

## TITLE: Transcriptome-wide association study of schizophrenia and chromatin activity yields mechanistic disease insights

AUTHORS: Alexander Gusev<sup>1,2,\*</sup>, Nick Mancuso<sup>3</sup>, Hilary K Finucane<sup>1,4</sup>, Yakir Reshef<sup>5</sup>, Lingyun Song<sup>6,7</sup>, Alexias Safi<sup>6,7</sup>, Edwin Oh<sup>8</sup>, Schizophrenia Working Group of the Psychiatric Genomics Consortium, Steven McCarroll<sup>9,10</sup>, Benjamin Neale<sup>2,10,11</sup>, Roel Ophoff<sup>12,13</sup>, Michael C O'Donovan<sup>14</sup>, Nicholas Katsanis<sup>8</sup>, Gregory E Crawford<sup>6,7</sup>, Patrick F Sullivan<sup>15,16</sup>, Bogdan Pasaniuc<sup>3,\*†</sup> and Alkes L Price<sup>1,2,\*†</sup>

<sup>1</sup>Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, Massachusetts, USA.

<sup>2</sup>Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA.

<sup>3</sup>David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, California, USA.

<sup>4</sup>Department of Mathematics, Massachusetts Institute of Technology, Cambridge, Massachusetts, USA.

<sup>5</sup>Department of Computer Science, Harvard University, Cambridge, Massachusetts, USA.

<sup>6</sup>Center for Genomic and Computational Biology, Duke University, Durham, North Carolina, USA.

<sup>7</sup>Department of Pediatrics, Division of Medical Genetics, Duke University Medical Center, Durham, North Carolina, USA.

<sup>8</sup>Center for Human Disease Modeling, Duke University Medical Center, Durham, North Carolina, United States.

<sup>9</sup>Department of Genetics, Harvard Medical School, Boston, Massachusetts, USA

<sup>10</sup>Stanley Center for Psychiatric Research, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

<sup>11</sup>Analytic and Translational Genetics Unit, Massachusetts General Hospital and Harvard Medical School, Boston, Massachusetts, USA.

<sup>12</sup>Center for Neurobehavioral Genetics, University of California, Los Angeles, Los Angeles, California, USA

<sup>13</sup>Department of Psychiatry, Brain Center Rudolf Magnus, University Medical Center Utrecht, Utrecht, The Netherlands

<sup>14</sup>MRC Centre for Psychiatric Genetics and Genomics, Cardiff University, Cardiff, UK

<sup>15</sup>Departments of Genetics and Psychiatry, University of North Carolina, Chapel Hill, North Carolina, USA

<sup>16</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm, Sweden

\*Corresponding authors

†Equal contribution

ABSTRACT: Genome-wide association studies (GWAS) have identified over 100 risk loci for schizophrenia, but the causal mechanisms remain largely unknown. We performed a transcriptome-wide association study (TWAS) integrating expression data from brain, blood, and adipose tissues across 3,693 individuals with schizophrenia GWAS of 79,845 individuals from the Psychiatric Genomics Consortium. We identified 157 genes with a transcriptome-wide significant association, of which 35 did not overlap a known GWAS locus; the largest number involved alternative splicing in brain. 42/157 genes were also associated to specific chromatin phenotypes measured in 121 independent samples (a 4-fold enrichment over background genes). This high-throughput connection of GWAS findings to specific genes, tissues, and regulatory mechanisms is an essential step toward understanding the biology of schizophrenia and moving towards therapeutic interventions.

## INTRODUCTION

1 Genome-wide association studies (GWAS) have yielded thousands of robustly associated variants for schizophre-  
2 nia (SCZ) and many other complex traits, but relatively few of these associations have implicated specific  
3 biological mechanisms<sup>1,2</sup>, as GWAS association signals often span many putative target genes, may affect  
4 gene expression through regulatory<sup>3</sup> or structural elements<sup>4</sup>, and may affect genes at considerable genomic  
5 distances via chromatin looping<sup>5,6</sup>. A growing body of research has demonstrated the enrichment of SCZ  
6 GWAS risk variants and heritability within regulatory elements identified through maps of chromatin mod-  
7 ifications and accessibility<sup>1,7-13</sup>. Since chromatin modifications are themselves under genetic control<sup>6,14-19</sup>,  
8 a causal mechanism for SCZ loci could lead from genetic variation to chromatin modifiers to gene expression  
9 and finally to disease risk. Indeed, QTLs for chromatin (and other molecular phenotypes) are enriched within  
10 GWAS associations, further supporting this hypothesis<sup>6,18,20,21</sup>.

11 In this work, we leveraged large gene expression cohorts from multiple tissues, as well as splice variants in  
12 brain, to perform a transcriptome-wide association study (TWAS)<sup>22-24</sup> in a large SCZ GWAS data set<sup>1</sup>  
13 to identify genes whose expression is associated with SCZ and mediated by genetics. We subsequently  
14 performed a TWAS for a diverse set of chromatin phenotypes to identify SCZ susceptibility genes that are  
15 also associated with specific regulatory elements. To our knowledge, this is the first TWAS to integrate  
16 analysis of gene expression, differential splicing, and chromatin variation, moving beyond top SNPs to  
17 implicate SCZ-associated molecular features across the regulatory cascade (Figure 1A).

## RESULTS

### 18 *TWAS for SCZ identifies new susceptibility genes*

19 We analyzed gene expression and genome-wide SNP array data in 3,693 individuals across four expression  
20 reference panels spanning three tissues: RNA-seq from the dorsolateral prefrontal cortex (PFC) of 621  
21 individuals - including SCZ and bipolar (BIP) cases and controls - collected by the CommonMind Consortium  
22 (CMC)<sup>25</sup> (see Web Resources), expression array data measured in peripheral blood from 1,245 unrelated  
23 individuals from the Netherlands Twin Registry (NTR)<sup>26</sup>, expression array data measured in blood from  
24 1,264 individuals from the Young Finns Study (YFS)<sup>23</sup>, and RNA-seq measured in adipose tissue from  
25 563 individuals from the Metabolic Syndrome in Men study (METSIM)<sup>23</sup>. The CMC/brain RNA-seq data  
26 further allowed the characterization of differentially spliced introns<sup>27</sup> (see Methods). Average cis and trans  
27 estimates of SNP-heritability of expression ( $h_g^2$ , see Methods) were highly significant in each panel, with  
28 nominally significant cis- $h_g^2$  ( $P < 0.01$ ) for a total of 18,084 genes summed across the four panels (10,819  
29 unique genes; Table S1), as well as an additional 9,009 differentially spliced introns in brain (in 3,908 unique  
30 genes; Table S1).

31 We performed a TWAS using each of the four gene expression reference panels and summary-level data from  
32 the PGC SCZ GWAS of 79,845 individuals<sup>1</sup> in order to identify genes associated to SCZ (Figure S1). Briefly,  
33 this approach integrates information from expression reference panels (SNP-expression correlation), GWAS  
34 summary statistics (SNP-SCZ correlation), and LD reference panels (SNP-SNP correlation) to assess the  
35 association between the cis-genetic component of expression and phenotype (expression-SCZ correlation)<sup>23</sup>  
36 (Figure 1A). In practice, the expression reference panel was used as the LD reference panel, and cis SNP-  
37 expression effect sizes were estimated using a sparse mixed linear model<sup>28</sup> (see Methods).

38 The TWAS identified 247 transcriptome-wide significant gene-SCZ and intron-SCZ associations (summed  
39 across expression reference panels) for a total of 157 unique genes, including 49 genes that were significant  
40 in more than one expression panel (Figure 2, Figure S1, Table 1, Table S2, S3). Of the 104 (non-HLA,  
41 autosomal) known PGC GWAS loci<sup>1</sup>, 47 loci overlapped with at least one TWAS gene locus (accounting  
42 for 122/157 genes) with the remaining 35/157 genes implicating novel loci. We excluded the MHC region  
43 (chr6:28-34MB) from our primary analyses due to its complex haplotype and LD structure. However, as  
44 a positive control we specifically tested the *C4A* gene recently fine-mapped for SCZ<sup>4</sup>, which lies inside the  
45 MHC, and confirmed a highly significant TWAS association between *C4A* expression in brain tissue and  
46 SCZ ( $P = 1.8 \times 10^{-18}$ ). Across all TWAS associations, the implicated gene was the nearest gene to the top  
47 SNP at the locus in only 56% of instances (using the 10,819 cis-heritable genes as background; decreasing

1 to 24% of instances when using all 26,469 known RefSeq genes), underscoring previous observations that  
2 the nearest gene to a GWAS hit is often not the most likely susceptibility gene when integrated with  
3 expression data<sup>23,24,29,30</sup>. Likewise, conditioning on the predicted expression of a TWAS-associated gene  
4 (using summary-level data<sup>31</sup>, see Methods) reduced the  $\chi^2$  of the lead GWAS SNP at the locus (including  
5 genome-wide significant and non-significant loci) from 42 to 10 on average, and explained more of the  
6 association signal than conditioning on the corresponding top expression-QTL (eQTL) (Table S4). For  
7 the 43 lead GWAS SNPs at genome-wide significant loci that were in LD ( $r^2 > 0.05$ ) with the predicted  
8 expression of at least one TWAS-significant gene (out of 47 overlapping index SNPs), joint conditioning  
9 on the predicted expression of all such genes reduced the median SNP P-value from  $P = 1.2 \times 10^{-10}$  to  
10  $P = 0.028$  (Table S5). Given that the TWAS typically captures only 60-80% of the cis component of gene  
11 expression at these expression panel sample sizes<sup>23</sup>, the complete elucidation of the cis component could  
12 potentially explain the entire GWAS signal at these loci.

13 For both total gene expression and differentially spliced introns, we frequently observed hotspots of multiple  
14 TWAS-associated genes in the same locus, a phenomenon previously observed for complex traits<sup>30</sup>. To quan-  
15 tify the total number of independently associated genes, we applied summary statistic-based approximate  
16 conditional and joint association methods<sup>31</sup> to identify genes/introns that had significant TWAS associations  
17 when analyzed jointly (see Methods). This yielded a set of 63 jointly significant associations, of which 17  
18 were at novel loci (Table 1, Table S3). 8 loci contained multiple significant associations (spanning 16 genes)  
19 in the joint analysis, indicative of allelic heterogeneity. Differentially spliced introns in CMC/brain accounted  
20 for more jointly significant associations than any other reference panel (Table 1), followed by gene expression  
21 in brain, emphasizing the importance of having both expression and splicing measured in a relevant tissue.

22 The differentially spliced introns accounted for 46 transcriptome-wide significant gene associations (of which  
23 10 were at novel loci), comparable to the 44 significant gene associations from brain (Table 1, Supplementary  
24 Materials), despite the fact that differentially spliced introns accounted for 30% fewer significantly cis-  
25 heritable genes than total expression (Table S1). Overall, 20/46 associations corresponded to genes that  
26 were not tested in the analysis of total gene expression due to non-significant expression heritability, and  
27 19 of the remaining 26 did not have a transcriptome-wide significant association for total gene expression.  
28 We identified multiple TWAS loci driven by specific splice-QTL (sQTL) SNPs that significantly explained  
29 a SCZ GWAS association independent of the eQTL effects (Figure S2, S3, S4; Supplementary Materials).  
30 This is consistent with the recent observation that sQTLs are typically independent of eQTLs at the same  
31 gene<sup>27</sup>. We note that, in contrast to total gene expression, effect direction for these associations is difficult  
32 to interpret because multiple alternatively spliced exons within an isoform yielded excised introns that were  
33 highly negatively correlated.

34 This SCZ GWAS data<sup>1</sup> was recently evaluated in a TWAS with gene expression using Summary-based  
35 Mendelian Randomization (SMR)<sup>24</sup>, identifying 16 transcriptome-wide significant associated genes (in con-  
36 trast to 157 identified here). This gap could be explained by several differences in the methods: SMR relies  
37 on individual eQTL significance and a first stage association at the top eQTL, which has been previously  
38 shown to have less power than our TWAS approach using all SNPs in the locus<sup>23</sup>; the SMR test uses an  
39 intentionally conservative second stage test for heterogeneity; and the SMR analysis used a different expres-  
40 sion panel from a meta-analysis of blood. Of the 16 gene associations identified by SMR, 12 were tested in  
41 our study in blood, all replicated at nominal  $P < 0.05$  (with consistent sign), and 9 were transcriptome-wide  
42 significant - a striking concordance given the different methods and independent expression panels used.

#### 43 *TWAS associations replicate in internal and external SCZ cohorts*

44 We first replicated the TWAS signal using internal cross-validation within the PGC SCZ GWAS cohorts<sup>1</sup>  
45 (we had permission to access individual genotypes for 58,246 of 79,845 samples). As a prerequisite step, we  
46 verified that TWAS using the raw GWAS data produced similar results as TWAS using summary statis-  
47 tics<sup>23</sup>, observing a correlation of Z-scores ranging from 0.85-0.90 despite the somewhat different set of GWAS  
48 samples and independent measures of LD (Table S6). Next, we down-sampled the PGC data into GWAS  
49 discovery samples of various sizes (10,000-50,000) to quantify the power and out-of-sample replication of the  
50 SCZ TWAS associations (where the expression panel size was always held constant; Figure 3A). We estimated  
51 TWAS effect-size precision as the slope of a regression of replication vs. discovery TWAS effect sizes across

1 transcriptome-wide significant associations (see Methods), where a slope below 1 represents over-estimated  
2 effect sizes in the discovery data due to winner's curse. Across random down-samples, average effect-size  
3 precision was 0.93 at a discovery sample size of 50,000, indicating minimal winner's curse and projecting  
4 highly accurate effect-size estimates in the full TWAS. The number of transcriptome-wide significant asso-  
5 ciated genes increased linearly with GWAS sample size and was similar across all four expression reference  
6 panels (with more associations in brain expression at the largest discovery size), consistent with a polygenic  
7 architecture and many undiscovered associations.

8 We next replicated the TWAS associations externally using case-control phenotypes from the CMC cohort  
9 (which consisted of SCZ+BIP cases and controls; all cases were included due to the high genetic correlation  
10 between these two diseases<sup>32</sup>). We observed significant replication across all four expression panels, with  
11 effect-size precision not significantly different from 1 (Figure 3B; Figure S5). Surprisingly, the TWAS gene  
12 effect-sizes achieved a higher effect-size precision than GWAS SNP effect-sizes (Figure S6), possibly due to  
13 TWAS aggregating heterogeneous effects in a locus more effectively than a single top SNP. We note that even  
14 though the same CMC samples were used for the TWAS brain expression reference panel and replication  
15 using case-control status, this is an independent replication because CMC case-control status was never used  
16 in the discovery TWAS.

#### 17 *Gene-based risk scores are more predictive than top SNPs*

18 Next, using all transcriptome-wide significant TWAS genes and differentially spliced introns identified in the  
19 PGC and their corresponding discovery effect sizes, we constructed gene-based risk scores (GeRS) from their  
20 predicted expression in the CMC (SCZ+BIP) case-control samples (see Methods). The GeRS from each  
21 expression panel were significantly associated to case-control status, with the strongest association coming  
22 from the CMC/brain expression GeRS, which explained  $1.4\times$  more SCZ variance (and was more significant;  
23  $P = 2.0 \times 10^{-4}$  vs.  $1.2 \times 10^{-3}$ ) than a genetic risk score computed using the 104 (non-HLA, autosomal)  
24 published genome-wide significant GWAS associations<sup>1</sup> (Figure 3C). The GeRS remained significant in a  
25 joint model with the published GWAS predictor ( $P = 0.01$ , Table S7) showing that the TWAS is prioritizing  
26 disease-relevant signal beyond the top SNPs.

27 We relaxed the significance threshold and found gene-based polygenic risk scores (GePRS) to be highly  
28 predictive across the full spectrum of TWAS association P-values (Figure 3D), as observed previously with  
29 SNP-based polygenic scores<sup>1,33,34</sup>. Although the prediction was significant in all tissues individually, there  
30 was evidence of increased effect in brain, with the prediction from brain (genes and introns) capturing 92%  
31 of the joint prediction from all tissues (Figure 3D; Figure S7). A GePRS from actual measured expression  
32 and differential splicing in brain was substantially less significant than the genetic GePRS (Figure S7), as  
33 expected if the non-genetic component of expression is independent of the TWAS signal. Based on polygenic  
34 theory<sup>35,36</sup>, the best TWAS GePRS was estimated to account for 26% of the total SCZ SNP-heritability (see  
35 Supplementary Material), indicating a substantial contribution from cis effects in these four tissues.

36 We sought to investigate temporal differences within the CMC/brain TWAS signal using individual tran-  
37 scriptomes collected by the BRAINSPAN study (see Web Resources) across developmental periods in the  
38 same brain sub-region (PFC) ranging from fetal to adult. For each of 19 developmental periods, we esti-  
39 mated differential expression relative to the other periods as an indicator of temporal specificity. We observed  
40 TWAS  $\chi^2$  statistics to be significantly positively correlated with differential expression from the mid-fetal  
41 developmental period ( $P < 0.05/19$ ), and corresponding significant negative correlation for differential ex-  
42 pression from post-fetal periods (Figure S8, S9). The effect was most significant in the CMC/brain TWAS,  
43 less significant in the two blood TWAS, and non-significant in the adipose TWAS. This further underscores  
44 tissue-specific differences in the TWAS associations, and prioritizes genes expressed in early development for  
45 relevance to SCZ.

#### 46 *Chromatin TWAS identifies specific regulatory mechanisms for SCZ-associated genes*

47 We next sought to identify relationships between the expression of associated genes from the SCZ TWAS  
48 and variation in cis-regulatory elements marked by chromatin (Figure 1A). We used population-level ChIP-  
49 Seq chromatin phenotypes measured in 76 HapMap YRI LCLs for H3k27ac (marking active enhancers),

1 H3k4me1 (enhancers), H3k4me3 (promoters), and DNase (open chromatin)<sup>6</sup>, and in 45 HapMap CEU LCLs  
2 for H3k27ac, H3k4me1, H3k4me3, PU1 (regulatory transcription factor) and RNA polymerase II (RPB2,  
3 associated with active transcription)<sup>18</sup>. For each of the nine chromatin phenotypes, sites with an excess of  
4 ChIP-Seq reads were categorized into local “peaks” corresponding to increased chromatin activity<sup>6,18</sup>. For  
5 each peak, the chromatin abundance across individuals was then treated as a single quantitative trait (with  
6 quality control mirroring the gene expression analyses; see Methods). Both cohorts additionally had gene  
7 expression measured by RNA-seq in the same samples, and we first used Haseman-Elston regression<sup>37</sup> to  
8 quantify the average genetic correlation between gene expression and all chromatin phenotype peaks in the  
9 cis locus, for all genes (see Methods). These cis genetic correlations were highly significant (and substantially  
10 higher than total correlation of measured phenotypes) across all chromatin phenotypes, persisting for peaks  
11 as far as 500kb from the TSS (Figure S10, S11, Table S8, see Methods). We also observed large and highly  
12 significant genetic correlations between peaks from different chromatin phenotypes (Figure S12), suggesting  
13 that such peaks may tag a single underlying biological feature.

14 Motivated by these findings, we applied individual-level TWAS methods<sup>23</sup> to predict expression from the  
15 much larger expression reference panels into samples with chromatin phenotypes and searched for expression-  
16 chromatin associations. Prediction was performed from expression to chromatin phenotype samples (instead  
17 of from chromatin phenotype to expression samples) due to improved numerical stability in the larger ex-  
18 pression panels, but we note that this choice was agnostic to the direction of causality (See Supplementary  
19 Materials). We confirmed by simulation that this chromatin TWAS strategy is well-calibrated (Figure S13)  
20 and much better powered to identify SNP → chromatin → expression associations compared to the con-  
21 ventional approach of testing each SNP for a significant association to both expression (eQTL) and nearby  
22 chromatin peaks (cQTL)<sup>6,18</sup> (Figure 4A, Figure S14).

23 We performed the chromatin TWAS for the 10,819 significantly heritable genes and 9,009 differentially  
24 spliced introns analyzed in the SCZ TWAS, together with all chromatin peaks in the ±500kb locus of each  
25 gene. Focusing on the 157 transcriptome-wide significant genes from the SCZ TWAS, we identified 42 genes  
26 (including 7 genes at novel loci) that also had Bonferroni significant chromatin TWAS associations (to a total  
27 of 78 individual chromatin peaks) in analyses using the same expression reference panel (Table 1, Table 2,  
28 Table S9, S10). Significant evidence of a chromatin-SCZ association was observed for the majority of genes  
29 using a separate TWAS-like test as well as the SMR<sup>24</sup> test using cQTLs, in spite of the low sample size  
30 (Supplementary Materials, Table S11). Individually, these 42 loci represent specific mechanistic hypotheses  
31 where the implicated chromatin phenotypes are disrupted and mediate the expression of susceptibility genes.  
32 Overall, there was a highly significant enrichment of gene-chromatin TWAS associations at SCZ TWAS genes  
33 relative to all heritable genes (OR=3.9;  $P = 1.4 \times 10^{-11}$  by Fisher’s exact test), suggesting that such mecha-  
34 nisms are particularly relevant for SCZ susceptibility genes. The enrichment was also individually significant  
35 ( $P < 0.05/5$ ) for all tissues except adipose (Table 1). Surprisingly, the SCZ-associated differentially spliced  
36 introns were also enriched for chromatin associations (OR=5.0,  $P = 8.1 \times 10^{-5}$ ), even though differentially  
37 spliced introns had many fewer chromatin associations overall (see below), suggesting that epigenetic reg-  
38 ulation of splicing may be particularly relevant for SCZ susceptibility genes. Only 8 of the 78 chromatin  
39 peaks underlying joint SCZ TWAS and chromatin TWAS associations were within the promoter (±2kb of  
40 the TSS) of their associated gene. This suggests that most regulatory elements affecting SCZ are distally  
41 located, as previously observed in other traits<sup>6,8,20</sup> and underscores the importance of searching broadly  
42 around the TSS.

43 We describe three specific examples of TWAS associations to SCZ and chromatin phenotypes. First, the  
44 expression of *SLC45A1* in CMC/brain was associated with SCZ ( $P = 3.5 \times 10^{-8}$ ; overlapping a significant  
45 GWAS locus) as well as a distal RPB2 peak (106kb from TSS;  $P = 1.5 \times 10^{-5}$ ), with significant marginal  
46 associations (QTLs) for both expression and chromatin (Figure 1B). TWAS prioritized this gene-peak com-  
47 bination from 15 genes and 60 peaks in the 1MB locus. Conditioning<sup>31</sup> on the predicted expression of  
48 *SLC45A1* explained all significant eQTLs/cQTLs and GWAS SNPs in the locus. Notably, *SLC45A1* was  
49 not the nearest gene to the top GWAS SNP nor to the associated chromatin peak, and the lead GWAS SNP  
50 was only nominally associated with the RPB2 peak. This highlights a chromatin association that could not  
51 have been identified using conventional approaches based on proximity or overlapping top hits. Second, the  
52 expression of *PPP2R3C* in NTR/blood was associated with SCZ ( $P = 3.4 \times 10^{-6}$ ) - despite no genome-wide  
53 significant SNPs at the locus - as well as two distal peaks for H3k4me1 (minimum  $P = 1.0 \times 10^{-9}$ ) and

1 two distal peaks for H3k27ac (minimum  $P = 4.1 \times 10^{-6}$ ) (Figure S15). The four chromatin peaks clustered  
2 together physically and suggest a coordinated regulatory effect on expression. Conditioning on the predicted  
3 expression of *PPP2R3C* again explained all significant marginal associations for the implicated phenotypes  
4 (Figure S15). *PPP2R3C* was the nearest gene to the most significantly associated SNP at the locus and  
5 to the implicated chromatin peaks. However, because the locus was not genome-wide significant, this as-  
6 sociation would not have been identified in a conventional analysis of known GWAS loci. *PPP2R3C* was  
7 also recently identified by SMR analysis of SCZ in an independent expression panel<sup>24</sup>; our findings identify  
8 specific chromatin features for experimental follow-up. Third, we describe the *MAPK3* locus in detail below.

### 9 *Allelic imbalance analyses localize TWAS association for MAPK3*

10 We highlight the SCZ TWAS association of *MAPK3* expression in CMC/brain ( $P = 1.3 \times 10^{-06}$ ), over-  
11 lapping a significant SCZ GWAS locus. Notably, *MAPK3* is located inside the 16p11 copy number variant  
12 that has been associated with SCZ and autism<sup>38-41</sup>, and has recently been linked with differential protein  
13 abundance in SCZ<sup>42</sup> and shown to respond to pharmacological targeting in cultured neurons<sup>43</sup>. In our anal-  
14 ysis, *MAPK3* was also nominally differentially expressed in the independent CMC (SCZ+BIP) case-control  
15 samples (Wilcoxon  $P = 0.03$ , over-expressed in controls) despite the small sample size. The chromatin  
16 TWAS identified associations between *MAPK3* and two peaks near the TSS (H3k27ac,  $P = 7 \times 10^{-6}$ ;  
17 RPB2,  $P = 1 \times 10^{-11}$ ). In the CEU chromatin phenotype samples, where *MAPK3* expression was also  
18 measured in LCLs, the H3k27ac and RPB2 peaks explained 36% ( $P = 7 \times 10^{-6}$ ) and 23% ( $P = 5 \times 10^{-4}$ ) of  
19 the variance in measured expression, respectively, with only the H3k27ac peak significant in a joint model.  
20 Because the chromatin TWAS associations were identified in LCLs, we examined additional epigenetic data  
21 from H3k27ac, H3k4me3 and ATAC-seq measured in brain tissues (including pre-frontal cortex) as part  
22 of the PsychENCODE project<sup>44</sup> and showed the presence of peaks across all three chromatin phenotypes  
23 (Figure S16, Supplementary Materials). Strikingly, both peaks were also nearly identical to two recently  
24 identified human-gained neuro-developmental enhancers in independent fetal cortex tissues<sup>45</sup> (Figure S16),  
25 providing compelling evidence of evolutionary importance. We then focused on two putatively functional  
26 SNPs within the ATAC-seq peak (rs28529403 and rs61764202, in tight LD with  $\rho = 0.7$ ) that were both  
27 recently classified as disrupting multiple transcription factor binding sites<sup>21</sup> and therefore plausible causal  
28 variants. Conditioning on either SNP accounted for all significant marginal associations across the corre-  
29 sponding expression and chromatin phenotypes (Figure S17, S18), consistent with these SNPs driving the  
30 local signal.

31 We sought to confirm the effect at these SNPs by looking at allelic imbalance, which naturally accounts for  
32 environmental differences across individuals<sup>46,47</sup>. First, we focused on allele-specific expression in the CMC  
33 RNA-seq data (which is statistically independent of the QTL-based signal used for TWAS). We phased  
34 the locus and, for each SNP, tested all heterozygous carriers for an imbalance in RNA-seq reads at the  
35 transcript SNP (see Methods). Both putative SNPs showed significant evidence of allele-specific expression  
36 (rs28529403,  $P = 2.5 \times 10^{-21}$ ; rs61764202,  $P = 2.0 \times 10^{-21}$ ), and were the most significantly imbalanced in  
37 the locus (Table S12, Figure S19).

38 Next, we collected additional ATAC-seq data at this locus, for a total of 314 CMC individuals<sup>44</sup>, and looked  
39 for allele-specific chromatin activity in the same manner. We restricted to 267 SNPs that were Bonferroni  
40 significant for allele-specific expression (above), and tested all heterozygous carriers for an imbalance in  
41 ATAC-seq reads using either of the two peak SNPs as anchors (see Methods). We observed a Bonferroni  
42 significant allelic imbalance at rs61764202 (94 heterozygous samples,  $P = 8.6 \times 10^{-5}$ ), which was more signif-  
43 icant than any other SNP in phase with rs61764202; we did not observe significant imbalance at rs28529403  
44 (100 samples,  $P = 0.15$ ) or any other SNP in phase with rs28529403 (Table S12). The imbalance was also  
45 individually significant in four samples but with opposite direction in one of the four (Figure S20, Table S13).  
46 For this sample, we cannot decisively rule out the possibility of artifact, a mis-imputed heterozygous call,  
47 confounding from an interaction involving nearby features, or tissue sub-type heterogeneity. However, such  
48 biases would not be expected to inflate the combined test across all individuals. Taken together, the allele-  
49 specific signal from both molecular phenotypes nominates the alternative allele at rs61764202 as disrupting  
50 chromatin activity, increasing *MAPK3* expression, which in turn decreases SCZ risk.

#### 1 *Chromatin TWAS associations explain the bulk of cis expression regulation*

2 We expanded the chromatin TWAS to all 10,819 heritable genes (not just SCZ-associated genes) in order  
3 to evaluate the properties of all chromatin TWAS associations. This yielded roughly  $7\times$  more Bonferroni  
4 significant expression-chromatin associations than using the conventional in-sample eQTL/cQTL overlap  
5 approach<sup>6,18</sup> (Figure 4B, Table S14), consistent with our previous simulations. Across all tissues, 806  
6 unique genes had a transcriptome-wide significant association (see Methods) with at least one chromatin  
7 phenotype (Figure 4B, Table S15), and 4,294 genes were significant at the 10% (per-phenotype) FDR used  
8 in previous studies<sup>6,18</sup> (Table S16). In contrast, only 224 of 9,009 differentially spliced introns in the CMC  
9 has a transcriptome-wide significant association with a nearby chromatin mark, corresponding to  $2\text{-}3\times$   
10 fewer associations than identified using total CMC gene expression (depending on the chromatin phenotype,  
11 Table S17). As expected, the distribution of associated chromatin peaks was centered at the TSS of the  
12 corresponding gene (Figure S21). However, additional associations were consistently observed when testing  
13 distal peaks ( $> 10\text{kb}$ ) after correcting for the additional tests performed, in contrast to the results from  
14 eQTL/cQTL overlap (Figure 4B, S22, S23). We evaluated these distal chromatin TWAS associations against  
15 genome-wide Hi-C measured in a reference LCL<sup>6</sup>, reasoning that truly interacting gene-peak pairs are likely  
16 to be in 3-dimensional chromatin contact. Across all chromatin phenotypes, we found distal ( $\geq 10\text{kb}$ ) gene-  
17 peak associations to be significantly enriched for Hi-C inferred loops relative to random (distance-matched)  
18 gene-peak pairs, with an odds ratio of 2.4 (95% CI 2.2-2.5;  $P < 1 \times 10^{-16}$ ) for pairs up to 500kb apart  
19 (Figure S24).

20 Our primary analyses leveraged external expression data in larger sample sizes, but we also analyzed directly  
21 measured expression in chromatin samples for validation. First, we confirmed that the predicted expression  
22 was significantly correlated with measured expression in the chromatin samples (Table S18). Even though the  
23 four expression reference panels all contained samples of European ancestry, this correlation was significant  
24 for both CEU and YRI target samples; with average  $R^2$  for YRI half as large as for CEU, providing an  
25 estimate of the power loss due to differences in linkage disequilibrium (LD) patterns across populations<sup>48</sup>.  
26 Next, for every TWAS-associated gene-chromatin peak pair, we measured the association between actual  
27 measured expression and the corresponding chromatin phenotype. If the chromatin TWAS association  
28 reflects true genetic correlation, we would expect the measured expression to also be highly correlated  
29 (barring tissue-specific differences between the expression panels). Indeed, across the 806 chromatin TWAS-  
30 associated genes, the correlation between measured expression and an associated chromatin phenotype was  
31 highly significant when compared against a distance-matched background null (Figure S10B). In CEU,  
32 where cross-population LD differences are minimized, the average individual TWAS-associated chromatin  
33 peak explained 21% of the variance in measured expression of its target gene for peaks within 2kb of the  
34 TSS, with the subset of distal peaks (2kb-500kb to the gene boundary, 32% of the peaks) explaining 18%  
35 (Figure S25, S26, S27, S28). This strikingly high variance explained is comparable to the total  $cis\text{-}h_g^2$   
36 of these genes (in blood, the tissue type most relevant to LCL; Table S19), implying that the  $cis$ -genetic effect  
37 on expression may be fully explained by individual chromatin TWAS peaks. We stress that the measured  
38 expression in chromatin samples was completely independent of the external expression data used in the  
39 chromatin TWAS, ensuring that these estimates were not biased by over-fitting or winner's curse. For  
40 the three chromatin phenotypes that were measured in both CEU and YRI, associated peaks identified in  
41 one population were still predictive of association with measured expression in the other (Figure S29, S30,  
42 Table S20), lending further support to previous observations that regulatory activity and the resulting impact  
43 on disease is stable across ethnicities<sup>16,49</sup> and supporting our use of measurements from multiple populations.

#### 44 *Chromatin modifications mediate SNP-expression associations that impact SCZ*

45 We next evaluated whether the SCZ and chromatin TWAS associated loci were more consistent with a  
46 chromatin mediating (SNP  $\rightarrow$  chromatin  $\rightarrow$  expression) or expression mediating (SNP  $\rightarrow$  expression  $\rightarrow$   
47 chromatin) causal model leading to SCZ susceptibility. We derived a statistic based on the ratio of  $cis$ -  
48 genetic covariance ( $cov_g$ ) between SCZ and the two molecular phenotypes (see Methods). Conceptually,  
49 the genetic effect of a given molecular phenotype on SCZ will be attenuated by environmental noise, which  
50 will manifest itself as lower  $cov_g$  to SCZ for phenotypes further along the molecular cascade. We estimated  
51 expression-SCZ and chromatin-SCZ  $cov_g$  in the CEU and YRI samples with both chromatin and expression

1 data, using cross-trait LD score regression<sup>50</sup>. Averaging across the 42 chromatin TWAS associations at  
2 genes identified in the SCZ TWAS, both the expression-SCZ and chromatin-SCZ  $cov_g$  were significantly  
3 higher than that of a random background of gene-peak pairs less than 500kb apart, with chromatin-SCZ  
4  $cov_g$  2.5× greater on average than expression-SCZ  $cov_g$  and more significantly different from the background  
5 (Figure S31, see Methods). This corresponds to the chromatin phenotype explaining 35% of the variance in  
6 expression under the model where it is the mediator, with the rest due to environmental or trans variance  
7 independent of disease (see Methods). Furthermore, regressing the chromatin phenotype out of expression  
8 led to non-significant estimates of expression-SCZ  $cov_g$ , but regressing the measured expression out of the  
9 chromatin phenotype did not significantly affect the chromatin-SCZ  $cov_g$  estimates (Figure S31B). Both  
10 observations strongly support a chromatin mediating model (SNP → chromatin → expression) previously  
11 hypothesized<sup>15-17</sup>, and extend it to SCZ etiology.

## DISCUSSION

12 The landmark PGC SCZ GWAS paper stated that “if most risk variants are regulatory, available eQTL  
13 catalogues do not yet provide power, cellular specificity, or developmental diversity to provide clear mech-  
14 anistic hypotheses for follow-up experiments”<sup>1</sup>. In this work, we apply cutting-edge methodology to data  
15 from expression and chromatin activity to provide mechanistic hypotheses. Applying the TWAS approach  
16 to SCZ, we identified 157 unique genes with transcriptome-wide significant associations, whose predicted  
17 expression explained the bulk of the corresponding GWAS SNP association. Genes below the transcriptome-  
18 wide significance threshold continued to be strongly associated with SCZ and exhibited clear preference for  
19 expression and splicing in the brain (though this can also reflect expression data quality). Indeed, alterna-  
20 tive splicing in the brain yielded the greatest number of independent TWAS associations, highlighting an  
21 important source of disease-relevant variation<sup>27</sup> with potential therapeutic implications<sup>51,52</sup>. 42 of the SCZ-  
22 associated genes were significantly associated with nearby chromatin variation (a 4-fold enrichment relative  
23 to non-SCZ genes), implicating specific regulatory features for functional follow-up. Our analyses strongly  
24 supported a model where chromatin variation (likely marking differences in transcription-factor binding)  
25 mediates the relationship between genetics and expression<sup>15-17</sup>, and connect this model with SCZ etiology.  
26 This may explain the modest overlap between eQTLs and GWAS hits previously reported<sup>29,30,46</sup>, where  
27 the downstream nature of the expression phenotype (relative to chromatin) would decrease power for such  
28 overlap analyses.

29 We note several limitations and future directions of this study. First, the discovery chromatin phenotypes  
30 analyzed here were measured in LCLs, preventing us from identifying regulatory elements that were brain-  
31 specific and of potentially greater relevance to SCZ. Second, although TWAS is not confounded by reverse-  
32 causality (disease → expression), instances where a SNP influences SCZ which in turn influences gene  
33 expression (or the same SNP influences SCZ and expression independently) are statistically indistinguishable  
34 from causal susceptibility genes. Therefore, although we can conclude that 26% of the SCZ  $h_g^2$  is explained by  
35 predicted expression of all genes analyzed here we cannot estimate the fraction that is truly causal without  
36 further mediation experiments. Finally, our simulations indicate that the chromatin TWAS approach will  
37 have substantially greater power as chromatin sample sizes increase into the hundreds<sup>44</sup>. As tissue acquisition  
38 may be the biggest hurdle for producing larger data sets, methods that do not depend on measurements  
39 from the same samples will remain critical. Above and beyond specific mechanistic findings for SCZ, our  
40 findings outline a systematic approach to identify biological mediators of complex disease.



## WEB RESOURCES

- 1 TWAS results:
- 2 <http://sashagusev.github.io/chromatinTWAS/>
- 3 TWAS methods:
- 4 <http://bogdan.bioinformatics.ucla.edu/software/twas/>
- 5 BRAINSPAN transcriptomes:
- 6 <http://www.brainspan.org/static/download.html>
- 7 CommonMind consortium:
- 8 <https://www.synapse.org/cmc>
- 9 Grubert et al data:
- 10 <http://chromovar3d.stanford.edu>
- 11 PGC summary data:
- 12 <https://www.med.unc.edu/pgc/downloads>
- 13 PLINK:
- 14 <https://www.cog-genomics.org/plink2>
- 15 PsychENCODE knowledge portal:
- 16 <https://www.synapse.org/#!Synapse:syn4921369/wiki/235539>
- 17 SNPWeights for principal component analysis:
- 18 <http://www.hsph.harvard.edu/alkes-price/software/>
- 19 Waszak et al data (provided by the authors):
- 20 <http://gardeux-vincent.eu/Cell2015/description.peaks.zip>,
- 21 <http://gardeux-vincent.eu/Cell2015/quantified.peaks.zip>,
- 22 <http://gardeux-vincent.eu/Cell2015/quantified.peaks.PEER.centered.zip>

## ACKNOWLEDGEMENTS

- 23 We would like to acknowledge Michael Gandal, Bryce van de Geijn, Arthur Ko, Po-Ru Loh, Luke O'Connor,
- 24 and Paivi Pajukanta for helpful discussions. This research was funded by NIH grants F32 GM106584,
- 25 R01 GM105857, R01 MH109978 and R01 MH107649. HKF was supported by the Fannie and John Hertz
- 26 Foundation. We are grateful to the CommonMind Consortium and the PsychENCODE Consortium for
- 27 making data publicly and readily available.
- 28 CMC: Data were generated as part of the CommonMind Consortium supported by funding from Takeda
- 29 Pharmaceuticals Company Limited, F. Hoffman-La Roche Ltd and NIH grants R01MH085542, R01MH093725,
- 30 P50MH066392, P50MH080405, R01MH097276, RO1-MH-075916, P50M096891, P50MH084053S1, R37MH057881
- 31 and R37MH057881S1, HHSN271201300031C, AG02219, AG05138 and MH06692. Brain tissue for the study
- 32 was obtained from the following brain bank collections: the Mount Sinai NIH Brain and Tissue Repository,
- 33 the University of Pennsylvania Alzheimers Disease Core Center, the University of Pittsburgh NeuroBioBank
- 34 and Brain and Tissue Repositories and the NIMH Human Brain Collection Core. CMC Leadership: Pamela
- 35 Sklar, Joseph Buxbaum (Icahn School of Medicine at Mount Sinai), Bernie Devlin, David Lewis (Univer-
- 36 sity of Pittsburgh), Raquel Gur, Chang-Gyu Hahn (University of Pennsylvania), Keisuke Hirai, Hiroyoshi
- 37 Toyoshiba (Takeda Pharmaceuticals Company Limited), Enrico Domenici, Laurent Essioux (F. Hoffman-La
- 38 Roche Ltd), Lara Mangravite, Mette Peters (Sage Bionetworks), Thomas Lehner, Barbara Lipska (NIMH).
- 39 The data reported in this paper are tabulated in the Supplementary Materials and archived at:
- 40 <http://sashagusev.github.io/chromatinTWAS/>

## METHODS

### 1 *Data and quality control*

2 Genotypes and expression from the NTR<sup>26</sup>, YFS<sup>23</sup>, and METSIM<sup>23</sup> were processed as described in ref<sup>23</sup>  
3 and the corresponding expression weights were downloaded directly from the TWAS web-site (see Web  
4 Resources). Genotypes and expression data from the CMC<sup>25</sup> were processed using the GTEx Consortium  
5 guidelines for eQTL analysis of RNA-seq data. Specifically, RNA-seq RPKM was quantile normalized across  
6 samples; genes having > 10 individuals with zero reads were removed; each gene was rank-normalized; 15  
7 PEER factors were computed; and the residual expression used.

8 We used the LeafCutter algorithm<sup>27,53</sup> to quantify de novo intron excision in the CMC RNA-seq data by  
9 clustering reads that spanned intron junctions. These clusters correspond to individual isoforms and enable  
10 an estimate of differential intron splicing computed from the ratio of reads spanning an intron relative to the  
11 total isoform read count. Splice variants were quantified using LeafCutter default parameters: a minimum  
12 of 50 reads per cluster, and a maximum intron length of 500kb. Based on the guidelines in ref.<sup>53</sup>, the  
13 following quality controls were applied to the inferred isoform clusters: clusters having > 10 individuals  
14 with zero reads were removed; clusters with < 100 individuals having > 20 reads were removed; and introns  
15 with < 5 individuals having non-zero counts were removed. The inferred per-sample abundance for each  
16 intron was then treated as a molecular phenotype, normalized, and PEER-corrected as with total expression  
17 above. This process identified 123,480 differentially spliced introns, of which 99,562 mapped to canonical  
18 gene introns. We treated the differential splicing of these 99,562 introns as quantitative traits in the same  
19 manner as total expression.

20 For genotype data, individuals failing a sex check or having 5% missing SNPs were removed. Additionally,  
21 SNPs were removed if they had > 5% missing calls;  $P < 0.05$  case-control missing association;  $P < 5 \times 10^{-6}$   
22 Hardy-Weinberg disequilibrium;  $P < 5 \times 10^{-3}$  association to batch;  $P < 5 \times 10^{-8}$  missing haplotype  
23 association; or frequency < 1%. Principal components (PCs) were computed using all samples for the NTR,  
24 YFS, and METSIM data directly and using SNPweights (v2.1)<sup>54</sup> for the CMC data, outliers were removed  
25 (samples > 6 standard deviations away the mean along any top component), and PCs included as fixed-  
26 effects in estimating  $h_g^2$ . For all datasets, related individuals with GRM values > 0.05 were also removed  
27 prior to estimating  $h_g^2$ .

### 28 $h_g^2$ estimation

29 Cis and trans  $h_g^2$  were estimated using variance-components, modeling the phenotype as a multi-variant  
30 Normal  $y \sim \sigma_{g,cis}^2 K_{cis} + \sigma_{g,trans}^2 K_{trans} + \sigma_e^2 I$  where  $K$  are the standard genetic relatedness matrices from  
31 SNPs in the cis locus ( $K_{cis}$ ) and in the rest of the genome ( $K_{trans}$ ). The  $\sigma^2$  parameters were fit for each  
32 gene using AI-REML as implemented in the GCTA software<sup>55</sup>, with principal components and sex included  
33 as fixed effects. For  $h_g^2$  of differentially spliced introns, the intron ratios condition out isoform abundance  
34 but total gene expression was also included as a covariate to account for any residual correlation. As in  
35 previous studies<sup>26</sup>, individual estimates outside the plausible 0-1 range were allowed in order to achieve  
36 unbiased mean estimates. The standard error of each estimate was approximated as the standard deviation  
37 divided by the square root of the number of genes tested; however, significant differences were confirmed by  
38 permutation tests (see below).

39 To evaluate the contribution of low-frequency variants, we imputed the NTR data to the Haplotype Reference  
40 Consortium reference, yielding high-quality imputed SNPs down to MAF of 0.001. On average, we did not  
41 observe a significantly non-zero contribution of imputed rare variants to cis- $h_g^2$ , nor did we see a significant  
42 change in common cis- $h_g^2$  due to denser imputation relative to array SNPs (Table S1). Though recent work  
43 has identified biases in estimates of  $h_g^2$  from rare variants<sup>56</sup>, we expect these biases to be small in the cis  
44 region and largely mitigated by the two-component model. We did not further evaluate the contribution of  
45 rare variants to trans- $h_g^2$ . No difference was observed when using dosages to construct the cis GRM.

46 In the CMC data, where SCZ/BIP and control status was also available, the average cis-genetic correlation  
47 of expression between (SCZ/BIP) cases and controls was 1.00 (s.e. 0.02), indicating consistent direction of  
48 eQTL effect sizes between cases and controls and motivating us to use the full cohort as a TWAS reference

1 panel (Table S21). We additionally performed multiple analyses of expression heritability associated with  
2 functional category (Supplementary Materials) demonstrating pervasive enrichment of chromatin marks  
3 near significantly heritable genes and underscoring the importance of chromatin variation to expression  
4 heritability<sup>57</sup>.

#### 5 *Schizophrenia TWAS*

6 We analyzed publicly available schizophrenia summary statistics from the PGC GWAS of 79,845 individuals<sup>1</sup>  
7 (see Web Resources). Summary statistic-based TWAS was performed as described previously<sup>23</sup> using cis  
8 SNP-expression effect sizes computed by the BSLMM software<sup>28</sup> for each of the four expression reference  
9 panels. Strand-ambiguous alleles (A/T, G/C) were removed from the summary data and all SNPs were set  
10 to the same strand. We evaluated TWAS predictions using either the SNPs genotyped in each expression  
11 reference panel, or imputed HapMap3 SNPs (which typically represent well-imputed SNPs). To account  
12 for multiple hypotheses, we applied Bonferroni correction within each expression panel that was used. This  
13 threshold was chosen so as to maximize consistency with previous published results and not penalize for  
14 additional (and often highly correlated) expression panels tested. Specifically, we report “transcriptome-  
15 wide” significance after correcting for the number of genes tested within each of the five reference panels  
16 (CMC, CMC-splicing, NTR, YFS, METSIM; 5,419 tests on average). This is consistent with the correction  
17 applied in previous TWAS results of multiple expression references<sup>23</sup>.

18 To ensure that results from the CMC data were not biased by ascertainment for SCZ/BIP cases, we performed  
19 a separate TWAS in which cis SNP-expression effect sizes were stratified on CMC case-control status and  
20 meta-analyzed; we observed no significant differences compared to using the full sample (Table S22).

21 Conditional and joint analysis was performed using the summary statistic-based method described in ref.<sup>31</sup>,  
22 which we applied to genes instead of SNPs. This approach requires marginal association statistics (i.e. the  
23 main TWAS results) and a correlation/LD matrix to evaluate the joint/conditional model. The correlation  
24 matrix was estimated by predicting the cis-genetic component of expression for each TWAS gene using the  
25 CMC genotypes and computing Pearson correlations across all pairs of genes. The 247 transcriptome-wide  
26 significant TWAS associations across four reference panels (spanning 157 unique genes) were then added  
27 to the model one at a time and retained if their conditional TWAS association remained significant after  
28 Bonferroni correction for 247 tests. The same procedure was used to perform TWAS joint/conditional  
29 analysis of marginal SNP/QTL associations.

#### 30 *Schizophrenia TWAS replication within PGC2*

31 We used the individual-level PGC data<sup>1</sup> to perform replication analyses and compare to summary-based  
32 results. Cis SNP-expression effect sizes were computed using the BSLMM software (as above) and pre-  
33 dicted into each PGC sub-cohort to produce individual-level gene expression predictions. To evaluate the  
34 importance of SNP platform, two separate analyses were performed for the CMC, using all genotyped SNPs  
35 in the expression reference and all imputed HapMap3 SNPs, respectively. Association between predicted  
36 expression and SCZ was assessed using logistic regression with sub-cohort label and 10 principal components  
37 included as covariates. For the sub-sampling analysis, the full study was randomly split into discovery and  
38 replication. The association was then measured in discovery samples only, and any significant genes again  
39 tested for association in the replication samples. For sub-cohort analysis, one cohort was removed at a time  
40 (marked as replication) and the same process repeated.

#### 41 *Gene-based risk scores*

42 We adopted SNP-based polygenic risk score analysis<sup>33</sup> to evaluate TWAS predictive accuracy and validation.  
43 Given a 1-by- $M$  vector  $z$  of signed association statistics in the discovery study and an  $N$ -by- $M$  matrix  $X$   
44 of genetic values for the corresponding  $M$  genes in the replication study, we constructed a gene-based risk  
45 score  $S = Xz$ . The  $M$  genes were either all transcriptome-wide significant genes (gene risk score - GeRS) or  
46 all genes passing relaxed p-value thresholds (gene polygenic risk score - GePRS). This risk score was then  
47 tested against case/control status by a standard linear model  $y \sim S + P + e$  where  $S$  is the risk score and  $P$   
48 is a matrix of principal components accounting for ancestry. The risk-score accuracy was then measured as

1 the  $R^2$  from the above model less the  $R^2$  from the model  $y \sim P + e$  to account for ancestry, and converted  
2 to the liability scale assuming a prevalence of 1%.

3 For the TWAS using METSIM, YFS, and NTR expression reference panels, the cis-genetic component of  
4 expression was predicted in CMC samples. For the TWAS using the CMC expression panel, either the total  
5 expression was used (Figure S7) or the cis-genetic component of expression was estimated directly using  
6 BSLMM (equivalent up to a scaling factor to estimating genetic values by dropping each individual in turn).  
7 We stress that the case/control label from the CMC data was never used to identify the TWAS associations,  
8 and that the GeRS or GePRS from the CMC expression panel were thus evaluated against an independent  
9 CMC case/control phenotype. Ascertaining cases in the CMC expression panel may increase the frequency of  
10 causal variants and make the prediction more accurate than using a randomly ascertained expression panel,  
11 however, we observed little difference when performing the TWAS using an expression panel consisting of  
12 CMC controls only (Table S22).

### 13 *Power simulations for chromatin TWAS*

14 Using real genotypes from the UK10K<sup>58</sup> study, we simulated a model where SNPs  $X$  are causal for a  
15 chromatin peak ( $C = X\beta_X + e_C$ ) and chromatin is causal for expression ( $E = C\beta_C + e_E$ ). 100 1MB loci  
16 were randomly selected across the genome, and causal SNPs in each locus were then randomly selected from  
17 common variants (MAF > 1%). Environmental noise was set such that SNPs explain 30% of the variance in a  
18 chromatin peak, and chromatin explains 30% of the variance in expression (consistent with our observations  
19 that expression  $h_g^2$  is  $\sim 10\%$ ). The TWAS was then performed using cis SNP-expression effect sizes computed  
20 by BSLMM either predicting the chromatin phenotype (with increasing sample size) into 1,000 independent  
21 individuals with expression or vice versa, and then performing an association between the predicted and  
22 measured phenotype. Separately, eQTL and cQTL association was computed for every common SNP in the  
23 locus using individuals with both expression and chromatin measured. To evaluate power at genome-wide  
24 significance, a TWAS association was reported as significant if it had  $P < 0.05$  after correcting for (average  
25 30 peaks per locus)  $\times$  (20,000 genes) = 600,000 tests. For the QTL-based approach, given  $M$  SNPs in  
26 the locus, a SNP was reported as significant if it had an eQTL  $P < 0.05$  after correcting for  $M \times 20,000$   
27 tests and a cQTL  $P < 0.05$  after correcting for  $M \times 30 \times 20,000$  tests. For a given chromatin sample  
28 size, the simulation was then performed at 100 random loci and 5 random seeds each, with the fraction  
29 of loci reported as significant by each method taken as the power (Figure S14). The TWAS simulation  
30 was separately performed under the null, using non-heritable chromatin and expression, and shown to be  
31 well-calibrated under the null (Figure s13).

### 32 *Individual-level chromatin TWAS*

33 We used cis SNP-expression effect sizes computed by BSLMM scores in the four expression reference panels  
34 (including differentially spliced introns) to predict individual-level expression in the 45 CEU<sup>18</sup> and 76<sup>6</sup> YRI  
35 individuals with measured chromatin phenotypes. We retained only SNPs that were typed in both studies  
36 and removed strand-ambiguous SNPs. We did not perform any additional QC of the functional features,  
37 which were all previously PEER-adjusted and normalized<sup>6,18</sup>. We note that even though the YRI target  
38 samples are of different ethnicity, this prediction does not require an LD-reference panel and is therefore  
39 only expected to suffer loss in power (but not increased type I error) due to the differences in LD. For each  
40 predicted gene, we identified all chromatin peaks within a given window of the TSS (primary results used  
41  $\pm 500kb$ ) and tested each mark for association to predicted expression by linear regression.

### 42 *Multiple hypothesis correction for chromatin TWAS*

43 The large number of correlated phenotypes analyzed - expression from five experiments and chromatin from  
44 nine experiments in two populations - allows for several approaches to multiple testing correction. For  
45 the chromatin TWAS, we corrected for the number of gene-peak pairs tested within a single expression  
46 reference and chromatin phenotype experiment (for example, number of gene-peak pairs when evaluating  
47 predicted CMC expression with the CEU:H3k27ac chromatin phenotype). This is directly comparable to  
48 the experiment-wide corrections applied in previous eQTL/cQTL analyses<sup>6,18</sup>. The same correction was

1 applied for the SCZ/chromatin TWAS overlap: for example, the 44 SCZ TWAS genes identified using  
2 CMC expression were within 500kb of 1,528 total peaks in the CEU:H3k27ac experiment and “overlap” was  
3 reported for any peak that had a chromatin TWAS association  $P < 0.05/1,528$ .

4 For comparison, we separately calculated the number of associations that were significant at 5% FDR across  
5 all molecular experiments. This yielded approximately  $3.5\times$  more chromatin TWAS associations and  $1.2\times$   
6 more SCZ and chromatin TWAS associations (Table S2), demonstrating that the above experiment-wide  
7 Bonferroni correction strategy corresponds to a conservative study-wide FDR.

#### 8 *eQTL/cQTL overlap analysis*

9 We compared the chromatin TWAS to the traditional approach of identifying SNPs that are significant  
10 both as cQTLs and eQTLs in real data (Figure 4B). For each population and given distance to TSS, we  
11 performed this analysis in two stages. Stage 1: We identified all eQTLs that were significant after Bonferroni  
12 correction for the total number of SNP-gene pairs tested. When distance to TSS was the maximal allowed  
13 (500kb), this resulted in 355 eQTLs in the YRI and 579 eQTLs in the CEU data. Stage 2: From the set  
14 of significant eQTLs, we then looked for those that were also significantly associated with peaks from a  
15 given chromatin phenotype (for peaks within the given distance to TSS), after Bonferroni correction for  
16 the number of eQTL-peak pairs tested. In both stages the tests were only counted for the given chromatin  
17 phenotype (e.g. H3K27ac in CEU). This was compared to the chromatin TWAS analysis where each gene  
18 was tested against any peak within the given distance, and number of significant results reported after  
19 Bonferroni correction for total number of gene-peak pairs tested.

#### 20 *Analysis of allele-specific expression*

21 For each molecular phenotype (RNA-seq/ATAC-seq), we phased the locus using EAGLE2<sup>59</sup> and evaluated  
22 haplotype allelic imbalance as follows. Given a “peak” SNP (for which there are RNA/ATAC-seq reads)  
23 and a “target” SNP (for which we want to evaluate allelic imbalance) we restricted to individuals that were  
24 heterozygous for both SNPs and counted the number of peak reads mapping to the REF/ALT haplotypes  
25 of the target SNP. We then assessed statistical significance for deviation from 50% balance using a binomial  
26 test. When the peak SNP and target SNP are identical, this is equivalent to a standard allelic-imbalance  
27 test across all heterozygous carriers in the peak. Duplicate reads were removed using samtools prior to all  
28 analyses and only variants with  $\geq 10$  reads at the peak SNP were tested. See Supplementary Materials for  
29 details on individual phenotypes analyzed.

#### 30 *Ratio of cis-genetic covariance between chromatin-SCZ and expression-SCZ*

31 To shed more light on the potential causal model, we sought to evaluate the evidence in support of two models  
32 of mediation:  $M_{CH}$ , where  $\text{SNP} \rightarrow \text{chromatin} \rightarrow \text{expression} \rightarrow \text{disease}$ ; and  $M_{EX}$ , where  $\text{SNP} \rightarrow \text{expression}$   
33  $\rightarrow \text{chromatin} \rightarrow \text{disease}$ . Under the assumption of linear, additive variance across molecular phenotypes, this  
34 can be estimated via the ratio of genetic covariance ( $cov_g$ ) between chromatin-SCZ and expression-SCZ. The  
35 fraction of environmental variance on expression ( $env_{EX}$ ) under each model of mediation can then computed  
36 from the following equation (see Supp. Note for derivation):

$$cov_{g,CH}/cov_{g,EX} = 1/\sqrt{1 - env_{EX}^2}$$

37 To compare these two models without bias from sample size or assay, we estimated the genetic covariance  
38 ( $cov_g$ ) using the PGC SCZ summary statistics and molecular data from the CEU/YRI individuals which  
39 had both chromatin and gene expression measured. For each SCZ/chromatin TWAS gene, we defined the  
40 target locus as  $\pm 50\text{kb}$  of the union of gene and peak boundary and estimated SCZ-chromatin and SCZ-  
41 expression  $cov_g$  using cross-trait LD score regression<sup>50</sup>, restricting to the well-imputed HapMap3 SNPs and  
42 using in-sample LD. To estimate significance, background  $cov_g$  was computed using the same procedure over  
43 200 randomly gene-peak pairs within 500kb of the TSS. The observed and background estimates were then  
44 compared using the non-parametric Kolmogorov-Smirnov test. We note that for the YRI samples LD in the  
45 GWAS summary statistics is expected to differ from LD in the eQTL/cQTL data, and confound the raw

1 estimate of  $cov_g$ ; however, because the random gene-peak pairs are also computed from the same population,  
2 we do not expect significance measured against this null to be inflated. We separately considered a partial  
3 correlation analysis, where each expression measurement was transformed to the residual of the associated  
4 chromatin peak in a standard linear regression (and likewise for each chromatin measurement). The two  
5 estimates of  $cov_g$  were again computed from these partial phenotypes as described above. We caution that  
6 the estimate of  $e$  in the above equation was computed from an average across all loci, and could also be  
7 consistent with confounding from different levels of measurement error for ChIP-seq and RNA-seq, a mixture  
8 of models  $M_{CH}$  and  $M_{EX}$  that favors model  $M_{CH}$ , or mediation by other unobserved molecular phenotypes.

# Bibliography

- 2 [1] Schizophrenia Working Group of the Psychiatric Genomics Consortium. (2014). Biological insights from  
3 108 schizophrenia-associated genetic loci. *Nature* *511*, 421–427.
- 4 [2] Price, A. L., Spencer, C. C., and Donnelly, P. (2015). Progress and promise in understanding the genetic  
5 basis of common diseases. In *Proc. R. Soc. B* volume 282 The Royal Society pp. 20151684.
- 6 [3] Soldner, F., Stelzer, Y., Shivalila, C. S., Abraham, B. J., Latourelle, J. C., Barrasa, M. I., Goldmann,  
7 J., Myers, R. H., Young, R. A., and Jaenisch, R. (2016). Parkinson-associated risk variant in distal  
8 enhancer of  $\alpha$ -synuclein modulates target gene expression. *Nature* *533*, 95–99.
- 9 [4] Sekar, A., Bialas, A. R., de Rivera, H., Davis, A., Hammond, T. R., Kamitaki, N., Tooley, K., Presumey,  
10 J., Baum, M., Van Doren, V., et al. (2016). Schizophrenia risk from complex variation of complement  
11 component 4. *Nature* *530*, 177–183.
- 12 [5] Claussnitzer, M., Dankel, S. N., Kim, K.-H., Quon, G., Meuleman, W., Haugen, C., Glunk, V., Sousa,  
13 I. S., Beaudry, J. L., Puvion, V., et al. (2015). Fto obesity variant circuitry and adipocyte browning  
14 in humans. *New England Journal of Medicine* *373*, 895–907.
- 15 [6] Grubert, F., Zaugg, J., Kasowski, M., Ursu, O., Spacek, D., Martin, A., Greenside, P., Srivas, R.,  
16 Phanstiel, D., Pekowska, A., et al. (2015). Genetic control of chromatin states in humans involves local  
17 and distal chromosomal interactions. *Cell* *162*, 1051–1065.
- 18 [7] Maurano, M. T., Humbert, R., Rynes, E., Thurman, R. E., Haugen, E., Wang, H., Reynolds, A. P.,  
19 Sandstrom, R., Qu, H., and Brody, J. (2012). Systematic localization of common disease-associated  
20 variation in regulatory dna. *Science* *337*, 1190–1195
- 21 [8] Trynka, G., Sandor, C., Han, B., Xu, H., Stranger, B. E., Liu, X. S., and Raychaudhuri, S. (2013).  
22 Chromatin marks identify critical cell types for fine mapping complex trait variants. *Nature Genetics*  
23 *45*, 124–130.
- 24 [9] Pickrell, J. K. (2014). Joint analysis of functional genomic data and genome-wide association studies  
25 of 18 human traits. *The American Journal of Human Genetics* *94*, 559–573.
- 26 [10] Gusev, A., Lee, S. H., Trynka, G., Finucane, H., Vilhjálmsson, B. J., Xu, H., Zang, C., Ripke, S.,  
27 Bulik-Sullivan, B., Stahl, E., et al. (2014). Partitioning heritability of regulatory and cell-type-specific  
28 variants across 11 common diseases. *The American Journal of Human Genetics* *95*, 535–552.
- 29 [11] Kichaev, G., Yang, W.-Y., Lindstrom, S., Hormozdiari, F., Eskin, E., Price, A. L., Kraft, P., and  
30 Pasaniuc, B. (2014). Integrating functional data to prioritize causal variants in statistical fine-mapping  
31 studies. *PLoS Genetics* *10*, e1004722
- 32 [12] Won, H.-H., Natarajan, P., Dobbyn, A., Jordan, D. M., Roussos, P., Lage, K., Raychaudhuri, S., Stahl,  
33 E., and Do, R. (2015). Disproportionate contributions of select genomic compartments and cell types  
34 to genetic risk for coronary artery disease. *PLoS Genetics* *11*, e1005622.
- 35 [13] Finucane, H. K., Bulik-Sullivan, B., Gusev, A., Trynka, G., Reshef, Y., Loh, P.-R., Anttila, V., Xu, H.,  
36 Zang, C., Farh, K., et al. (2015). Partitioning heritability by functional annotation using genome-wide  
37 association summary statistics. *Nature Genetics* *47*, 1228–1235.

- 1 [14] Degner, J. F., Pai, A. A., Pique-Regi, R., Veyrieras, J.-B., Gaffney, D. J., Pickrell, J. K., De Leon, S.,  
2 Michelini, K., Lewellen, N., Crawford, G. E., et al. (2012). Dnase [thinsp] i sensitivity qtls are a major  
3 determinant of human expression variation. *Nature* *482*, 390–394.
- 4 [15] McVicker, G., van de Geijn, B., Degner, J. F., Cain, C. E., Banovich, N. E., Raj, A., Lewellen, N.,  
5 Myrthil, M., Gilad, Y., and Pritchard, J. K. (2013). Identification of genetic variants that affect histone  
6 modifications in human cells. *Science* *342*, 747–749.
- 7 [16] Kasowski, M., Kyriazopoulou-Panagiotopoulou, S., Grubert, F., Zaugg, J. B., Kundaje, A., Liu, Y.,  
8 Boyle, A. P., Zhang, Q. C., Zakharia, F., Spacek, D. V., et al. (2013). Extensive variation in chromatin  
9 states across humans. *Science* *342*, 750–752.
- 10 [17] Kilpinen, H., Waszak, S. M., Gschwind, A. R., Raghav, S. K., Witwicki, R. M., Orioli, A., Migliavacca,  
11 E., Wiederkehr, M., Gutierrez-Arcelus, M., Panousis, N. I., et al. (2013). Coordinated effects of sequence  
12 variation on dna binding, chromatin structure, and transcription. *Science* *342*, 744–747.
- 13 [18] Waszak, S., Delaneau, O., Gschwind, A., Kilpinen, H., Raghav, S., Witwicki, R., Orioli, A., Wiederkehr,  
14 M., Panousis, N., Yurovsky, A., et al. (2015). Population variation and genetic control of modular  
15 chromatin architecture in humans. *Cell* *162*, 1039–1050.
- 16 [19] Taudt, A., Colome-Tatche, M., and Johannes, F. (2016). Genetic sources of population epigenomic  
17 variation. *Nat Rev Genet* *advance online publication*, –.
- 18 [20] Farh, K. K.-H., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W. J., Beik, S., Shores, N., Whitton,  
19 H., Ryan, R. J., Shishkin, A. A., et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune  
20 disease variants. *Nature* *518*, 337–343.
- 21 [21] Moyerbrailean, G. A., Kalita, C. A., Harvey, C. T., Wen, X., Luca, F., and Pique-Regi, R. (2016).  
22 Which genetics variants in dnase-seq footprints are more likely to alter binding? *PLoS Genetics* *12*,  
23 e1005875.
- 24 [22] Gamazon, E. R., Wheeler, H. E., Shah, K. P., Mozaffari, S. V., Aquino-Michaels, K., Carroll, R. J.,  
25 Eyler, A. E., Denny, J. C., Nicolae, D. L., Cox, N. J., et al. (2015). A gene-based association method  
26 for mapping traits using reference transcriptome data. *Nature Genetics* *47*, 1091–1098.
- 27 [23] Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B., Jansen, R., de Geus, E., Boomsma, D.,  
28 Wright, F., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies.  
29 *Nature Genetics* *48*, 245–252.
- 30 [24] Zhu, Z., Zhang, F., Hu, H., Bakshi, A., Robinson, M., Powell, J., Montgomery, G., Goddard, M., Wray,  
31 N., Visscher, P., et al. (2016). Integration of summary data from gwas and eqtl studies predicts complex  
32 trait gene targets. *Nature Genetics*.
- 33 [25] Fromer, M., Roussos, P., Sieberts, S. K., Johnson, J. S., Kavanagh, D. H., Perumal, T. M., Ruderfer,  
34 D. M., Oh, E. C., Topol, A., Shah, H. R., et al. (2016). Gene expression elucidates functional impact  
35 of polygenic risk for schizophrenia. *bioRxiv* pp. 052209.
- 36 [26] Wright, F., Sullivan, P., Brooks, A., Zou, F., Sun, W., Xia, K., Madar, V., Jansen, R., Chung, W.,  
37 Zhou, Y.-H., et al. (2014). Heritability and genomics of gene expression in peripheral blood. *Nature*  
38 *Genetics* *46*, 430–437.
- 39 [27] Li, Y., van+de+Geijn, B., Raj, A., Knowles, D., Petti, A., Golan, D., Gilad, Y., and Pritchard, J.  
40 (2016). Rna splicing is a primary link between genetic variation and disease. *Science* *352*, 600–604.
- 41 [28] Zhou, X., Carbonetto, P., and Stephens, M. (2013). Polygenic modeling with bayesian sparse linear  
42 mixed models. *PLoS Genetics* *9*, e1003264.
- 43 [29] Nica, A. C., Montgomery, S. B., Dimas, A. S., Stranger, B. E., Beazley, C., Barroso, I., and Dermitzakis,  
44 E. T. (2010). Candidate causal regulatory effects by integration of expression qtls with complex trait  
45 genetic associations. *PLoS Genetics* *6*, e1000895.



- 1 [30] Nicolae, D. L., Gamazon, E., Zhang, W., Duan, S., Dolan, M. E., and Cox, N. J. (2010). Trait-associated  
2 snps are more likely to be eqtls: annotation to enhance discovery from gwas. *PLoS Genetics* *6*, e1000888.
- 3 [31] Yang, J., Ferreira, T., Morris, A. P., Medland, S. E., Madden, P. A., Heath, A. C., Martin, N. G.,  
4 Montgomery, G. W., Weedon, M. N., Loos, R. J., et al. (2012). Conditional and joint multiple-snp  
5 analysis of gwas summary statistics identifies additional variants influencing complex traits. *Nature*  
6 *Genetics* *44*, 369–375.
- 7 [32] Cross-Disorder Group of the Psychiatric Genomics Consortium. (2013). Genetic relationship between  
8 five psychiatric disorders estimated from genome-wide snps. *Nature Genetics* *45*, 984–994.
- 9 [33] Purcell, S. M., Wray, N. R., Stone, J. L., Visscher, P. M., O'Donovan, M. C., Sullivan, P. F., Sklar, P.,  
10 Ruderfer, D. M., McQuillin, A., Morris, D. W., et al. (2009). Common polygenic variation contributes  
11 to risk of schizophrenia and bipolar disorder. *Nature* *460*, 748–752.
- 12 [34] Vilhjálmsson, B. J., Yang, J., Finucane, H. K., Gusev, A., Lindström, S., Ripke, S., Genovese, G., Loh,  
13 P.-R., Bhatia, G., Do, R., et al. (2015). Modeling linkage disequilibrium increases accuracy of polygenic  
14 risk scores. *The American Journal of Human Genetics* *97*, 576–592.
- 15 [35] Palla, L. and Dudbridge, F. (2015). A fast method that uses polygenic scores to estimate the variance  
16 explained by genome-wide marker panels and the proportion of variants affecting a trait. *The American*  
17 *Journal of Human Genetics* *97*, 250–259.
- 18 [36] Daetwyler, H. D., Villanueva, B., and Woolliams, J. A. (2008). Accuracy of predicting the genetic risk  
19 of disease using a genome-wide approach. *PLoS one* *3*, e3395.
- 20 [37] Haseman, J. and Elston, R. (1972). The investigation of linkage between a quantitative trait and a  
21 marker locus. *Behavior Genetics* *2*, 3–19.
- 22 [38] McCarthy, S. E., Makarov, V., Kirov, G., Addington, A. M., McClellan, J., Yoon, S., Perkins, D. O.,  
23 Dickel, D. E., Kusenda, M., Krastoshevsky, O., et al. (2009). Microduplications of 16p11. 2 are associated  
24 with schizophrenia. *Nature Genetics* *41*, 1223–1227.
- 25 [39] Golzio, C., Willer, J., Talkowski, M. E., Oh, E. C., Taniguchi, Y., Jacquemont, S., Reymond, A., Sun,  
26 M., Sawa, A., Gusella, J. F., et al. (2012). *Kctd13* is a major driver of mirrored neuroanatomical  
27 phenotypes of the 16p11. 2 copy number variant. *Nature* *485*, 363–367.
- 28 [40] Migliavacca, E., Golzio, C., Männik, K., Blumenthal, I., Oh, E. C., Harewood, L., Kosmicki, J. A.,  
29 Loviglio, M. N., Giannuzzi, G., Hippolyte, L., et al. (2015). A potential contributory role for ciliary  
30 dysfunction in the 16p11. 2 600 kb bp4-bp5 pathology. *The American Journal of Human Genetics* *96*,  
31 784–796.
- 32 [41] Maillard, A., Ruef, A., Pizzagalli, F., Migliavacca, E., Hippolyte, L., Adaszewski, S., Dukart, J., Ferrari,  
33 C., Conus, P., Männik, K., et al. (2015). The 16p11. 2 locus modulates brain structures common to  
34 autism, schizophrenia and obesity. *Molecular psychiatry* *20*, 140–147.
- 35 [42] Föcking, M., Lopez, L., English, J., Dicker, P., Wolff, A., Brindley, E., Wynne, K., Cagney, G., and  
36 Cotter, D. (2015). Proteomic and genomic evidence implicates the postsynaptic density in schizophrenia.  
37 *Molecular psychiatry* *20*, 424–432.
- 38 [43] Blizinsky, K. D., Diaz-Castro, B., Forrest, M. P., Schürmann, B., Bach, A. P., Martin-de Saavedra,  
39 M. D., Wang, L., Csernansky, J. G., Duan, J., and Penzes, P. (2016). Reversal of dendritic phenotypes  
40 in 16p11. 2 microduplication mouse model neurons by pharmacological targeting of a network hub.  
41 *Proceedings of the National Academy of Sciences* pp. 201607014.
- 42 [44] Akbarian, S., Liu, C., Knowles, J. A., Vaccarino, F. M., Farnham, P. J., Crawford, G. E., Jaffe, A. E.,  
43 Pinto, D., Dracheva, S., Geschwind, D. H., et al. (2015). The psychencode project. *Nature neuroscience*  
44 *18*, 1707–1712.

- 1 [45] Reilly, S. K., Yin, J., Ayoub, A. E., Emera, D., Leng, J., Cotney, J., Sarro, R., Rakic, P., and Noonan,  
2 J. P. (2015). Evolutionary changes in promoter and enhancer activity during human corticogenesis.  
3 *Science* *347*, 1155–1159.
- 4 [46] Consortium, G. et al. (2015). The genotype-tissue expression (gtex) pilot analysis: Multitissue gene  
5 regulation in humans. *Science* *348*, 648–660.
- 6 [47] Pickrell, J. K., Marioni, J. C., Pai, A. A., Degner, J. F., Engelhardt, B. E., Nkadori, E., Veyrieras, J.-B.,  
7 Stephens, M., Gilad, Y., and Pritchard, J. K. (2010). Understanding mechanisms underlying human  
8 gene expression variation with rna sequencing. *Nature* *464*, 768–772.
- 9 [48] Rosenberg, N. A., Huang, L., Jewett, E. M., Szpiech, Z. A., Jankovic, I., and Boehnke, M. (2010).  
10 Genome-wide association studies in diverse populations. *Nature Reviews Genetics* *11*, 356–366.
- 11 [49] Kichaev, G. and Pasaniuc, B. (2015). Leveraging functional-annotation data in trans-ethnic fine-  
12 mapping studies. *The American Journal of Human Genetics* *97*, 260–271.
- 13 [50] Bulik-Sullivan, B., Finucane, H. K., Anttila, V., Gusev, A., Day, F. R., Loh, P.-R., Duncan, L., Perry,  
14 J. R., Patterson, N., Robinson, E. B., et al. (2015). An atlas of genetic correlations across human  
15 diseases and traits. *Nature Genetics*.
- 16 [51] Sibley, C. R., Blazquez, L., and Ule, J. (2016). Lessons from non-canonical splicing. *Nature Reviews*  
17 *Genetics*.
- 18 [52] Nelson, C. E., Hakim, C. H., Ousterout, D. G., Thakore, P. I., Moreb, E. A., Rivera, R. M. C.,  
19 Madhavan, S., Pan, X., Ran, F. A., Yan, W. X., et al. (2016). In vivo genome editing improves muscle  
20 function in a mouse model of duchenne muscular dystrophy. *Science* *351*, 403–407.
- 21 [53] Li, Y. I., Knowles, D. A., and Pritchard, J. K. (2016). Leafcutter: Annotation-free quantification of  
22 rna splicing. *bioRxiv* pp. 044107.
- 23 [54] Chen, C.-Y., Pollack, S., Hunter, D. J., Hirschhorn, J. N., Kraft, P., and Price, A. L. (2013). Improved  
24 ancestry inference using weights from external reference panels. *Bioinformatics* pp. btt144.
- 25 [55] Yang, J., Lee, S. H., Goddard, M. E., and Visscher, P. M. (2011). Gcta: a tool for genome-wide complex  
26 trait analysis. *The American Journal of Human Genetics* *88*, 76–82.
- 27 [56] Yang, J., Bakshi, A., Zhu, Z., Hemani, G., Vinkhuyzen, A. A., Lee, S. H., Robinson, M. R., Perry, J. R.,  
28 Nolte, I. M., van Vliet-Ostaptchouk, J. V., et al. (2015). Genetic variance estimation with imputed  
29 variants finds negligible missing heritability for human height and body mass index. *Nature Genetics*.
- 30 [57] Liu, X., Finucane, H. K., Gusev, A., Bhatia, G., Gazal, S., O'Connor, L., Bulik-Sullivan, B., Wright,  
31 F., Sullivan, P., Neale, B., et al. (2016). Functional partitioning of local and distal gene expression  
32 regulation in multiple human tissues. *bioRxiv*.
- 33 [58] Consortium, U. et al. (2015). The uk10k project identifies rare variants in health and disease. *Nature*  
34 *526*, 82–90.
- 35 [59] Loh, P.-R., Danecek, P., Palamara, P. F., Fuchsberger, C., Reshef, Y. A., Finucane, H. K., Schoenherr,  
36 S., Forer, L., McCarthy, S., Abecasis, G. R., et al. (2016). Reference-based phasing using the haplotype  
37 reference consortium panel. *bioRxiv* pp. 052308.

#### LIST OF SUPPLEMENTARY MATERIALS

- 1 Supplementary Text
- 2 Table S1-S24
- 3 Figure S1-S38

Table 1: Number of TWAS-associated genes across all phenotypes and tissues.

	CMC/brain introns*	CMC/brain	NTR/blood	YFS/blood	MET/adipose	Total**
Heritable	(9,009) 3,890	5,514	2,743	5,418	4,654	11,749
SCZ associated	(80) 46	44	35	48	39	157
SCZ associated (joint†)	(21) 20	13	12	10	9	63
SCZ associated (novel)	(12) 10	9	6	6	7	35
SCZ associated (novel joint†)	(4) 4	4	3	4	2	17
Chromatin associated	(224) 125	244	182	346	232	806
SCZ and chromatin associated	(10) 8	11	10	13	7	42
Enrichment P-value	8.1e-05	2.0e-05	4.4e-04	6.5e-05	2.9e-02	1.4e-11

\* Number of differentially spliced introns shown in parenthesis.

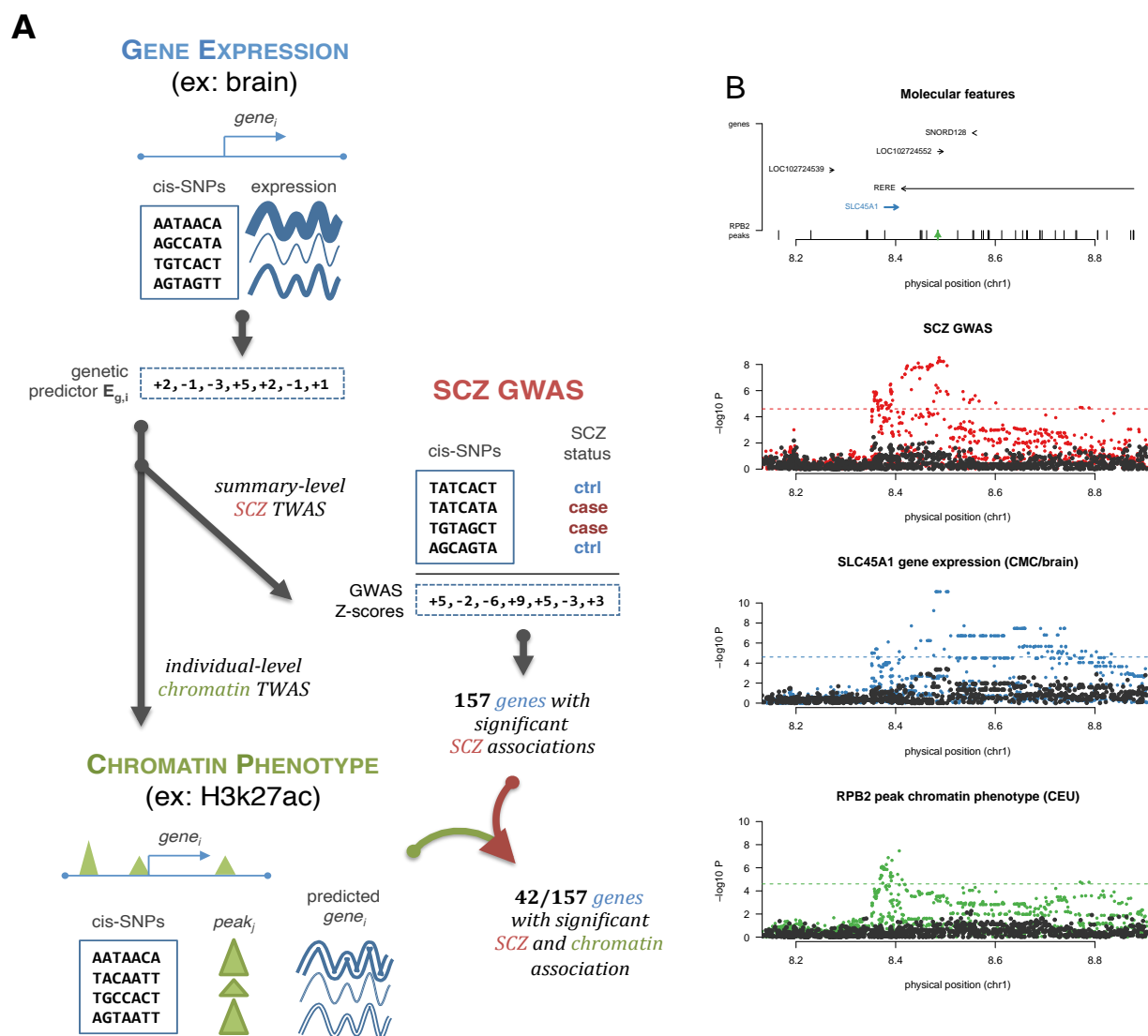
\*\* Total number of unique gene associations.

† Significant in a summary-based joint analysis.

Table 2: **TWAS genes with association to schizophrenia and chromatin phenotypes.** 42 genes (including 7 genes at novel loci, highlighted with a [\*]) had a significant TWAS association with SCZ and chromatin phenotypes. For each significant TWAS association with SCZ, the number of significant gene-chromatin associations (FWER 5% among TWAS gene-mark associations, by Bonferroni correction) are reported. In the middle columns ‘.’ represents genes that were not heritable in the study and therefore not TWAS-associated. In the right columns ‘.’ represents no identified association; genes with no chromatin associations are not shown. Top panel shows results from genes, with TSS listed as position; bottom panel shows results from differentially spliced introns in CMC with exon-exon junction listed as position (details in Table S9).

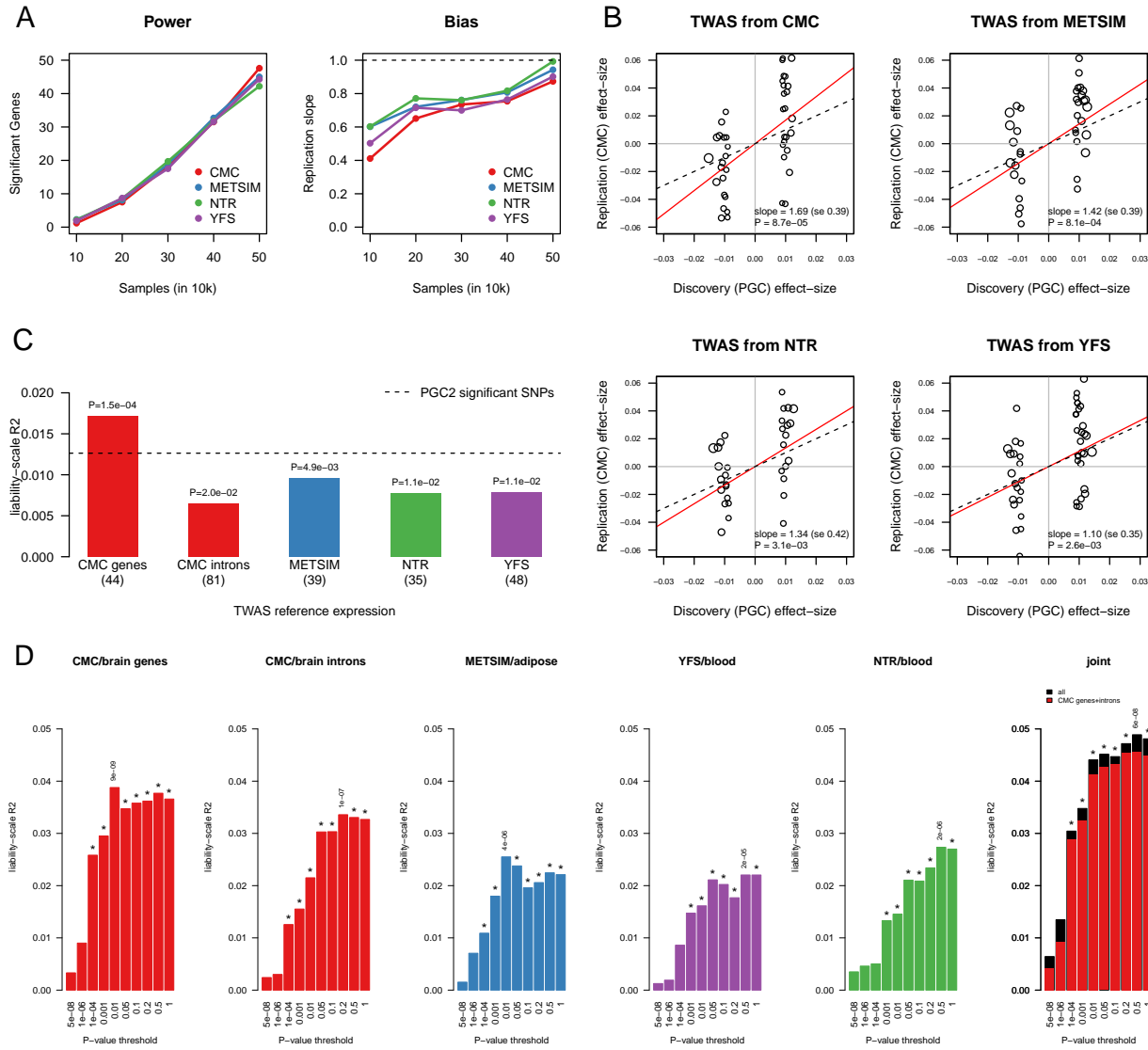
Gene	Chr	Position	TWAS P-value				Associated marks				
			YFS/blood	MET/adipose	NTR/blood	CMC/brain	DHS	H3K27AC	H3K4ME1	H3K4ME3	PU1
RERE	1	8,483,747	4e-07	2e-06	2e-06	.	.	.	.	.	3
SLC45A1	1	8,378,144	.	.	.	4e-08	.	.	.	.	1
MAP7D1*	1	36,621,565	6e-04	.	1e-06	.	.	.	.	.	1
MED8	1	43,855,483	5e-01	.	.	2e-06	.	.	1	.	.
ANP32E	1	150,207,026	.	.	1e-08	.	.	.	1	.	.
MRPS21	1	150,266,261	3e-06	3e-03	6e-03	2e-02	.	.	1	.	.
RFWD2*	1	176,176,380	4e-06	.	.	.	.	1	.	.	.
C2orf69	2	200,775,978	.	6e-10	.	.	.	.	1	.	.
GLT8D1	3	52,737,714	.	.	5e-08	3e-08	.	1	.	.	.
GLYCTK	3	52,321,835	2e-08	.	.	.	.	1	.	.	.
GNL3	3	52,719,935	7e-09	6e-07	.	5e-02	.	.	1	.	.
NEK4	3	52,804,965	.	.	.	2e-09	.	.	.	1	.
NT5DC2	3	52,567,793	6e-06	6e-06	.	7e-01	.	1	.	1	.
PPM1M	3	52,279,808	2e-07	2e-07	.	2e-03	.	1	.	.	.
TMEM110	3	52,931,597	1e-02	4e-01	1e-08	6e-06	.	1	1	2	.
PCCB	3	135,969,166	1e-08	1e-10	.	3e-10	1	.	3	.	.
RP11-53O19.3	5	44,826,178	.	6e-06	.	.	.	1	.	.	.
DND1	5	140,053,171	.	8e-07	1e-02	.	.	1	.	1	.
IK	5	140,027,383	4e-06	1e-06	.	5e-05	.	1	.	1	.
NDUFA2	5	140,027,370	2e-06	.	.	4e-06	.	1	.	2	.
PCDHA2	5	140,174,443	.	.	.	7e-06	.	1	.	1	.
ZMAT2	5	140,080,031	5e-06	1e-03	.	3e-06	.	.	.	1	.
AS3MT	10	104,629,209	.	6e-08	7e-09	1e-05	.	1	.	.	.
MPHOSPH9	12	123,717,785	4e-09	1e-05	.	2e-08	.	1	.	.	1
KIAA0391*	14	35,591,526	7e-01	2e-07	5e-01	.	.	2	1	.	.
PPP2R3C*	14	35,591,748	6e-05	1e-01	3e-06	2e-02	.	2	2	.	.
MAPK3	16	30,134,630	5e-05	.	.	1e-06	.	1	.	.	1
GFOD2	16	67,753,273	.	.	6e-07	2e-05	.	1	.	2	.
TSNAXIP1	16	67,840,780	.	.	.	2e-06	.	.	1	2	.
DUS2L	16	68,038,024	1e-06	.	3e-06	4e-04	.	.	.	2	.
PRMT7	16	68,344,876	1e-05	8e-04	.	8e-06	.	.	1	1	.
GRAP*	17	18,950,336	.	.	5e-07	.	.	.	.	.	1
RNF112*	17	19,314,490	8e-06	.	.	.	.	.	.	.	1
ACTR5	20	37,377,096	2e-07	2e-04	.	7e-01	1	.	1	.	.
CBR3	21	37,507,262	6e-03	2e-03	2e-06	5e-04	1	.	2	.	.
CMC/brain splicing											
TBC1D5	3	17,255,862 - 17,279,655	.	.	.	3e-06	.	.	1	.	.
NEK4	3	52,800,010 - 52,800,194	.	.	.	1e-06	.	.	1	.	.
CCDC90B	11	82,985,783 - 82,991,184	.	.	.	3e-07	.	1	.	.	.
SBNO1	12	123,821,038 - 123,825,535	.	.	.	4e-10	.	.	.	1	.
KLC1	14	104,145,855 - 104,151,323	.	.	.	7e-12	.	.	1	1	.
RTN1*	14	60,074,210 - 60,193,637	.	.	.	1e-06	.	.	1	.	.
TAOK2	16	29,997,825 - 29,998,165	.	.	.	4e-06	.	.	.	.	1
PPP4C	16	30,094,168 - 30,094,715	.	.	.	2e-06	.	.	.	.	1

\* Novel, not overlapping 108 PGC SCZ GWAS loci



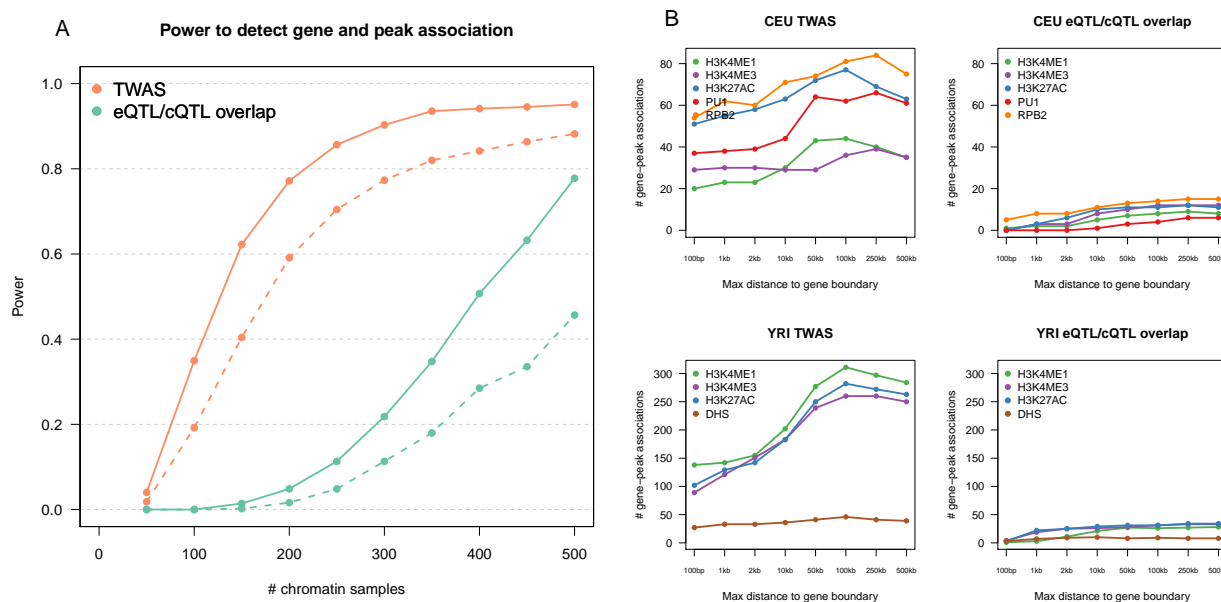
**Figure 1: Schematic and example of TWAS approach.** (left) Illustration of the TWAS approach: genetic predictor of gene expression ( $E_g$  is learned in a reference panel (top); integrated with SCZ GWAS association statistics to infer SCZ- $E_g$  association (middle); further integrated with individual-level chromatin phenotypes to infer genes with SCZ and chromatin- $E_g$  associations (bottom). (right) Example association of *SLC45A1* gene expression and SCZ, as well as *SLC45A1* expression and a distal RPB2 chromatin peak. Top panel shows locus schematic with all nearby genes and chromatin peaks and corresponding associations highlighted. Three lower panels show Manhattan plots of marginal association statistics before and after conditioning on the TWAS predicted expression (colored/dark dots, respectively). Dashed line shows local significance threshold after Bonferroni correction for number of SNPs. The full 1MB cis locus harbors 15 genes and 60 RPB2 peaks.





**Figure 3: Replication of TWAS associations.** (a) The PGC data was randomly split into increasingly large discovery samples (size on the x-axis) and TWAS statistics were estimated from each reference panel. Left panel reports the number of significant genes (after 5% FWER correction) for a given GWAS sample size. Right panel reports the slope from a regression of  $\beta_{\text{replication}} \sim \beta_{\text{discovery}}$  for significant genes identified at each sample size (where all  $56,000 - x$  remaining samples were used as replication). (b) TWAS effect-sizes for association to schizophrenia identified in the PGC (x-axis) compared to corresponding estimates in the CMC (y-axis). Dotted line corresponds to  $y = x$ ; red line corresponds to the slope from a (Z-score weighted) regression of CMC  $\sim$  PGC, with estimate and p-value shown in bottom right. (c) Schizophrenia risk prediction  $R^2$  shown for risk scores constructed from significant TWAS genes (bars) and PGC2 GWAS SNPs (dashed line, comparable to the 1% – 3% reported in ref<sup>1</sup>). Number of genes used in each score reported in parenthesis. Linear-regression  $R^2$  of phenotype on predictor (after subtracting  $R^2$  from jointly fit ancestry PCs) was transformed to the liability scale assuming schizophrenia prevalence of 1%, a linear transformation consistent across all predictors. (d) Schizophrenia risk prediction  $R^2$  for polygenic gene risk scores across multiple significance thresholds. Significant correlations (after Bonferroni correction for number of thresholds tested) are indicated with a (\*) and most significant P-value reported. Right-most panel shows prediction from all tissues jointly (black) and from CMC/brain genes + differentially expressed introns jointly (red).





**Figure 4: Power and detection of significant gene-mark associations.** (a) Molecular phenotypes were simulated under the SNP  $\rightarrow$  chromatin  $\rightarrow$  expression model and two methods to detect gene-mark associations evaluated. (TWAS, orange) corresponds to predicting expression from a held-out reference panel with 1,000 individuals and testing each proximal chromatin peak for association. (eQTL/cQTL, green) corresponds to identifying SNPs that are significantly associated with both chromatin and expression at the locus. For a given chromatin phenotype sample size (x-axis), power was measured as the number of instances where locus was deemed significant after accounting for number of gene-mark pairs tested (TWAS) or number of SNPs and gene-mark pairs tested (eQTL/cQTL). Solid (dashed) lines correspond to 1 (2) chromatin-causing variants in the simulation. (b) Genes significantly associated with a chromatin phenotype peak in YRI/CEU. For a given distance from the gene (x-axis), the number of unique genes is reported after experiment-wide Bonferroni correction for tests across all chromatin phenotypes. Results from TWAS prediction-based approach shown on left; results from overlapping QTL approach on the right (see Methods).