

Genetic control of RNA splicing and its distinctive role in complex trait variation

Jian Yang (✉ jian.yang@westlake.edu.cn)

Westlake University <https://orcid.org/0000-0003-2001-2474>

Ting Qi

Westlake University

Yang Wu

The University of Queensland <https://orcid.org/0000-0002-0128-7280>

Futao Zhang

The University of Queensland

Jian Zeng

University of Queensland <https://orcid.org/0000-0001-8801-5220>

Article

Keywords: genome-wide association studies (GWAS), gene regulation

Posted Date: February 10th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-155233/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Genetics on August 18th, 2022. See the published version at <https://doi.org/10.1038/s41588-022-01154-4>.

1 Genetic control of RNA splicing and its distinctive role in complex trait variation

2
3 Ting Qi^{1,2,3}, Yang Wu³, Futao Zhang^{3,4,5}, Shouye Liu^{1,6}, Jian Zeng³, Jian Yang^{1,2,3,*}

4
5 ¹School of Life Sciences, Westlake University, Hangzhou, Zhejiang 310024, China

6 ²Westlake Laboratory of Life Sciences and Biomedicine, Hangzhou, Zhejiang 310024, China

7 ³Institute for Molecular Bioscience, The University of Queensland, Brisbane, Queensland 4072,
8 Australia

9 ⁴Neuroscience Research Australia, Sydney, New South Wales 2031, Australia

10 ⁵Clinical Genetics and Genomics, NSW Health Pathology Randwick, Sydney, New South Wales
11 2031, Australia

12 ⁶Institute for Advanced Research, Wenzhou Medical University, Wenzhou, Zhejiang 325027,
13 China

14 *Correspondence: Jian Yang (jian.yang@westlake.edu.cn)

15 16 Abstract

17 Most genetic variants identified from genome-wide association studies (GWAS) in humans are
18 noncoding, indicating their role in gene regulation. Prior studies have shown considerable links
19 of GWAS signals to expression quantitative trait loci (eQTLs), but the links to other genetic
20 regulatory mechanisms such as splicing QTLs (sQTLs) are underexplored. Here, we introduce a
21 transcript-based sQTL method (named THISTLE) with improved power for sQTL detection.
22 Applying THISTLE along with LeafCutter, an event-based sQTL method, to brain transcriptomic
23 data ($n = 1,073$), we identified 7,491 genes with sQTLs with $P < 5 \times 10^{-8}$ (the largest brain cis-
24 sQTL collection to date), ~68% of which were distinct from eQTLs. Integrating the sQTL data into
25 GWAS for ten brain-related complex traits (including diseases), we identified 107 genes
26 associated with the traits through the sQTLs, ~68% of which could not be discovered using eQTL
27 data. Our study demonstrates the distinctive role of most sQTLs in genetic regulation of
28 transcription and complex trait variation.

30 **Introduction**

31 Genome-wide association studies (GWAS) have led to the discovery of tens of thousands of
32 genetic variants associated with human complex traits (including diseases)^{1,2}. However, the
33 majority of the trait-associated variants are of uncharacterised function, and the mechanisms
34 through which the genetic variants exert their effects on the traits remain largely elusive.
35 Considering that most of the GWAS signals are located in non-coding regions of the genome, one
36 hypothesis is that the genetic variants affect the traits through genetic regulation of gene
37 expression^{3,4}. The effects of genetic variation on mRNA abundance (also known as expression
38 quantitative trait loci or eQTLs) have been studied for more than a decade⁵⁻⁷, and nearly all genes
39 have been found with one or more genetic variants associated with their mRNA abundance⁶⁻⁸.
40 These advances have propelled the development of methods⁹⁻¹³ to integrate eQTL data with
41 GWAS to prioritize genes responsible for the GWAS signals. However, only a moderate proportion
42 of the GWAS signals have been attributed to cis-eQTLs¹⁴⁻¹⁸, likely because of various reasons
43 including limited power, spatiotemporal eQTL effects, the focus on genomic regions in cis, and
44 mechanisms beyond genetic control of mRNA abundance.

45

46 Genetic control of pre-mRNA splicing (also called splicing quantitative trait locus or sQTL) is
47 another fundamental mechanism of gene regulation but remains heavily underexplored
48 compared to eQTL. Currently, there is no consensus in the literature regarding the relationship
49 between eQTL and sQTL. For example, there are observations showing that sQTLs are largely
50 independent of eQTLs^{7,19,20} and hypothesized to be one of the major contributors to genetic risk
51 of disease²⁰, whereas a recent study shows that sQTLs do not have statistically significant
52 contribution to trait heritability conditional on eQTLs²¹. These results motivated us to investigate
53 the role of sQTLs in complex traits using a larger data set with a rich collection of sQTLs and a
54 broader spectrum of splicing events in a single tissue. We focused our study on the brain because
55 of data availability (see below) and the considerable links of sQTLs to neurodegenerative diseases
56 such as Alzheimer's disease²² and Parkinson's disease²³ reported recently.

57

58 There are two major classes of methods for sQTL mapping²⁴, transcript-based methods²⁵⁻²⁹ and
59 event-based methods³⁰⁻³³, each favouring different types of splicing events. The transcript-based
60 methods typically involve the generation of a splicing phenotype such as isoform ratio (i.e., a ratio
61 of mRNA abundance of an isoform to the overall mRNA abundance of the gene), and the test of
62 associations between genetic variants and the splicing phenotype²⁷⁻²⁹. sQTLseeker²⁹ is one of the
63 most commonly used transcript-based methods. It can detect sQTLs arising from complex
64 splicing events because it is sensitive to splicing events involving main isoform changes. However,
65 sQTLseeker often reports a modest number of sQTLs and suffers from computational burden. On

66 the other hand, event-based methods³⁰⁻³³ such as LeafCutter³³ typically involve the quantification
 67 of splicing events by counting RNA-seq reads aligned to specific exons or splice junctions, the
 68 estimation of a splicing phenotype (e.g., exon-inclusion level), and the association tests for the
 69 splicing phenotype. More specifically, LeafCutter³³, which does not require read assembly or
 70 inference of isoforms, allows the identification and quantification of novel splicing events by
 71 focusing on intron excision. LeafCutter can detect sQTLs arising from subtle splicing patterns (e.g.,
 72 exon skipping) because it is sensitive to local splicing events.

73

74 In this study, we aim to investigate the genetic control of RNA splicing by generating the largest
 75 collection of sQTLs to date in brain and to understand the role of sQTLs in complex trait variation.
 76 To harvest as many sQTLs as possible, we used both the transcript- and event-based methods for
 77 sQTL analysis. The limited number of sQTLs detected by sQTLseekeR motivated us to develop a
 78 more powerful and efficient transcript-based sQTL method (dubbed THISTLE) that requires only
 79 isoform-level eQTL (isoform-eQTL) summary data. We show by simulation that THISTLE has
 80 well-calibrated test-statistics and improves power over sQTLseekeR. We applied THISTLE
 81 together with LeafCutter to the PsychENCODE data^{34,35} (the largest publicly available brain
 82 transcriptomic data, $n = 1,073$), integrated the sQTL data into GWAS for ten brain-related traits
 83 (including diseases) of large sample sizes (n ranging from 46,350 to 766,345), and compared the
 84 results with those for eQTLs identified using the same data.

85

86 **Results**

87 **Overview of THISTLE**

88 Details of the THISTLE method can be found in the Methods section. For ease of understanding,
 89 let us take a gene with two transcript isoforms as an example. If a genetic variant is associated
 90 with a splicing event (e.g., exon skipping), the variant is expected to be associated with the
 91 difference in mature mRNA abundance between the two transcript isoforms (**Fig. 1**). Hence, an
 92 sQTL test can be performed by assessing the association of the variant with the difference in
 93 mRNA abundance between the isoforms, which is equivalent to a test of the difference in variant
 94 effect on mRNA abundance (i.e., isoform-eQTL effect) between the isoforms. In other words, it is
 95 a test of heterogeneity in isoform-eQTL effect between isoforms. Without loss of generality, let m
 96 be the number of transcript isoforms for a gene and $\hat{\mathbf{b}} = \{\hat{b}_1, \dots, \hat{b}_j, \dots, \hat{b}_m\}$, where \hat{b}_j is the
 97 estimated isoform-eQTL effect for isoform j . We have $\hat{\mathbf{b}} \sim MVN(\mathbf{b}, \mathbf{S})$ with \mathbf{S} being the variance-
 98 covariance matrix of $\hat{\mathbf{b}}$. The difference between b_j and b_k ($j \neq k$) can be estimated as

99 $\hat{d}_{jk} = \hat{b}_j - \hat{b}_k$. Then, if we let $\hat{\mathbf{d}} = \{\hat{d}_{jk}\}_{j,k \in \{1, \dots, m\}, j < k}$, we have $\hat{\mathbf{d}} \sim MVN(\mathbf{d}, \mathbf{V})$ with \mathbf{V} being the
 100 variance-covariance matrix of $\hat{\mathbf{d}}$. The covariance between \hat{d}_{gh} and \hat{d}_{jk} is $cov(\hat{d}_{gh}, \hat{d}_{jk}) =$

101 $cov(\hat{b}_g, \hat{b}_j) - cov(\hat{b}_n, \hat{b}_j) - cov(\hat{b}_g, \hat{b}_k) + cov(\hat{b}_n, \hat{b}_k)$, and the covariance between \hat{b}_g and \hat{b}_j is
102 $cov(\hat{b}_g, \hat{b}_j) = \theta_{gj} S_g S_j$, where S_g^2 and S_j^2 are the variance of \hat{b}_g and \hat{b}_j , respectively, and θ_{gj} is the
103 correlation between \hat{b}_g and \hat{b}_j . We have shown previously^{11,36} that all the parameters in \mathbf{V} can be
104 estimated from summary-level data if individual-level data are unavailable (**Methods**). Under the
105 null hypothesis of no sQTL effect, $\mathbf{d} = \mathbf{0}$. We use an approximate multivariate approach, similar
106 to the heterogeneity in dependent instruments (HEIDI) test¹¹, to test against the null hypothesis
107 given \mathbf{V} . We name the method THISTLE, which stands for testing for heterogeneity between
108 isoform-eQTL effects, and have implemented it in the OSCA³⁷ software.

109

110 **Calibration of THISTLE**

111 We first used simulations to calibrate THISTLE (**Supplementary Note**) and compare it with the
112 competing method, sQTLseeker. The simulation results showed that there was no inflation in
113 false-positive rate (FPR) under the null of no sQTL effect for both THISTLE and sQTLseeker, in
114 either the absence or the presence of eQTL effect (**Supplementary Fig. 1**). With respect to
115 statistical power, the area under the receiver operating characteristic curve (AUC) for THISTLE
116 was larger than that for sQTLseeker, and the difference in power between the two methods
117 increased with the increase of $-\log_{10}(\text{p-value})$ threshold (**Supplementary Fig. 2**). The difference
118 in power was more prominent in the analysis of the PsychENCODE data (see below for more
119 details), with 699,445 sQTLs discovered by THISTLE for 8,241 genes vs. 178,961 sQTLs
120 discovered by sQTLseeker for 4,887 genes using a false discovery rate (FDR) threshold of 0.01
121 (**Supplementary Note**), predominantly owing to a substantial number of genes filtered out by
122 sQTLseeker as being non-informative. Nevertheless, for the genes tested in both methods, 70%
123 of the sQTLs discovered by THISTLE with $P_{\text{sQTL}} < 5 \times 10^{-8}$ were replicated in sQTLseeker with
124 $P_{\text{sQTL}} < 0.05$ (**Supplementary Fig. 3**). Note that the reciprocal analysis (selecting sQTLs from
125 sQTLseeker for replication with THISTLE) was infeasible because the sQTLseeker p-values were
126 capped at 1×10^{-6} owing to the use of permutation to compute the p-values. In addition, the
127 performance of THISTLE using individual-level SNP genotype and RNA-seq data is similar to that
128 using summary-level isoform-eQTL data (**Supplementary Fig. 4**), demonstrating the flexibility
129 of the method.

130

131 **Identifying cis-sQTLs in brain**

132 We applied THISTLE along with LeafCutter to a brain transcriptomic data set from the
133 PsychENCODE Consortium. We processed the PsychENCODE data from six cohorts: UCLA-ASD,
134 BrainGVEX, the Lieber Institute for Brain Development (LIBD), the CommonMind Consortium
135 (CMC), the CMC's NIMH Human Brain Collection Core (HBCC), and the Bipseq of Bipolar cohorts
136 (Bipseq). Generation and processing of the genotype and RNA-seq data have been detailed

137 elsewhere^{34,35}. Further quality control (QC) in each cohort and the combined data set has been
138 summarized in the Methods section (**Supplementary Figs. 5 & 6**). After QC, we included in the
139 analysis ~3.8 million genotyped or imputed common variants (minor allele frequency, MAF >
140 0.01) on 1,073 unrelated individuals of European ancestry, including 348 subjects with
141 schizophrenia (SCZ), 164 subjects with bipolar disorder (BIP), 17 subjects with autism spectrum
142 disorder (ASD), and 544 controls. For THISTLE analysis, we excluded isoforms with low
143 expression level (i.e., Transcripts Per Million (TPM) < 0.1 in more than 20% of individuals) or
144 genes with only one isoform, retaining 77,511 isoforms of 13,366 genes (**Methods**). To remove
145 confounding factors, VariancePartition³⁸ was applied to decompose the variation in isoform
146 abundance into components attributable to multiple known biological and technical factors such
147 as study, bank, RNA quality (RIN), and age at death (**Supplementary Fig. 7**). To account for other
148 unknown biological and technical factors such as hidden batch effects, Probabilistic Estimation of
149 Expression Residuals (PEER) method³⁹ was applied to obtain a set of latent covariates. The
150 adjusted isoform abundances for each gene were used for an isoform-eQTL analysis in OSCA³⁷
151 with the first 5 genetic principal components (PCs) and 50 PEER factors fitted as covariates,
152 followed by an sQTL analysis using THISTLE (**Supplementary Fig. 8**). Recognising the likely
153 insufficient power to detect trans-sQTLs and for ease of computing, we limited the sQTL test to
154 SNPs within 2 Mb of each gene on either side. In total, we identified 610,039 significant cis-sQTLs
155 for 4,881 genes with $P_{sQTL} < 5 \times 10^{-8}$, comprising 433,220 unique SNPs (**Supplementary Fig.**
156 **9a**). The use of a p-value threshold of 5×10^{-8} is because it is the genome-wide significance level
157 used in GWAS and also the default threshold used in SMR to select instrument SNPs¹¹ (see below).
158 Note that if we use a less stringent p-value threshold (e.g., FDR < 0.01), there were 1,114,064
159 significant cis-sQTLs for 8,241 genes, comprising 699,445 unique SNPs. Of the 4,881 genes with
160 at least one sQTL with $P_{sQTL} < 5 \times 10^{-8}$, 87 (1.8%) were lowly expressed in brain with median
161 Fragments Per Kilobase of transcript per Million mapped reads (FPKM) ranging from 0.2 to 1
162 (**Supplementary Fig. 10a**).

163

164 We also used LeafCutter, the state-of-the-art event-based sQTL method, to identify sQTLs for
165 splicing events that may not be well captured by a transcript-based method like THISTLE. Using
166 LeafCutter, we quantified excision ratios of 202,344 introns in 46,868 intronic excision clusters
167 (**Methods**) with 46,199 clusters uniquely mapped to 16,518 genes and the remaining clusters
168 either mapped to more than one gene or not located in the gene body. In the following analyses,
169 we focused only on the 46,199 uniquely mapped clusters. To define the set of covariates for
170 adjustment, VariancePartition³⁸ was applied to decompose the variance in intron excision ratio
171 into components attributable to biological and technical factors (**Methods**). The results showed
172 that a substantial proportion of variance in gene expression level could be explained by multiple

173 biological and technical factors such as study, bank, library batch, and RIN (**Supplementary Fig.**
174 **7b**). In contrast, variance in intron excision ratio explained by these factors were, on average,
175 minor (**Supplementary Fig. 7a**), likely because these factors affect both the numerator and
176 denominator of the intron excision ratio so that the effects are mostly cancelled out each other.
177 Each of the intron excision ratio was then adjusted for the selected covariates (i.e., bank, study,
178 RIN, age of death, PMI, library batch, platform, sex, and diagnosis) and used as a quantitative trait
179 for association analysis with SNPs within 2 Mb of each intron using FastQTL⁴⁰ for sQTL test, with
180 the first 5 genetic PCs and 15 PEER factors (estimated from the adjusted intron excision ratio)
181 fitted as covariates (**Supplementary Fig. 8**). In total, we identified 1,860,296 cis-sQTLs for 7,451
182 LeafCutter intron clusters in 5,629 genes with $P_{sQTL} < 5 \times 10^{-8}$, comprising 494,283 unique SNPs
183 (**Supplementary Fig. 9b**). Of note, if we use an FDR threshold of 0.01, there were 3,479,048
184 significant cis-sQTLs for 16,518 intron clusters in 10,186 genes, comprising 787,277 unique SNPs.
185 Of the 5,629 genes with at least one sQTL with $P_{sQTL} < 5 \times 10^{-8}$, 786 genes (14.0%) were lowly
186 expressed in brain (**Supplementary Fig. 10c**).

187

188 We define genes with at least one significant sQTL with $P_{sQTL} < 5 \times 10^{-8}$ as sGenes. Combining the
189 sQTL results from THISTLE and LeafCutter, there were a total of 632,481 sQTLs associated with
190 7,491 sGenes, of which 1,862 (25%) and 2,610 (35%) sGenes were detected only by THISTLE and
191 LeafCutter, respectively, with 3,019 sGenes (40%) detected by both methods (**Supplementary**
192 **Fig. 11a**). Although the proportion of sGenes over the tested genes for THISTLE ($4,881/13,366 =$
193 36.5%) was slightly higher than that for LeafCutter ($5,629/16,518 = 34.1\%$), LeafCutter
194 identified more sGenes (5,629 vs. 4,881) than THISTLE (**Supplementary Fig. 11a**), partly
195 because THISTLE focuses only on genes with more than one annotated isoform while LeafCutter
196 does not have such limitation. For example, 602 LeafCutter sGenes were filtered out in the
197 THISTLE analysis because they have only one annotated isoform (**Supplementary Fig. 10d**). In
198 addition, we have developed an online R shiny application to query the sQTLs with $FDR < 0.01$
199 (<https://yanglab.westlake.edu.cn/qi-sqtl/>).

200

201 Of the 3,019 sGenes detected by both methods, there were 1,954 sGenes for which the THISTLE
202 sQTL signal was shared with the LeafCutter sQTL signal as indicated by a COLOC⁹ PP4 value of $>$
203 0.8 , in line with the observation that the THISTLE top sQTL SNP was in high linkage
204 disequilibrium (LD), $r^2 > 0.8$, with the LeafCutter top sQTL SNP for 1,723 sGenes (**Supplementary**
205 **Fig. 11b**). Note that PP4 is the posterior probability for the 4th hypothesis of COLOC that a trait is
206 associated with another trait through a shared genetic variant⁹. These results mean that apart
207 from the 1,862 THISTLE-specific sGenes and 2,610 LeafCutter-specific sGenes, there were
208 additionally 1,065 sGenes for which the sQTL signals detected by the two methods were distinct,

209 demonstrating the benefit of using a combination of transcript- and event-based methods for
210 sQTL detection.

211

212 **Quantifying the relationship between sQTLs and eQTLs**

213 To assess the relationship between sQTLs and eQTLs, we performed an eQTL analysis using the
214 same data as used in the sQTL analysis above and identified 663,326 cis-eQTLs with $P_{\text{eQTL}} <$
215 5×10^{-8} for 9,524 genes. Similar as above, we define eGenes as genes with at least one significant
216 eQTL with $P_{\text{eQTL}} < 5 \times 10^{-8}$. We found that 65% (4,893/7,491) of the sGenes were also eGenes
217 (**Fig. 2a**), a proportion higher than that (50%) reported in a recent study using transcriptomic
218 data from fetal brain¹⁹. We hypothesize that the difference is due to the larger sample size (n)
219 used in our study than that used in the fetal brain study ($n = 233$). To test this hypothesis, we used
220 a down-sampling strategy to assess the sGene-eGene overlap in several subsets of data with n
221 ranging from 100 to 1,073. The result showed that the power of either sGene or eGene discovery
222 was proportional to n and that the difference in discovery power between sGene and eGene
223 increased with increasing n (**Fig. 2c**), in line with the observation from prior work⁷. We also found
224 that the sGene-eGene overlap was positively correlated with n (**Fig. 2e**), which is expected if most
225 genes have both eQTL and sQTL. Of the 4,893 overlapping genes, there were 2,418 genes (32% of
226 all sGenes) for which the sQTL signal was shared with the eQTL signal as indicated by a COLOC
227 PP4 value of > 0.8 , in line with the observation that LD r^2 between the top sQTL SNP and the top
228 eQTL SNP was > 0.8 for only 2,094 genes (**Fig. 2b**). In summary, our results suggest that although
229 a large proportion of sGenes are expected to be eGenes with large n and this proportion will
230 increase with increasing n , most sQTLs (an estimate of $\sim 68\%$ in this study) are distinct from
231 eQTLs.

232

233 **sQTLs are enriched for splicing sites**

234 Having shown that sQTLs were largely distinct from eQTLs (**Fig. 2**), we then asked whether the
235 sQTLs and eQTLs show different patterns of functional enrichment. To do this, we annotated the
236 7,187 top sQTL SNPs and 9,106 top eQTL SNPs (only one top sQTL or eQTL per gene) using
237 SnpEff⁴¹ alongside the functional annotation data from the brain samples of the Roadmap
238 Epigenomics Mapping Consortium (REMC)⁴² (**Methods**). The fold enrichment was computed as
239 the proportion of sQTLs or eQTLs in a functional category divided by the mean of a null
240 distribution generated by resampling “control” cis-SNPs with MAF matched. It is reassuring that
241 sQTLs were more enriched in “splice regions” (i.e., splice acceptor, splice donor, and splice region)
242 than eQTLs (on average, 26.7-fold of enrichment for sQTLs vs. 9.9-fold for eQTLs; **Fig. 3a & 3b**).
243 eQTLs were more enriched in transcription start site (TSS) than sQTLs (11.5-fold, s.e. = 0.25, for
244 eQTLs vs. 7.8-fold, s.e. = 0.28, for sQTLs; **Fig. 3c & 3d**), consistent with our observation that eQTLs

245 were much closer to TSS than sQTLs (**Fig. 3e**). The results remained largely unchanged if we
246 further matched the control SNPs with the top cis-sQTL (or cis-eQTL) SNPs by distance to TSS
247 (**Supplementary Fig. 12**). It is of note that the enrichment of eQTLs in “splice regions” and the
248 enrichment of sQTLs in TSS are not unexpected because of the shared sQTL and eQTL signals for
249 about one third of the sGenes.

250

251 **Enrichment of the sQTLs for heritability of complex traits**

252 As described above, an sQTL is a genetic variant that affects transcript abundance because of its
253 effect on splicing event. Hence, if an sGene is functionally related to a trait, the sQTL of the gene
254 would exert an effect on the trait. We tested whether the brain cis-sQTLs are enriched for genetic
255 variants associated with complex traits and disorders related to the brain. We acquired GWAS
256 summary statistics for ten brain-related traits, i.e., fluid intelligence score (IQ), educational
257 attainment (EA), smoking initiation (SmkInit), schizophrenia (SCZ), Alzheimer’s disease (AD),
258 Parkinson’s disease (PD), autism spectrum disorder (ASD), insomnia, bipolar disorder (BIP), and
259 major depression (MD) from published studies⁴³⁻⁵² (**Methods & Supplementary Table 1**). Both
260 the sQTLs and eQTLs showed more inflated GWAS test-statistics compared to the other SNPs for
261 all the ten traits, and the levels of inflation were indistinguishable between the sQTLs and the
262 eQTLs (**Supplementary Fig. 13**). Nevertheless, if a phenotype is polygenic, the GWAS test-
263 statistics are expected to be inflated even in the absence of confounding factors, and the level of
264 inflation depends on sample size, trait heritability, and LD between SNPs and causal variants^{53,54}.
265 To account for the LD, we performed a stratified LD score regression (S-LDSC)^{55,56} analysis to
266 quantify the enrichment of the top cis-sQTL SNPs for SNP-based heritability (h_{SNP}^2) in comparison
267 with that of the top cis-eQTL SNPs when fitted together with 53 other functional categories in the
268 “baselineLD model”⁵⁵ (**Methods**). The result showed that both the sQTLs and eQTLs were highly
269 enriched for heritability of all the traits except AD, and the fold enrichment for sQTLs (median
270 enrichment = 65) was similar to that of eQTLs (median enrichment = 55) on average across traits
271 (**Fig. 4a**). It is of note that the top SNPs were ascertained in the cis-regions, and SNPs in or near
272 genes are known to explain disproportionately more trait variation than those in intergenic
273 regions⁵⁷. To avoid such an ascertainment bias, we re-computed the enrichment by comparing
274 the per-SNP heritability of the top sQTL (or eQTL) SNPs to a set of MAF-matched “control SNPs”
275 resampled from the SNPs used in the sQTL (or eQTL) analysis (i.e., the SNPs in cis). In this case,
276 the overall levels of enrichment decreased, but the fold enrichment of sQTLs (median = 1.85)
277 remained similar to that of eQTLs (median = 1.73) on average across traits (**Fig. 4b**). We further
278 quantified the heritability enrichment at all the significant cis-sQTL (or cis-eQTL) SNPs, clumped
279 cis-sQTL (or cis-eQTL) SNPs, or fine-mapped cis-sQTL (or cis-eQTL) SNPs with maximum causal
280 posterior probabilities²¹, and the conclusion that sQTLs have a similar level of heritability

281 enrichment to eQTLs still held (**Supplementary Figs. 14a, 15a & 16a**). In addition, to account
282 for overlap between the sQTLs and eQTLs, we estimated the contributions of the top cis-sQTL and
283 the top cis-eQTL SNPs to trait heritability in a joint model and determined that the annotation
284 effect size (τ) remained statistically significant for both the sQTLs and eQTLs for most traits (**Fig.**
285 **4c, Supplementary Figs. 14b, 15b & 16b**).

286
287 We have shown that both the sQTLs and the eQTLs were enriched for trait heritability. However,
288 it is unclear whether or not the enrichment is driven by mediated genetic effects on trait through
289 pre-mRNA splicing (or mRNA abundance). To address this question, we applied a recently
290 developed method, mediated expression score regression (MESRC)¹⁸, to quantify the proportion of
291 heritability mediated (h_{med}^2/h_g^2) by the cis-sQTLs and the cis-eQTLs. We included only the sQTLs
292 discovered by LeafCutter because the MESRC analysis requires estimated sQTL effects which are
293 not provided by THISTLE. Across the ten traits, we observed a median estimate of h_{med}^2/h_g^2 of
294 0.13 and 0.11 from expression scores and splicing scores, respectively (**Fig. 4d**), where the
295 estimate for expression scores is similar to that reported by a previous study¹⁸. All the results
296 from the heritability enrichment and mediation analyses suggest that sQTLs play a unique role in
297 explaining complex trait variation, at a level that is comparable to that for eQTLs (see below for
298 more discussion).

300 **Identifying complex trait genes using sQTL data**

301 To leverage the cis-sQTLs to prioritize functional genes for the ten brain-related traits, we applied
302 the summary-data-based Mendelian Randomization (SMR) approach¹¹ to test if an sGene is
303 associated with a trait through sQTL and the COLOC approach⁹ to assess whether the sGene-trait
304 association is driven by the same set of causal variants. We used a combination of SMR and COLOC
305 for this analysis because while SMR is more powerful than COLOC in detecting sGene-trait
306 associations¹¹, the HEIDI test in SMR (for filtering out sGene-trait associations not driven by the
307 same set of causal variants) is inapplicable to the THISTLE sQTL summary data owing to the
308 unavailability of sQTL effects. We identified 212 sGenes in total for the ten traits at a genome-
309 wide significance level ($P_{SMR} < 2.1 \times 10^{-6}$), correcting for ~23,400 SMR tests for each trait, 107
310 of which showed a COLOC PP4 value of > 0.8 (**Fig. 4e & Supplementary Table 2**), consistent with
311 a plausible mechanism that genetic variants affect a trait through genetic control of splicing.
312 *MPHOSPH9* is a notable example where *MPHOSPH9* was identified to be associated with SCZ using
313 sQTL summary data from either THISTLE or LeafCutter (**Supplementary Fig. 17**), and the top
314 cis-sQTL SNP, rs1727307, for the intronic excision event, 12:123707674:123710835:clu_71282_-,
315 is located in the introns that are actively retained in mature mRNAs (intron retention). We also
316 included the eQTL data in the SMR analysis for the ten traits and identified 317 eGenes ($P_{SMR} <$

317 5.3×10^{-6}), 149 of which reached $PP4 > 0.8$ in COLOC (**Fig. 4e** & **Supplementary Table 2**).
318 Between the two sets of trait-associated genes discovered through sQTLs and eQTLs respectively,
319 34 genes were in common, meaning that we identified 73 more genes through sQTLs on top of
320 the 149 genes identified through eQTLs, a ~49% increase (**Fig. 4e**). One of the examples is *DGKZ*
321 (**Fig. 5**), the functional relevance of which to SCZ has been implicated in prior work using gene
322 expression data from blood⁵⁸ or multiple tissues⁵⁹. In this study, *DGKZ* was identified to be
323 associated with SCZ using the sQTL data, implying a possible mechanism that the SNP effects on
324 SCZ are mediated by genetic regulation of RNA splicing of *DGKZ*. Notably, no genome-wide
325 significant eQTL was identified to be associated with the overall expression level of *DGKZ* in
326 PsychENCODE (**Fig. 5b**), meaning that *DGKZ* would have been missed if we analyzed the eQTL
327 data only. In summary, ~68% (73/107) of the trait-associated genes identified through sQTLs
328 were not detected using eQTL data, which again suggests the distinctive role of sQTL in complex
329 traits.

330

331 For each of the 34 trait-associated genes that could be identified by using either the sQTL or the
332 eQTL data, we further performed a COLOC analysis to test whether the sQTL and eQTL signals are
333 driven by the same or distinct causal variants (i.e., sharing vs. linkage). We found evidence for
334 sharing ($PP4 > 0.80$) for 25 genes (23% of all trait-associated sGenes). One typical example is
335 *SETD6*, for which the COLOC analysis ($PP4 = 0.94$) suggested that the sQTL, eQTL and SCZ GWAS
336 signals are driven by the same causal variant(s) (**Supplementary Fig. 18**). However, this result
337 also means that there were 9 genes, each of which was associated with a trait through two
338 distinctive genetic regulatory mechanisms (see **Supplementary Fig. 19** for a typical example).
339 Taken all together, our results show the distinctiveness of most sQTLs in mediating the polygenic
340 effects for complex traits and demonstrate the substantial gain of power in gene discovery for
341 complex traits by integrating sQTL data into GWAS.

342

343 **Discussion**

344 In this study, we have generated, to date, the largest catalog of genetic variants associated with a
345 broad spectrum of alternative splicing events in human brain, significantly expanding our
346 understanding of genetic control of RNA splicing. By comparing the sQTLs with the eQTLs
347 identified in this study, we showed that ~68% of the sQTLs are distinct from the eQTLs,
348 propelling research in the future to increase sample sizes for sQTL detection. By integrating
349 sQTLs with GWAS for ten brain-related traits, ~68% of the trait-associated genes identified
350 through sQTLs could not be discovered through eQTLs, demonstrating the unique role of sQTLs
351 in dissecting the genetic architecture underpinning complex trait variation. Moreover, the trait-
352 associated genes identified through sQTLs in this study provide important leads for further

353 mechanistic work to elucidate their functions in the development of the brain-related traits and
354 disorders.

355

356 We applied both THISTLE and LeafCutter to the PsychENCODE data for sQTL discovery. These
357 two methods complement to each other in several ways. First, deriving transcript abundance
358 from short-read RNA-seq data in a transcript-based method like THISTLE is more challenging
359 than directly quantifying intron excision ratios from split reads in an event-based method like
360 LeafCutter, because each sequencing read only represents a small part of the transcript isoform,
361 and alternative transcripts often have substantial overlaps. Second, local splicing events may not
362 be able to capture the complete landscape of gene regulation but have more ability to detect
363 sQTLs. Third, although LeafCutter identified a larger number of sGenes than THISTLE (5,629 vs.
364 4,881; **Supplementary Fig. 11a**), the number of genome-wide significant trait-associated genes
365 discovered by SMR through the THISTLE sQTLs is larger than that through the LeafCutter sQTLs
366 (103 vs. 74; **Supplementary Table 2 & Supplementary Fig. 21**). This is partly due to that the
367 Bonferroni correction in the SMR analysis using the LeafCutter sQTL summary data ($m = \sim 18,500$
368 tests for each trait) was more conservative than that using the THISTLE sQTL summary data (m
369 $= \sim 4,800$). These results implicate that THISTLE and LeafCutter are complementary and should
370 be used together for sQTL detection in practice. Of note that THISTLE requires only summary-
371 level isoform-eQTL data, a feature that can facilitate sQTL analyses using data from multiple
372 cohorts when sharing individual-level data is prohibitive.

373

374 With individual-level bulk RNA-seq and SNP genotype data on 1,073 prefrontal cortex samples
375 from the PsychENCODE consortium, we identified 632,481 sQTLs associated with splicing levels
376 of 7,491 genes and 663,326 eQTLs associated with expression levels of 9,524 genes at p -value $<$
377 5×10^{-8} . We have developed an online tool to query the sQTLs and eQTLs with $FDR < 0.01$. The
378 full sQTL and eQTL summary data without p -value thresholding will also be publicly available
379 upon publication. These data sets may be useful for future studies to understand the molecular
380 mechanisms underpinning the genetic regulation of splicing in brain, to identify functional genes
381 and variants for other brain-related phenotypes (e.g., cognitive and neuroimaging phenotypes
382 and neuropsychiatric disorders), and to improve genomic risk prediction, etc. Our study also
383 informs future analyses to quantify the relationship between sGenes and eGenes, between sQTLs
384 and eQTLs, or more generally between any two types of molecular QTLs. We show that the low
385 to moderate sGene-eGene overlap as observed in previous studies^{7,20} are due to small n because
386 the overlap is a function of n (**Fig. 2**). For example, only $\sim 8\%$ of sGenes are eGenes when $n = 100$,
387 and this proportion increases to 64% as n increases to 1,073. If every gene has sQTL and eQTL
388 (the latter has almost been proved to be true⁶), we predict that this proportion will eventually

389 approach to 1 with ever increasing n . Nevertheless, even if this is the case, the underlying sQTL
390 and eQTL causal variants are largely distinct. In this study, we estimated that the sQTL and eQTL
391 causal variants are shared for only $\sim 32\%$ of the sGenes (**Fig. 2**). This 32% may even be an
392 overestimation because the limited power of COLOC in distinguishing close linkage (a scenario
393 where the sQTL and eQTL causal variants are distinct but in high LD) from sharing (a scenario
394 where the sQTL and eQTL share the same causal variant(s)), especially when n is not sufficiently
395 large. In an extreme scenario where the sQTL and eQTL causal variants are in perfect LD, there is
396 no power to distinguish linkage from sharing.

397

398 By integrating the sQTLs with GWAS data, we confirmed that sQTLs were enriched for genetic
399 variants associated with complex traits, as in previous studies^{19,20,22,23,60}. We note that a previous
400 study by Hormozdiari et al.²¹ shows that sQTLs do not have a significant contribution to disease
401 heritability in a joint analysis of five BLUEPRINT molecular QTLs, namely eQTL, sQTL, H3K27ac
402 histone QTL (hQTL), H3K4me1 hQTL, and DNA methylation QTL (mQTL)²¹, which does not seem
403 consistent with our results. The discrepancy is likely due to the small sample sizes of the
404 molecular QTL data in Hormozdiari et al. ($n = \sim 200$ per cell type). In the present study, we have
405 provided multiple lines of evidence that the role of sQTLs in complex trait variation is largely
406 distinct from that of eQTLs. First, as discussed above, only up to 32% of the sQTL and eQTL signals
407 are shared (**Fig. 2**). Second, sQTLs contributed significantly to trait heritability conditional on
408 eQTLs (**Fig. 4c**). Third, $\sim 68\%$ of the trait-associated genes detected by integrating GWAS data
409 with the sQTLs were not detected by using the eQTLs (**Fig. 4e**). Our results also imply that the
410 contribution of sQTLs in mediating polygenic effects is comparable to that of eQTLs. For example,
411 we observed that the inflation in GWAS test-statistics at the sQTLs was indistinguishable from
412 that at the eQTLs (**Supplementary Fig. 13**), that heritability enrichment of the sQTLs was similar
413 to that of the eQTLs (median fold enrichment of 1.85 for sQTLs vs. 1.73 for eQTLs across the ten
414 traits), that the proportion of mediated heritability for the traits through cis-sQTLs was on a par
415 with that through cis-eQTLs (median h_{med}^2/h_g^2 of 0.11 for sQTLs vs. 0.13 for eQTLs across the ten
416 traits), and that the number of trait-associated genes identified through the sQTLs was of the
417 same magnitude as that through the eQTLs (107 vs. 149). Hence, large-scale sQTL studies in blood
418 and other tissues or specific cell types are urgently needed to make more discovery of sQTLs to
419 improve our understanding of genetic regulation of RNA splicing and to facilitate the translation
420 of GWAS maps into mechanisms.

421

422 Our study has several limitations. First, the THISTLE method relies on accurate assembly and
423 quantification of mRNA isoforms, which are limited by the short-read RNA-seq technology but
424 could be improved in the future by long-read RNA-seq that allows direct identification of isoforms.

425 For example, the Pacific Biosciences (PacBio) isoform sequencing has successfully identified
426 many novel transcripts and alternative splicing events in a range of tissues and cell types with
427 well-characterized transcriptomes⁶¹. Second, the eQTLs and sQTLs identified in this study only
428 explain a small fraction of the GWAS loci for the ten traits, and the majority of the loci remain to
429 be unexplained, which is likely due to the insufficient sample size of the transcriptome data and
430 mechanisms beyond genetic control of transcriptional regulation. Thus, a comprehensive set of
431 genetic variants associated with other molecular phenotypes (e.g., chromatin accessibility, DNA
432 methylation, histone modification, and protein abundance) from studies with very large sample
433 sizes is essential to unveil the molecular basis of polygenic trait variation. Third, we have
434 identified a large number of genetic variants associated with mRNA abundance and pre-mRNA
435 splicing, but the causative functional variants are mostly unknown. Fine-mapping analyses and
436 functional experiments are required to pinpoint the causal variants and elucidate the
437 mechanisms behind genetic control of transcriptional regulation and the following cascades of
438 events that lead to phenotypic changes. Last but not least, we focus this study only on brain
439 transcriptomic data and brain-related traits. Although it is likely that our conclusions can be
440 generalised to other tissues and complex traits given the strong correlation of eQTL effects across
441 tissues^{7,36,62} and the extremely similar pattern of genetic architecture of brain phenotypes to
442 other complex traits, the generalisation of our conclusions needs to be confirmed by well-
443 powered studies in the future. Despite these limitations, our study develops a powerful and
444 flexible transcript-based sQTL method, generates the largest set of brain cis-sQTL summary data
445 to date (with an online tool for data query), demonstrates an analysis paradigm to assess the
446 relationship between two types of molecular QTLs, and provides multiple lines of evidence that
447 most sQTLs are distinct from eQTLs.

448

449 **Methods**

450 **Testing for heterogeneity between isoform-eQTL effects (THISTLE)**

451 Let m be the number of transcript isoforms for a gene, \mathbf{y}_j be a vector of mRNA abundances across
452 n individuals for isoform j , and \mathbf{x} be a vector of genotypes of a variant. If we define $\hat{\mathbf{b}} =$
453 $\{\hat{b}_1, \dots, \hat{b}_j, \dots, \hat{b}_m\}$ with \hat{b}_j being the estimated isoform-eQTL effect for isoform j , we have $\hat{\mathbf{b}} \sim N(\mathbf{b}, \mathbf{S})$
454 with \mathbf{S} being the variance-covariance matrix of $\hat{\mathbf{b}}$. The difference between \hat{b}_j and \hat{b}_k ($j \neq k$) can
455 be estimated as:

$$456 \hat{d}_{jk} = \hat{b}_j - \hat{b}_k$$

457 If we define $\hat{\mathbf{d}} = \{\hat{d}_{jk}\}_{j,k \in \{1, \dots, m\}, j < k}$, we have $\hat{\mathbf{d}} \sim N(\mathbf{d}, \mathbf{V})$ with \mathbf{V} being the variance-covariance
458 matrix of $\hat{\mathbf{d}}$. The variance of \hat{d}_{jk} over repeated experiments (i.e., a diagonal element of \mathbf{V}) can be
459 written as

460
$$var(\hat{d}_{jk}) = var(\hat{b}_j - \hat{b}_k) = var(\hat{b}_j) + var(\hat{b}_k) - 2cov(\hat{b}_j, \hat{b}_k)$$

461 The covariance between \hat{d}_{jk} and \hat{d}_{gh} over repeated experiments (i.e., an off-diagonal element of
462 \mathbf{V}) can be written as

463
$$cov(\hat{d}_{jk}, \hat{d}_{gh}) = cov(\hat{b}_j - \hat{b}_k, \hat{b}_g - \hat{b}_h) = cov(\hat{b}_j, \hat{b}_g) - cov(\hat{b}_j, \hat{b}_h) - cov(\hat{b}_k, \hat{b}_g) + cov(\hat{b}_k, \hat{b}_h)$$

464 The covariance between \hat{b}_j and \hat{b}_k can be estimated as $cov(\hat{b}_j, \hat{b}_k) = \theta_{jk}S_jS_k$, where S_j^2 and S_k^2

465 are the variances of \hat{b}_j and \hat{b}_k , respectively, and θ_{jk} is the correlation between \hat{b}_j and \hat{b}_k . Since the

466 isoform abundances are measured on the same set of individuals, \hat{b}_j and \hat{b}_k are likely to be

467 correlated (i.e., $\theta_{jk} \neq 0$). We have shown in our previous studies^{11,36} that $\theta_{jk} \approx r_p\rho$, where $\rho =$

468 $\frac{n_s}{\sqrt{n_j n_k}}$ measures the sample overlap with n_j and n_k being the sample sizes of isoforms j and k ,

469 respectively, and n_s being the number of overlapping individuals between two isoforms, and r_p

470 is the Pearson correlation of mRNA abundances between two isoforms in the overlapping sample.

471 If the individual-level data are unavailable, θ_{jk} can be approximated by the Pearson correlation

472 of the estimated isoform-eQTL effects between two isoforms across the “null” SNPs (e.g., $P_{\text{isoform-eQTL}} > 0.01$)^{11,36}. Under the null hypothesis of no sQTL effect (i.e., $\mathbf{d} = \mathbf{0}$), there is no heterogeneity

473 in isoform-eQTL effect between the isoforms. We have a vector of standard normal variables $\mathbf{z}_d =$

474 $\{z_{d(jk)}\}_{j,k \in \{1, \dots, m\}, j < k}$ with $z_{d(jk)} = \hat{d}_{jk} / \sqrt{var(\hat{d}_{jk})}$ and $\mathbf{z}_d \sim N(\mathbf{0}, \mathbf{R})$, where \mathbf{R} is the correlation

475 matrix with $cor(z_{d(jk)}, z_{d(gh)}) = cov(\hat{d}_{jk}, \hat{d}_{gh}) / \sqrt{var(\hat{d}_{jk})var(\hat{d}_{gh})}$. To test against the null

476 hypothesis ($\mathbf{d} = \mathbf{0}$), we use an approximate multivariate approach, similar to the heterogeneity

477 in dependent instruments (HEIDI) test¹¹, to construct a THISTLE test statistic $T_{THISTLE} = \mathbf{z}_d \mathbf{I} \mathbf{z}_d^T$,

478 with \mathbf{I} being an identify matrix. This distribution can be approximated by the Satterthwaite or

479 Saddlepoint method^{63,64}. We named the method THISTLE and has been implemented it in the

480 OSCA³⁷ software package (<https://yanglab.westlake.edu.cn/software/osca/#THISTLE>).

481

482

483 **Data used in this study**

484 This study is approved by the Ethics Committee of Westlake University (approval number:

485 20200722YJ001). We used the genotype and RNA-seq data in dorsolateral prefrontal cortex

486 (DLPFC) tissue from the PsychENCODE consortium. We processed the PsychENCODE data from

487 six cohorts: UCLA-ASD, BrainGVEX, the Lieber Institute for Brain Development (LIBD), the

488 CommonMind Consortium (CMC), the CMC’s NIMH Human Brain Collection Core (HBCC), and

489 Bipseq of Bipolar cohorts (Bipseq). Generation of the genotype and RNA-seq data and imputation

490 of the genotype data have been detailed elsewhere^{34,35}. RNA-seq data were filtered with a

491 conventional QC process, e.g., retaining individuals with more than 10 million total reads and RNA

492 integrity number (RIN) > 5.5. Genotyped and imputed SNP data were filtered with standard QC

493 criteria in each cohort separately using PLINK2⁶⁵, i.e., genotyping rate > 0.95, Hardy-Weinberg
494 equilibrium test p-value > 1×10^{-6} , MAF > 0.01, and imputation information score > 0.3.
495 Individuals with duplicated sample ID across six cohorts, individuals with non-European ancestry
496 as inferred from principal component analysis (**Supplementary Note; Supplementary Fig. 5**),
497 or individuals with SNP-derived genetic relatedness > 0.05 were further removed. The resulting
498 merged consensus data from six cohorts was filtered again according to the QC criteria described
499 above. Finally, 1,073 unrelated individuals of European ancestry including 348 subjects with
500 schizophrenia (SCZ), 164 subjects with bipolar disorder (BIP), 17 subjects with autism spectrum
501 disorder (ASD), and 544 controls on 3,777,685 genotyped or imputed common SNPs were
502 retained for further analysis (**Supplementary Fig. 6**).

503

504 We included in this study ten brain-related phenotypes, i.e., intelligence (IQ), educational
505 attainment (EA), smoking initiation (SmkInit), schizophrenia (SCZ), Alzheimer's disease (AD),
506 Parkinson's disease (PD), autism spectrum disorder (ASD), insomnia, bipolar disorder (BIP), and
507 major depression (MD). GWAS summary statistics for IQ ($n = 269,867$)⁴³, EA ($n = 766,345$)⁴⁴,
508 SmkInit (311,629 cases and 321,173 controls)⁴⁵, SCZ (40,675 cases and 64,643 controls)⁴⁶, AD
509 (71,880 cases and 315,120 controls)⁴⁷, PD (33,674 cases and 449,056 controls)⁵², ASD (18,381
510 cases and 27,969 controls)⁴⁸, insomnia ($n = 386,533$)⁴⁹, BIP (20,352 cases and 31,358 controls)⁵⁰,
511 and MD (170,756 cases and 329,443 controls)⁶⁶ were from the latest published studies
512 respectively (**Supplementary Table 1**).

513

514 **Identification of cis-sQTLs using THISTLE and LeafCutter**

515 To identify sQTLs using THISTLE, we first quantified the isoform abundance from short-read
516 RNA-seq data. Detailed procedures of isoform quantification are available elsewhere^{34,35}. In brief,
517 RNA-seq reads of each of the 1,073 PsychENCODE samples were aligned to GRCh37 genome
518 assembly by STAR v2.4.0f1⁶⁷, and isoform abundances, measured by transcripts per million
519 (TPM), were quantified using the RSEM software package⁶⁸. After filtering out isoforms with low
520 expression (e.g., isoform-level TPM < 0.1 in more than 20% of the subjects) or genes with less
521 than two annotated isoforms, a total of 77,511 isoforms of 13,366 genes were retained.
522 VariancePartition³⁸ was then used to determine which known covariates (both biological and
523 technical) should be used to adjust isoform expression using mixed-effects modeling.
524 Probabilistic Estimation of Expression Residuals (PEER) method³⁹ was used to compute a set of
525 latent covariates to capture hidden batch effects and other unknown technical and biological
526 factors. We then performed isoform-eQTL analysis in OSCA³⁷ using the adjusted isoform
527 expression data with the top 5 genetic principal components (PCs) and 50 PEER factors fitted as
528 covariates, which was then followed by the THISTLE analysis to test for sQTLs. It is of note that

529 we used the same data for cis-eQTL analysis, which was essentially a meta-analysis of isoform-
530 eQTL summary data across isoforms for each gene.

531

532 To identify sQTLs using LeafCutter³³, the alignments from STAR⁶⁷ were used as input into
533 LeafCutter to identify intron clusters with at least one differentially excised intron by jointly
534 modelling intron clusters using a Dirichlet-multinomial generalized linear model³³. Default
535 LeafCutter parameters were used to detect clusters, i.e., 50 split reads across all individuals were
536 required to support each cluster, introns up to 500 kb were included, and introns used in less
537 than 40% of individuals or with almost no variation were filtered out. The intron excision ratio
538 (the ratio of intron defining reads to the total number of reads of an intron cluster) was
539 standardized across individuals for each intron and quantile-normalized across introns. In total,
540 202,344 introns in 46,868 intronic excision clusters were identified with 46,199 clusters uniquely
541 mapped to 16,518 genes and the remaining clusters either mapped to more than one gene or not
542 located in the gene body. For covariate adjustment, we first identified known biological and
543 technical covariates that were significantly associated with intron excision ratio by
544 variancePartition³⁸ and then adjusted intron excision ratio by a mixed linear model, i.e., intron
545 excision ratios \sim covariates + (1|Donor), where each intron was regressed on the covariates with
546 donor (individual) fitted as random effects. PEER³⁹ was used to compute a set of latent covariates
547 to capture unknown technical and biological factors. We then used FastQTL⁴⁰ to perform an
548 association test for sQTL based on the following model: adjusted intron excision ratio \sim genotype
549 + PC₁ + ... + PC₅ + PEER₁ + ... + PEER₁₅ (**Supplementary Fig. 8**), where PC₁ + ... + PC₅ represent
550 the top 5 genetic PCs and PEER₁ + ... + PEER₁₅ represent 15 PEER factors.

551

552 **Enrichment of the sQTLs and eQTLs for functional annotations**

553 To test if the sQTLs and eQTLs are functionally enriched, we annotated the top cis-sQTL or cis-
554 eQTL SNPs using the annotations from SnpEff⁴¹ (e.g., splice region, intronic, and upstream) and
555 14 main functional annotations (e.g., promoter and enhancer) in 10 brain samples from the NIH
556 Roadmap Epigenomics Mapping Consortium (REMC)⁴². More specifically, we annotated 9,106
557 unique top-associated cis-eQTL SNPs and 7,187 unique top-associated cis-sQTL SNPs in different
558 functional annotations based on their physical positions and quantified the proportion of eQTLs
559 or sQTLs in each functional category. To ameliorate ascertainment bias, we sampled at random
560 the same number of cis-SNPs (i.e., SNPs included in the sQTL or eQTL analysis) as “control SNPs”,
561 with their MAF matched with the top sQTL (or eQTL) SNPs. This sampling procedure was
562 repeated 1,000 times. We computed the fold enrichment in each functional annotation category
563 as the ratio of the proportion of sQTL (or eQTL) SNPs in a functional category over the mean
564 proportion of the control SNPs in the category across 1,000 replicates. The sampling variance of

565 the fold enrichment can be calculated approximately by the Delta method⁶⁹ (**Supplementary**
566 **Note**).

567

568 **Enrichment of the sQTLs and eQTLs for trait heritability**

569 The stratified LD score regression (S-LDSC)^{55,56} was used to quantify the enrichment of
570 heritability attributable to the sQTLs and eQTLs (when fitted together with 53 other functional
571 categories in the “baseline model”) for the ten brain-related traits. Details of the “baseline model”
572 that consists of 53 functional categories can be found elsewhere⁵⁵. We created a binary
573 annotation for sQTLs and eQTLs respectively. In brief, we assigned an annotation value of 1 for
574 the most significant sQTL (or eQTL) SNP for each gene and a zero value for the remaining SNPs.
575 This resulted in an sQTL annotation category with 6,547 SNPs and an eQTL annotation category
576 with 8,192 SNPs with an overlap of 1,486 SNPs. LD scores of the SNPs were computed using SNP
577 genotype data of the individuals of European ancestry from the 1000 Genomes Project (Phase
578 3)⁷⁰ with a window size of 1 cM. Heritability enrichment of a category was computed as the
579 proportion of heritability explained by the category divided by the proportion of SNPs in the
580 category. Considering that SNPs in or near genes explain disproportionately more trait variation
581 than intergenic SNPs⁵⁷, we also computed the fold enrichment of heritability as the per-SNP
582 heritability for the top cis-sQTL (or cis-eQTL) SNPs divided by a mean of a distribution generated
583 by resampling MAF-matched cis-SNPs (i.e., the control SNPs mentioned above). Sampling
584 variance of the fold enrichment of heritability can be calculated approximately by the Delta
585 method⁶⁹ (**Supplementary Note**). Note that both the per-SNP heritability and annotation effect
586 size (τ) reported by S-LDSC are used to quantify the relevance of a functional category to the trait
587 variation⁵⁶, and the main difference between the two metrics lies in how the overlapping
588 annotations are dealt with.

589

590 **Author Contributions**

591 JY conceived and supervised the study. TQ and JY designed the experiment and developed the
592 method. TQ performed all the simulations and data analyses under the assistance or guidance
593 from YW, JZ, and JY. TQ, FZ and SL developed the application tools. TQ and JY wrote the
594 manuscript with the participation of all authors.

595

596 **Acknowledgements**

597 This research was partly supported by the Westlake Education Foundation, the Australian
598 National Health and Medical Research Council (1107258 and 1113400) and the Australian
599 Research Council (DP160101343 and FT180100186). We thank the Westlake University High-
600 Performance Computing Center for assistance in computing. This study uses data from the

601 PsychENCODE consortium (Synapse accession: syn4921369). A full list of acknowledgments to
602 this data set can be found in the Supplementary Note.

603

604 **Data Availability**

605 PsychENCODE data: <https://www.synapse.org/#!Synapse:syn4921369/files/>. Online tool for
606 querying the sQTLs and eQTLs with FDR < 0.01: <https://yanglab.westlake.edu.cn/qi-sqtl/>. The
607 full sQTL and eQTL summary data will be publicly available upon publication.

608

609 **Code Availability**

610 The computer code and documentation of THISTLE are available at

611 <https://yanglab.westlake.edu.cn/software/osca/#THISTLE>.

612

613 **URLs**

614 OSCA: <https://yanglab.westlake.edu.cn/software/osca>.

615 LeafCutter: <https://davidaknowles.github.io/leafcutter/index.html>.

616 FastQTL: <http://fastqtl.sourceforge.net/>

617 S-LDSC: <https://github.com/bulik/ldsc>.

618 SMR: <https://cnsgenomics.com/software/smr/#Overview>.

619 COLOC: <https://cran.r-project.org/web/packages/coloc/index.html>.

620 MESOC: <https://github.com/douglasyao/mesc>.

621

622 **References**

623 1 Visscher, P. M. *et al.* 10 years of GWAS discovery: biology, function, and translation. *The*
624 *American Journal of Human Genetics* **101**, 5-22 (2017).

625 2 Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association
626 studies, targeted arrays and summary statistics 2019. *Nucleic acids research* **47**, D1005-
627 D1012 (2018).

628 3 Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in
629 regulatory DNA. *Science* **337**, 1190-1195 (2012).

630 4 Ward, L. D. & Kellis, M. Interpreting noncoding genetic variation in complex traits and
631 human disease. *Nature biotechnology* **30**, 1095 (2012).

632 5 Morley, M. *et al.* Genetic analysis of genome-wide variation in human gene expression.
633 *Nature* **430**, 743-747 (2004).

634 6 Vösa, U. *et al.* Unraveling the polygenic architecture of complex traits using blood eQTL
635 meta-analysis. *bioRxiv*, 447367 (2018).

636 7 GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human
637 tissues. *Science* **369**, 1318-1330 (2020).

638 8 Liu, B., Gloudemans, M. J., Rao, A. S., Ingelsson, E. & Montgomery, S. B. Abundant
639 associations with gene expression complicate GWAS follow-up. *Nature genetics* **51**, 768-
640 769 (2019).

641 9 Giambartolomei, C. *et al.* Bayesian test for colocalisation between pairs of genetic
642 association studies using summary statistics. *PLoS genetics* **10**, e1004383 (2014).

643 10 Gusev, A. *et al.* Integrative approaches for large-scale transcriptome-wide association
644 studies. *Nature genetics* **48**, 245-252 (2016).

645 11 Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex
646 trait gene targets. *Nature genetics* **48**, 481-487 (2016).

647 12 Hormozdiari, F. *et al.* Colocalization of GWAS and eQTL signals detects target genes. *The
648 American Journal of Human Genetics* **99**, 1245-1260 (2016).

649 13 Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene
650 expression variation inferred from GWAS summary statistics. *Nature communications* **9**,
651 1-20 (2018).

652 14 Pavlides, J. M. W. *et al.* Predicting gene targets from integrative analyses of summary data
653 from GWAS and eQTL studies for 28 human complex traits. *Genome medicine* **8**, 84 (2016).

654 15 Mancuso, N. *et al.* Integrating gene expression with summary association statistics to
655 identify genes associated with 30 complex traits. *The American Journal of Human Genetics*
656 **100**, 473-487 (2017).

657 16 Huckins, L. M. *et al.* Gene expression imputation across multiple brain regions provides
658 insights into schizophrenia risk. *Nature genetics* **51**, 659-674 (2019).

659 17 Porcu, E. *et al.* Mendelian randomization integrating GWAS and eQTL data reveals genetic
660 determinants of complex and clinical traits. *Nature communications* **10**, 1-12 (2019).

661 18 Yao, D. W., O'Connor, L. J., Price, A. L. & Gusev, A. Quantifying genetic effects on disease
662 mediated by assayed gene expression levels. *Nature Genetics*, 1-8 (2020).

663 19 Walker, R. L. *et al.* Genetic Control of Expression and Splicing in Developing Human Brain
664 Informs Disease Mechanisms. *Cell* **179**, 750-771. e722 (2019).

665 20 Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease. *Science*
666 **352**, 600-604 (2016).

667 21 Hormozdiari, F. *et al.* Leveraging molecular quantitative trait loci to understand the
668 genetic architecture of diseases and complex traits. *Nature genetics* **50**, 1041-1047 (2018).

669 22 Raj, T. *et al.* Integrative transcriptome analyses of the aging brain implicate altered
670 splicing in Alzheimer's disease susceptibility. *Nature genetics* **50**, 1584 (2018).

- 671 23 Li, Y. I., Wong, G., Humphrey, J. & Raj, T. Prioritizing Parkinson's disease genes using
672 population-scale transcriptomic data. *Nature communications* **10**, 1-10 (2019).
- 673 24 Park, E., Pan, Z., Zhang, Z., Lin, L. & Xing, Y. The expanding landscape of alternative splicing
674 variation in human populations. *The American Journal of Human Genetics* **102**, 11-26
675 (2018).
- 676 25 Katz, Y., Wang, E. T., Airoidi, E. M. & Burge, C. B. Analysis and design of RNA sequencing
677 experiments for identifying isoform regulation. *Nature methods* **7**, 1009 (2010).
- 678 26 Trapnell, C. *et al.* Differential analysis of gene regulation at transcript resolution with
679 RNA-seq. *Nature biotechnology* **31**, 46 (2013).
- 680 27 Lappalainen, T. *et al.* Transcriptome and genome sequencing uncovers functional
681 variation in humans. *Nature* **501**, 506-511 (2013).
- 682 28 Battle, A. *et al.* Characterizing the genetic basis of transcriptome diversity through RNA-
683 sequencing of 922 individuals. *Genome research* **24**, 14-24 (2014).
- 684 29 Monlong, J., Calvo, M., Ferreira, P. G. & Guigó, R. Identification of genetic variants
685 associated with alternative splicing using sQTLseekeR. *Nature communications* **5**, 4698
686 (2014).
- 687 30 Shen, S. *et al.* rMATS: robust and flexible detection of differential alternative splicing from
688 replicate RNA-Seq data. *Proceedings of the National Academy of Sciences* **111**, E5593-
689 E5601 (2014).
- 690 31 Ongen, H. & Dermitzakis, E. T. Alternative splicing QTLs in European and African
691 populations. *The American Journal of Human Genetics* **97**, 567-575 (2015).
- 692 32 Vaquero-Garcia, J. *et al.* A new view of transcriptome complexity and regulation through
693 the lens of local splicing variations. *elife* **5**, e11752 (2016).
- 694 33 Li, Y. I. *et al.* Annotation-free quantification of RNA splicing using LeafCutter. *Nature*
695 *genetics* **50**, 151 (2018).
- 696 34 Gandal, M. J. *et al.* Transcriptome-wide isoform-level dysregulation in ASD, schizophrenia,
697 and bipolar disorder. *Science* **362** (2018).
- 698 35 Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the
699 human brain. *Science* **362** (2018).
- 700 36 Qi, T. *et al.* Identifying gene targets for brain-related traits using transcriptomic and
701 methylomic data from blood. *Nature communications* **9**, 2282 (2018).
- 702 37 Zhang, F. *et al.* OSCA: a tool for omic-data-based complex trait analysis. *Genome biology*
703 **20**, 107 (2019).
- 704 38 Hoffman, G. E. & Schadt, E. E. variancePartition: interpreting drivers of variation in
705 complex gene expression studies. *BMC bioinformatics* **17**, 483 (2016).

706 39 Stegle, O., Parts, L., Piipari, M., Winn, J. & Durbin, R. Using probabilistic estimation of
707 expression residuals (PEER) to obtain increased power and interpretability of gene
708 expression analyses. *Nature protocols* **7**, 500 (2012).

709 40 Ongen, H., Buil, A., Brown, A. A., Dermitzakis, E. T. & Delaneau, O. Fast and efficient QTL
710 mapper for thousands of molecular phenotypes. *Bioinformatics* **32**, 1479-1485 (2015).

711 41 Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide
712 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118;
713 iso-2; iso-3. *Fly* **6**, 80-92 (2012).

714 42 Kundaje, A. *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**,
715 317 (2015).

716 43 Savage, J. E. *et al.* Genome-wide association meta-analysis in 269,867 individuals
717 identifies new genetic and functional links to intelligence. *Nature genetics* **50**, 912 (2018).

718 44 Lee, J. J. *et al.* Gene discovery and polygenic prediction from a genome-wide association
719 study of educational attainment in 1.1 million individuals. *Nature genetics* **50**, 1112
720 (2018).

721 45 Liu, M. *et al.* Association studies of up to 1.2 million individuals yield new insights into the
722 genetic etiology of tobacco and alcohol use. *Nature genetics* **51**, 237 (2019).

723 46 Pardiñas, A. F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant
724 genes and in regions under strong background selection. *Nature genetics* **50**, 381 (2018).

725 47 Jansen, I. E. *et al.* Genome-wide meta-analysis identifies new loci and functional pathways
726 influencing Alzheimer's disease risk. (2019).

727 48 Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder.
728 *Nature genetics* **51**, 431 (2019).

729 49 Jansen, P. R. *et al.* Genome-wide analysis of insomnia in 1,331,010 individuals identifies
730 new risk loci and functional pathways. *Nature genetics* **51**, 394 (2019).

731 50 Stahl, E. A. *et al.* Genome-wide association study identifies 30 loci associated with bipolar
732 disorder. *Nature genetics* **51**, 793 (2019).

733 51 Howard, D. M. *et al.* Genome-wide meta-analysis of depression identifies 102 independent
734 variants and highlights the importance of the prefrontal brain regions. *Nature*
735 *neuroscience* **22**, 343 (2019).

736 52 Nalls, M. A. *et al.* Identification of novel risk loci, causal insights, and heritable risk for
737 Parkinson's disease: a meta-analysis of genome-wide association studies. *The Lancet*
738 *Neurology* **18**, 1091-1102 (2019).

739 53 Yang, J. *et al.* Genomic inflation factors under polygenic inheritance. *European Journal of*
740 *Human Genetics* **19**, 807-812, doi:10.1038/ejhg.2011.39 (2011).

741 54 Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from
742 polygenicity in genome-wide association studies. *Nature genetics* **47**, 291 (2015).

743 55 Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide
744 association summary statistics. *Nature genetics* **47**, 1228 (2015).

745 56 Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits
746 shows action of negative selection. *Nature genetics* **49**, 1421 (2017).

747 57 Yang, J. *et al.* Genome partitioning of genetic variation for complex traits using common
748 SNPs. *Nature genetics* **43**, 519 (2011).

749 58 Alinaghi, S. *et al.* Expression analysis and genotyping of DGKZ: a GWAS-derived risk gene
750 for schizophrenia. *Molecular biology reports*, 1-7 (2019).

751 59 Pers, T. H. *et al.* Comprehensive analysis of schizophrenia-associated loci highlights ion
752 channel pathways and biologically plausible candidate causal genes. *Human molecular*
753 *genetics* **25**, 1247-1254 (2016).

754 60 Takata, A., Matsumoto, N. & Kato, T. Genome-wide identification of splicing QTLs in the
755 human brain and their enrichment among schizophrenia-associated loci. *Nature*
756 *communications* **8**, 14519 (2017).

757 61 Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves variant
758 detection and assembly of a human genome. *Nature biotechnology*, 1-8 (2019).

759 62 Liu, X. *et al.* Functional architectures of local and distal regulation of gene expression in
760 multiple human tissues. *The American Journal of Human Genetics* **100**, 605-616 (2017).

761 63 Kuonen, D. Miscellanea. Saddlepoint approximations for distributions of quadratic forms
762 in normal variables. *Biometrika* **86**, 929-935 (1999).

763 64 Davies, R. B. Numerical inversion of a characteristic function. *Biometrika* **60**, 415-417
764 (1973).

765 65 Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
766 datasets. *Gigascience* **4**, s13742-13015-10047-13748 (2015).

767 66 Wray, N. R. *et al.* Genome-wide association analyses identify 44 risk variants and refine
768 the genetic architecture of major depression. *Nature genetics* **50**, 668 (2018).

769 67 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).

770 68 Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or
771 without a reference genome. *BMC bioinformatics* **12**, 323 (2011).

772 69 Lynch, M. & Walsh, B. *Genetics and analysis of quantitative traits*. Vol. 1 (Sinauer
773 Sunderland, MA, 1998).

774 70 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092
775 human genomes. *Nature* **491**, 56 (2012).

776

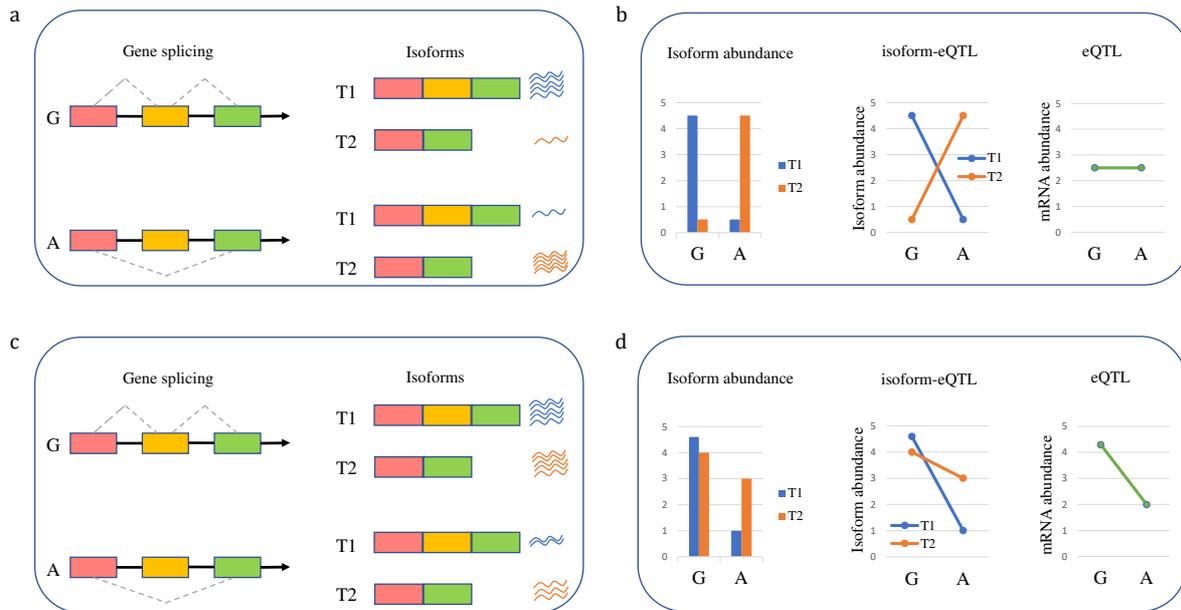


Figure 1. Schematic of THISTLE sQTL analysis. In this toy example, a genetic variant with two alleles, G and A, is associated with a splicing event (e.g., exon skipping) in a gene with two isoforms, T_1 and T_2 . Panels **a**) and **b**): a schematic of THISTLE sQTL analysis in the absence of eQTL effect. In this scenario, individuals with G allele show higher mean transcription abundance for T_1 than T_2 , and individuals with A allele show higher mean abundance for T_2 than T_1 (**a**), meaning that the genetic variant is associated with the difference in mRNA abundance between the isoforms. In other words, there is a difference in the isoform-eQTL effect between isoforms T_1 and T_2 (**b**). However, there is no difference in overall gene expression between individuals with different alleles, meaning that this genetic variant is not an eQTL. Panels **c**) and **d**): a schematic of THISTLE sQTL analysis in the presence of eQTL effect. In this scenario, individuals with G allele show similar transcription abundance between T_1 and T_2 , and individuals with A allele show lower abundance for T_1 than T_2 (**c**). The isoform-eQTL effect for T_1 is different from that for T_2 albeit in the same direction (**d**). In this case, there is a difference in overall gene expression between alleles G and A, meaning that this genetic variant is also an eQTL.

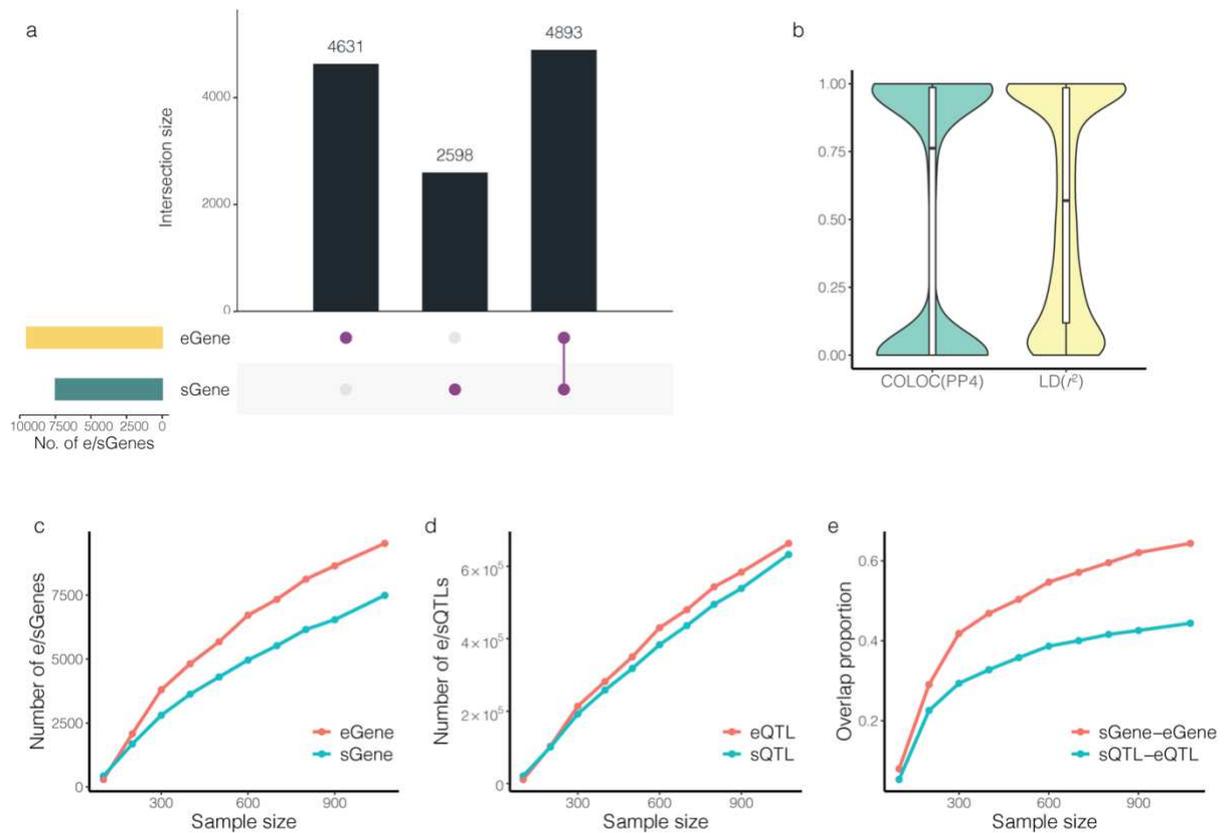


Figure 2. Relationship between sQTLs and eQTLs. **a)** Overlap between sGenes and eGenes in the PsychENCODE data. **b)** LD r^2 or COLOC PP4 between the top cis-sQTL SNP and the top cis-eQTL SNP for the 4,893 overlapping genes. **c)** The number of eGenes (or sGenes) discovered as a function of sample size. **d)** The number of eQTLs (or sQTLs) discovered as a function of sample size. **e)** The overlap between sGenes and eGenes (or between sQTLs and eQTLs) as a function of sample size. sQTL-eQTL overlap is defined as the proportion of sGenes for which the top sQTL is a significant eQTL SNP for the same gene.

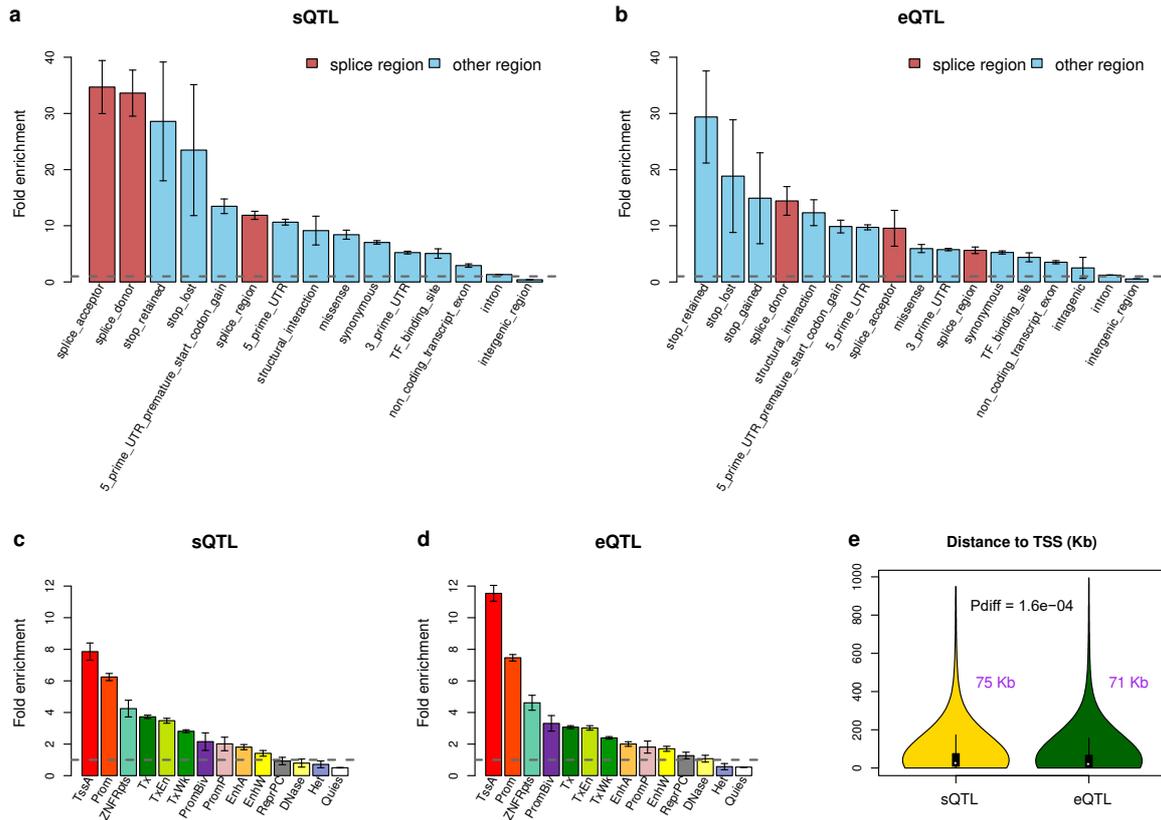


Figure 3. Enrichment of the top cis-sQTL or cis-eQTL SNPs for functional annotations from SnpEff (**a, b**) and REMC (**c, d**). Fold enrichment is computed by comparing the top cis-sQTL (or cis-eQTL) SNPs in a functional category with the control SNPs. Each error bar represents the 95% confidence interval around an estimate. The grey dashed line represents no enrichment. **e**) Distance of the top cis-sQTL (or cis-eQTL) SNP to TSS of the gene. The text in purple represents the median across genes, and P_{diff} is computed by a t -test for a mean difference between the two distributions.

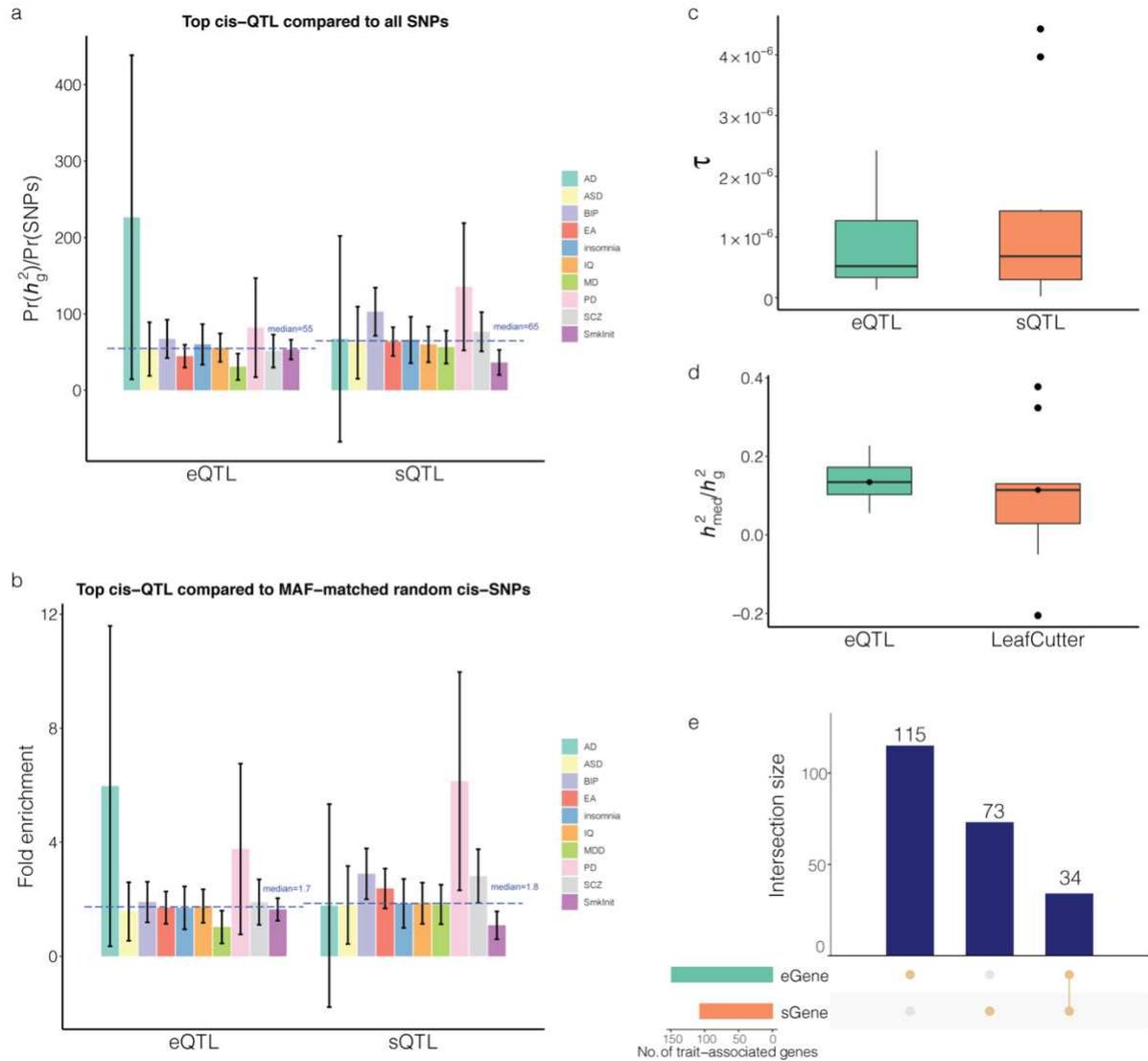


Figure 4. Enrichment of the top cis-sQTL or cis-eQTL SNPs for heritability of the ten brain-related traits. In panel **a**), heritability enrichment is defined as a ratio of the proportion of heritability explained by the SNPs in query (estimated by S-LDSC) to the proportion of the SNPs. In panel **b**), heritability enrichment is defined as a ratio of the proportion of heritability explained by the SNPs in query to that by the control SNPs. In panel **c**), annotation effect size (τ) is used to assess the contribution of the top cis-sQTL SNPs to heritability when fitted jointly with the top cis-eQTL SNPs. Panel **d**) shows the proportion of mediated heritability (h_{med}^2/h_g^2) by the cis-sQTL (or cis-eQTL) SNPs, estimated by MESC. In panels **a**) and **b**), each error bar represents the 95% confidence interval of an estimate, and the blue dashed line represents the median value across traits. Only the sQTLs discovered by LeafCutter were included in the MESC analysis. Panel **e**) shows the overlap of the trait-associated sGenes identified by SMR and COLOC with the trait-associated eGenes.

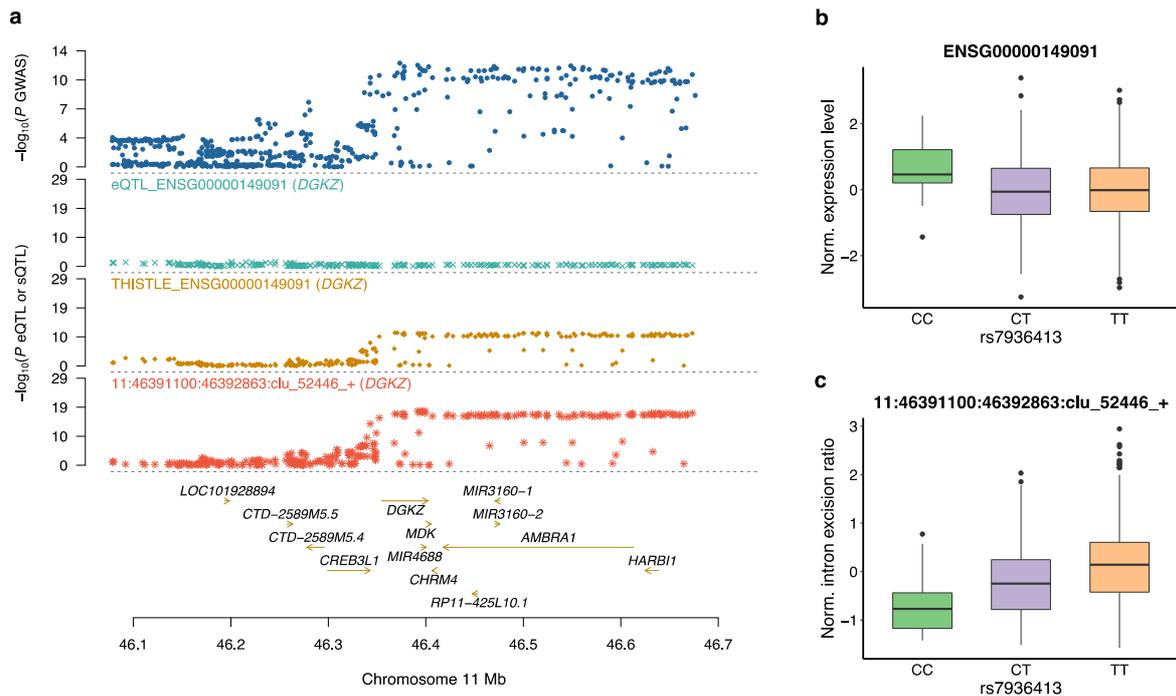


Figure 5. Association of *DGKZ* with SCZ through sQTLs but not eQTLs. **a)** The GWAS, sQTL, and eQTL p-values. The top plot shows $-\log_{10}(p\text{-values})$ of SNPs from the GWAS meta-analysis for SCZ. The second, third, and fourth plots show $-\log_{10}(p\text{-values})$ from the eQTL analysis, THISTLE sQTL analysis, and LeafCutter sQTL analysis (intron 11:46391100:46392863:clu_52446_+) for *DGKZ*, respectively. **b)** Association of rs7936413 (the top sQTL SNP) with the overall mRNA abundance of *DGKZ*. Each dot represents mRNA abundance of an individual. **c)** Association of rs7936413 with intron excision ratio of 11:46391100:46392863:clu_52446_+. Each dot represents intron excision ratio of an individual.

Figures

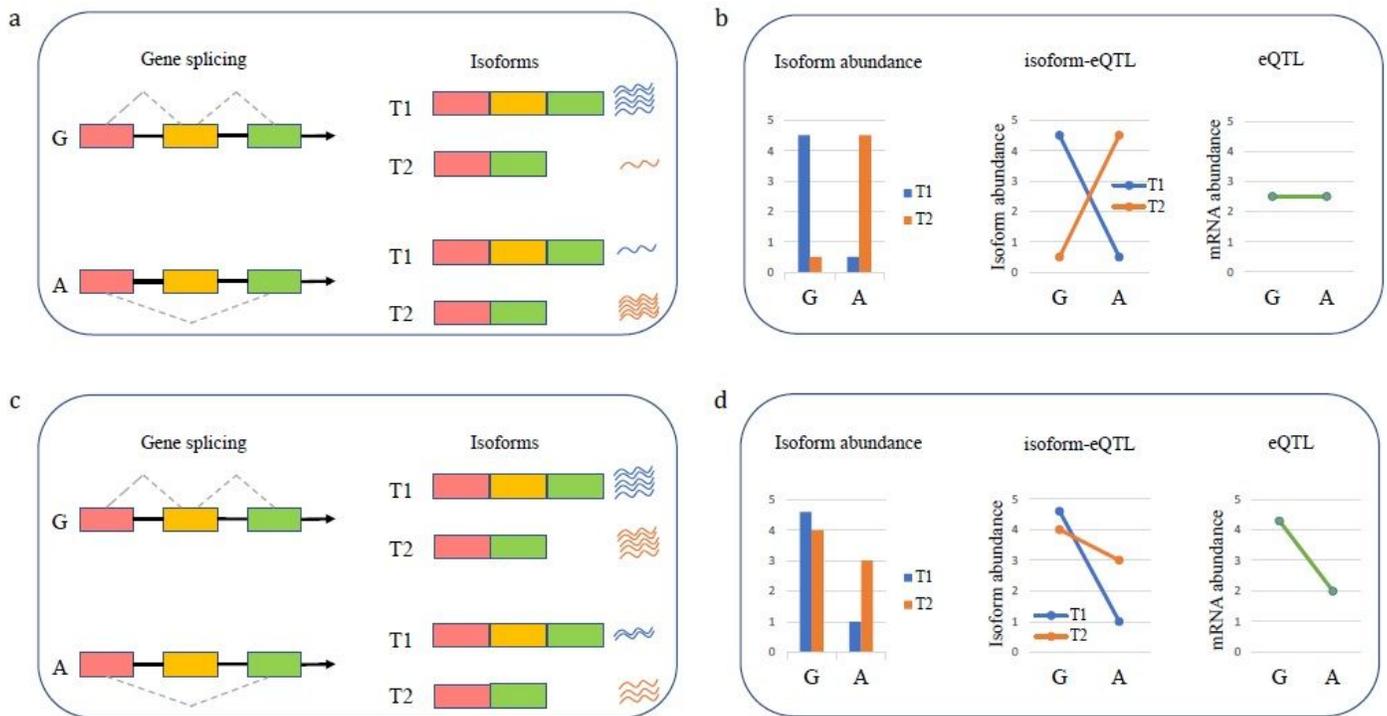


Figure 1

Schematic of THISTLE sQTL analysis. In this toy example, a genetic variant with two alleles, G and A, is associated with a splicing event (e.g., exon skipping) in a gene with two isoforms, T1 and T2. Panels a) and b): a schematic of THISTLE sQTL analysis in the absence of eQTL effect. In this scenario, individuals with G allele show higher mean transcription abundance for T1 than T2, and individuals with A allele show higher mean abundance for T2 than T1 (a), meaning that the genetic variant is associated with the difference in mRNA abundance between the isoforms. In other words, there is a difference in the isoform-eQTL effect between isoforms T1 and T2 (b). However, there is no difference in overall gene expression between individuals with different alleles, meaning that this genetic variant is not an eQTL. Panels c) and d): a schematic of THISTLE sQTL analysis in the presence of eQTL effect. In this scenario, individuals with G allele show similar transcription abundance between T1 and T2, and individuals with A allele show lower abundance for T1 than T2 (c). The isoform-eQTL effect for T1 is different from that for T2 albeit in the same direction (d). In this case, there is a difference in overall gene expression between alleles G and A, meaning that this genetic variant is also an eQTL.

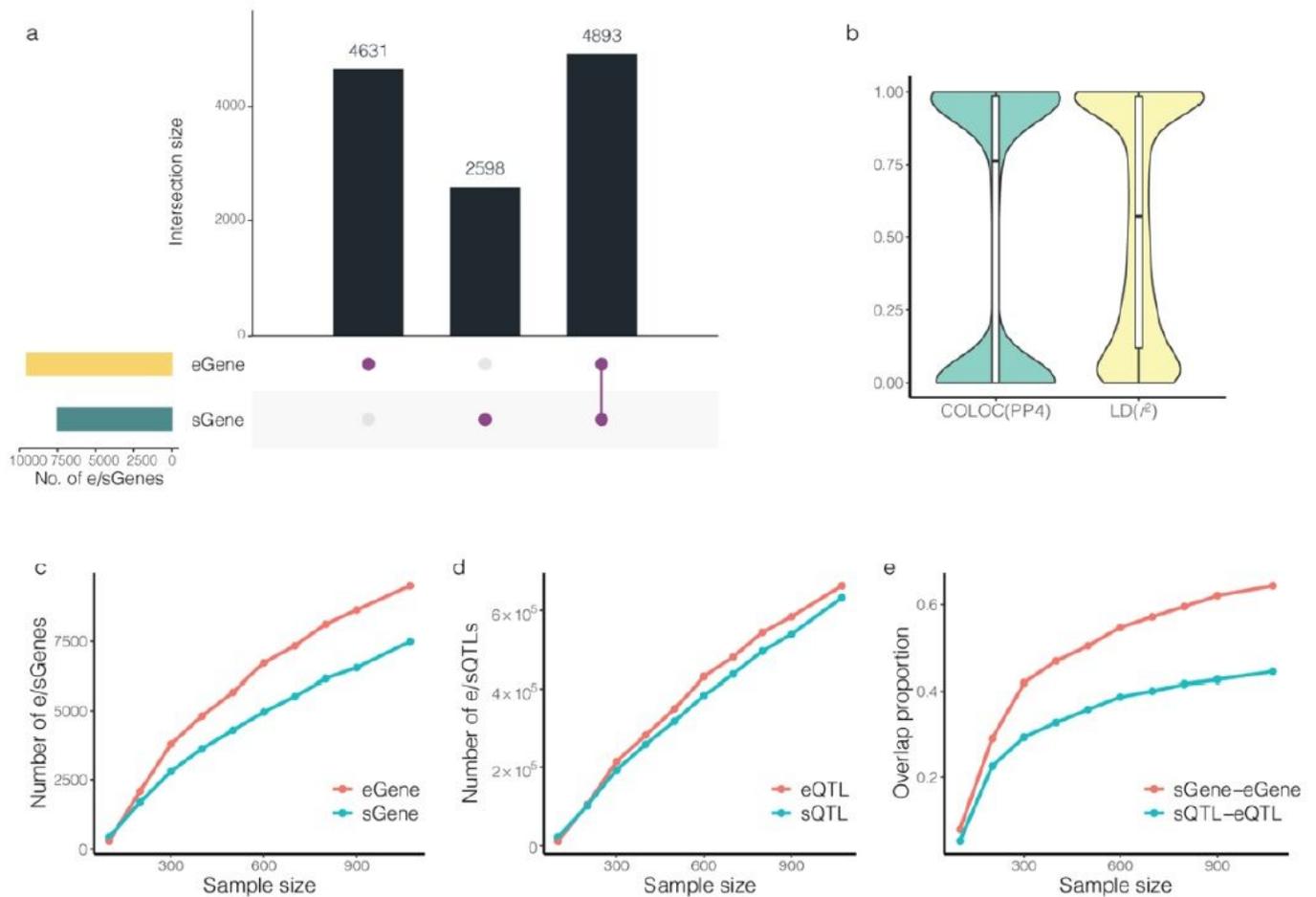


Figure 2

Relationship between sQTLs and eQTLs. a) Overlap between sGenes and eGenes in the PsychENCODE data. b) LD r^2 or COLOC PP4 between the top cis-sQTL SNP and the top cis-eQTL SNP for the 4,893 overlapping genes. c) The number of eGenes (or sGenes) discovered as a function of sample size. d) The number of eQTLs (or sQTLs) discovered as a function of sample size. e) The overlap between sGenes and eGenes (or between sQTLs and eQTLs) as a function of sample size. sQTL-eQTL overlap is defined as the proportion of sGenes for which the top sQTL is a significant eQTL SNP for the same gene.

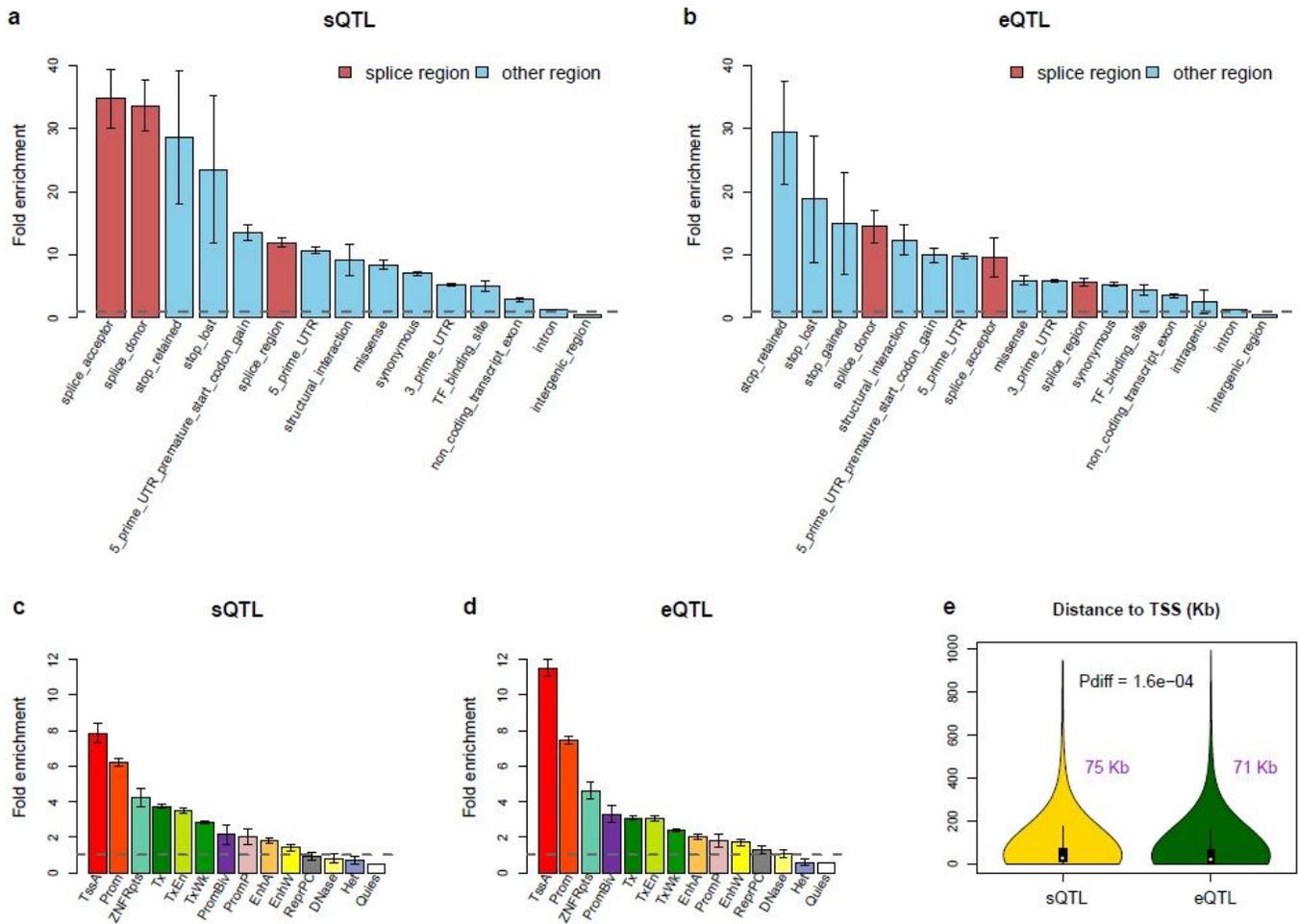


Figure 3

Enrichment of the top cis-sQTL or cis-eQTL SNPs for functional annotations from SnpEff (a, b) and REMC (c, d). Fold enrichment is computed by comparing the top cis-sQTL (or cis-eQTL) SNPs in a functional category with the control SNPs. Each error bar represents the 95% confidence interval around an estimate. The grey dashed line represents no enrichment. e) Distance of the top cis-sQTL (or cis-eQTL) SNP to TSS of the gene. The text in purple represents the median across genes, and Pdiff is computed by a t-test for a mean difference between the two distributions.

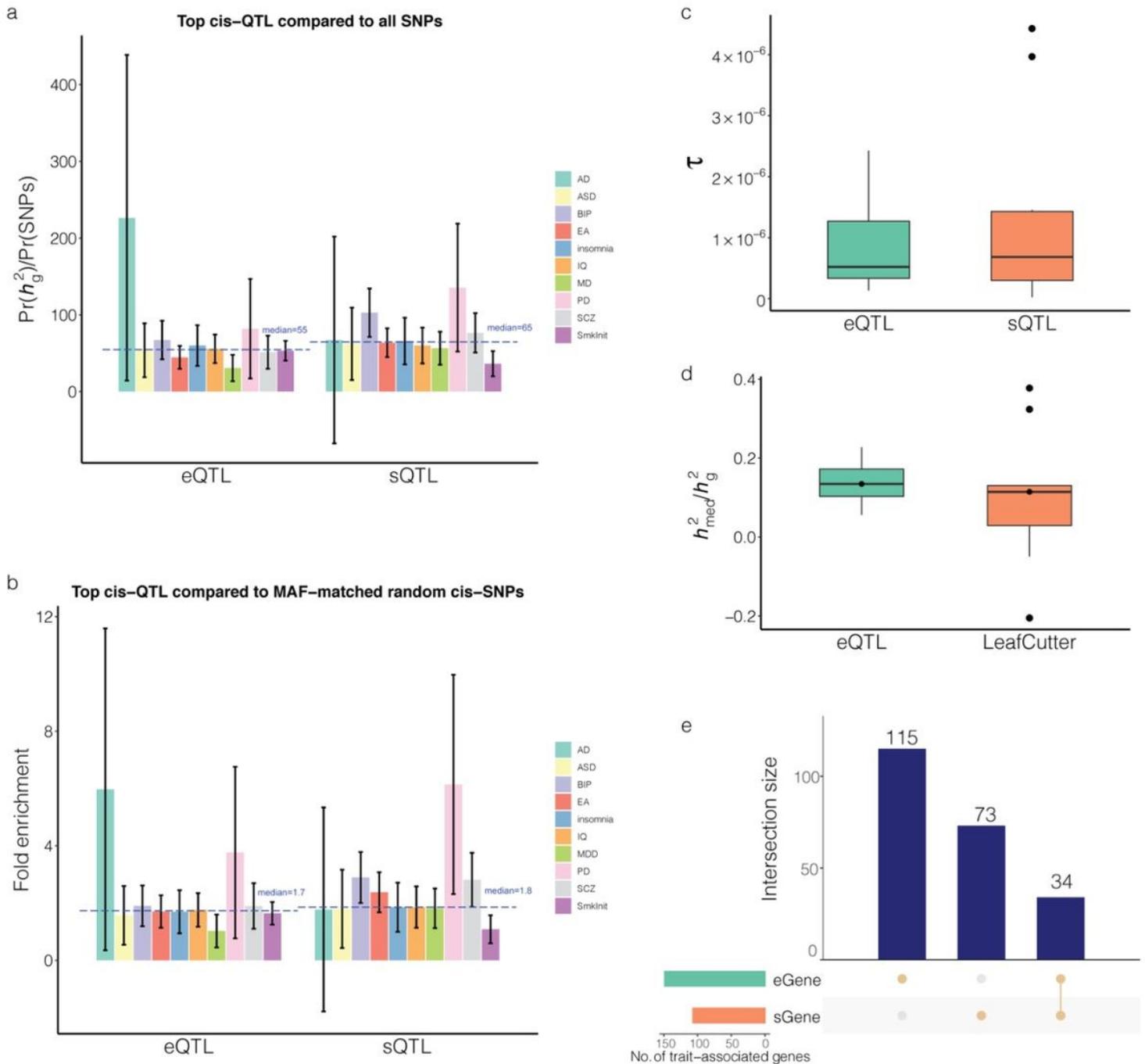


Figure 4

Enrichment of the top cis-sQTL or cis-eQTL SNPs for heritability of the ten brain-related traits. In panel a), heritability enrichment is defined as a ratio of the proportion of heritability explained by the SNPs in query (estimated by S-LDSC) to the proportion of the SNPs. In panel b), heritability enrichment is defined as a ratio of the proportion of heritability explained by the SNPs in query to that by the control SNPs. In panel c), annotation effect size (τ) is used to assess the contribution of the top cis-sQTL SNPs to heritability when fitted jointly with the top cis-eQTL SNPs. Panel d) shows the proportion of mediated heritability (h_{med}^2/h_g^2) by the cis-sQTL (or cis-eQTL) SNPs, estimated by MESC. In panels a) and b), each error bar represents the 95% confidence interval of an estimate, and the blue dashed line represents the median

value across traits. Only the sQTLs discovered by LeafCutter were included in the MESC analysis. Panel e) shows the overlap of the trait-associated sGenes identified by SMR and COLOC with the trait-associated eGenes.

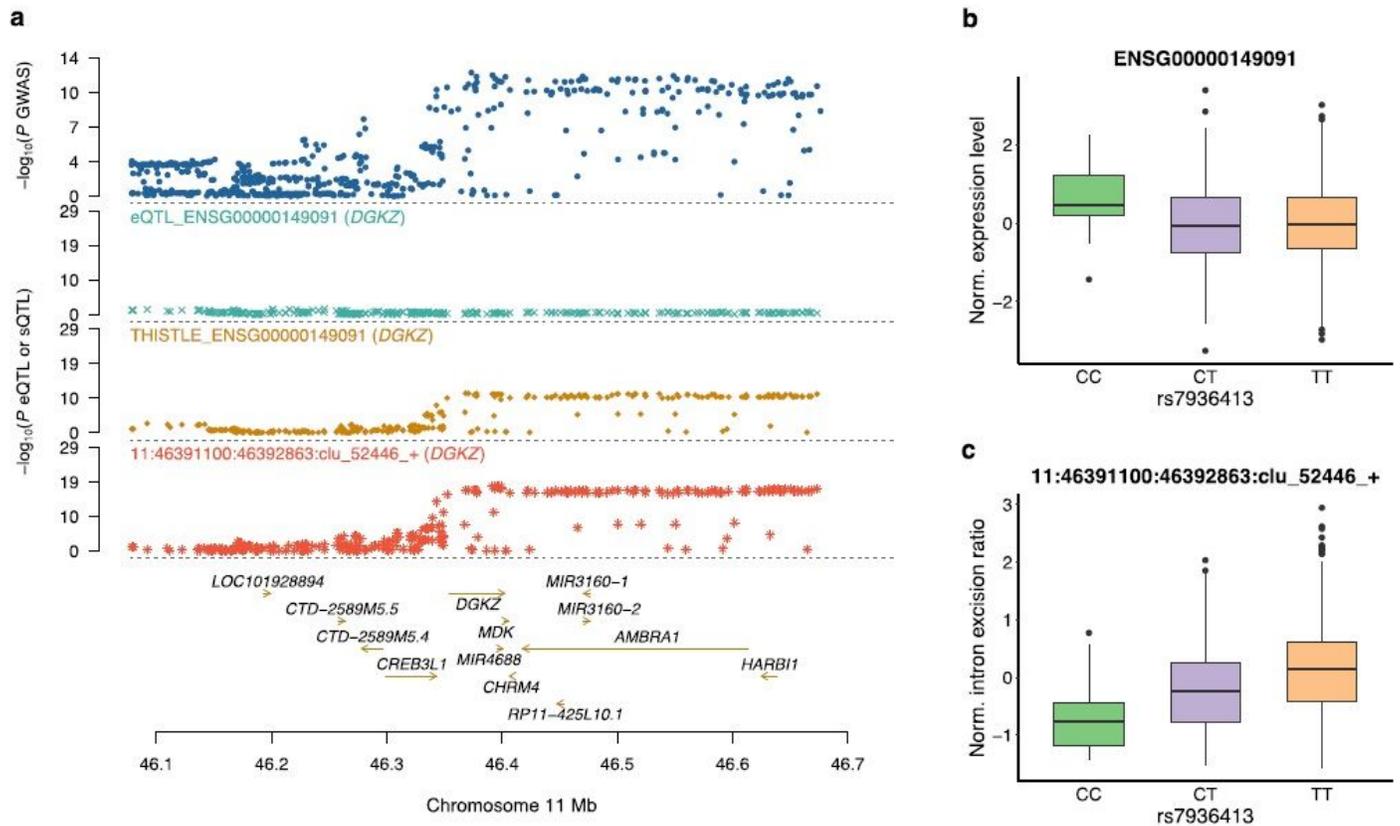


Figure 5

Association of DGKZ with SCZ through sQTLs but not eQTLs. a) The GWAS, sQTL, and eQTL p-values. The top plot shows $-\log_{10}(\text{p-values})$ of SNPs from the GWAS meta-analysis for SCZ. The second, third, and fourth plots show $-\log_{10}(\text{p-values})$ from the eQTL analysis, THISTLE sQTL analysis, and LeafCutter sQTL analysis (intron 11:46391100:46392863:clu_52446_+) for DGKZ, respectively. b) Association of rs7936413 (the top sQTL SNP) with the overall mRNA abundance of DGKZ. Each dot represents mRNA abundance of an individual. c) Association of rs7936413 with intron excision ratio of 11:46391100:46392863:clu_52446_+. Each dot represents intron excision ratio of an individual.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [sQTLSoM25Jan2021.pdf](#)