genes with more conditional eQTLs tend to be under lower evolutionary constraint, as measured by the probability of loss intolerance (pLI) calculated from large-scale exome sequencing (Lek et al., 2016). While 35% of genes with no detectable eQTLs are highly constrained ( $pLI > 0.9$ ), only 10% of genes with 4 eQTLs exceed this cutoff (**Supplementary Figure 2C**).



**Figure 4: Properties of brain eQTL meta-analysis.** **A)** Number of genes having a significant primary or conditional eQTL for degree up to 12. Inset shows number of genes for eQTL degree 4 to 12. **B)** Replication rate measured by Storey's  $\pi_1$  in the current study and PsychENCODE for eQTLs discovered in the granule cell layer of the dentate gyrus enriched for excitatory neurons (Jaffe et al., 2020), and purified microglia (Kosoy, in preparation). Error bar indicates standard error from 100 bootstrap samplings. P-value indicates one-sided z-test. **C)** Enrichment of variants in the 95% causal sets for each gene in open chromatin regions assayed in each of 4 cell populations. Results are shown for eQTL degree 1 to 4. Error bars indicate standard deviation, '#' indicates Bonferroni adjusted p-value < 0.05 and '\*' indicates nominal p-value < 0.05.

## Variants in credible sets are enriched for risk for brain-related traits

Integration of variants in the 95% credible set for primary and conditional eQTLs with large-scale GWAS summary statistics using stratified linkage disequilibrium scores regression (Finucane et al., 2015) finds significant enrichments across 22 complex traits after accounting for baseline annotations (**Figure 5**). Variants in the 95% credible set for primary eQTLs were enriched for 21 traits, including 8 neuropsychiatric and behavioral traits, and 4 neurodegenerative diseases. Meanwhile, the enrichment for conditional eQTLs was limited to AD, BD and alcohol use. These enrichments indicate that our meta-analysis and statistical fine-mapping captures risk variants for brain-related phenotypes.



**Figure 5. Heritability enrichment of variants in the 95% causal set for 22 complex traits.** Linkage disequilibrium score regression (LDSC) enrichments are shown for variants in the 95% causal set for primary and secondary eQTLs. Error bars indicate standard errors. ‘#’ indicates p-value passes 5% Bonferroni cutoff for 44 tests and ‘\*’ indicates p-value < 0.05. See Supplementary Table 1 for trait abbreviations and references.

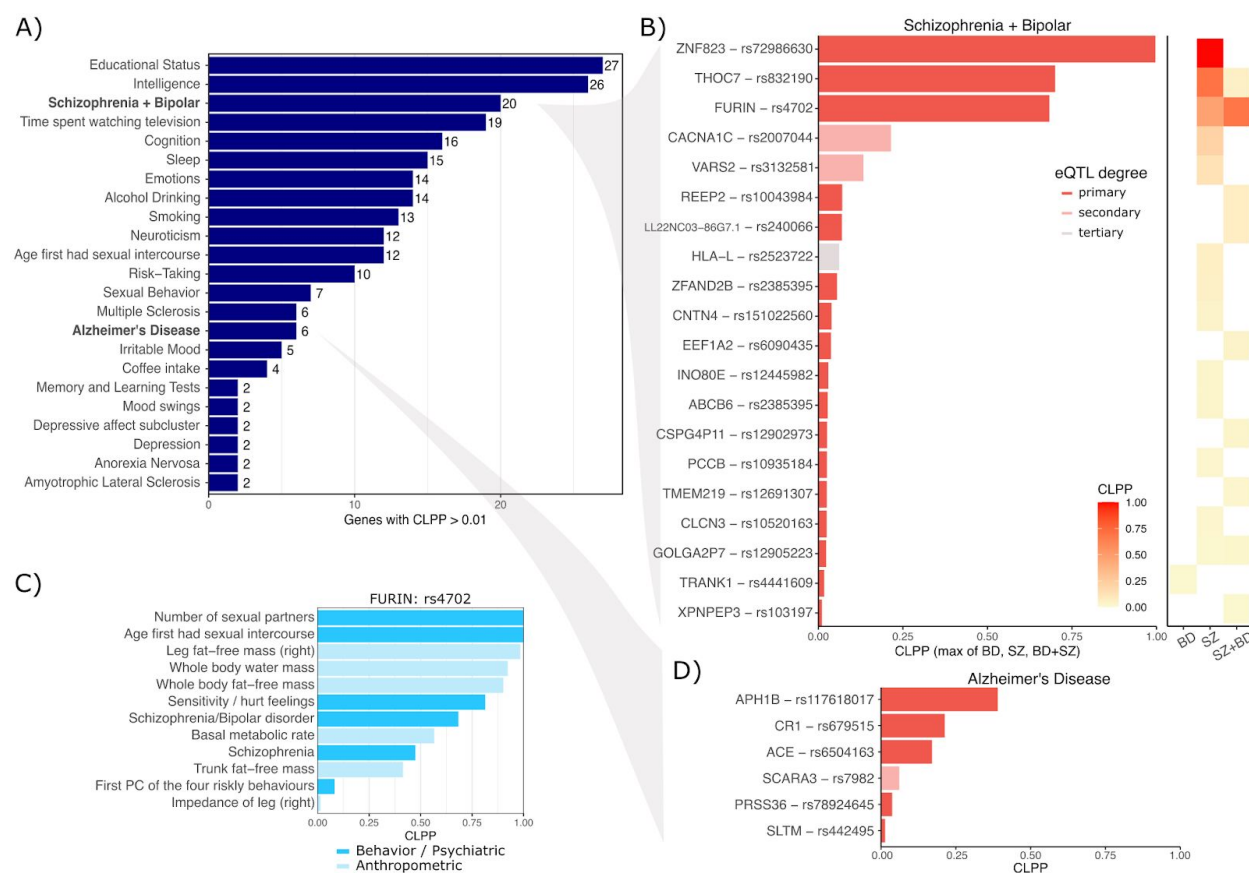
## Fine mapping identifies candidate causal variants conveying risk for brain-related traits

Integrating our eQTL fine-mapping results with candidate causal variants from large-scale GWAS (Wang et al., 2020) using a joint fine-mapping approach (Hormozdiari et al., 2016) identifies 6,978 variant-trait pairs ( $CLPP > 0.01$ ) including 2,048 unique candidate causal variants and 1626 unique genes among 683 complex traits (**Supplementary Figure 3**). These results include 301 variant-trait pairs for 23 brain-related traits for 189 and 179 unique candidate causal variants and genes, respectively (**Figure 6A**). Analysis of SZ and BD, two neuropsychiatric diseases with high genetic co-heritability (Bipolar Disorder and Schizophrenia Working Group of the Psychiatric Genomics Consortium, 2018; Cross-Disorder Group of the Psychiatric Genomics Consortium., 2019; Pardiñas et al., 2018) identified candidate causal variants for 20 genes predicted to confer risk for one or both diseases (**Figure 6B**). The top genes with  $CLPP > 0.5$  for either of these diseases include *ZNF823*, *THOC7* and *FURIN*. While

these genes have been implicated in SZ or BP previously, and in fact the candidate causal variant for *FURIN*, rs4702, has been validated experimentally (Schrode et al., 2019), candidate causal variants for the other two genes have not been previously identified. Moreover, integrating results from analysis of SZ, BP and SZ+BP versus controls indicates the specificity of these candidates causal variants. *ZNF823* is predicted to confer risk to SZ, but not BD. *THOC7* has a substantially larger CLPP score for SZ compared to the joint SZ+BP GWAS. Conversely, *FURIN* has a higher CLPP for the joint SZ+BP GWAS than for SZ alone. Notably, the candidate causal variants driving the colocalization with SZ and BD for *CACNA1C* and *VAR2* are in fact secondary eQTLs, emphasizing the importance of including conditional eQTL analysis. The candidate causal variant for *CACNA1C*, rs2007044, which is within an intronic enhancer, has been previously shown to affect transcription due to reduced promoter interaction (Roussos et al., 2014).

In addition, analysis of candidate causal variables across many phenotypes enables insight into pleiotropy. *FURIN* and rs4702 are also implicated in the number of sexual partners, age at first sexual intercourse, risk taking behavior, and emotional sensitivity / hurt feelings, and multiple anthropometric traits (**Figure 6C, Supplementary Figure 4**). Sharing of a candidate causal variant and gene between SZ+BP and these risk-taking behavior traits is particularly interesting given that impulsiveness is a clinical feature of both SZ and BD (Najt et al., 2007; Ouzir, 2013), and is associated with more severe psychiatric symptoms and decreased level of functioning (Cerimele and Katon, 2013).

Analysis of AD identified candidate causal variants for 6 genes. While these genes have been highlighted previously (Jansen et al., 2019), our analysis highlights variants and their mechanistic link to disease (**Figure 6D**).



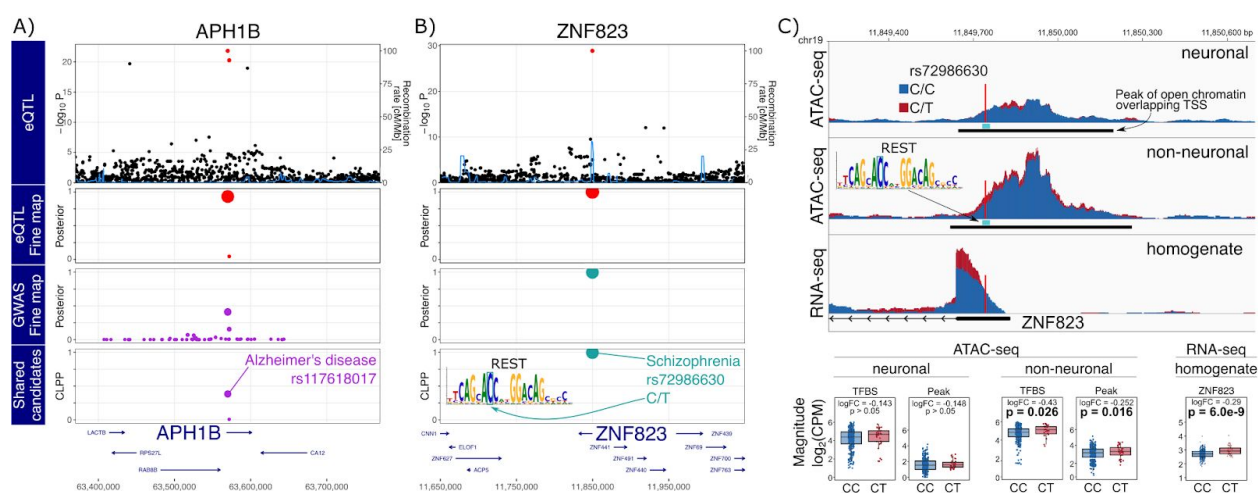
**Figure 6: Summary of joint fine-mapping colocalization with brain-related traits** **A)** Number of genes with colocalization posterior probability (CLPP) > 0.01 for **B)** Genes with CLPP > 0.01 for Schizophrenia (SZ) and Bipolar Disorder (BP) and a joint GWAS of SZ+BP versus controls. For each gene, the max CLPP across SZ, BP and SZ+BP is shown. Right Panel shows CLPP for BP, SCZ and SZ+BP compared to controls. **C)** A validated casual variant, rs4702, that affects expression of *FURIN* is predicted to affect risk for multiple complex behavioral, psychiatric and anthropometric traits. **D)** Genes with CLPP > 0.01 for AD.

## Candidate causal variants elucidate potential molecular mechanisms

rs117618017 is the top causal variant for AD and drives the expression of *APH1B*, a subunit of the gamma-secretase complex, which includes multiple AD risk genes as components (**Figure 7A**). This missense coding variant was identified in a GWAS meta-analysis for AD (Jansen et al., 2019), but an attempt to experimentally validate a functional effect from this single amino acid change yielded only negative results (Zhang et al., 2020). Yet our analysis indicates an alternative molecular mechanism, whereby, instead of acting by changing protein sequence, the

minor allele of rs117618017 increases AD risk by directly increasing gene expression of *APH1B*.

The top hit for SZ is rs72986630, which is predicted to drive expression of *ZNF823*, a zinc finger protein with little additional annotation (**Figure 7B**). This C/T SNP is located in the 5' UTR of the gene and the minor allele, T (MAF ~6%), is protective against SZ. This variant is predicted to disrupt a binding site for the RE1 silencing transcription factor (REST), also known as neuron-restrictive silencing factor. REST is upregulated during neurogenesis and in adult non-neuronal cells, and acts by silencing neuron specific genes (Hwang and Zukin, 2018; Schoenherr and Anderson, 1995). Analysis of chromatin accessibility in this region using a large-scale ATAC-seq dataset from purified neuronal and non-neuronal nuclei from the anterior cingulate cortex (ACC) of post mortem brains of 368 donors elucidated the molecular mechanism (Bendl, et al. in preparation) (**Figure 7C**). In non-neuronal cells, but not in neuronal cells, individuals heterozygous at this site have higher chromatin accessibility at both the 644 bp ATAC-seq peak ( $p = 0.016$ ) and the 21 bp motif ( $p = 0.026$ ), and this corresponds to decreased binding of REST at this site. Since REST is a transcriptional silencer, decreased binding of REST should lead to increased expression of *ZNF823*. Querying RNA-seq data from brain homogenate from these samples confirms that heterozygous individuals have increased expression of *ZNF823* ( $p=6.01 \times 10^{-9}$ ).



**Figure 7: GWAS-eQTL colocalization by joint fine-mapping. A,B)** Starting from the top, the plot shows  $-\log_{10} p$ -values from eQTL analysis, poster probabilities from statistical fine-mapping of eQTL results, poster probabilities from statistical fine-mapping of GWAS results, and colocalization posterior probabilities (CLPP) for combining eQTL and GWAS fine-mapping. **A)**

Expression of *APH1B* and AD risk share rs117618017 as a candidate causal variant. **B)** Expression of *ZNF823* and SZ risk share rs72986630 as a candidate causal variant. This variant is predicted to disrupt a REST binding site motif. **C)** Individuals heterozygous for rs72986630 have increased chromatin accessibility at the peak and REST binding site in non-neuronal cells. Genome-plot shows chromatin accessibility for neuronal (top) and non-neuronal (middle) nuclei, and gene expression from brain homogenate bottom. The lower panel shows boxplots comparing chromatin accessibility and gene expression between individuals with two reference alleles (i.e. CC) compared to CT heterozygotes.

## Discussion

Integration of eQTL and GWAS is a powerful method to understand the molecular mechanism influencing complex traits. While transcript-wide association studies aim to identify genes underlying a complex trait, correlated expression and co-regulation can be challenging to overcome (Mancuso et al., 2019; Wainberg et al., 2019). Joint fine-mapping focuses instead on identifying variants that drive both gene expression and a downstream trait (Hormozdiari et al., 2016). Despite recent successes, fine-mapping is often limited by statistical power and linkage disequilibrium (Hormozdiari et al., 2016; Schaid et al., 2018). Our mmQTL workflow addresses both of these issues by performing a trans-ethnic eQTL meta-analysis of 3,188 RNA-seq samples from 2,029 donors, with an effective sample size of 2,974, to produce the largest resource to date characterizing the genetics of gene expression in the human brain. This analysis has substantially boosted the catalog of genes with detected conditional eQTLs, while increasing the resolution of statistical fine-mapping.

Despite being performed on bulk RNA-seq data, our analysis is able to replicate eQTLs discovered in purified microglia (Kosoy, in preparation) and neurons (Jaffe et al., 2020), and the replication rate is substantially higher than for PsychENCODE (Wang et al., 2018a). Moreover, we identify candidate causal variants enriched in cell type specific open chromatin regions. While much recent work has pursued generating eQTLs from purified cell populations (Jaffe et al., 2020; de Paiva Lopes et al., 2020; Young et al., 2019), and eQTL discovery from single cell/nucleus RNA-seq is becoming tractable (Mandric et al., 2020; van der Wijst et al., 2020), our eQTL meta-analysis from bulk tissue illustrates that large sample size and sophisticated statistical modelling has substantial power to replicate eQTLs from smaller studies of purified cell types.

While the number of genes with detectable eQTLs approaches saturation, there is substantial value in increasing sample size. Here, we use individuals of diverse ancestry paired with a linear mixed model in our mmQTL workflow to increase the resolution of statistical fine-mapping. Moreover, we perform conditional eQTL analysis to identify genes with up to 12 independent eQTLs. These conditional eQTLs tend to have smaller effect sizes, be farther from transcription start sites, and affect genes that are more cell type specific. The number of genes with secondary and tertiary eQTL does not appear close to saturation, underscoring the regulatory variation that remains to be identified.

Integrating statistical fine-mapping for eQTLs and GWAS across hundreds of complex traits enabled insight into candidate causal variants, mechanisms of disease genetics and pleiotropy. Focusing on regulatory mechanisms for genes underlying brain-related traits, we identified 20 genes and candidate causal variants predicted to drive risk for SZ and BD, plus another 6 for AD. While other methods focus on discovering disease genes, here we focus on discovering gene-variant pairs underlying disease risk in order to elucidate the molecular mechanisms that convey risk.

Here we highlighted two examples. The SNP rs117618017 is a candidate causal variant causing a single amino acid change in *APH1B*. While experimental results of the impact of this amino acid change were negative (Zhang et al., 2020), our analysis instead supports a mechanism where this variant increases disease risk by increasing expression of *APH1B*. Our analysis predicts that rs72986630 drives expression of *ZNF823* and is protective against SZ. By integrating chromatin accessibility data from post mortem brains, we traced the predicted chain of causality and found that the minor allele disrupts binding of REST in non-neuronal cells, which then increases expression of *ZNF823*. The lack of an effect in neuronal nuclei is consistent with the higher expression of REST in non-neuronal cells during adulthood, silencing neuron-specific genes (Hwang and Zukin, 2018; Schoenherr and Anderson, 1995).

While we focused on regulatory mechanisms for genes underlying SZ, BD and AD, all results are available from the Brain eQTL meta-analysis (BREMA) resource ([icahn.mssm.edu/brema](https://icahn.mssm.edu/brema)).

Further integration of multi-omics data with trans-ethnic fine-mapping and large-scale GWAS promises to yield further insight into the molecular mechanisms underlying disease risk. Future

studies are poised to perform multiple genomic assays, namely RNA-seq and ATAC-seq, on multiple tissues or brain regions, and target multiple cell types either by sorting or single cell/nucleus methods (Mandric et al., 2020; van der Wijst et al., 2020). Moreover, these studies will increasingly include individuals of diverse ancestry (Wojcik et al., 2019). Our mmQTL method will enable the field to take advantage of these repeated measures datasets while modeling effect size heterogeneity and controlling the false positive rate. Efforts to trace the chain of causality from variants and molecular mechanisms to pleiotropy across complex phenotypes are poised to yield insight into novel therapeutic targets.

## Resources

Brain eQTL meta-analysis resource: <http://icahn.mssm.edu/brema>

mmQTL: <https://github.com/jxzb1988/mmQTL>

## Contributions

B.Z., G.E.H. and P.R. conceived and designed the study; B.Z. designed and implemented the statistical method; B.Z. and G.E.H. performed the analysis; J.F.F. generated the cell type specific expression and chromatin accessibility data; J.B. and R.K. preprocessed and analyzed cell type specific expression and chromatin accessibility data; J.F.F. and P.R. supervised the data generation; G.E.H. and P.R. supervised the data analysis; G.E.H., B.Z. and P.R. wrote the manuscript with the help of all authors.

## Funding

The project was supported by the National Institute of Mental Health, NIH grants R01-MH109677, U01-MH116442, R01-MH125246 and R01-MH109897, the National Institute on Aging, NIH grants R01-AG050986, R01-AG067025 and R01-AG065582, and the Veterans Affairs Merit BX004189 to P.R.. G.E.H. was supported in part by NARSAD Young Investigator Grant 26313 from the Brain & Behavior Research Foundation. J.B. was supported in part by NARSAD Young Investigator Grant 27209 from the Brain & Behavior Research Foundation. Research reported in this paper was supported by the Office of Research Infrastructure of the National Institutes of Health under award numbers S10OD018522 and S10OD026880. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

## Methods

### Obtaining and processing of RNA-seq datasets

Imputed genotypes from GTEx v7 were downloaded from dbGAP (Accession [phs000424.v7.p2](#)). For ROSMAP, the imputed genotypes were downloaded from the Synapse website (id: [syn3157329](#)). The imputed genotype for each cohort in the PsychENCODE study was downloaded from Synapse website (id: [syn21052530](#)), and was then filtered to retain variants with imputation quality  $\geq 0.3$ . Filtered genotypes from each cohort were merged and variants with MAF  $\geq 1\%$  and Hardy-Weinberg Equilibrium p-value  $\geq 1e-6$  were retained.

The original PsychENCODE analysis performed eQTL detection using 1,387 individuals (Wang et al., 2018a). In the current work we exclude a small fraction of these individuals. First, the full PsychENCODE dataset contains GTEx samples which we excluded in order to avoid redundancy with our separate GTEx analysis. Second, the original analysis used ~5 million imputed SNPs. Since accurate statistical fine-mapping depends on including the true causal variant in the analysis, we included additional well-imputed SNPs at the cost of excluding a small set of samples. Excluding samples with  $< 8$  million well-imputed (info score  $> 0.3$ ) variants yielded 1,215 individuals used in this study.

The normalized gene expression of GTEx v7 was downloaded from GTEx Portal (GTEx\_Analysis\_2016-01-15\_v7\_RNASeQCv1.1.8\_gene\_tpm.gct.gz, <https://gtexportal.org/>), and we regressed out covariates from the companion file GTEx\_Analysis\_v7\_eQTL\_covariates.tar.gz with linear regression. Normalized data from PsychENCODE (DER-01\_PEC\_Gene\_expression\_matrix\_normalized.txt) was downloaded from <http://resource.psychencode.org/>, and as the downloaded gene expression is already normalized regressing out the effect of covariates, no further normalization was taken. Data from ROSMAP (syn3388564, ROSMAP\_RNAseq\_FPKM\_gene.tsv) was downloaded from <https://adknowledgeportal.synapse.org>. The provided FPKM abundance values were quantile normalized, log2 transformed, standardized to normal distribution, and 20 principal components of the gene expression matrix were regressed out.

### Linear mixed model eQTL analysis

Given expression abundance of a gene measured in  $n$  tissues from same set of individuals, the gene expression in tissue  $t$  can be modeled as:

$$y_{i,t} = x_{i,j}\beta_{j,t} + \sum_{k=1}^m x_{i,k}\alpha_{k,t} + \varepsilon_{i,t} \quad (1)$$

where,  $y_{i,t}$  is the measured gene expression value for individual  $i$  in tissue which has been normalized so that it has mean 0 and variance 1,  $x_{i,j}$  is the genotype dosage for individual  $i$  at variant  $j$  normalized so that it has mean 0 and variance 1,  $\beta_{j,t}$  is effect size for variant  $j$  and tissue  $t$ . The next term models the polygenic background across  $m$  variants where  $x_{i,k}$  is the genotype dosage value for individual  $i$  at variant  $k$  and  $\alpha_{k,t}$  is the effect size for variant  $k$  and tissue  $t$  with distribution  $\mathcal{N}(0, \sigma_{gt}^2)$ , where  $\sigma_{gt}^2$  is the tissue-specific parameter for genetic background. Finally,  $\varepsilon_{i,t}$  is the normally distributed error variance for individual  $i$  and tissue  $t$  with distribution  $\mathcal{N}(0, \sigma_{\varepsilon_t}^2)$ , where  $\sigma_{\varepsilon_t}^2$  is the tissue-specific parameter for random noise.

This linear mixed model can be transformed for practical estimation the effect size  $\beta_{j,t}$ . Equation (1) can be rewritten as

$$y_{i,t} = x_{i,j}\beta_{j,t} + \hat{\varepsilon}_{i,t} \quad (2)$$

where  $\hat{\varepsilon}_{i,t} = \sum_{k=1}^m x_{i,k}\alpha_{k,t} + \varepsilon_{i,t}$  and has a distribution  $\mathcal{N}(0, K\sigma_{gt}^2 + \sigma_{\varepsilon_t}^2)$ , where  $K$  is a genetic relatedness matrix estimated based on genome-wide genotypes.

Considering that the phenotype was collected among  $l$  individuals, we can write formula (2) into a vector format:

$$Y_t = X_j\beta_{j,t} + \hat{\varepsilon}_t \quad (3)$$

where  $Y_t$ ,  $X_j$  and  $\hat{\varepsilon}_t$  are  $l$ -dimensional vectors, and contain normalized phenotype, normalized genotype of variant  $j$ , and noise, respectively.

From Equation (3),  $\beta_{j,t}$  can be estimated as

$$\hat{\beta}_{j,t} = (X_j^T V^{-1} X_j)^{-1} (X_j^T V^{-1} Y_t) \quad (4)$$

where  $V = K\sigma_{gt}^2 + \sigma_{\varepsilon_t}^2$  and produces an unbiased estimator since

$$\begin{aligned} E(\hat{\beta}_{j,t}) &= E((X_j^T V^{-1} X_j)^{-1} (X_j^T V^{-1} (X_j\beta_{j,t} + \hat{\varepsilon}_t))) \\ &= \beta_{j,t} + E((X_j^T V^{-1} X_j)^{-1} X_j^T V^{-1} \hat{\varepsilon}_t) = \beta_{j,t} + E(\tilde{\varepsilon}_t) = \beta_{j,t}. \end{aligned}$$

## Modeling covariance across tissues

While standard meta-analysis assumes that effect size estimates are statistically independent, analysis of multiple tissues from the same set of subjects produces covariance between the coefficient estimates. Here we explicitly model this covariance in order to control the false positive rate.

Denote the estimate for variant  $j$  across all tissues as the vector  $\hat{\beta}_j = [\hat{\beta}_{j,1}, \hat{\beta}_{j,2}, \dots, \hat{\beta}_{j,l}]$ . Since individuals overlap across the multiple tissues, the entries of  $\hat{\beta}_j$  will be correlated. Estimating coefficients for tissues 1 and 2 using Equation 4 gives

$$\hat{\beta}_{j,1} = (X_1^T V_1^{-1} X_1)^{-1} (X_1^T V_1^{-1} Y_1) \quad (5)$$

$$\hat{\beta}_{j,2} = (X_2^T V_2^{-1} X_2)^{-1} (X_2^T V_2^{-1} Y_2) \quad (6)$$

where an index is added to distinguish the two tissues which may have partial sample overlapping. These estimates are not statistically independent since

$$\begin{aligned} E(\hat{\beta}_{j,1} \hat{\beta}_{j,2}) &= E((X_1^T V_1^{-1} X_1)^{-1} (X_1^T V_1^{-1} Y_1) (X_2^T V_2^{-1} X_2)^{-1} (X_2^T V_2^{-1} Y_2)) \\ &= E((X_1^T V_1^{-1} X_1)^{-1} (X_1^T V_1^{-1} (X_1 \beta_{j,1} + \hat{\epsilon}_1)) (X_2^T V_2^{-1} X_2)^{-1} (X_2^T V_2^{-1} (X_2 \beta_{j,2} + \hat{\epsilon}_2))) \\ &= E(\beta_{j,1} \beta_{j,2}) + E(C \hat{\epsilon}_1 \hat{\epsilon}_2) \end{aligned}$$

where  $C = (X_1^T V_1^{-1} X_1)^{-1} X_1^T V_1^{-1} V_2^{-1} X_2 (X_2^T V_2^{-1} X_2)^{-1}$  is only involved with transformed genotypes projected by a covariance matrix. Noting that  $\hat{\epsilon}_1$  and  $\hat{\epsilon}_2$  are the summed contribution from polygenic background and noise, if there are sample overlapping and the phenotypes share causal variants in two tissues, then

$$E(C \hat{\epsilon}_1 \hat{\epsilon}_2) = C N_{shared} \sigma_{g,1} \sigma_{g,2} \neq 0,$$

where  $N_{shared}$  is the number of shared individuals, and  $\sigma_{g,1}$  and  $\sigma_{g,2}$  are the genetic component for polygenic background in tissue 1 and tissue 2. Finally we note that

$$cov(\hat{\beta}_{j,1}, \hat{\beta}_{j,2}) = E((\hat{\beta}_{j,1} - \beta_{j,1})(\hat{\beta}_{j,2} - \beta_{j,2})) = E(C \hat{\epsilon}_1 \hat{\epsilon}_2),$$

explicitly indicating that there is nonzero covariance between estimators. Our mmQTL method estimates the covariance matrix among  $n$  tissues based on the non-significant z-score in tissues, and set it to be  $\hat{C}$ . This matrix is defined so that the covariance between tissues  $i$  and  $j$  is estimated by the covariance between z-scores from non-significant variants ( $p > 0.05$ ) according to:

$$\hat{C} = cov(Z_i, Z_j),$$

where  $Z_i$  and  $Z_j$  are vectors containing statistical Z-scores.

## Fixed-and random effects meta-analysis

The results from multiple analyses are aggregated using either a fixed or random effects meta-analysis. The true effects sized are assumed to be drawn from a normal distribution  $\mathcal{N}(\beta, \sigma_\beta^2)$  centered at the true effect size  $\beta$  with variance  $\sigma_\beta^2$ . For a fixed effect model, the true effect size is fixed at a constant value which is equivalent to setting  $\sigma_\beta^2 = 0$  and for the random effects model  $\sigma_\beta^2 \geq 0$ . From this hierarchical framework, we obtain estimators for variant  $j$  among tissues, denoted as a vector  $\hat{\beta}_j = \beta_j \mathbf{1} + \tilde{\epsilon}_j$  which has a distribution  $\mathcal{N}(\beta_j \mathbf{1}, \sigma_{g_j}^2 I + \hat{C})$ . We applied the Brent-method implemented in C++ Boost library (<https://www.boost.org>) to estimate  $\beta_j$  and  $\sigma_{g_j}^2$ . To test the difference with null hypothesis, we applied the random-effect model2 (Han and Eskin, 2011; Han et al., 2016) to obtain a p-value.

## Detection of conditional eQTLs

We applied a stepwise selection strategy explore cis-region and identify conditionally independent eQTL associations. An iterative strategy is applied to find conditional independent eQTL: previously detected eQTL signals are regressed out and another round of eQTL detection was initiated. If one or more variants with p-value less than  $10^{-6}$ , the variant with the smallest p-value is added to the list of conditionally independent effects. The process is repeated until no addition variant has a p-value  $< 10^{-6}$ . If a high-order eQTL is in high LD with low-order eQTL ( $r^2 \geq 0.3$ ), the high-order eQTL will be excluded in order to avoid attenuating the estimated effect size of low-order eQTL.

Importantly, we demonstrate statistically that the order in which conditional eQTLs are detected is biologically meaningful: large-effect eQTL shared among tissues are likely to be detected first, while and small-effect eQTL or tissue-specific effect will be detected as higher-order eQTL.

Consider two true causal variants  $i$  and  $j$  where the estimated effect has the distribution around the true value according to  $\hat{\beta}_i \sim \mathcal{N}(\beta_i, \sigma_i^2 I + \hat{C}_i)$ , where  $\hat{C}_i$  is defined above. The non-centrality parameter (NCP) reflecting the statistical power for this variant is

$$NCP_i = \frac{\beta_i}{\sqrt{(N * tr((\sigma_i^2 I + \hat{C}_i)^{-1}))^{-1}}}.$$

The ratio between the NPC of variant  $i$  and the second variant  $j$  is

$$\frac{NCP_i}{NCP_j} = \frac{\beta_i}{\beta_j} / \sqrt{\frac{tr((\sigma_j^2 I + \hat{C}_j)^{-1})}{tr((\sigma_i^2 I + \hat{C}_i)^{-1})}} \quad (7)$$

From empirical observation that the effect size (with the genotype and response normalized) of primary eQTL is much larger than that of non-primary eQTL,  $\hat{C}_i \approx \hat{C}_j$  and both are positive definite, and can be decomposed into  $U\Sigma U$ , in which  $U$  consists of the eigenvectors, and  $\Sigma$  is a diagonal matrix with elements being eigenvalues, denoted as  $diag(\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_N)$ .  $\sigma_i^2 I + \hat{C}_i$  can be decomposed into  $U(\Sigma + \sigma_i^2)U$ , and  $\sigma_j^2 I + \hat{C}_j$  to be  $U(\Sigma + \sigma_j^2)U$ . Therefore, Equation (7) can be rewritten as

$$\frac{NCP_i}{NCP_j} = \frac{\frac{\beta_i}{\beta_j}}{\frac{\sum_k^N (\sigma_j^2 + \lambda_k)^{-1/2}}{\sum_k^N (\sigma_i^2 + \lambda_k)^{-1/2}}}.$$

Based on this, it is apparent that the difference in statistical power for variants  $i$  and  $j$  is mainly determined by effect size, and its variance. For a variant with larger effect size, and smaller variance, it has a higher statistical power, which is consistent with the empirical that primary eQTL has a much larger normalized effect size ( $\beta_i > \beta_j$ ), and smaller variance because of the

sharing among tissues ( $\sigma_j^2 > \sigma_i^2$ ), so  $\frac{NCP_i}{NCP_j} > 1$ . It follows that on average, mmQTL will pick the independent eQTL signal in a biologically meaningful manner, so that eQTL with a larger influence on expression abundance among conditions tends to be selected first.

## Multiple testing

Multiple testing correction is performed at the level of the locus as well as genome-wide. Empirically, we performed the locus-level control applying Bonferroni correction, which is a most conservative strategy, and Benjamini-Hochberg method (Benjamini and Hochberg, 1995) on genome-wide correction, and we found that a p-value  $10^{-6}$  is enough for two-level multiple test correction. While studies often use more liberal multiple testing cutoffs because of the limited

statistical power, the statistical fine-mapping that is the focus of this analysis can perform poorly on genes that only pass a liberal cutoff (Hormozdiari et al., 2016, 2018).

## Computing empirical effective sample size

In linear regression model for QTL analysis, given that both phenotype and genotype were normalized, the estimator for the allelic effect size is  $\hat{\beta} = (X^T X)^{-1} X Y$ , and its variance is  $var(\hat{\beta}) = (X^T X)^{-1} \sigma_e^2$ . Letting  $R_i^2$  be the variance explained by the explored variant and  $N_i$  be the (effective) sample size for study  $i$ , the variance of the effect size estimate is  $var(\hat{\beta}_i) = \frac{1 - R_i^2}{N_i}$ . Consider two studies, where the (effective) sample size of the first study is easy to estimate just by using the number of samples, and the second has some confounding factors such as repeat measurements or population structure. Assuming that the effect size of a given causal variant is constant in the two studies, the ratio of the variances is determined only by  $N_1$  and  $N_2$ :

$$\frac{var(\hat{\beta}_2)}{var(\hat{\beta}_1)} = \frac{N_1}{N_2}.$$

Therefore the effective sample size,  $N_2$ , can be computed from known values by

$$N_2 = \frac{N_1 var(\hat{\beta}_1)}{var(\hat{\beta}_2)}.$$

We used individual brain tissue in GTEx dataset as study 1 to define  $N_1$ , and eQTL results from fixed-effect meta-analysis as study 2. The genome-wide variance ratio was set to be the median ratio of variances based on all variants with  $abs(z\text{-score}) \geq 10$  in the fixed-effect meta analysis. When evaluating the effective size of a meta-analysis, the effective sample size was computed by treating each brain tissue in GTEx as the baseline study and then taking the mean estimate effective sample size over 13 brain regions.

## Replication of eQTLs from purified cell types

In order to assess the replication of eQTLs discovered in independent datasets, we considered the lead SNP for each gene with a genome-wide significant eQTL in the granule cell layer of the dentate gyrus enriched for excitatory neurons (Jaffe et al., 2020) and purified microglia (Kosoy, et al., in preparation). For the set of lead SNPs from each dataset, the p-values were extracted

from the current eQTL analysis as well as the PsychENCODE analysis (Wang et al., 2018a) and Storey's  $\pi_1$  was evaluated using qvalue (Storey and Tibshirani, 2003). The PsychENCODE p-values were obtained from [http://resource.psychencode.org/Datasets/Derived/QTLs/Full\\_hg19\\_cis-eQTL.txt.gz](http://resource.psychencode.org/Datasets/Derived/QTLs/Full_hg19_cis-eQTL.txt.gz).

Uncertainty in  $\pi_1$  estimates were evaluated using 100 bootstraps where SNPs were sampled with replacement and  $\pi_1$  was recomputed each time. A p-value comparing the replication rate for the current and PsychENCODE analysis was computed using a one-sided z-test using the estimated  $\pi_1$  values and their bootstrap variances.

## Simulation pipeline to evaluate mmQTL performance

Genotype and gene expressions data simulated to compare to empirical performance of eQTL analysis using a linear model with 5 genotype principal components compared to a linear mixed model. Results from eQTL analysis of 5 simulated tissues were then aggregated using either Sidak correction, or a fixed- or random-effects meta-analysis.

Biologically realistic genotype data reflecting real human populations was simulated with a sampling-based simulation package, hapgen2 (Su et al., 2011), and haplotype information for European, African, and Asian populations from the 1000 Genomes Project ([https://mathgen.stats.ox.ac.uk/impute/data\\_download\\_1000G\\_2010\\_interim.html](https://mathgen.stats.ox.ac.uk/impute/data_download_1000G_2010_interim.html)). We simulated 500 individuals for each population, and merged these individuals into a single trans-ancestry dataset with sample size 1,500. We also simulated 1,500 individuals solely based on European haplotype information.

Based on these genotypes, we adapted the phenotypesimulator pipeline (Meyer and Birney, 2018) to perform 800 simulations for each scenario, simulating one gene expression trait for each simulation. For each gene a single eQTL was simulated to affect expression abundance explaining 1% phenotypic variance, and the contribution due to polygenic background was set to be 30%. We applied phenotypesimulator's simulating strategy to account for shared environmental factors, measurement noise, and polygenic background to create correlated phenotypes. In the simulation, we simulated phenotypes in 5 tissues, and set the number of tissues that the causal genetic variant affects to be 1, 2, 3, 4, 5. To demonstrate the robustness of mmQTL to control for population structure and batch effect, we set two different levels of phenotype correlation, a low level,  $r=0.12$ , and the high level  $r=0.45$ . For power analysis, any

simulated causal variants located in high LD ( $r^2 \geq 0.8$ ) with a variant passing the multiple testing cutoff was considered to be detected.

We also performed a null simulation with no true causal variants where all effect sizes were set to zero. Results from 50 simulations were aggregated and we used genomic inflation factor (Devlin and Roeder, 1999) and QQ plots to assess the false positive rate.

Comparison of fine-mapping resolution between a European and trans-ethnic population was performed using simulated pure 1,500 European individuals and trans-ethnic 1,500 individuals with 500 individuals in each of European, African and Asian population. Gene expression phenotype was simulated in a single tissue, and a causal variant was randomly chosen to explain 2% phenotypic variance. For 1,500 European individuals, we applied standard linear regression model to detect eQTL and then fine-mapping was conducted to obtain a 95% credible set candidate for causal variants, while for 1,500 trans-ethnic individuals, we used mixed linear model to detect eQTL and then fine-mapping was taken to find a 95% credible set. The size of the 95% credible set was used to compare the fine-mapping resolution, a smaller number indicating a higher fine-mapping resolution.

## Integration with ATAC-seq data

Variants in the 95% credible set were overlaid with open chromatin regions from four distinct populations of cells (glutamatergic neurons, GABAergic neurons, oligodendrocytes, and a mixture of microglia/astrocytes) identified by ATAC-seq (Hauberg et al., 2020). In order to reduce the influence of the low fine-mapping resolution of conditional eQTL, if the size of the 95% credible set for a single gene contained  $>10$  variants, only the 10 variants with highest PIP were included. Enrichment of variants within open chromatin regions was evaluated using a Fisher's exact test implemented in QTLTools (Delaneau et al., 2017).

## Evaluating GWAS enrichments for variants in credible sets

We applied a strategy developed in Hormozdiari et al. (Hormozdiari et al., 2018): for each eQTL, we performed fine-mapping and compute the causal posterior probability (CPP) of each cis-SNP and only variants in the fine-mapped 95% credible set are left for following analysis. For each SNP in cis-regions, we assign an annotation value based on the maximum value of CPP across all molecular phenotypes; SNPs that do not belong to any 95% Credible Set are assigned an annotation value of 0, which is referred as MaxCPP in (Hormozdiari et al., 2018). Stratified

linkage disequilibrium (LD) score regression (S-LDSC) (Finucane et al., 2015) was then used to partition trait heritability using the constructed functional annotations, and the estimated enrichment was used to measure the importance of each eQTL category on human complex traits or diseases. To rule out the potential influences of the correlation among eQTL categories, we aggregate the baselineLD model, which includes a set of 75 functional annotations, and functional annotations for eQTL categories and run S-LDSR simultaneously.

GWAS summary statistics were obtained for 22 human complex traits or diseases, which contain both brain traits and non-brain traits (**Supplementary Table 1**). The summary results were firstly converted into the required format for LDSR by the provided “`munge_sumstats.py`” command in the LDSC package (<https://github.com/bulik/ldsc>).

## eQTL detection in cell type specific datasets

Microglia from fresh human brain specimens (101 samples, including 27 non-Europeans) were prepared using the Adult Brain Dissociation Kit (Miltenyi Biotech). Tissue homogenates were incubated in antibody (CD45: BD Pharmingen, Clone HI30 and CD11b: BD Pharmingen, Clone ICRF44) at 1:500 for 1 hour in the dark at 4°C with end-over-end rotation. Prior to FACS, DAPI (Thermoscientific) was added at 1:1000 to facilitate identification of dead cells. Viable (DAPI negative) CD45<sup>+</sup>/CD11b<sup>+</sup> cells were isolated by FACS using a FACS Aria flow cytometer (BD Biosciences) (Kosoy et al., in preparation). RNA was extracted from FACS sorted cells (Arcturus PicoPure RNA Isolation Kit, Life Technologies) and sequencing libraries generated using the SMARTer Stranded RNA-seq kit (Clontech), according to manufacturer’s instructions. Variants with MAF > 5%, and Hardy-Weinberg equilibrium test p-value > 10<sup>-6</sup> were retained and analyzed using a linear mixed model implemented in mmQTL. Gene expression was normalized using log2 CPM and eQTL analysis was performed on residuals after regression out 15 principal components of the gene expression. For each gene, a Benjamini-Hochberg (BH) FDR correction was applied across all variants tested in the cis regulatory region to obtain the minimum q-value. Then, the minimum q-values across all genes are adjusted again by the BH FDR method to compute the genome-wide FDR. Limited by the small sample size, we chose a less conservative FDR cutoff of 10%.

## Statistical fine-mapping

For each detected eQTL, we conducted a fine-mapping analysis applying the CAVIAR method (Hormozdiari et al., 2014) implemented in mmQTL to find a 95% credible set for causal variants.

Briefly, meta-analysis p-value based on a random-effect model in each round of conditional eQTL detection was firstly converted to z-score, which was then used as input for fine-mapping. CAVIAR will calculate the posterior inclusion probability (PIP) of each variant to causal, and a set of variants prioritized by PIP score were outputted with summed PIP equal to 0.95.

## Detecting colocalization between eQTL and GWAS signals

Joint statistical fine-mapping of eQTL's and GWAS signals (Hormozdiari et al., 2016) was performed by multiplying the estimated posterior inclusion probability (PIP) for a given variant from the eQTL analysis by the PIP for this variant from GWAS of traits compiled in CausalDB (Wang et al., 2020) to obtain a colocalization posterior probability (CLPP). A gene is considered to share a candidate causal variant with a GWAS trait if at least one variant has a  $CLPP > 0.01$  (Hormozdiari et al., 2016).

## Trait classification

CausalDB (Wang et al., 2020) provided the MeSH Category for each GWAS trait. However, brain-related traits fall in multiple MeSH categories and there is no single criterion to identify such traits. We performed manual inspection of traits in CausalDB that could be considered neuropsychiatric, neurodegenerative or behavioral and termed them 'brain related'.

## Validation of rs72986630 effect in chromatin accessibility and gene expression data

To further investigate one prioritized functional variant rs72986630 that reside in REST TF binding site overlapping TSS of *ZNF823*, we queried our unpublished ATAC-seq data set (Bendl et al., in preparation) of neuronal and non-neuronal samples from ACC brain region generated on postmortem human brains from CommonMind cohort (Hoffman et al., 2019). This dataset consists of samples from 370 donors (114 SZ cases, 64 BD cases, 64 controls) with rs72986630 MAF of 6.0%. Since only two donors carry the ALT/ALT (i.e. T/T) genotype, we excluded them for further analysis.

To generate ATAC-seq data set, neuronal and non-neuronal cell populations were isolated from postmortem tissue by fluorescence-activated nuclear sorting using anti-NeuN antibody. ATAC-seq libraries were created using an established protocol (Buenrostro et al., 2015). Raw sequencing reads were mapped to human genome hg38 using STAR aligner (Dobin et al.,

2013). The samples of the same cell type (neuron / non-neuron) and genotype at rs72986630 (CC / CT) were subsampled and merged, creating 4 BAM files with a uniform depth of 1 billion pair-end reads. Subsampling ratios were calculated per each sample individually within those four respective groups (genotype by cell type) to ensure that each of them contributed the same number of reads, regardless of their per-sample read counts. Using these BAM files, bigWig files were created and peaks were called by the MACS (v2.1) with the same parameters as described in (Hauberg et al., 2020), but using an FDR threshold of 0.01. After removing peaks overlapping the blacklisted genomic regions and peaks not being sufficiently accessible (CPM>1 in at least 10% of samples was required), 498,183 peaks remained. The final read count matrix of 664 samples by 498,183 peaks was normalized using the trimmed mean of M-values (TMM) method. The following covariates were selected by Bayesian Information Criterion (BIC) method to be added to the base covariates, i.e. genotype by cell type: mean GC content, fraction of reads with GC content 0-19%, 20-39%, 40-59%, fraction of reads in peaks, fraction of unmapped reads, AT dropout, and mean insert size. Since our dataset contains up to two samples per individual, we ran differential analysis to get differentially accessible peaks between CC and CT carriers using dream (v1.17.9) (Hoffman and Roussos, 2020) that accounts for correlation structure in repeated measures. As an alternative approach, instead of quantifying changes between all open chromatin regions, we performed differential analysis between CC and CT carriers on TF binding sites of REST motif. We used footprinting to narrow down our focus only to 31,534 REST TF binding sites that are bound in at least one set of samples (out of 4 sets, i.e. genotype by cell type) as predicted by TOBIAS (Bentsen et al., 2020).

The analysis of differential gene expression between REF/REF (C/C) and REF/ALT (C/T) genotype at rs72986630 followed the same approach as the analysis of chromatin accessibility. We used a subset of 338 homogenate RNA-seq samples of ACC brain region from CommonMind Consortium (Hoffman et al., 2019) that originate from the same donors as ATAC-seq samples. We performed differential analysis only for sufficiently expressed protein-coding genes (CPM>1 in at least 30% of samples was required), i.e. we start the analysis with a count matrix of 338 samples by 14,893 genes that were normalized by trimmed mean of M-values (TMM) method. The following technical covariates were selected by BIC method: institution, expression profiling efficiency, intronic rate, intragenic rate, fraction of reads with GC content 20-39%, 40-59%, and AT dropout.

# References

- Asimit, J.L., Hatzikotoulas, K., McCarthy, M., Morris, A.P., and Zeggini, E. (2016). Trans-ethnic study design approaches for fine-mapping. *Eur. J. Hum. Genet.* **24**, 1330–1336.
- Benjamini, Y., and Hochberg, Y. (1995). Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society: Series B (Methodological)* **57**, 289–300.
- Benner, C., Spencer, C.C.A., Havulinna, A.S., Salomaa, V., Ripatti, S., and Pirinen, M. (2016). FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* **32**, 1493–1501.
- Bennett, D.A., Buchman, A.S., Boyle, P.A., Barnes, L.L., Wilson, R.S., and Schneider, J.A. (2018). Religious orders study and rush memory and aging project. *J Alzheimers Dis* **64**, S161–S189.
- Bentsen, M., Goymann, P., Schultheis, H., Klee, K., Petrova, A., Wiegandt, R., Fust, A., Preussner, J., Kuenne, C., Braun, T., et al. (2020). ATAC-seq footprinting unravels kinetics of transcription factor binding during zygotic genome activation. *Nat. Commun.* **11**, 4267.
- Bipolar Disorder and Schizophrenia Working Group of the Psychiatric Genomics Consortium (2018). Genomic Dissection of Bipolar Disorder and Schizophrenia, Including 28 Subphenotypes. *Cell* **173**, 1705–1715.
- Buenrostro, J.D., Wu, B., Chang, H.Y., and Greenleaf, W.J. (2015). ATAC-seq: A Method for Assaying Chromatin Accessibility Genome-Wide. *Curr. Protoc. Mol. Biol.* **109**, 21.29.1–21.29.9.
- Cao, J., O'Day, D.R., Pliner, H.A., Kingsley, P.D., Deng, M., Daza, R.M., Zager, M.A., Aldinger, K.A., Blecher-Gonen, R., Zhang, F., et al. (2020). A human cell atlas of fetal gene expression. *Science* **370**.
- Cerimele, J.M., and Katon, W.J. (2013). Associations between health risk behaviors and symptoms of schizophrenia and bipolar disorder: a systematic review. *Gen. Hosp. Psychiatry* **35**, 16–22.
- Cross-Disorder Group of the Psychiatric Genomics Consortium. (2019). Genomic Relationships, Novel Loci, and Pleiotropic Mechanisms across Eight Psychiatric Disorders. *Cell* **179**, 1469–1482.e11.
- Darmanis, S., Sloan, S.A., Zhang, Y., Enge, M., Caneda, C., Shuer, L.M., Hayden Gephart, M.G., Barres, B.A., and Quake, S.R. (2015). A survey of human brain transcriptome diversity at the single cell level. *Proc Natl Acad Sci USA* **112**, 7285–7290.
- Delaneau, O., Ongen, H., Brown, A.A., Fort, A., Panousis, N.I., and Dermitzakis, E.T. (2017). A complete tool set for molecular QTL discovery and analysis. *Nat. Commun.* **8**, 15452.
- Devlin, B., and Roeder, K. (1999). Genomic control for association studies. *Biometrics* **55**, 997–1004.

- Dobbyn, A., Huckins, L.M., Boocock, J., Sloofman, L.G., Glicksberg, B.S., Giambartolomei, C., Hoffman, G.E., Perumal, T.M., Girdhar, K., Jiang, Y., et al. (2018). Landscape of Conditional eQTL in Dorsolateral Prefrontal Cortex and Co-localization with Schizophrenia GWAS. *Am. J. Hum. Genet.* *102*, 1169–1184.
- Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* *29*, 15–21.
- Fairfax, B.P., Humburg, P., Makino, S., Naranbhai, V., Wong, D., Lau, E., Jostins, L., Plant, K., Andrews, R., McGee, C., et al. (2014). Innate immune activity conditions the effect of regulatory variants upon monocyte gene expression. *Science* *343*, 1246949.
- 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. (2015). Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.* *47*, 1228–1235.
- Finucane, H.K., Reshef, Y.A., Anttila, V., Slowikowski, K., Gusev, A., Byrnes, A., Gazal, S., Loh, P.-R., Lareau, C., Shores, N., et al. (2018). Heritability enrichment of specifically expressed genes identifies disease-relevant tissues and cell types. *Nat. Genet.* *50*, 621–629.
- Fromer, M., Roussos, P., Sieberts, S.K., Johnson, J.S., Kavanagh, D.H., Perumal, T.M., Ruderfer, D.M., Oh, E.C., Topol, A., Shah, H.R., et al. (2016). Gene expression elucidates functional impact of polygenic risk for schizophrenia. *Nat. Neurosci.* *19*, 1442–1453.
- Gallagher, M.D., and Chen-Plotkin, A.S. (2018). The Post-GWAS Era: From Association to Function. *Am. J. Hum. Genet.* *102*, 717–730.
- Glassberg, E.C., Gao, Z., Harpak, A., Lan, X., and Pritchard, J.K. (2019). Evidence for weak selective constraint on human gene expression. *Genetics* *211*, 757–772.
- GTEx Consortium (2020). The GTEx Consortium atlas of genetic regulatory effects across human tissues. *Science* *369*, 1318–1330.
- GTEx Consortium, Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group, Statistical Methods groups—Analysis Working Group, Enhancing GTEx (eGTEx) groups, NIH Common Fund, NIH/NCI, NIH/NHGRI, NIH/NIMH, NIH/NIDA, Biospecimen Collection Source Site—NDRI, et al. (2017). Genetic effects on gene expression across human tissues. *Nature* *550*, 204–213.
- Habib, N., Avraham-David, I., Basu, A., Burks, T., Shekhar, K., Hofree, M., Choudhury, S.R., Aguet, F., Gelfand, E., Ardlie, K., et al. (2017). Massively parallel single-nucleus RNA-seq with DroNc-seq. *Nat. Methods* *14*, 955–958.
- Han, B., and Eskin, E. (2011). Random-effects model aimed at discovering associations in meta-analysis of genome-wide association studies. *Am. J. Hum. Genet.* *88*, 586–598.
- Han, B., Duong, D., Sul, J.H., de Bakker, P.I.W., Eskin, E., and Raychaudhuri, S. (2016). A

general framework for meta-analyzing dependent studies with overlapping subjects in association mapping. *Hum. Mol. Genet.* 25, 1857–1866.

Hauberg, M.E., Creus-Muncunill, J., Bendl, J., Kozlenkov, A., Zeng, B., Corwin, C., Chowdhury, S., Kranz, H., Hurd, Y.L., Wegner, M., et al. (2020). Common schizophrenia risk variants are enriched in open chromatin regions of human glutamatergic neurons. *Nat. Commun.* 11, 5581.

Hoffman, G.E., and Roussos, P. (2020). dream: Powerful differential expression analysis for repeated measures designs. *Bioinformatics*.

Hoffman, G.E., Bendl, J., Voloudakis, G., Montgomery, K.S., Sloofman, L., Wang, Y.-C., Shah, H.R., Hauberg, M.E., Johnson, J.S., Girdhar, K., et al. (2019). CommonMind Consortium provides transcriptomic and epigenomic data for Schizophrenia and Bipolar Disorder. *Sci. Data* 6, 180.

Hormozdiari, F., Kostem, E., Kang, E.Y., Pasaniuc, B., and Eskin, E. (2014). Identifying causal variants at loci with multiple signals of association. *Genetics* 198, 497–508.

Hormozdiari, F., van de Bunt, M., Segrè, A.V., Li, X., Joo, J.W.J., Bilow, M., Sul, J.H., Sankararaman, S., Pasaniuc, B., and Eskin, E. (2016). Colocalization of GWAS and eQTL Signals Detects Target Genes. *Am. J. Hum. Genet.* 99, 1245–1260.

Hormozdiari, F., Gazal, S., van de Geijn, B., Finucane, H.K., Ju, C.J.-T., Loh, P.-R., Schoech, A., Reshef, Y., Liu, X., O'Connor, L., et al. (2018). Leveraging molecular quantitative trait loci to understand the genetic architecture of diseases and complex traits. *Nat. Genet.* 50, 1041–1047.

Hwang, J.-Y., and Zukin, R.S. (2018). REST, a master transcriptional regulator in neurodegenerative disease. *Curr. Opin. Neurobiol.* 48, 193–200.

Jaffe, A.E., Straub, R.E., Shin, J.H., Tao, R., Gao, Y., Collado-Torres, L., Kam-Thong, T., Xi, H.S., Quan, J., Chen, Q., et al. (2018). Developmental and genetic regulation of the human cortex transcriptome illuminate schizophrenia pathogenesis. *Nat. Neurosci.* 21, 1117–1125.

Jaffe, A.E., Hoeppner, D.J., Saito, T., Blanpain, L., Ukaigwe, J., Burke, E.E., Collado-Torres, L., Tao, R., Tajinda, K., Maynard, K.R., et al. (2020). Profiling gene expression in the human dentate gyrus granule cell layer reveals insights into schizophrenia and its genetic risk. *Nat. Neurosci.* 23, 510–519.

Jansen, I.E., Savage, J.E., Watanabe, K., Bryois, J., Williams, D.M., Steinberg, S., Sealock, J., Karlsson, I.K., Hägg, S., Athanasiu, L., et al. (2019). Genome-wide meta-analysis identifies new loci and functional pathways influencing Alzheimer's disease risk. *Nat. Genet.* 51, 404–413.

Jansen, R., Hottenga, J.-J., Nivard, M.G., Abdellaoui, A., Laport, B., de Geus, E.J., Wright, F.A., Penninx, B.W.J.H., and Boomsma, D.I. (2017). Conditional eQTL analysis reveals allelic heterogeneity of gene expression. *Hum. Mol. Genet.* 26, 1444–1451.

Kim-Hellmuth, S., Aguet, F., Oliva, M., Muñoz-Aguirre, M., Kasela, S., Wucher, V., Castel, S.E., Hamel, A.R., Viñuela, A., Roberts, A.L., et al. (2020). Cell type-specific genetic regulation of gene expression across human tissues. *Science* 369.

Lake, B.B., Chen, S., Sos, B.C., Fan, J., Kaeser, G.E., Yung, Y.C., Duong, T.E., Gao, D., Chun, J., Kharchenko, P.V., et al. (2018). Integrative single-cell analysis of transcriptional and

epigenetic states in the human adult brain. *Nat. Biotechnol.* 36, 70–80.

Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al. (2016). Analysis of protein-coding genetic variation in 60,706 humans. *Nature* 536, 285–291.

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.

Mandric, I., Schwarz, T., Majumdar, A., Hou, K., Briscoe, L., Perez, R., Subramaniam, M., Hafemeister, C., Satija, R., Ye, C.J., et al. (2020). Optimized design of single-cell RNA sequencing experiments for cell-type-specific eQTL analysis. *Nat. Commun.* 11, 5504.

Meyer, H.V., and Birney, E. (2018). PhenotypeSimulator: A comprehensive framework for simulating multi-trait, multi-locus genotype to phenotype relationships. *Bioinformatics* 34, 2951–2956.

Morris, A.P. (2011). Transethnic meta-analysis of genomewide association studies. *Genet. Epidemiol.* 35, 809–822.

Najt, P., Perez, J., Sanches, M., Peluso, M.A.M., Glahn, D., and Soares, J.C. (2007). Impulsivity and bipolar disorder. *Eur. Neuropsychopharmacol.* 17, 313–320.

Nalls, M.A., Blauwendraat, C., Vallerga, C.L., Heilbron, K., Bandres-Ciga, S., Chang, D., Tan, M., Kia, D.A., Noyce, A.J., Xue, A., et al. (2019). Identification of novel risk loci, causal insights, and heritable risk for Parkinson's disease: a meta-analysis of genome-wide association studies. *Lancet Neurol.* 18, 1091–1102.

Ng, B., White, C.C., Klein, H.-U., Sieberts, S.K., McCabe, C., Patrick, E., Xu, J., Yu, L., Gaiteri, C., Bennett, D.A., et al. (2017). An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome. *Nat. Neurosci.* 20, 1418–1426.

Nott, A., Holtman, I.R., Coufal, N.G., Schlachetzki, J.C.M., Yu, M., Hu, R., Han, C.Z., Pena, M., Xiao, J., Wu, Y., et al. (2019). Brain cell type-specific enhancer-promoter interactome maps and disease-risk association. *Science* 366, 1134–1139.

Ouzir, M. (2013). Impulsivity in schizophrenia: A comprehensive update. *Aggress. Violent Behav.* 18, 247–254.

de Paiva Lopes, K., Snijders, G.J.L., Humphrey, J., Allan, A., Sneeboer, M., Navarro, E., Schilder, B.M., Vialle, R.A., Parks, M., Missall, R., et al. (2020). Atlas of genetic effects in human microglia transcriptome across brain regions, aging and disease pathologies. *BioRxiv*.

Pardiñas, A.F., Holmans, P., Pocklington, A.J., Escott-Price, V., Ripke, S., Carrera, N., Legge, S.E., Bishop, S., Cameron, D., Hamshere, M.L., et al. (2018). Common schizophrenia alleles are enriched in mutation-intolerant genes and in regions under strong background selection. *Nat. Genet.* 50, 381–389.

Raj, T., Rothamel, K., Mostafavi, S., Ye, C., Lee, M.N., Replogle, J.M., Feng, T., Lee, M., Asinowski, N., Frohlich, I., et al. (2014). Polarization of the effects of autoimmune and neurodegenerative risk alleles in leukocytes. *Science* 344, 519–523.

- Roussos, P., Mitchell, A.C., Voloudakis, G., Fullard, J.F., Pothula, V.M., Tsang, J., Stahl, E.A., Georgakopoulos, A., Ruderfer, D.M., Charney, A., et al. (2014). A role for noncoding variation in schizophrenia. *Cell Rep.* 9, 1417–1429.
- 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.
- Schizophrenia Working Group of the Psychiatric Genomics Consortium (2014). Biological insights from 108 schizophrenia-associated genetic loci. *Nature* 511, 421–427.
- Schoenherr, C.J., and Anderson, D.J. (1995). The neuron-restrictive silencer factor (NRSF): a coordinate repressor of multiple neuron-specific genes. *Science* 267, 1360–1363.
- Schrode, N., Ho, S.-M., Yamamuro, K., Dobbyn, A., Huckins, L., Matos, M.R., Cheng, E., Deans, P.J.M., Flaherty, E., Barretto, N., et al. (2019). Synergistic effects of common schizophrenia risk variants. *Nat. Genet.* 51, 1475–1485.
- Sekar, A., Bialas, A.R., de Rivera, H., Davis, A., Hammond, T.R., Kamitaki, N., Tooley, K., Presumey, J., Baum, M., Van Doren, V., et al. (2016). Schizophrenia risk from complex variation of complement component 4. *Nature* 530, 177–183.
- Storey, J.D., and Tibshirani, R. (2003). Statistical significance for genomewide studies. *Proc Natl Acad Sci USA* 100, 9440–9445.
- Sul, J.H., Martin, L.S., and Eskin, E. (2018). Population structure in genetic studies: Confounding factors and mixed models. *PLoS Genet.* 14, e1007309.
- Su, Z., Marchini, J., and Donnelly, P. (2011). HAPGEN2: simulation of multiple disease SNPs. *Bioinformatics* 27, 2304–2305.
- Visscher, P.M., Wray, N.R., Zhang, Q., Sklar, P., McCarthy, M.I., Brown, M.A., and Yang, J. (2017). 10 years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet.* 101, 5–22.
- 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.
- Wang, D., Liu, S., Warrell, J., Won, H., Shi, X., Navarro, F.C.P., Clarke, D., Gu, M., Emani, P., Yang, Y.T., et al. (2018a). Comprehensive functional genomic resource and integrative model for the human brain. *Science* 362.
- Wang, J., Huang, D., Zhou, Y., Yao, H., Liu, H., Zhai, S., Wu, C., Zheng, Z., Zhao, K., Wang, Z., et al. (2020). CAUSALdb: a database for disease/trait causal variants identified using summary statistics of genome-wide association studies. *Nucleic Acids Res.* 48, D807–D816.
- Wang, M., Beckmann, N.D., Roussos, P., Wang, E., Zhou, X., Wang, Q., Ming, C., Neff, R., Ma, W., Fullard, J.F., et al. (2018b). The Mount Sinai cohort of large-scale genomic, transcriptomic and proteomic data in Alzheimer’s disease. *Sci. Data* 5, 180185.
- van der Wijst, M.G.P., Brugge, H., de Vries, D.H., Deelen, P., Swertz, M.A., LifeLines Cohort Study, BIOS Consortium, and Franke, L. (2018). Single-cell RNA sequencing identifies celltype-specific cis-eQTLs and co-expression QTLs. *Nat. Genet.* 50, 493–497.

van der Wijst, M., de Vries, D.H., Groot, H.E., Trynka, G., Hon, C.C., Bonder, M.J., Stegle, O., Nawijn, M.C., Idaghdour, Y., van der Harst, P., et al. (2020). The single-cell eQTLGen consortium. *Elife* 9.

Wojcik, G.L., Graff, M., Nishimura, K.K., Tao, R., Haessler, J., Gignoux, C.R., Highland, H.M., Patel, Y.M., Sorokin, E.P., Avery, C.L., et al. (2019). Genetic analyses of diverse populations improves discovery for complex traits. *Nature* 570, 514–518.

Wray, N.R., Ripke, S., Mattheisen, M., Trzaskowski, M., Byrne, E.M., Abdellaoui, A., Adams, M.J., Agerbo, E., Air, T.M., Andlauer, T.M.F., et al. (2018). Genome-wide association analyses identify 44 risk variants and refine the genetic architecture of major depression. *Nat. Genet.* 50, 668–681.

Yang, J., Zaitlen, N.A., Goddard, M.E., Visscher, P.M., and Price, A.L. (2014). Advantages and pitfalls in the application of mixed-model association methods. *Nat. Genet.* 46, 100–106.

Young, A., Kumasaka, N., Calvert, F., Hammond, T.R., Knights, A.J., Panousis, N., Schwartzentruber, J., Liu, J., Kundu, K., Segel, M., et al. (2019). A map of transcriptional heterogeneity and regulatory variation in human microglia. *BioRxiv*.

Zaitlen, N., Paşaniuc, B., Gur, T., Ziv, E., and Halperin, E. (2010). Leveraging genetic variability across populations for the identification of causal variants. *Am. J. Hum. Genet.* 86, 23–33.

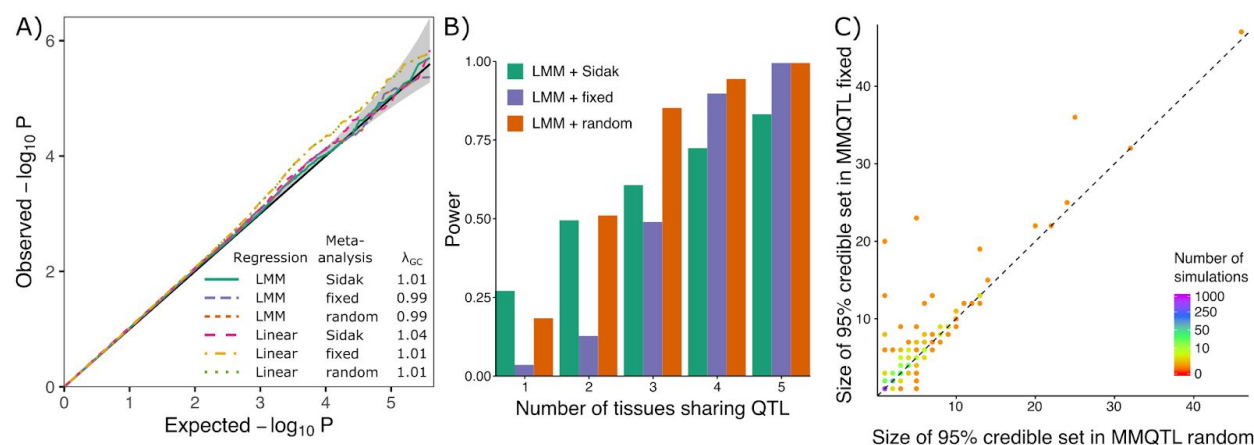
Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., La Manno, G., Juréus, A., Marques, S., Munguba, H., He, L., Betsholtz, C., et al. (2015). Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* 347, 1138–1142.

Zhang, X., Zhang, C., Prokopenko, D., Liang, Y., Han, W., Tanzi, R.E., and Sisodia, S.S. (2020). Negative evidence for a role of APOE ε4 variant in Alzheimer's disease. *Hum. Mol. Genet.* 29, 955–966.

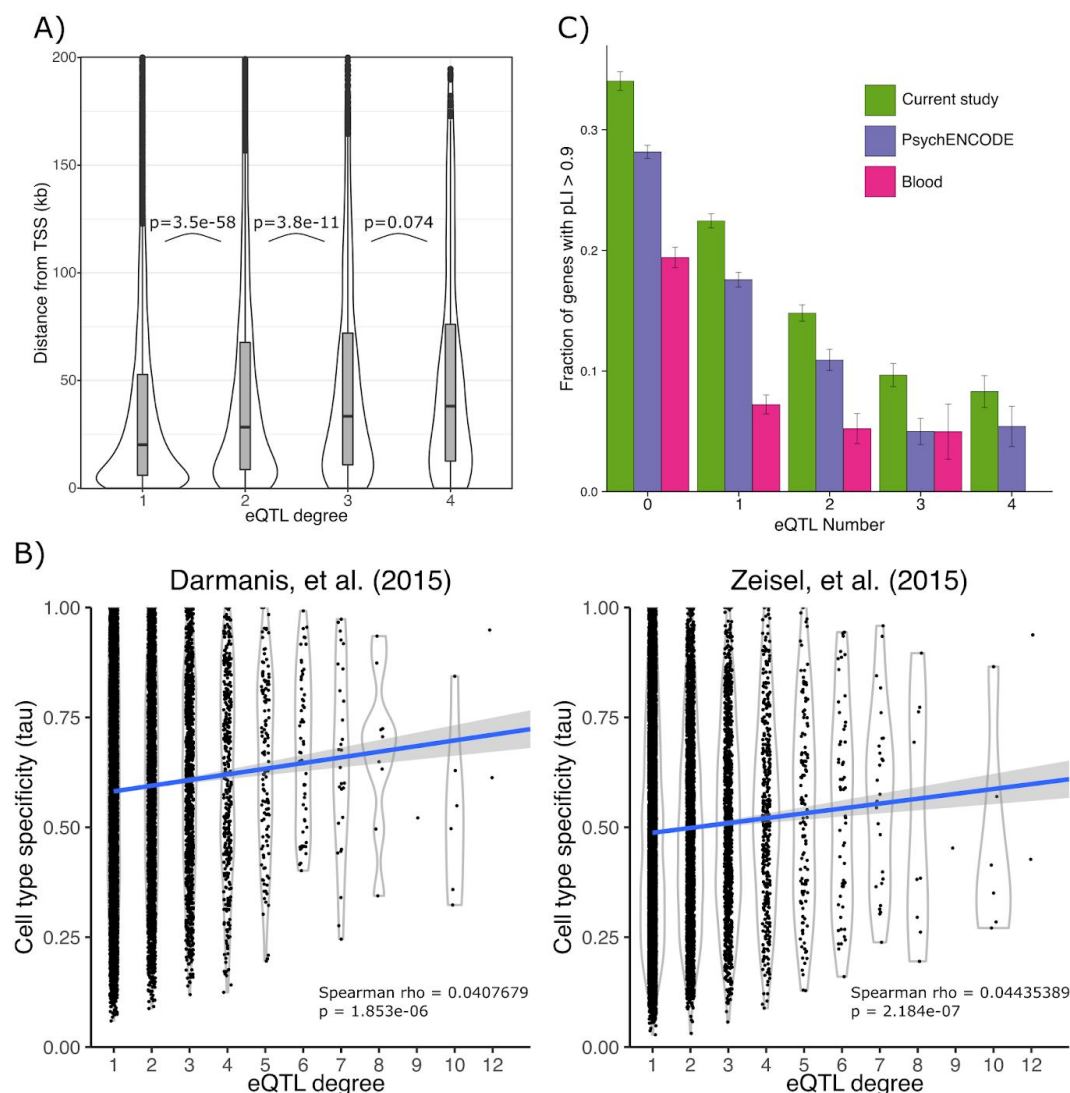
Zhernakova, D.V., Deelen, P., Vermaat, M., van Iterson, M., van Galen, M., Arindrarto, W., van 't Hof, P., Mei, H., van Dijk, F., Westra, H.-J., et al. (2017). Identification of context-dependent expression quantitative trait loci in whole blood. *Nat. Genet.* 49, 139–145.

Zhou, X., and Stephens, M. (2012). Genome-wide efficient mixed-model analysis for association studies. *Nat. Genet.* 44, 821–824.

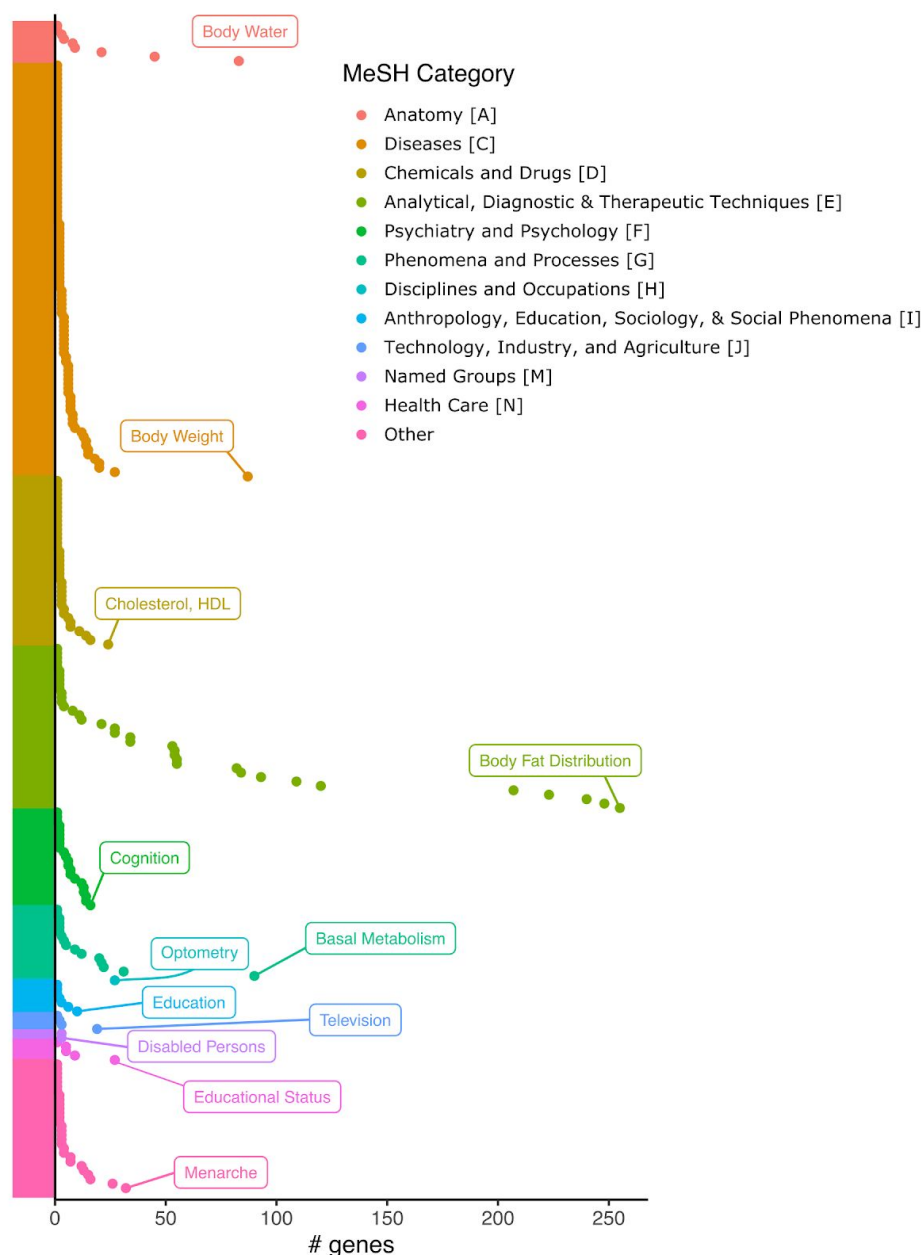
## Supplementary Figures



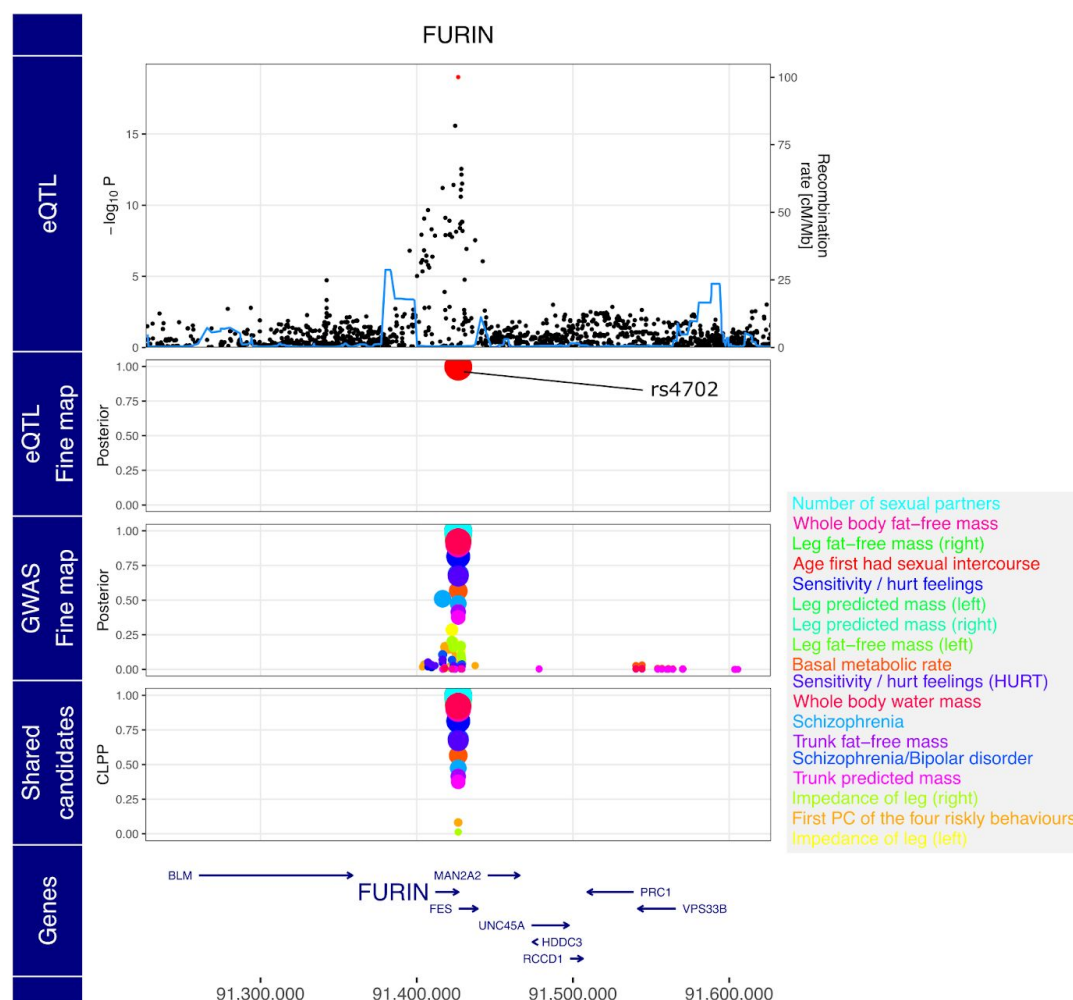
**Supplementary Figure 1: Biologically motivated simulations demonstrate performance of mmQTL workflow: low correlation scenario. A)** QQ plot of results from null simulation shows that the linear mixed model (LMM) with fixed or random effect meta-analysis accurately controls the false positive rate for, while linear regression with 5 genotype principal components did not. The Sidak method was very conservative in both cases.  $\lambda_{GC}$  indicates the genomic control inflation factor. **B)** Power from LMM followed by 3 types of meta-analysis versus the number of tissues sharing an eQTL. **C)** Size of the 95% credible sets from fixed- (y-axis) and random- (x-axis) effects meta-analysis from simulations in Figure 2C.



**Supplementary Figure 2: Properties of conditional eQTLs.** **A)** The distribution of the distance to the transcription start site is shown for the lead variant for eQTL analysis of increasing degree. P-values indicate significance of one-sided Wilcoxon test between adjacent groups. Box plot indicates median, interquartile range (IQR) and 1.5\*IQR. **B)** Cell type specificity metric  $\tau$  plotted against the number of independent eQTLs discovered for each gene. **C)** The fraction of genes with high evolutionary constraint ( $pLI > 0.9$ ) is shown increasing eQTL degree for the current study, PsychENCODE (Wang et al., 2018a), and whole blood (Glassberg et al., 2019). Error bars indicate standard error based on asymptotic estimate of binomial proportion.



**Supplementary Figure 3. Number of genes colocalizing for each MeSH category with CLPP > 0.01.** The phenotype with the highest number of colocalized genes for each MeSH category is indicated.



**Supplementary Figure 4. Expression of *FURIN* and risk for multiple complex traits share rs4702 as a candidate causal variant.** Starting from the top, the plot shows  $-\log_{10}$  p-values from eQTL analysis, poster probabilities from statistical fine-mapping of eQTL results, poster probabilities from statistical fine-mapping of GWAS results, and colocalization posterior probabilities (CLPP) for combining eQTL and GWAS fine-mapping. Traits are shown in the box on the right in decreasing order to CLPP value.

## Supplementary Tables

| Abbreviation | Trait/disease                            | Reference (doi)               |
|--------------|------------------------------------------|-------------------------------|
| ADHD         | Attention deficit hyperactivity disorder | 10.1016/j.jaac.2010.06.008    |
| ALS          | Amyotrophic lateral sclerosis            | 10.1016/j.neuron.2018.02.027  |
| PD           | Parkinson's disease                      | 10.1016/S1474-4422(19)30320-5 |
| RA           | Rheumatoid arthritis                     | 10.1038/nature12873           |
| BMI          | Body mass index                          | 10.1038/nature14177           |
| EduYear      | Educational years                        | 10.1038/nature17671           |
| HEIGHT       | Height                                   | 10.1038/ng.3097               |
| CD           | Crohn's disease                          | 10.1038/ng.3359               |
| IBD          | Inflammatory bowel disease               | 10.1038/ng.3359               |
| UC           | Ulcerative colitis                       | 10.1038/ng.3359               |
| CAD          | Cardiovascular disease                   | 10.1038/ng.3396               |
| DS           | Depression symptom                       | 10.1038/ng.3552               |
| Neu          | Neuroticism                              | 10.1038/ng.3552               |
| SLE          | Systemic lupus erythematosus             | 10.1038/ng.3603               |
| T2D          | Type 2 Diabetes                          | 10.1038/s41588-018-0241-6     |
| DRINKING     | Alcohol Assumption                       | 10.1038/s41588-018-0309-3     |
| AD           | Alzheimer's disease                      | 10.1038/s41588-018-0311-9     |
| ASD          | Autism spectrum disorder                 | 10.1038/s41588-019-0344-8     |
| BD           | Bipolar disorder                         | 10.1038/s41588-019-0397-8     |
| MDD          | Major depression disorder                | 10.1038/s41593-018-0326-7     |
| SZ           | Schizophrenia                            | 10.1101/2020.09.12.20192922   |
| MS           | Multiple sclerosis                       | 10.1126/science.aav7188       |

**Supplementary Table 1. Trait/Disease, abbreviation and reference for GWAS included in LD-score regression analysis**