

Supplementary Information

Data acquisition, quality control, and normalization for known technical factors

Data from the ROSMAP study were used in this work, and are available at: RADC Research Resource Sharing Hub at www.radc.rush.edu, and <https://www.synapse.org/#!/Synapse:syn3219045>. The ROSMAP study is a longitudinal study where all participants are healthy at enrolment. The sampled subjects thus represent a relatively random set of older individuals. By the time of death, 58% and 38% of participants are diagnosed with pathological and clinical AD, respectively (**Table S1**). These percentages are consistent with AD population prevalence. We note that a single person performed all of the dissections of the frozen tissues in isolating the gray matter for gene expression, DNA methylation, and histone modification data generation to minimize technical variability in sample preparation.

*Genotype data*¹. Genotyping of the ROS and MAP subjects was performed on the Affymetrix Genome-Wide HumanSNP Array6.0 ($n=1709$) and the Illumina OmniQuad Express platform ($n=384$). DNA was extracted from whole blood, lymphocytes, or frozen brain tissue, as previously described¹. To minimize population admixture, only self-declared non-Hispanic Caucasians were genotyped. At the sample level, samples with genotyping success rate $< 95\%$, discordant genetically inferred and reported gender, or excess inter/intra-heterozygosity were excluded. At the probe level, genotyping data from both platforms were processed with same quality-control (QC) metrics: Hardy-Weinberg equilibrium $p > 0.001$, genotype call rate < 0.95 , misshap test $< 1 \times 10^{-9}$. QC was performed using version 1.08p of the PLINK software². EIGENSTRAT³ was used with the default setting to remove population outliers and to generate a genotype covariance matrix. The resultant datasets include 729,463 SNPs for 1,709 individuals (Affy) and 624,668 SNPs for 384 individuals (Omni). Dosages for all SNPs (>35 million) on the 1000 Genomes reference were imputed using version 3.3.2 version of the BEAGLE software⁴ (1000 Genomes Project Consortium interim phase I haplotypes, 2011 Phase 1b data freeze(verify) data freeze). Imputed SNPs were filtered based on minor allele frequency (MAF) > 0.01 and imputation INFO score > 0.3 , resulting in 7,321,515 SNPs available for analysis.

*Gene expression data*⁵. Gene expression data were generated using RNA-sequencing from Dorsolateral Prefrontal Cortex (DLPFC) of 540 individuals, at an average sequence depth of 90M reads. Detailed description of data generation and processing was previously described (Mostafavi, Gaiteri et al., under review) and summarized here.

Samples were submitted to the Broad Institute's Genomics Platform for transcriptome analysis following the dUTP protocol with Poly(A) selection developed by Levin and colleagues⁶. All samples were chosen to pass two initial quality filters: RNA integrity (RIN) score >5 and quantity threshold of 5 ug (and were selected from a larger set of 724 samples). Sequencing was performed on the Illumina HiSeq with 101bp paired-end reads and achieved coverage of 150M reads of the first 12 samples. These 12 samples will serve as a deep coverage reference and included 2 males and 2 females of non-impaired, mild cognitive impaired, and Alzheimer's cases (Figure S3A). The remaining samples were sequenced with a target coverage of 50M reads; the mean coverage for the samples passing QC is 95 million reads (median 90 million reads) (Figure S3A). The libraries were constructed and pooled according to the RIN scores such that similar RIN scores would be pooled together (Figure S3B). Varying RIN scores results in a larger spread of insert sizes during library construction and leads to uneven coverage distribution throughout the pool.

RNA-seq data were processed by our parallelized pipeline. This pipeline included trimming the beginning and ending bases from each read, identifying and trimming adapter sequences from reads, detecting and removing rRNA reads, aligning reads to reference genome (using Bowtie⁷) and quantification of transcript expression levels (using RSEM⁸). Specifically, RNA-Seq reads in FASTQ format were inspected using FASTQC program. Barcode and adapter contamination, low quality regions (8bp at beginning and 7bp at ending of each fastq reads) were trimmed using FASTX-toolkit. To remove rRNA contamination, we aligned trimmed reads to rRNA reference (rRNA genes were downloaded from UCSC genome browser selecting the RepeatMask table) by BWA then extracted only paired unmapped reads for transcriptome alignment. rRNA depleted reads were then mapped to transcriptome reference (gencode v14) using Trinity package with RSEM as output option. Gene expression FPKM values were estimated by "rsem-calculate-expression" from RSEM.

Samples from 494 individuals were used in the eQTL analysis, which include those that had QC'd genotype and pass the expression outlier test (a D-statistic below 0.9⁹). To quantify the contribution of experimental and other confounding factors to the overall expression profiles, we performed a PCA analysis on log transformed FPKM values in all samples, and computed the correlation between the top 10 PCs and experimental factors. As shown in **Figure S4**, we observed significant correlations between many of these technical and confounding factors and top expression PCs. Thus, with the log transformed FPKM data, we used the COMBAT algorithm¹¹ to account for the effect of batch and linear regression to remove the effects of RIN, post-mortem interval (PMI), sequencing depth, study index (ROS sample or MAP sample), genotyping PCs, age at death, and sex. (i.e. all factors shown in **Figure S3**). Finally, only highly expressed

genes were kept (mean expression $>2 \log_2$ -FPKM), resulting in 13,484 expressed genes for eQTL analysis. This FPKM-based threshold was determined based on visual inspection of histogram of mean expression values to approximately define two expression distributions: a) no expression or very low expression and b) moderate to high expression.

*DNA methylation data*¹². DNA methylation data were generated using the 450K Illumina array from DLPFC of 740 individuals. Detailed description of data acquisition and QC are previously published¹². Briefly, methylation probes that coincided with common polymorphic sites were removed. Initial normalization of CpG probes to account for differences between type I and type II probes, was performed using the BMIQ algorithm from the Watermelon package¹³ and beta-values were extracted for further analysis. The SNM approach¹⁴ was then used to regress out the effects of batch, PMI, sex, age at death, and a previously published estimate of proportion of neurons present in each sample¹². In this study, samples from 468 individuals were analyzed for which gene expression data was also available. As described below, this decision was made to enable using gene expression data to estimate the proportions of the five major brain cell types. This correction for cell type proportions was done in addition to the regression approach for removing the effect of generic neuronal proportions based on DNAm marks¹².

*Histone modification data*⁵. Histone modification data were generated using H3K9Ac ChIP-sequencing from DLPFC of 714 individuals. Single-end reads were aligned by the BWA algorithm¹¹, and peaks were detected in each sample separately using the MACS2 algorithm¹² (using the broad peak option and a q-value cutoff of 0.001). A series of QC steps were employed to identify and remove low quality reads¹³, and samples that did not reach (i) $\geq 15 \times 10^6$ unique reads, (ii) non-redundant fraction ≥ 0.3 , (iii) cross correlation ≥ 0.03 , (iv) fraction of reads in peaks ≥ 0.05 and (v) ≥ 6000 peaks were removed. Cross correlation was defined as the maximum Pearson's correlation between the read coverage on the negative and positive strand after binning reads into 10bp bins¹⁵. Cross correlation was calculated after shifting the reads on the negative strand by s base pairs for $s = 0, 10, 20, \dots, 1000$, and the maximum cross correlation was reported. In total, 669 samples passed quality control. Distribution of these QC metrics across the samples are reported in Figure S4.

H3K9Ac domains were defined by calculating all genomic regions that were detected as a peak in at least 100 of the 669 samples (15%). Regions within 100bp from each other were merged and very small regions of less than 100bp were removed. Reads were then extended towards the 3' end to the fragment size of the respective sample. The fragment size was estimated by the shift s_{max} that maximized the cross correlation

(mean $s_{max} = 271\text{bp}$). Finally, the number of extended reads in each H3K9Ac region was determined for each sample. Only uniquely mapped distinct reads were considered. Quantified histone acetylation data were quantile normalized to account for variability in sequencing depth across individuals. Samples from 433 individuals for which gene expression data were available were used in our analysis.

Additional removal of known and hidden confounding factors

In addition to the data-specific QC and normalization described above, the effect of ancestry, cell type composition, and “hidden factors” were regressed out from the gene expression, DNA methylation, and histone acetylation data. Variables representing ancestry were defined using the top three principal components of the genotype data. Cell type composition was estimated using gene expression levels of markers of major brain cell types: neurons (ENO2), oligodendrocytes (OLIG2, MBP, CNP), astrocytes (GFAP), microglia (CD68), and endothelial cells (CD34). Hidden confounding factors included top N PCs from the gene expression, DNA methylation, and histone modification data (separately). PCA-based hidden factors typically capture variation in cell type proportions across individuals and other unmeasured confounding factors^{19,20}. Following previous studies²¹, for each molecular phenotype data, we varied N from 1 to 30 at a \log_{10} scale, and defined “optimal” N as the value at which the number of significant hits in chromosome 18 saturated, **Figure S1**. We chose to assess performance on only chromosome 18 as opposed to all chromosomes to avoid overfitting. The optimal N was found to be approximately 10 for all three data types.

xQTL Association Analysis

Spearman’s rank correlation was used to estimate the association strength between the alleles of each SNP and the three molecular phenotypes measured. For eQTL analysis, we used SNPs that are up to 1MB upstream or downstream from the TSS of each gene. For mQTL analysis, we used SNPs that are within 5KB of each methylation site. For haQTL analysis, we used SNPs that are within 1MB of each acetylation peak. The window sizes are informed by prior studies²²⁻²⁴. For each xQTL type, we declare an association as significant if its p-value is less than 0.05 after Bonferroni correction. Bonferroni threshold was determined separately for eQTL ($p < 8 \times 10^{-8}$), mQTL ($p < 5 \times 10^{-9}$), and haQTL ($p < 4 \times 10^{-10}$) analysis based on the number of tested associations.

Replication estimation with π_1 statistics

Replication analysis for eQTLs and mQTLs was performed using previous brain-based studies²⁵⁻²⁷, blood-based studies^{28,29}, and the GTEx study³⁰ to evaluate cross sample

replication and cross tissue replication. Replication rates were estimated using the π_1 statistics³¹, which provides an estimate of the proportion of xQTLs that are significant based on their p-value distribution. Only associations comprising the top SNP for each eQTL gene and each mQTL probe were included in the π_1 estimation, to avoid including many SNPs in LD with each other in this analysis. For π_1 estimation, we used p-values from this study restricted to the eQTLs and mQTLs found in previous studies. That is, we used an existing reference eQTL or mQTL list, and assessed the replication of those reported xQTLs in our dataset. When possible, we also estimated π_1 in the other direction. Specifically, we assessed the replication rate of our eQTLs in a large whole-blood dataset²⁸. To determine if the replication rate is higher than chance level, we generated empirical null distributions by computing π_1 for 10^4 random p-value subsets of size m , where m is the number of eQTLs or mQTLs. Only p-values of associations that do not overlap with the eQTLs and mQTLs are used for null estimation.

Genomic annotations

To examine if the xQTL SNPs are enriched in specific gene regions, we used genomic annotations from the ChomHMM resource³², which comprise 15 categories: 1. Active TSS (TssA), 2. Flanking Active TSS (TssAFlnk), 3. Transcription at gene 5' and 3' (TxFlnk), 4. Strong transcription (Tx), 5. Weak transcription (TxWk), 6. Genic enhancers (EnhG), 7. Enhancers (Enh), 8. ZNF genes & repeats (ZNF/Rpts), 9. Heterochromatin (Het), 10. Bivalent/Poised TSS (TssBiv), 11. Flanking Bivalent TSS/Enhancers (BivFlnk), 12. Bivalent Enhancers (EnhBiv), 13. Repressed PolyComb (ReprPC), 14. Weak Repressed PolyComb (ReprPCWk), and 15. Quiescent/Low (Quies). We also used the knownGene table (GRCh37/hg19 assembly) provided on the UCSC genome browser website³³ to examine if the xQTL SNPs are enriched in exons and introns. For each xQTL type, we computed the odds ratio of the xQTL SNPs being in each of the gene regions versus all other tested SNPs, i.e. those within predefined windows from the molecular features (1Mb, 5Kb, and 1Mb for eQTL, mQTL, and haQTL analyses). We further estimated the probability of observing an xQTL SNP at a certain distance away from the TSS of its respective gene(s) by computing the number of xQTL SNPs at different distances away from TSS and dividing that by the number of tested SNPs to account for sampling biases.

Estimation of xQTL SNP sharing across molecular phenotypes

The π_1 statistics was employed to estimate the sharing of xQTL SNPs across molecular phenotypes. Using sharing between mQTLs and eQTLs as an example with methylation and gene expression being the “discovery” and “test” phenotypes, respectively, we computed π_1 with p-values of the tested SNP-expression associations that consist of

mQTL SNPs. This π_1 analysis provides an estimate of the proportion of SNP-expression associations that are significant when we restrict to associations comprising mQTL SNPs, i.e. if majority of mQTL SNPs also drive gene expression, then the corresponding π_1 would be high. Importantly, since an mQTL SNP might be tested for association with expression levels of multiple genes, a decision had to be made regarding which associations to include in the π_1 estimation. A lenient strategy would be to retain only the strongest association for each mQTL SNP, and a more stringent strategy would be to include all tested associations. With the lenient strategy, we estimated a cross-phenotype sharing (π_1) of ~ 0.83 - 0.97 for different pairs of phenotypes. With the more stringent strategy, we estimated a cross-phenotype sharing (π_1) of ~ 0.1 - 0.35 . For the more stringent strategy, we note that as we decreased the allowable genomic distance between a “discovery” SNP and a “tested” feature, which by construction shrinks the coverage of xQTL SNPs, the cross-phenotype sharing increased (**Figure S2**).

The lenient strategy likely provides an over-estimate of π_1 , since the retained associations were selected by their strength, i.e. the smaller p-values kept. To tighten up our assessment of xQTL SNP sharing while not being too stringent, we examined the distance between each pair of “discovery” SNP and “test” feature, which we found to be a prime determinant of cross-phenotype sharing. For example, the strongest associated eQTL gene for each mQTL SNP is often the gene closest to the mQTL SNP (**Figure 3C**). We also observed similar trends for other cross-phenotype comparisons (results not shown). Based on this observation, we modified our analysis to only consider the closest feature to each xQTL SNP (**Figure 3D**).

Mediation analysis

We applied causal inference test (CIT)³⁴ to investigate whether the effect of a regulatory *cis* eQTL SNP is propagated through its impact on DNA methylation and/or histone modification (causal model) as well as whether the effect of an eQTL SNP on DNA methylation and/or histone modification is mediated through gene expression (reactive model). In brief, for the causal model, applying CIT involves testing the following four associations: 1) an eQTL SNP is associated with the first PC of its associated histone acetylation peaks and methylation probes (i.e. epigenome PC), 2) this eQTL SNP is associated with expression of a gene, 3) this eQTL SNP is associated with the epigenome PC conditioned on gene expression, and 4) this eQTL SNP is independent of gene expression given epigenome PC. Testing the reactive model involves reversing the role of gene expression and epigenome PC. A p-value was assigned to each set of associations using the Intersection-Union test³⁴. Bonferroni correction was applied to account for the number of tested association sets, m . We declared an association set as conforming to the causal model (or epigenetic mediation model) if $p_{\text{Causal}} < 0.05/m$

and $p_{\text{React}} > 0.05/m$ and conforming to the reactive model (or transcription mediation model) if $p_{\text{Causal}} > 0.05/m$ and $p_{\text{React}} < 0.05/m$. An association set was declared as conforming to the independent model if $p_{\text{Causal}} > 0.05/m$ and $p_{\text{React}} > 0.05/m$. The remaining association sets were considered unclassified. The above analysis was performed on 20,916 association sets for the 10,897 xQTL SNPs that are associated with all three molecular phenotypes. We restricted analysis to these shared xQTL SNPs since only these SNPs would fulfill conditions 1 and 2. The same analysis was also performed to assess the mediation of the shared xQTL SNPs through DNA methylation and histone acetylation separately. In this analysis, when multiple CpG probes (or acetylation peaks) were associated with a given xQTL SNP, we used their first PC to summarize their combination.

Disease enrichment analysis

We performed enrichment analysis on reported p-values of 16 GWAS datasets downloaded from the Psychiatric Genomics Consortium website: <https://www.med.unc.edu/pgc/results-and-downloads>. Data from the following GWAS studies were used in the analysis:

- Attention deficit hyperactivity disorder (ADHD)³⁶
- Alzheimer's disease (stage 1 data from IGAP)³⁷
- Anxiety (case vs. control and factor score from ANGST)³⁸
- Autism³⁹
- Bipolar disorder^{40,41}
- Major depressive disorder (MDD)⁴²
- Schizophrenia^{41,43}
- Body mass index (BMI)⁴⁴
- Height⁴⁵
- Crohn's disease⁴⁶
- Ulcerative colitis⁴⁷
- Inflammatory bowel disease⁴⁸
- Diabetes⁴⁹

Enrichment was assessed using stratified LD score regression (LDSR) to estimate partitioned heritability⁵⁰. For example, we labeled all xQTL SNPs as one category and SNPs in the LDSR baseline model as background. Significant enrichment was declared at an α of 0.05.

Cell type specificity analysis

We used a previous approach to estimate the cell-specificity of an eQTL SNP, based on a statistical model that tests for an interaction effect between the SNP genotype and proportion of a cell type of interest⁵². Proportions of neurons, astrocytes, microglia, oligodendrocytes, and endothelial cells were estimated with known cell type markers for

these cells. Specifically, ENO2 was used as the marker for neurons, CD68 for microglia, OLIG2 for oligodendrocytes, GFAP for astrocytes, and CD34 for endothelial cells. To reduce the number of tests, we only tested for cell-specificity of the lead eQTL SNPs. That is, we tested for cell-specificity of leads SNPs that impacted the expression levels of 3,388 genes with at least one significant eQTL SNP. In this analysis, we only corrected for known confounding factors, since regressing out the effect of hidden confounding factors would remove the effect of cell-specific expression⁵².

xQTL-weighted GWAS

We used the weighted Bonferroni procedure⁵³ to prioritize xQTL SNPs in GWAS analysis. This procedure involves weighting p-values (or summary statistics) from a GWAS study by their potential relevance. Provided that the weights are non-negative and average to one, strong control on family-wise error rate is guaranteed⁵³. We used this approach with a simple binary weighting scheme on the 16 GWAS datasets listed in the previous section, as well as GWAS datasets pertaining to systolic and diastolic blood pressure⁵⁴ and BHRadj BMI, which were excluded from the LDSR enrichment analysis due to unavailability of some of the required summary statistics. Specifically, p-values of xQTL SNPs are weighted by w_1 and all other SNPs are weighted by w_0 , where $w_1 = s/(1+(s-1)n_1/n)$ and $w_0 = 1/(1+(s-1)n_1/n)$ with $s = w_1/w_0$ ranging from 1 to 100. n_1 is the number of xQTL SNPs in our list and n is the number of SNPs in a study.

When only summary statistics are available, i.e. GWAS p-values, selecting the optimal s is nontrivial, since w_1 and w_0 do not depend on the p-values, i.e. w_1 and w_0 only depend on s , the number of SNPs, and the number of xQTL SNPs. Hence, we cannot “train” w_1 and w_0 based on p-values. We thus instead proposed to divide the list of p-values into random half splits and use the following criterion: $J(s) = (D^1(s)/\pi_1^1 + D^2(s)/\pi_1^2) / |D^1(s)/\pi_1^1 - D^2(s)/\pi_1^2|$, where $D^i(s)$ is the number of SNPs in half split i with weighted p-values $< 5 \times 10^{-8}$, and π_1^i is the estimated proportion of SNPs in half split i that are significant based on their non-weighted p-values. The rationale behind the proposed criterion, $J(s)$, is two-fold. First, if a given s is generalizable, then the detection rate should be similar for two half splits of randomly selected SNPs, as opposed to being large for one half but not the other. Second, among the s values that provide high “reproducibility” between splits, we should select the one that maximizes detection rate. We note that $|D^1(s) - D^2(s)|$ would not reflect reproducibility if the ground truth number of significant SNPs are different between the two splits. This complication is alleviated by dividing $D^i(s)$ by π_1^i .

To determine the number of independent significant SNPs, we applied PLINK1.9⁵⁶ pairwise LD pruning function ($r^2 = 0.2$) on the 1000 Genomes phase 1 data⁵⁷.

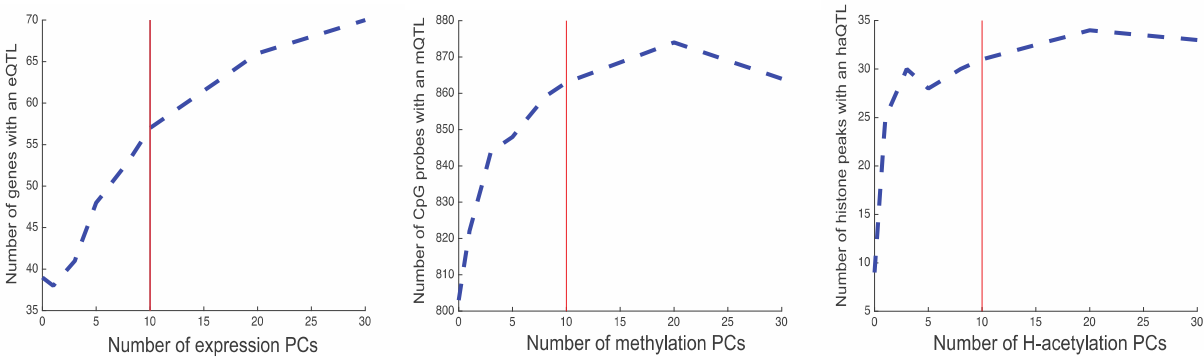


Figure S1. Figure shows the number of “features” (genes, methylation probes, histone peaks) detected to have a significant xQTL as the number of PCs (hidden confounds) was increased from 0 to 30. The optimal number of PCs to account for was deemed to be 10 for all three data types (data for gene expression shown in the right panel, DNA methylation shown in the in the middle panel, histone acetylation shown in the left panel). To avoid overfitting, this analysis was performed on features that reside on chromosome 18.

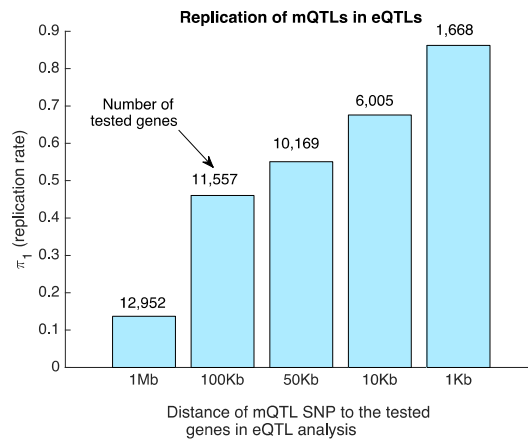


Figure S2. Figure shows the π_1 statistics for replication rate of mQTLs (“discovery”) in eQTLs (“replication”). As the window size centered at an mQTL for defining the “replication set” from eQTL data is decreased, the π_1 statistic also decreases but so does the coverage of testable eQTLs (SNPs and genes).

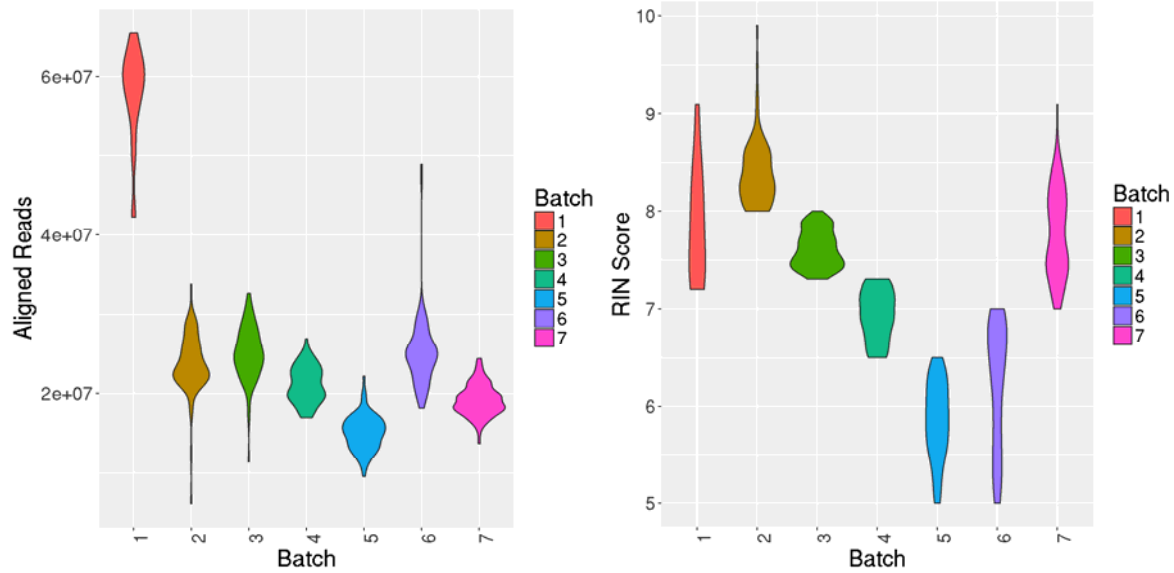


Figure S3. Figure on the right shows the number of aligned reads for RNA-seq data per batch. Figure on the left shows the distribution of RIN scores per RNA-seq batch.

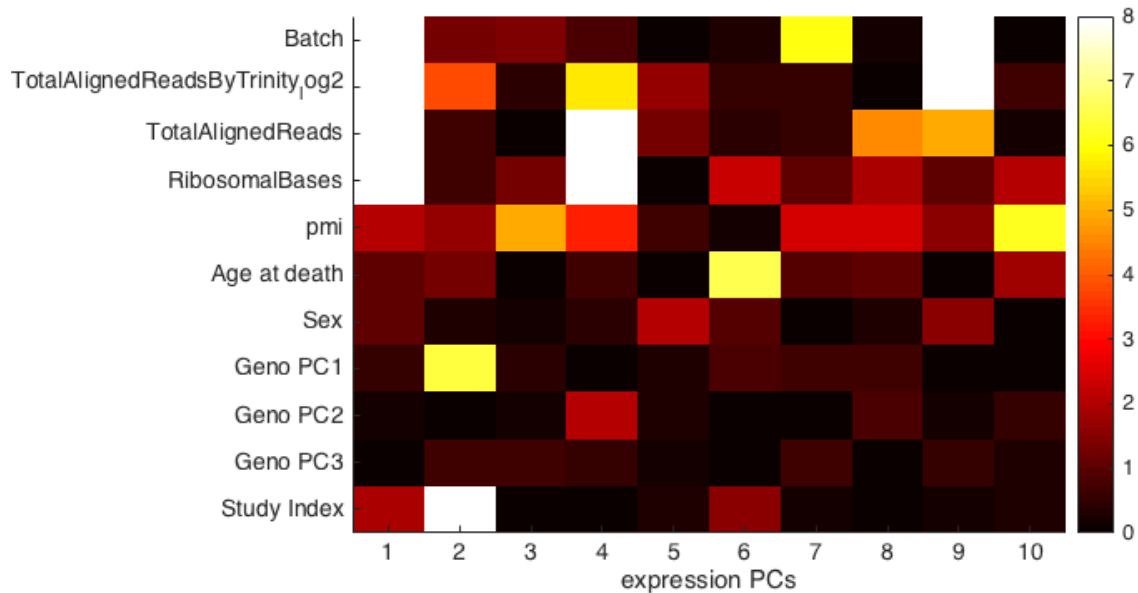


Figure S4. Figure shows the strength of the association between top expression PCs and 11 technical and biological confounding factors. Batch refers to the date of RNA preparation. PMI refers to the postmortem interval. Genotype PCs were computed as the top 3 PCs of genotype data. Study index refers to RUSH vs MAP samples. The heatmap depicts the log₁₀ p-values for correlation coefficient.

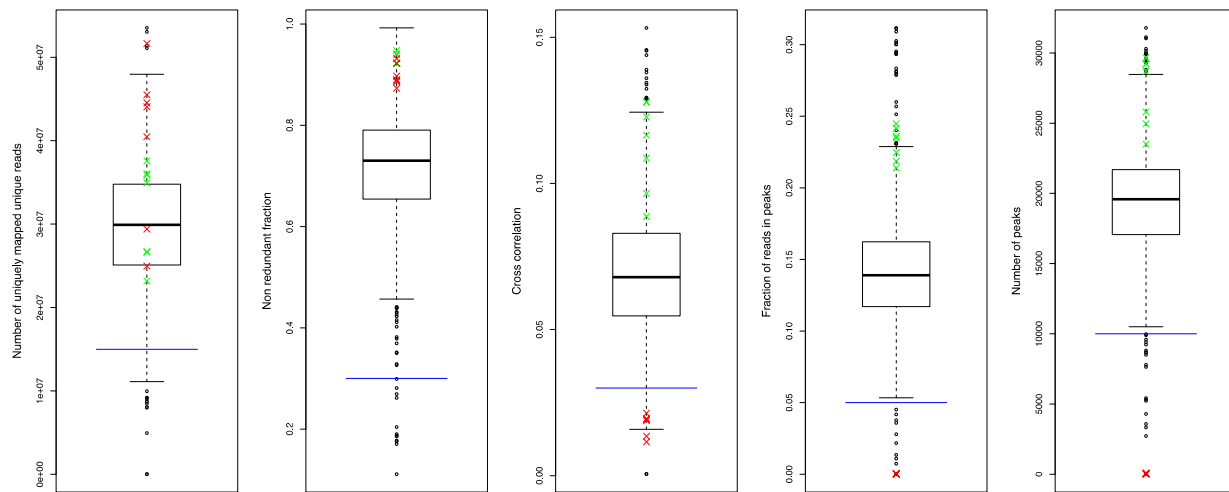


Figure S5. Figure shows several quality control metrics for H3K9Ac acetylation ChIP-Seq dataset.

References

- 1 De Jager, P. L. *et al.* A genome-wide scan for common variants affecting the rate of age-related cognitive decline. *Neurobiol Aging* **33**, 1017 e1011-1015, doi:10.1016/j.neurobiolaging.2011.09.033 (2012).
- 2 Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* **81**, 559-575, doi:10.1086/519795 (2007).
- 3 Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics* **38**, 904-909 (2006).
- 4 Browning, B. L. & Browning, S. R. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am J Hum Genet* **84**, 210-223, doi:10.1016/j.ajhg.2009.01.005 (2009).
- 5 Lim, A. S. *et al.* Diurnal and seasonal molecular rhythms in human neocortex and their relation to Alzheimer's disease. *Nat Commun* **8**, 14931, doi:10.1038/ncomms14931 (2017).
- 6 Levin, J. Z. *et al.* Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nat Methods* **7**, 709-715, doi:10.1038/nmeth.1491 (2010).
- 7 Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25, doi:10.1186/gb-2009-10-3-r25 (2009).
- 8 Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* **12**, 323, doi:10.1186/1471-2105-12-323 (2011).
- 9 Consortium, G. T. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648-660, doi:10.1126/science.1262110 (2015).
- 10 Trapnell, C. *et al.* Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc* **7**, 562-578, doi:10.1038/nprot.2012.016 (2012).

- 11 Johnson, W. E., Li, C. & Rabinovic, A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics* **8**, 118-127 (2007).
- 12 De Jager, P. L. *et al.* Alzheimer's disease: early alterations in brain DNA methylation at ANK1, BIN1, RHBDF2 and other loci. *Nature Neuroscience* **17**, 1156-1163 (2014).
- 13 Teschendorff, A. E. *et al.* A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data. *Bioinformatics* **29**, 189-196 (2013).
- 14 Mecham, B. H., Nelson, P. S. & Storey, J. D. Supervised normalization of microarrays. *Bioinformatics* **26**, 1308-1315, doi:10.1093/bioinformatics/btq118 (2010).
- 15 Kharchenko, P. V., Tolstorukov, M. Y. & Park, P. J. Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat Biotechnol* **26**, 1351-1359, doi:10.1038/nbt.1508 (2008).
- 16 Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).
- 17 Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137, doi:10.1186/gb-2008-9-9-r137 (2008).
- 18 Landt, S. G. *et al.* ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res* **22**, 1813-1831, doi:10.1101/gr.136184.111 (2012).
- 19 Mostafavi, S. *et al.* Normalizing RNA-sequencing data by modeling hidden covariates with prior knowledge. *PLoS One* **8**, e68141, doi:10.1371/journal.pone.0068141 (2013).
- 20 Stegle, O., Parts, L., Durbin, R. & Winn, J. A Bayesian framework to account for complex non-genetic factors in gene expression levels greatly increases power in eQTL studies. *PLoS Comput Biol* **6**, e1000770, doi:10.1371/journal.pcbi.1000770 (2010).
- 21 Battle, A. *et al.* Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome research*, doi:10.1101/gr.155192.113 (2013).
- 22 Banovich, N. E. *et al.* Methylation QTLs are associated with coordinated changes in transcription factor binding, histone modifications, and gene expression levels. *PLoS Genet* **10**, e1004663, doi:10.1371/journal.pgen.1004663 (2014).
- 23 Gutierrez-Arcelus, M. *et al.* Passive and active DNA methylation and the interplay with genetic variation in gene regulation. *Elife* **2**, e00523 (2013).
- 24 Do, C. *et al.* Mechanisms and Disease Associations of Haplotype-Dependent Allele-Specific DNA Methylation. *Am J Hum Genet* **98**, 934-955, doi:10.1016/j.ajhg.2016.03.027 (2016).
- 25 Ramasamy, A. *et al.* Genetic variability in the regulation of gene expression in ten regions of the human brain. *Nat Neurosci* **17**, 1418-1428, doi:10.1038/nn.3801 (2014).
- 26 Hannon, E. *et al.* Methylation QTLs in the developing brain and their enrichment in schizophrenia risk loci. *Nat Neurosci* **19**, 48-54, doi:10.1038/nn.4182 (2016).
- 27 Fromer, M. *et al.* Gene expression elucidates functional impact of polygenic risk for schizophrenia. *Nat Neurosci* **19**, 1442-1453, doi:10.1038/nn.4399 (2016).
- 28 Battle, A. *et al.* Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome research* **24**, 14-24 (2014).
- 29 Raj, T. *et al.* Polarization of the effects of autoimmune and neurodegenerative risk alleles in leukocytes. *Science* **344**, 519-523, doi:10.1126/science.1249547 (2014).
- 30 The Genotype-Tissue Expression (GTEx) project. *Nat Genet* **45**, 580-585, doi:10.1038/ng.2653 (2013).
- 31 Storey, J. D. & Tibshirani, R. Statistical significance for genomewide studies. *Proc Natl Acad Sci U S A* **100**, 9440-9445, doi:10.1073/pnas.1530509100 (2003).
- 32 Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and characterization. *Nat Methods* **9**, 215-216, doi:10.1038/nmeth.1906 (2012).
- 33 Rosenbloom, K. R. *et al.* The UCSC Genome Browser database: 2015 update. *Nucleic Acids Res* **43**, D670-681, doi:10.1093/nar/gku1177 (2015).

- 34 Millstein, J., Zhang, B., Zhu, J. & Schadt, E. E. Disentangling molecular relationships with a causal inference test. *BMC Genet* **10**, 23, doi:10.1186/1471-2156-10-23 (2009).
- 35 Baron, R. M. & Kenny, D. A. The moderator-mediator variable distinction in social psychological research: conceptual, strategic, and statistical considerations. *J Pers Soc Psychol* **51**, 1173-1182 (1986).
- 36 Neale, B. M. *et al.* Meta-analysis of genome-wide association studies of attention-deficit/hyperactivity disorder. *J Am Acad Child Adolesc Psychiatry* **49**, 884-897, doi:10.1016/j.jaac.2010.06.008 (2010).
- 37 Lambert, J. C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat Genet* **45**, 1452-1458, doi:10.1038/ng.2802 (2013).
- 38 Otowa, T. *et al.* Meta-analysis of genome-wide association studies of anxiety disorders. *Mol Psychiatry* **21**, 1391-1399, doi:10.1038/mp.2015.197 (2016).
- 39 Robinson, E. B. *et al.* Genetic risk for autism spectrum disorders and neuropsychiatric variation in the general population. *Nat Genet* **48**, 552-555, doi:10.1038/ng.3529 (2016).
- 40 Large-scale genome-wide association analysis of bipolar disorder identifies a new susceptibility locus near ODZ4. *Nat Genet* **43**, 977-983, doi:10.1038/ng.943 (2011).
- 41 Ruderfer, D. M. *et al.* Polygenic dissection of diagnosis and clinical dimensions of bipolar disorder and schizophrenia. *Mol Psychiatry* **19**, 1017-1024, doi:10.1038/mp.2013.138 (2014).
- 42 Sparse whole-genome sequencing identifies two loci for major depressive disorder. *Nature* **523**, 588-591, doi:10.1038/nature14659 (2015).
- 43 Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-427, doi:10.1038/nature13595 (2014).
- 44 Speliotes, E. K. *et al.* Association analyses of 249,796 individuals reveal 18 new loci associated with body mass index. *Nat Genet* **42**, 937-948, doi:10.1038/ng.686 (2010).
- 45 Lango Allen, H. *et al.* Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* **467**, 832-838, doi:10.1038/nature09410 (2010).
- 46 Franke, A. *et al.* Genome-wide meta-analysis increases to 71 the number of confirmed Crohn's disease susceptibility loci. *Nat Genet* **42**, 1118-1125, doi:10.1038/ng.717 (2010).
- 47 Anderson, C. A. *et al.* Meta-analysis identifies 29 additional ulcerative colitis risk loci, increasing the number of confirmed associations to 47. *Nat Genet* **43**, 246-252, doi:10.1038/ng.764 (2011).
- 48 Liu, J. Z. *et al.* Association analyses identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk across populations. *Nat Genet* **47**, 979-986, doi:10.1038/ng.3359 (2015).
- 49 Gaulton, K. J. *et al.* Genetic fine mapping and genomic annotation defines causal mechanisms at type 2 diabetes susceptibility loci. *Nat Genet* **47**, 1415-1425, doi:10.1038/ng.3437 (2015).
- 50 Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat Genet* **47**, 1228-1235, doi:10.1038/ng.3404 (2015).
- 51 Pers, T. H., Timshel, P. & Hirschhorn, J. N. SNPsnap: a Web-based tool for identification and annotation of matched SNPs. *Bioinformatics* **31**, 418-420, doi:10.1093/bioinformatics/btu655 (2015).
- 52 Westra, H. J. *et al.* Cell Specific eQTL Analysis without Sorting Cells. *PLoS Genet* **11**, e1005223, doi:10.1371/journal.pgen.1005223 (2015).
- 53 Roeder, K., Devlin, B. & Wasserman, L. Improving power in genome-wide association studies: weights tip the scale. *Genet Epidemiol* **31**, 741-747, doi:10.1002/gepi.20237 (2007).
- 54 Ehret, G. B. *et al.* Genetic variants in novel pathways influence blood pressure and cardiovascular disease risk. *Nature* **478**, 103-109, doi:10.1038/nature10405 (2011).

- 55 Heid, I. M. *et al.* Meta-analysis identifies 13 new loci associated with waist-hip ratio and reveals sexual dimorphism in the genetic basis of fat distribution. *Nat Genet* **42**, 949-960, doi:10.1038/ng.685 (2010).
- 56 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* **4**, 7, doi:10.1186/s13742-015-0047-8 (2015).
- 57 Abecasis, G. R. *et al.* An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56-65, doi:10.1038/nature11632 (2012).