1    Developmental and genetic regulation of the human cortex transcriptome in schizophrenia

2

3    Andrew E Jaffe[1,2,3,4,#], Richard E Straub[1], Joo Heon Shin[1], Ran Tao[1], Yuan Gao[1], Leonardo
4    Collado Torres[1,3,4], Tony Kam-Thong[5], Hualin S Xi[6], Jie Quan[6], Qiang Chen[1], Carlo
5    Colantuoni[1,7,8], Bill Ulrich[1], Brady J. Maher[1,9], Amy Deep-Soboslay[1],  The BrainSeq Consortium,
6    Alan Cross[10], Nicholas J. Brandon[10], Jeffrey T Leek[3,4], Thomas M. Hyde[1,7,9], Joel E. Kleinman[1,7],
7    Daniel R Weinberger[1,8,9,11,&]

8    1.  Lieber Institute for Brain Development, Johns Hopkins Medical Campus, Baltimore, MD
9        21205, USA
10   2.  Department of Mental Health, Johns Hopkins Bloomberg School of Public Health, Baltimore,
11       MD, 21205, USA
12   3.  Department of Biostatistics, Johns Hopkins Bloomberg School of Public Health, Baltimore,
13       MD, 21205, USA
14   4.  Center for Computational Biology, Johns Hopkins University, Baltimore MD 21205 USA
15   5.  Roche Pharma Research and Early Development, Pharmaceutical Sciences, Roche
16       Innovation Center Basel, F. Hoffmann-La Roche Ltd, Grenzacherstrasse 124, 4070 Basel,
17       Switzerland
18   6.  Computational Sciences, Pfizer Inc, Cambridge, MA 02140
19   7.  Department of Neurology, Johns Hopkins School of Medicine, Baltimore, MD 21205, USA
20   8.  Department of Neuroscience, Johns Hopkins School of Medicine, Baltimore, MD 21205,
21       USA
22   9.  Department of Psychiatry and Behavioral Sciences, Johns Hopkins School of Medicine,
23       Baltimore, MD 21205, USA
24   10. AstraZeneca Neuroscience iMed, IMED Biotech Unit, Cambridge, MA, 02139 USA
25   11. McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins School of Medicine,
26       Baltimore, MD 21205, USA

27   # andrew.jaffe@libd.org (co-corresponding)

28   & drweinberger@libd.org (co-corresponding)

29

30

**Summary:**

GWAS have identified over 108 loci that confer risk for schizophrenia, but risk mechanisms for individual loci are largely unknown. Using developmental, genetic, and illness-based RNA sequencing expression analysis, we characterized the human brain transcriptome around these loci and found enrichment for developmentally regulated genes with novel examples of shifting isoform usage across pre- and post-natal life. Within patients and controls, we implemented a novel algorithm for RNA quality adjustment, and identified 237 genes significantly associated with diagnosis that replicated in an independent case-control dataset. These genes implicated synaptic processes and were strongly regulated in early development ($p < 10^{-20}$). Lastly, we found 42.5% of risk variants associate with nearby genes and diverse transcript features that converge on developmental regulation and subsequent dysregulation in illness and 34 loci show convergent directionality with illness association implicating specific causative transcripts. These data offer new targets for modeling schizophrenia risk in cellular systems.

44

45

**Key Words:** schizophrenia, functional genomics, RNA sequencing, human postmortem brain, differential expression analysis, RNA degradation

48

**Introduction:**

Schizophrenia (SCZD) is a prevalent neuropsychiatric disorder with a combination of genetic and environmental risk factors. Research over the last several decades has suggested that SZCD is a neurodevelopmental disorder arising through altered connectivity and plasticity in relevant neural circuits. However, discovering the causative mechanisms of these putatively developmental deficits has been very challenging[1]. The most consistent evidence of etiologic mechanisms related to SCZD has come from a recent genome-wide association study (GWAS) in which over a hundred independent single nucleotide polymorphisms (SNPs) were identified having an allele frequency difference between patients with schizophrenia and unaffected controls[2]. While these findings have identified regions in the genome harboring genetic risk variants, almost all of the associated SNPs are non-coding, located in intronic or intergenic sequence, and hypothesized to have some role in regulating expression[3]. However, the exact gene(s) and transcript(s) potentially regulated by risk-associated genetic variation are uncertain, as most of these genomic regions contain multiple genes. In principle, the effects of non-coding genetic variation, by whatever mechanisms (e.g. promoter, enhancer, splicing, noncoding RNA, epigenetics, etc), should be observed in the transcriptome. Therefore, to better understand how these regions of genetic risk and their underlying genotypes may confer risk of schizophrenia and to better characterize the molecular biology of the disease state, we sequenced the polyA+ transcriptomes from the prefrontal cortex of 495 individuals with ages across the lifespan, ranging from the second trimester of fetal life to 85 years of age (see Table S1), including 175 patients with schizophrenia (see Figure S1).

Here we identify novel expression associations with genetic risk and with illness state and explore developmentally regulated features, including a subset of genes with previously uncharacterized isoform shifts in expression patterns across the fetal-postnatal developmental transition. We further identify many more expression quantitative loci (eQTLs) in schizophrenia risk regions then previously observed by surveying the full spectrum of associated expression features to generate potential molecular mechanisms underlying genetic risk. We also explore differential gene expression associated with the state of illness in a comparison of the postmortem brains of patients with schizophrenia with non-psychiatric controls. By incorporating a novel, experiment-based algorithm to account for RNA quality differences which have not been adequately controlled in earlier studies, we report a high degree of replication across independent case-control gene expression datasets[4,5]. By combining genetic risk at the population-level with eQTLs and case-control differences, we identify putative human frontal cortex mechanisms underlying risk for schizophrenia and replicable molecular features of the illness state.

84

**Results**

We performed deep polyA+ RNA-sequencing of 495 individuals, ranging in age from the second trimester of fetal life to 85 years old (see Table S1), including 175 patients with

88    schizophrenia (see Figure S1). We quantified expression across multiple transcript features,
89    including: annotated 1) genes and 2) exons, 3) annotation-guided transcripts[4] as well as
90    alignment-based 4) exon-exon splice junctions[5] and 5) expressed regions (ERs)[6]. These last
91    two expression features were selected to reduce reliance on the potentially incomplete
92    annotation of the brain transcriptome[7] (Results S1). We find a large number of moderately
93    expressed and previously unannotated splice junctions that tag potential transcripts with
94    alternative exonic boundaries or exon skipping (Figure S2), 95% of which are also found in
95    other large RNA-seq datasets, including a subset that were brain-specific (Table 1). Similarly,
96    we find that only 56.1% of ERs were annotated to strictly exonic sequence – many ERs
97    annotated to strictly intronic (22.3%) or intergenic (8.5%) sequence, or were transcribed beyond
98    existing annotation (e.g. extended UTRs, extended exonic sequence).

99

### Developmental regulation of transcription and shifting isoform usage

101    Characterizing expression changes in unaffected individuals, particularly across brain
102    development beginning with prenatal life, has previously offered disease-relevant insights into
103    particular genomic loci [8-12]. Specifically, we and others [7,13,14] have shown that genomic risk loci
104    associated with neurodevelopmental disorders including schizophrenia are enriched for
105    transcript features showing differential expression between fetal and postnatal brains. Here too,
106    among the 320 control samples, the strongest component of expression change corresponded
107    to large expression changes in the contrast of pre-natal and early postnatal life, in line with
108    previous data [7] (Figure 1A). We further defined a developmental regulation statistic for each
109    expressed feature using a generalized additive model (see Methods) and found widespread
110    developmental regulation of these expressed features (Results S2, Table S3, Figure S3),
111    including of previously unannotated sequence (Table S4). Motivated in part by previous reports
112    of preferential fetal isoform use among schizophrenia candidate risk genes[10,11] (e.g predominant
113    fetal versus predominant postnatal isoforms), we next formally identified the subset of genes
114    showing alternative isoform expression patterns across  fetal and postnatal life using those
115    exons, junctions, transcripts, and ERs that meet the statistical criteria for developmental
116    regulation (, i.e. those genes with at least one developmentally changing feature, see Methods).
117    There were 9,161 Ensembl genes (28% of the set of developmentally regulated genes) with
118    both positive and negative expression features having genome-wide significant correlations to
119    age (each with  $p_{bonf}<0.05$, Figure 1B, Table S5, Figure S4). In other words, these represent
120    alternate transcript isoforms of the same gene that show opposite patterns of expression across
121    the prenatal-postnatal transition.

122    We performed gene set analyses of those genes with shifting isoform usage compared
123    to the larger set of genes with at least one developmentally regulated feature but without shifting
124    isoform usage to identify more specific biological functions of this unique form of developmental
125    regulation (Table S6). The set of developmentally shifting isoforms was relatively enriched for
126    localization, catalytic activity, signaling-related processes, including synaptic transmission and
127    cell communication, and neuronal development, among many others. Interestingly, genes
128    identified with shifting isoforms across development based exclusively on junction counts were

129    enriched for both dopaminergic (FDR=$5.11 \times 10^{-7}$) and glutamatergic (FDR=$4.02 \times 10^{-4}$) synapse
130    KEGG pathways (Figure 1C), the two neurotransmitter systems most prominently implicated in
131    schizophrenia pathogenesis and treatment.

132

133    *Schizophrenia risk is associated with shifting isoform usage across brain development*

134    Based on the KEGG analysis, we hypothesized that the genes with developmentally
135    regulated isoform shifts may relate to risk for schizophrenia. Indeed, genes within the SZCD
136    GWAS risk loci were more likely to harbor these isoform shifts occurring in the fetal-postnatal
137    developmental transition compared with the rest of the expressed transcriptome (Figure 2D).
138    For example, genes with developmental isoform shifts identified by exon and junction counts
139    were 62% (p=$4.1 \times 10^{-5}$) and 69% (p=$5.6 \times 10^{-7}$) more likely to lie within the PGC2 risk regions
140    (with permutation-based p = 0.048 and 0.02 respectively, see Methods) than developmentally
141    regulated genes without isoforms shifts (Table S7). These results further underscore the role of
142    changes in the regulation of transcription and splicing in early brain development in
143    schizophrenia risk.

144

145    *Expression associations with chronic schizophrenia illness*

146    We next explored the expression landscape of the prefrontal cortex of the schizophrenia
147    illness state and its potential link with developmental regulation and genetic risk. We performed
148    differential expression modeling using 351 higher quality adult samples (196 controls, 155
149    cases), and found extensive bias by RNA degradation within both univariate analysis (where
150    12,686 genes were differentially expressed at FDR<5%) and even after adjusting for measured
151    levels of RNA quality typical of all prior studies (Figure S5). We therefore implemented a
152    statistical framework based on an independent molecular degradation experiment (see
153    Methods, Results S3), called "quality surrogate variable analysis" (qSVA, see Methods)[15]. We
154    further utilized replication RNA-seq data from the CommonMind Consortium (CMC) dataset,
155    using a subset of age range-matched 159 schizophrenia patients and 172 controls.
156    Interestingly, adjusting for observed factors related to RNA quality that characterize all earlier
157    studies of gene expression in schizophrenic brain, the proportion of genes with differentially
158    expressed features at genome wide significant FDR < 5% that replicate (with directionality and
159    marginal significance) in the CMC dataset was very small (only 11.0%, 244/2,215, Figure S6). In
160    contrast, using our new statistical qSVA approach, 40.1% of differentially expressed genes at
161    FDR < 5% (N=75/183) replicate in the CMC dataset.  At genome-wide significant FDR<10%
162    (see Methods), we identified 237 genes with 556 DE features that replicated in the CMC dataset
163    (33.6% gene-level replication rate, Table S8, Table S9).

164    The differences in expression levels between cases and controls of these DE features
165    were generally small in both our discovery and the replication datasets (Figure 2A, Figure S7),
166    perhaps a direct result of the clinical and molecular heterogeneity of this disorder [13,16]. Gene
167    ontology analysis implicated transporter- and channel-related signaling as significantly

168   consistently downregulated in patients compared to controls across genes annotated in all three
169   expression summarizations (Figure 2B, Table S10). These results suggested decreased
170   signaling in patients with schizophrenia, but could raise the possibility that these replicated
171   expression differences between patients and controls relate to the epiphenomena of illness,
172   such as treatment with antipsychotics which affect signaling in the brain[14], as the majority of
173   patients were on anti-psychotics at the time of death (64%, Table S1). Only two genes (*KLC1*
174   and *PPP2R3A*) in the 108 significant schizophrenia GWAS loci were significantly differentially
175   expressed. However, in an exploratory analysis, we found that overall the differential expression
176   statistics within the loci were significantly different than those features outside the loci (Results
177   S4, Figure S8, Table S11).  We also investigated the relationships between transcription and
178   genomic risk for schizophrenia using genome wide Risk Profile Scores (RPS) from each subject
179   calculated as previously described[2] (see Methods). Using the subset of 209 Caucasian
180   samples, we largely found a lack of association between RPS and expression of individual
181   expression features. We further found a lack of enrichment of RPS on expression comparing the
182   differentially expressed and replicated case-control features to the rest of the transcriptome, as
183   well as lack of directionally consistency between RPS- and diagnosis-associated statistics
184   among expressed features (Table S12).These results further suggest that the significant case-
185   control expression differences show little overlap with genetic risk for the disorder.

186       In an earlier study of the epigenetic landscape of frontal cortex of patients with
187   schizophrenia, we showed that DNA methylation levels in patients were closer to fetal
188   methylation levels than to those of adult control samples[17]. Here we tested for analogous effects
189   in the RNA-seq data related to the illness state. Every significant gene with differentially
190   expressed features in the adult case-control analysis and replicated in the independent dataset
191   showed evidence for developmental regulation across at least two expression feature types. We
192   further found that the expression features normally more highly expressed in postnatal life
193   tended to be more lowly expressed in patients compared to controls (max: $p=3.24 \times 10^{-11}$, min:
194   $p=1.05 \times 10^{-70}$, Figure 2C) and features more highly expressed in fetal life tended to be more
195   highly expressed in patients with schizophrenia compared to controls regardless of
196   summarization type (max: $p=6.86 \times 10^{-33}$, min: $p < 10^{-100}$, Figure 2D). Analogous analyses for
197   developmental regulation of schizophrenia-associated features without adjusting for the RNA
198   quality qSVs were significant in the opposite directions, namely that schizophrenia-associated
199   changes were further from, rather than closer to, fetal expression levels, which would be
200   predicted as an artifact of residual RNA quality confounding (as the quality of the samples rank
201   as fetal > adult control > adult SZ, see Table S1).These results further converge on a role for
202   genes changing during brain development and maturation in schizophrenia, specifically that
203   both DNA methylation and expression levels in adult patients appear to reflect levels in the
204   developing brain more strongly than do those of unaffected individuals.  These results also
205   underscore the risk of spurious findings based on uncorrected RNA quality confounding.

206

207   *Clinical enrichment of eQTL associations for schizophrenia*

208    In order to elucidate the RNA features associated with schizophrenia risk variants
209    themselves, rather than LD regions, we first performed a genome-wide *cis* (<500kb) expression
210    quantitative trait loci (eQTL) analysis within the 412 post-adolescent subjects (see Methods).
211    We then replicated SNP-feature eQTL pairs in two independent datasets (CMC and GTEx) to
212    form a significant and replicated set of high-confidence eQTLs ("core eQTLs") to interrogate for
213    clinical risk. Overall, we confirm widespread genetic control of expression in human brain
214    (Figure 3A, Table S13, Table S14) which can be accessed by the research community via a
215    user friendly browser at eqtl.brainseq.org. The majority of these core eQTLs were largely gene-
216    specific, consistent across both interrogated races (Results S5, Figure 3B), very proximal to the
217    TSS of genes (Figure 3C), and expression of the majority of genes was associated with more
218    than 1 LD-independent SNP. We further found the largest effect sizes on average for the
219    junction and ER eQTLs (Figure 3D), and interestingly a large proportion of these eQTLs
220    corresponded to unannotated transcriptional activity. While the majority of genes with features
221    as eQTLs have multiple possible Ensembl transcript isoforms, 67.0% of exon eQTLs and 31.1%
222    of junction eQTLs were specific to a single Ensembl transcript when multiple transcripts existed
223    for a given gene. Together these eQTL results, based on analyses across multiple feature
224    summarizations from RNA sequencing data, demonstrate more widespread genetic regulation
225    of expression than previously reported, including extensive transcript specificity and regulation
226    of unannotated sequence in the human brain.  These observations offer novel insights into cis
227    mechanisms of common variation in brain.

228    We next explored the landscape of eQTLs associated with genetic risk using the set of
229    core eQTLs among the PGC2-significant risk SNPs. There were 46 (42.5%) risk SNPs
230    significantly associated with expression levels of 764 nearby features at FDR < 1% significance
231    (Table S15, of the 108/128 PGC2 "index" SNPs that were present in our data).  This proportion
232    of SNPs was highly significant (p=$4.6 \times 10^{-9}$) relative to the 19.7% of all tested common variants
233    with eQTL signal. These risk eQTL features mapped back to 134 Ensembl genes (and 116 gene
234    symbols). Interestingly, many of these adult-identified eQTLs show similar genetic regulation in
235    fetal brain, analogous to observations made among methylation quantitative trait loci (meQTLs)
236    at these risk SNPs[17,18]. In addition, the majority of these genetic risk-associated expression
237    features were developmentally regulated across the lifespan in our control samples (N=528,
238    69.1%), suggesting that different developmental time periods may have different magnitudes of
239    eQTL effects among these risk-associated expressed features. Even more strikingly, 57 genes
240    (42.5%) having risk SNP eQTLs had shifting isoform usage across pre- and post-natal life, a
241    highly enriched proportion relative to all expressed genes (9,205/43,046 genes with isoform
242    shifts, OR=2.73, p=$4.21 \times 10^{-9}$), further highlighting developmental relevance for schizophrenia
243    genetic risk.

244

245    *Identifying risk transcripts through clinical and molecular convergence*

246    In principle, transcripts with convergence of directionality of association with both genetic
247    risk and with clinical illness are especially attractive as putative molecular mechanisms of illness
248    and as candidates for designing pathogenic models.  We therefore used allele information and

249 the corresponding effect estimates in the GWAS (odds ratio) and eQTL analyses (change in
250 expression per allele copy) to test which of the expressed features that differentiated patients
251 and controls were directionally consistent with the allelic directionality of genetic risk. In other
252 words, if a given allele at a risk associated SNP is associated with increased expression of a
253 nearby feature via the eQTL analysis, we tested whether expression of the same feature was
254 more likely to be higher in patients compared to controls (Results S6, Figure 4A). Indeed, we
255 found two-fold enrichment of this directional consistency between the GWAS, case-control, and
256 FDR-significant eQTLs ($p=4.2x10^{-6}$). In order to reduce the chances of off-target eQTLs being
257 driven by correlated expression features, we lastly retained SNP-feature associations that were
258 at least half as significant on the $\log_{10}$ scale as the best eQTL for each GWAS index SNP (e.g. if
259 the best eQTL association for a risk SNP was $p < 1x10^{-10}$, then eQTLs with $p > 1x10^{-5}$ were
260 excluded).

261      This approach resulted in core eQTL signal across 34 GWAS loci with convergence of
262 allele and illness state expression directionality (encompassing 59 genes across 231 features,
263 Table S16). These are transcript features and genes that show consistent directionality of
264 association of genetic risk alleles and of illness state, implicating these as particularly likely
265 pathogenic factors underlying the GWAS loci. In the majority of these loci (N=22, 64.7%) the
266 eQTLs represented a single Ensembl gene, nominating it as the most well supported "risk
267 gene", based on schizophrenia genetic risk coupled to expression features in the polyA+
268 fraction of postmortem DLPFC (Figure 4 and Figure S9). Of the remaining 12 GWAS loci that
269 did not associate with features in only a single gene, 6 (17.6%) were similarly associated with
270 the expression of two nearby genes, which may indicate different candidate genes whose
271 expression is driven by different chromosomal haplotypes (Table 2). The effects of the
272 remaining risk SNPs were more difficult to untangle, each associating with similar significance
273 levels with three (N=2), four (N=3) or seven (N=1) genes. Finally, there were transcript-specific
274 signals for many of the loci, including individual isoforms among multi-gene risk eQTLs. For
275 example, there were 40 genes with transcript annotation (excluding gene- and region-level
276 summarizations) and 25 genes contained features tagging single Ensembl transcripts (Table
277 S17). This transcript specificity cautions that many gene count- or microarray-based approaches
278 may miss important risk gene signals in eQTL analyses (see Discussion). These eQTL results
279 highlight significant and independently-replicated risk-associated schizophrenia candidate
280 genes and their specific transcripts that comprise links in the causative chain of schizophrenia in
281 the human brain.

282

## Discussion

284      We have explored the diverse landscape of expression correlates of schizophrenia risk
285 and illness state in the postmortem human frontal cortex. Using deep RNA sequencing to define
286 convergent measures of gene expression and early brain development, we identified
287 widespread developmental regulation of transcription, including novel discoveries related to
288 preferential isoform usage across brain development. These isoform "shifts" were associated
289 with genetic risk for schizophrenia, and the directionality of dysregulation of developmentally

290    regulated features suggest a more fetal-like expression profile in patients with schizophrenia
291    compared to controls. Our approach to transcript characterization, which included extensive
292    characterization of unannotated sequence, revealed that many more schizophrenia risk
293    associated SNPs are brain eQTLs than previously reported - many risk SNPs only associate
294    with a single gene, or even a single transcript, and many of these adult-identified eQTLs show
295    similar genetic regulation in the fetal brain. Lastly, we identified significant and replicated genes
296    differentially expressed in patients with schizophrenia compared to unaffected controls using a
297    new experiment-based statistical framework to estimate and reduce the effects of latent RNA
298    degradation bias. Without this new approach to RNA quality adjustment, replication across
299    datasets is markedly limited if not negligible, and the directionality of the association with
300    developmental isoform shifts is anomalous. These data suggest a convergence of
301    developmental regulation and genetic risk for schizophrenia that appears relatively stable in
302    patients ascertained at death, following decades of illness after diagnosis.  We previously
303    observed analogous stability of epigenetic marks highlighting prenatal life in adult patients with
304    schizophrenia[17], suggesting that both genetic and environmental risk factors implicated in
305    schizophrenia illness involve early developmental events which are still observable in the brain
306    tissue of adult individuals who have been ill for many years.

307         While our approach utilizing convergent expression features – genes, exons, transcripts,
308    junctions, and expressed regions – results in more complicated data processing and analysis, it
309    can potentially cast a wider net in the search for valid biological signals in RNA sequencing
310    datasets. Using all convergent features overcomes the limitations related to any given feature
311    summarization, including the inability to measure and interrogate unannotated or novel
312    transcribed sequence using gene and exon counts, and the difficulties in full transcript assembly
313    from short sequencing reads [19]. We note that both quantifying and analyzing splice junctions,
314    and also transcript-level data, rely on junction-spanning reads for statistical power. In our data,
315    there were approximately 3 times (IQR: 2.86-3.24) more reads available by gene/exon counting
316    approaches than those that contain splice junctions, likely explaining why gene counts
317    discovered more differentially expressed genes in the schizophrenia diagnosis analyses. Two
318    relatively new approaches utilized here – direct quantification and statistical analyses of splice
319    junction counts and expressed regions – can identify differential expression signal when it is
320    outside of the annotated transcriptome. The junction-level approach can also identify previously
321    uncharacterized novel transcribed sequence, which we replicated in other large publicly
322    available datasets, as well as delineate individual transcripts or classes of transcripts that share
323    a particular splice junction, and as read lengths increase, the proportion of reads containing
324    splice junctions will increase, making junction- and transcript-based approaches even more
325    powerful, including those recently developed to identify splicing QTLs [20].

326         Our analysis of RNA-seq data identified widespread shifts in preferential isoform use
327    across brain development, which would have been impossible to identify using only gene-level
328    data and incomplete with only exon-level data (Figure 4). The genes with these isoform shifts
329    were significantly enriched for neurodevelopmental and cellular signaling processes, and as well
330    as for genes in regions of genetic risk for schizophrenia. A prevalent hypothesis suggests that
331    schizophrenia is a neurodevelopmental disorder that arises because of altered connectivity and
332    plasticity in the early assembly of relevant neural circuits[1], and the potential convergence of

333    genetic risk with developing signaling processes across human brain development should  point
334    to specific candidate molecular disruptions occurring during the wiring of the fetal brain. Indeed,
335    inefficient or disrupted signaling and tuning is thought to underlie the expression of illness in the
336    adult brain [1], and the most successful therapeutics work through improving these processes[14].
337    Consistent with this hypothesis, we find evidence for differences in the expression of genes
338    coding for subunits of ion channels in the cortices of patients with schizophrenia compared to
339    controls. We observed significant differential expression of both voltage-gated (*KCNA1*, *KCNC3*,
340    KCNK1, *KCNN1*, *SCN9A*) and ligand gated ion channels (*GRIN3A*, *GABRA5*, *GABRB3*),
341    transporters (*SLC16A2*, *ALC25A33*, *SLC26A11*, *SLC35F2*, *SLC7A3*), and ion channel auxiliary
342    subunits (*KCNIP3*, *SCN1B*), supporting other evidence that the clinical phenomenology of
343    schizophrenia is associated with altered neuronal excitability [21].  While these findings implicating
344    basic mechanisms of cortical circuit dynamics may underlie fundamental aspects of the clinical
345    disorder, the possibility that they are driven by the effects of pharmacological treatment and are
346    thus state dependent epiphenomena cannot be excluded.

347          Our eQTL analyses are among the largest and most comprehensive to date in human
348    brain tissue, requiring stringent genome-wide significance and independent replication, and offer
349    additional insights into the genetic regulation of RNA expression levels. We show that over half
350    of the genes with eQTL signals are associated with three or more nearby LD-independent
351    SNPs, suggesting a potential mechanism to explain how expression levels at the population-
352    level are relatively normally distributed while so many genes show distinct eQTL signals. Our
353    data also suggest more widespread regulation of specific transcript isoforms, which we were
354    able to identify using exon- and junction-level analyses. This transcript-specific genetic
355    regulation was particularly prevalent among schizophrenia risk variants, where 62.5% of loci
356    containing multiple transcripts showed clinically- and molecularly-consistent eQTL signal to a
357    single Ensembl transcript isoform. Overall, we have identified many more eQTLs to genome-
358    wide significant schizophrenia risk variants – 42.5% - than previously reported, experimentally
359    implicating far more potential "risk" genes within these loci than previously characterized.

360          These eQTL associations within the genome-wide significant schizophrenia loci
361    generate novel putative biological mechanisms underlying risk for the disorder. We highlight 34
362    GWAS loci that contain significant eQTLs that show convergence of consistent altered
363    expression across both genetic risk and case status in the human brain. These loci often point
364    to individual "risk" genes or even more specific "risk" transcripts that can represent targetable
365    entry points for more focused cellular assays and model organism work to better characterize
366    schizophrenia risk mechanisms. Moreover, observation of convergence in directionality of
367    molecular association with genetic risk and clinical state identifies a compelling strategy and
368    directionality for target rescue, specifically to increase or decrease the function of the target
369    transcript(s) and downstream effectors. Focusing solely on increased or decreased expression
370    in brains of patients compared to controls, without considering genetic risk variants and their
371    regulation of local gene expression, will likely predominantly highlight molecular changes
372    resulting from the schizophrenia illness state, as we suggest with consistent down-regulation of
373    ion channels.

374    We stress the priority of identifying the most relevant cellular consequences of genetic
375    risk, which we view as production of particular isoforms with predicted directionality, rather than
376    trying to identify "causal" mutations tagged by "marker" risk SNPs from the GWAS.  We suggest
377    that identifying convergence between genetic risk and potential molecular consequences of the
378    disorder is likely to result in better, or at least more consistent support for, targets for drug
379    discovery efforts.

380

381

382    **Author Contributions**

383    A.E.J – performed primary data processing and analyses, led the writing of the manuscript

384    R.E.S – contributed to data analysis and writing of the manuscript

385    J.H.S., R.T. , Y.G. – performed RNA sequencing data generation (RNA extraction, library
386    preparation, and sequencing) and QC analyses

387    L.C.T.,J.T.L – performed region-level data generation and assisted in data analysis and
388    interpretation

389    T.K.T.,S.X.,J.Q.,C.C.,B.J.M., A.C.,N.B.,BrainSeq – provided feedback on manuscript and
390    contributed to data analyses and interpretations on eQTL analyses.

391    B.U. – created user-friendly database of eQTLs

392    A.D.S., T.M.H.,J.E.K.- collected, consented, characterized, and dissected human brain tissue

393    D.R.W. – designed and oversaw the research project, wrote the manuscript

394    Tony Kam-Thong is employed by F. Hoffmann-La Roche

395    Hualin S Xi and Jie Quan are employees of Pfizer Inc.

396    Alan Cross, and Nicholas J.Brandon were full time employees and shareholders in AstraZeneca
397    at the time these studies were conducted.

398    The remaining authors declare no competing financial interests.

399

400    **Data Availability:** sequencing reads and genotype data are available through SRA and dbGaP
401    at accession numbers: [TBD] following publication**.**

402

403    **Acknowledgements:**

404    We thank Dr. Ronald Zielke, Robert D. Vigorito, and Robert M. Johnson of the National Institute
405    of Child Health and Human Development Brain and Tissue Bank for Developmental Disorders

447    Maura Furey, Wayne C. Drevets, Joo Heon Shin, Andrew E. Jaffe, Yankai Jia, Richard E.
448    Straub, Amy Deep-Soboslay, Thomas M. Hyde, Joel E. Kleinman, Daniel R. Weinberger

449

## Figure Legends

451    **Figure 1:** Developmental regulation of expression. (A) Principal component #1 of the gene-level
452    expression data versus age; PCW: post-conception weeks, remaining ages are in years. (B)
453    Expression features fall into two main development regulation signatures, increasing in
454    expression from fetal to postnatal life (orange) or decreasing from fetal to postnatal life (blue). Y-
455    axis is Z-scaled expression (to standard normal), dark lines represent median expression levels,
456    and confidence bands represent $25^{th}$-$75^{th}$ percentiles of expression levels for each class of
457    features. (C) KEGG pathways enriched for genes with isoform shifts, stratified by which feature
458    type identified the gene as having a switch. Coloring/scaling represents -log10(FDR) for gene
459    set enrichment. Analogous data for GO gene sets (biological processes, BP, and molecular
460    function, MF) are available in Table S6. DER: differentially expressed region. Enrichment
461    analyses for isoform shift genes among PGC2 schizophrenia GWAS risk loci with exon and
462    junction counts using both (D) parametric p-values) and (E) permutation-based p-values. OR:
463    odds ratio.

464    **Figure 2:** Differential expression comparing patients with schizophrenia to controls. (A)
465    Histogram of fold changes of the diagnosis effect of those features that were significant and
466    independently replicated, colored by feature type. (B) Gene set analyses of genes with
467    decreased expression in patients compared to controls by feature type. Coloring/scaling
468    represents -log10(FDR) for gene set enrichment. Significant directional effects of developmental
469    regulation among diagnosis-associated features for those features that (C) increased and (D)
470    decreased across development (i.e. those features shown in Figure 1B). P-values provided for
471    Wilcoxon rank sign test for those features developmentally regulated among case-control
472    differences to those not developmentally regulated.

473    **Figure 3:** Expression quantitative trait loci (eQTL) in the human brain. (A) Table summarizing
474    eQTL results across the five feature summarization types (gene, exon, transcript, junction and
475    expressed region). Bonf: Bonferroni-adjusted p-value, Ens: Ensembl, Sym: Symbol, IQR:
476    interquartile range, $25^{th}$-$75^{th}$ percentiles, bp = base pairs, kb = kilobases. (B) Similar
477    standardized effect sizes (T-statistics) for eQTLs identified in full sample within just Caucasian
478    (CAUC) and African American (AA) computed separately, here using exon-level eQTL results.
479    (C) Histogram of the distance of the most significant SNP to each eQTL expression feature. (D)
480    Boxplots of the effect sizes (absolute fold changes per minor allele copy) stratified by feature
481    type.

482    **Figure 4:** Clinical enrichment of schizophrenia risk. (A) Image describing enforcement of
483    directional consistency between GWAS, eQTL, and case-control differences. We highlight
484    example eQTLs to GWAS-positive variants among annotated exons in (B) *SRR* and (C) *STAT6*,
485    as well as unannotated expressed features, for example an intronic region in (D) *HSPD1*
486    */HSPE1*, shown as a star in (E), and a junction linked to (F) only one annotated exon of

487 IMMP2L, shown as a bar in (G). In panels E and G, dark blue: exon, light blue: intron;
488 coordinates relative to hg19.

489 **Tables**

| | | Transcript Class | | | |
|---|---|---|---|---|---|
| **Number** | **Datasets** | **In Ensembl** | **Exon Skip** | **Alt Boundary** | **Novel** |
| 251557 | Neither | 2680(1.3%) | 118(1.4%) | 315(1.6%) | 4613(24.8%) |
| | Geuvadis | 134(0.1%) | 3(0%) | 14(0.1%) | 116(0.6%) |
| | GTEX | 7883(3.8%) | 1257(14.8%) | 4312(22.3%) | 6768(36.3%) |
| | Both | 194342(94.8%) | 7136(83.8%) | 14730(76%) | 7136(38.3%) |

490

491 **Table 1:** splice junction annotation and characterization in GTEx and GEUVADIS for moderately
492 expressed junctions (mean reads per 80M mapped reads, RP80M > 1). Each column
493 represents a 2x2 table for presence of identified junctions in 495 DLPFC samples in two
494 independent polyA+ datasets.

495

496

| GWAS Rank | Risk Direction | Gene | Feature Type | Best p-value |
|---|---|---|---|---|
| 3 | UP | *AS3MT* | E | 6.73E-40 |
| 3 | UP | *BORCS7*[+] | G,E,R | 9.23E-34 |
| 3 | UP | *AS3MT** | J | 9.57E-25 |
| 7 | DOWN | *FTSJ2* | E,J,G,R | 2.64E-07 |
| 8 | DOWN | *GNL3LP1* | E | 3.11E-10 |
| 8 | DOWN | *ERCC8* | E,R | 6.22E-10 |
| 9 | DOWN | *ARL6IP4* | J | 4.12E-36 |
| 16 | UP | *IMMP2L* | J,T | 6.04E-10 |
| 17 | UP | *SNX19* | J | 9.77E-07 |
| 23 | UP | *C2orf82* | T | 6.33E-06 |
| 24 | DOWN | *NRGN* | R | 3.5E-05 |
| 24 | DOWN | *VSIG2* | R | 8.6E-05 |
| 31 | UP | *GOLGA2P7* | E,R | 5.12E-30 |
| 31 | UP | *GOLGA6L5P* | R | 8.13E-24 |
| 33 | UP | *HSPD1* | R | 3.34E-14 |
| 34 | UP | *XPNPEP3* | E,J,R | 3.63E-08 |
| 42 | DOWN | *PCCB* | G,E,R | 2.95E-11 |
| 47 | DOWN | *SRR* | E,G,R,J | 2.1E-12 |
| 51 | DOWN | *GATAD2A* | R | 2.51E-08 |
| 51 | DOWN | *NDUFA13* | J | 9.36E-06 |
| 52 | UP | *VPS45* | R | 1.88E-16 |
| 57 | UP | *RERE* | E | 7.7E-07 |
| 73 | UP | *CTNNA1* | E | 5.18E-07 |
| 84 | UP | *GPM6A* | R | 3.12E-06 |
| 94 | UP | *PRMT7* | R,J,E | 5.37E-09 |
| 94 | UP | *TSNAXIP1* | G,R | 8.78E-06 |
| 96 | DOWN | *ATPAF2* | E | 4.58E-07 |
| 103 | UP | *STAT6* | E,J | 1.71E-07 |
| 106 | DOWN | *SOX2-OT* | R | 8.18E-12 |
| 106 | DOWN | *RP11-275H4.1* | T | 2.62E-06 |
| 117 | DOWN | *FANCL* | E,R | 2.45E-17 |
| 120 | DOWN | *BRINP2* | R | 6.7E-06 |
| 121 | UP | *CD46* | E,R,G,J | 1.57E-37 |
| 125 | UP | *IRF3* | J,E | 4.51E-10 |

497  **Table 2:** eQTL and clinical directional consistent signal around the GWAS risk variants for
498  schizophrenia. GWAS rank is relative to Supplementary Table 2 in Ripke et al [2]. Direction refers
499  to whether schizophrenia risk represents increased or decreased expression of the particular
500  features. Feature type indicates which of (G)enes, (E)xons, (J)unctions, (T)ranscripts, or
501  expressed (R)egions showed eQTL association to the risk SNP. The best p-value for the risk
502  SNP to features in each gene are recorded in the last column. [+]*C10orf32* in Ensembl v75 but
503  updated in later versions; * unannotated junction found to tag a new transcript of the gene,
504  AS3MT[d2d3], that has been more fully characterized by Li et al[22].

505

506

## References

1    Weinberger, D. R. & Levitt, P. in *Schizophrenia*    393-412 (Wiley-Blackwell, 2011).

2    Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-+, doi:10.1038/nature13595 (2014).

3    Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190-1195, doi:10.1126/science.1222794 (2012).

4    Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from RNA-seq reads. *Nature biotechnology* **33**, 290-295, doi:10.1038/nbt.3122 (2015).

5    Nellore, A. *et al.* Human splicing diversity across the Sequence Read Archive. *bioRxiv*, doi:10.1101/038224 (2016).

6    Collado Torres, L. *et al.* Flexible expressed region analysis for RNA-seq with derfinder. *Nucleic Acids Research* **In Press.**, doi:10.1101/015370 (2016).

7    Jaffe, A. E. *et al.* Developmental regulation of human cortex transcription and its clinical relevance at single base resolution. *Nat Neurosci* **18**, 154-161, doi:10.1038/nn.3898 (2015).

8    Tan, W. *et al.* Molecular cloning of a brain-specific, developmentally regulated neuregulin 1 (NRG1) isoform and identification of a functional promoter variant associated with schizophrenia. *The Journal of biological chemistry* **282**, 24343-24351, doi:10.1074/jbc.M702953200 (2007).

9    Kao, W. T. *et al.* Common genetic variation in Neuregulin 3 (NRG3) influences risk for schizophrenia and impacts NRG3 expression in human brain. *Proceedings of the National Academy of Sciences of the United States of America* **107**, 15619-15624, doi:10.1073/pnas.1005410107 (2010).

10   Tao, R. *et al.* Expression of ZNF804A in human brain and alterations in schizophrenia, bipolar disorder, and major depressive disorder: a novel transcript fetally regulated by the psychosis risk variant rs1344706. *JAMA psychiatry* **71**, 1112-1120, doi:10.1001/jamapsychiatry.2014.1079 (2014).

11   Hyde, T. M. *et al.* Expression of GABA signaling molecules KCC2, NKCC1, and GAD1 in cortical development and schizophrenia. *The Journal of neuroscience : the official journal of the Society for Neuroscience* **31**, 11088-11095, doi:10.1523/JNEUROSCI.1234-11.2011 (2011).

12   Birnbaum, R., Jaffe, A. E., Hyde, T. M., Kleinman, J. E. & Weinberger, D. R. Prenatal expression patterns of genes associated with neuropsychiatric disorders. *The American journal of psychiatry* **171**, 758-767, doi:10.1176/appi.ajp.2014.13111452 (2014).

13   Buchanan, R. W. & Carpenter, W. T. Domains of psychopathology: an approach to the reduction of heterogeneity in schizophrenia. *The Journal of nervous and mental disease* **182**, 193-204 (1994).

14   Winterer, G. & Weinberger, D. R. Genes, dopamine and cortical signal-to-noise ratio in schizophrenia. *Trends in neurosciences* **27**, 683-690, doi:10.1016/j.tins.2004.08.002 (2004).

15   Jaffe, A. E. *et al.* A framework for RNA quality correction in differential expression analysis. *bioRxiv*, doi:10.1101/074245 (2016).

16   Schizophrenia Working Group of the Psychiatric Genomics, C. Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-427, doi:10.1038/nature13595 (2014).

17   Jaffe, A. E. *et al.* Mapping DNA methylation across development, genotype and schizophrenia in the human frontal cortex. *Nature neuroscience* **19**, 40-47, doi:10.1038/nn.4181 (2016).

555　18　Hannon, E. *et al.* Methylation QTLs in the developing brain and their enrichment in
556　　　schizophrenia risk loci. *Nature neuroscience* **19**, 48-54, doi:10.1038/nn.4182 (2016).
557　19　Steijger, T. *et al.* Assessment of transcript reconstruction methods for RNA-seq. *Nature*
558　　　*methods* **10**, 1177-1184, doi:10.1038/nmeth.2714 (2013).
559　20　Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease.
560　　　*Science* **352**, 600-604, doi:10.1126/science.aad9417 (2016).
561　21　Uhlhaas, P. J. & Singer, W. Abnormal neural oscillations and synchrony in
562　　　schizophrenia. *Nature reviews. Neuroscience* **11**, 100-113, doi:10.1038/nrn2774 (2010).
563　22　Li, M. *et al.* A human-specific AS3MT isoform and BORCS7 are molecular risk factors in
564　　　the 10q24.32 schizophrenia-associated locus. *Nature medicine*, doi:10.1038/nm.4096
565　　　(2016).

566

567

568　**Methods**

569

570　<u>Postmortem brain samples</u>
571　Post-mortem human brain tissue was obtained by autopsy primarily from the Offices of the Chief
572　Medical Examiner of the District of Columbia, and of the Commonwealth of Virginia, Northern
573　District, all with informed consent from the legal next of kin (protocol 90-M-0142 approved by the
574　NIMH/NIH Institutional Review Board). Additional post-mortem fetal, infant, child and adolescent
575　brain tissue samples were provided by the National Institute of Child Health and Human
576　Development Brain and Tissue Bank for Developmental Disorders (http://www.BTBank.org)
577　under contracts NO1-HD-4-3368 and NO1-HD-4-3383. The Institutional Review Board of the
578　University of Maryland at Baltimore and the State of Maryland approved the protocol, and the
579　tissue was donated to the Lieber Institute for Brain Development under the terms of a Material
580　Transfer Agreement. Clinical characterization, diagnoses, and macro- and microscopic
581　neuropathological examinations were performed on all samples using a standardized paradigm,
582　and subjects with evidence of macro- or microscopic neuropathology were excluded. Details of
583　tissue acquisition, handling, processing, dissection, clinical characterization, diagnoses,
584　neuropathological examinations, RNA extraction and quality control measures were described
585　previously in Lipska, et al. [23]. The Brain and Tissue Bank cases were handled in a similar
586　fashion (http://medschool.umaryland.edu/BTBank/ProtocolMethods.html). Antipsychotic use
587　was measured using toxicology at time of death.
588

589　<u>RNA extraction and sequencing</u>
590　Post-mortem tissue homogenates of dorsolateral prefrontal cortex grey matter (DLPFC)
591　approximating BA46/9 in postnatal samples and the corresponding region of PFC in fetal
592　samples were obtained from all subjects. Total RNA was extracted from ~100 mg of tissue using
593　the RNeasy kit (Qiagen) according to the manufacturer's protocol. The poly-A containing RNA
594　molecules were purified from 1 µg DNAse treated total RNA and sequencing libraries were
595　constructed using the Illumina TruSeq© RNA Sample Preparation v2 kit. Sequencing
596　indices/barcodes were inserted into Illumina adapters allowing samples to be multiplexed in

597 across lanes in each flow cell. These products were then purified and enriched with PCR to
598 create the final cDNA library for high throughput sequencing using an Illumina HiSeq 2000 with
599 paired end 2x100bp reads.

600

601 <u>RNA sequencing data processing</u>

602

603 The Illumina Real Time Analysis (RTA) module performed image analysis, base calling, and the
604 BCL Converter (CASAVA v1.8.2), generating FASTQ files containing the sequencing reads.
605 These reads were aligned to the human genome (UCSC hg19 build) using the spliced-read
606 mapper TopHat (v2.0.4) using the reference transcriptome to initially guide alignment, based on
607 known transcripts of the previous Ensembl build GRCh37.67 (the "–G" argument in the
608 software) [24]. We achieved a median of 85.3 million (IQR: 71.7M-111.2M) aligned reads per
609 sample (see Table S1).

610

611 We characterized the transcriptomes of these 495 samples using five convergent
612 measurements of expression ("feature summarizations")– (1) gene and (2) exon counts, and (3)
613 transcript-level quantifications that rely on existing gene annotation, and two annotation-
614 agnostic approaches we have developed that are determined solely from the read alignments –
615 (4) read coverage supporting exon-exon splice junctions (e.g. coordinates of potentially intronic
616 sequence that are spliced out of mature transcripts captured by a single read) and (5) read
617 coverage overlapping each base in each sample which we have summarized into contiguous
618 "expressed regions" (ERs, see Methods, Figure S1). These last three measurements generate
619 expression for features of interest that can "tag" elements of transcripts in the data that are not
620 constrained by limitations or incompleteness of existing annotation, and the counts for these
621 features can then be directly used for differential expression analysis.
622 1. Gene counts were generated using the featureCounts tool[25] based on the more recent
623    Ensembl v75, which was the last stable release for the hg19 genome build, using single end
624    read counting [featureCounts –a $GTF –o $OUT $BAM]. We converted counts to RPKM
625    values using the total number of aligned reads across the autosomal and sex chromosomes
626    (dropping reads mapping to the mitochondria chromosome).
627 2. Exon counts were also generated using the featureCounts tool[25] based on the more recent
628    Ensembl v75, using single end read counting, and allowing reads to be assigned to multiple
629    exons (e.g. those with splice junctions) [featureCounts –O –f –a $GTF –o $OUT $BAM]. We
630    converted counts to RPKM values using the total number of aligned reads across the
631    autosomal and sex chromosomes (dropping reads mapping to the mitochondria
632    chromosome).
633 3. Junction counts were generated by first filtering the TopHat BAM file to primary alignments
634    only [samtools view -bh -F 0x100 $BAM > $NEWBAM ] and regtools [26] was used to extract
635    analogous junction information (coordinates and number of reads supporting) as the TopHat
636    output. We found that native TopHat output (junctions.bed) was based on both primary and
637    secondary alignments, which could influence the degree of potentially novel splice junctions.
638    We used a modified version of TopHat's "bed_to_juncs" program to retain the number of
639    supporting reads (in addition to returning the coordinates of the spliced sequence, rather

640     than the maximum fragment range), and used R code (see Supplementary Code) to

641     combine and annotate these junctions across all samples. We annotated identified splice

642     junctions using Ensembl v75 – while the initial alignment was guided by Ensembl v67, novel

643     junctions, by definition, are identified in the second genome alignment, rather than the initial

644     guided transcriptome alignment step. We converted counts to "RP80M" values, or "reads per

645     80 million mapped" using the total number of aligned reads across the autosomal and sex

646     chromosomes (dropping reads mapping to the mitochondria chromosome), which can be

647     interpreted as the number of reads supporting the junction in an average library size (we

648     were targeting 80M reads in the sequencing). Most junctions were lowly expressed in our

649     homogenate tissue, with fewer than 1 average normalized supporting read (N=3,330,642;

650     92.98%) including approximately half unique to a single individual (N= 1,779,241, 49.67%).

651 4.  Transcripts were assembled using StringTie[4] (version 1.1.2) guided by Ensembl v75

652     annotation within each sample [stringtie $BAM –o $OUT –G $GTF]. We then used

653     "CuffMerge" [27] to merge all assembled transcriptomes across all samples, and then re-

654     quantified the expression of each transcript isoform in each sample again using StringTie to

655     this global set of transcripts [stringtie $BAM –B –e –o $OUT –G $GTF_ALL] to have

656     expression measurements on the same transcripts across all samples. We then used the

657     "ballgown" tool[28] to merge all assembled and quantified transcripts across all samples (N=

658     733,339), and used liberal filtering to remove lowly or uniquely expressed transcripts (mean

659     FPKM > 0.025), resulting in 188,578 transcripts across the 495 samples.

660 5.  Expressed regions (ERs) were calculated using the "derfinder" R Bioconductor package[6]

661     using a cutoff of 5 normalized (to 80M reads) read coverage, which identified 389,797 ERs.

662     We retained the 275,885 ERs that were at least 12 basepairs, and annotated the ERs to

663     Ensembl v75.

664

665 <u>Genotype data processing</u>

666 SNP genotyping with HumanHap650Y_V3 (N=135), Human 1M-Duo_V3 (N=357), and Omni5

667 (N=3) BeadChips (Illumina, San Diego, CA) was carried out according to the manufacturer's

668 instructions with DNA extracted from cerebellar tissue. Genotype data were processed and

669 normalized with the crlmm R/Bioconductor package[29] separately by platform. Genotype

670 imputation was performed on high-quality observed genotypes (removing low quality and rare

671 variants) using the prephasing/imputation stepwise approach implemented in IMPUTE2[30] and

672 Shape-IT[31], with the imputation reference set from the full 1000 Human Genomes Project Phase

673 3 data set, separately by platform. We retained common SNPs (MAF > 5%) that were present in

674 the majority of samples (missingness < 10%) that were in Hardy Weinberg equilibrium (at p >

675 $1x10^{-6}$) using the Plink[32] version 1.9 tool kit [`plink --bfile $BFILE --geno 0.1 --maf 0.05 --hwe

676 0.000001`]. We then identified linkage disequilibrium (LD)-independent SNPs to use in genome-

677 wide clustering of samples and in the number of independent eQTL tests performed [`plink –

678 bfile $BFILE --indep 100 10 1.25`]. Multidimensional scaling (MDS) was performed on the

679 autosomal LD-independent construct genomic ancestry components on each sample, which can

680 be interpreted as quantitative levels of ethnicity – the first component separated the Caucasian

681 and African American samples. This processing and quality control steps resulted in 7,421,423

682 common variants in this dataset of 495 subjects.

683

684    Risk polygene score (RPS) analysis

685    We computed risk polygene scores (RPS) using the allelic dosage files following imputation
686    described above and the SNPs from provided by the PGC to the Lieber Institute that did not
687    contain completely different clinical subjects used in the GWAS[2]. We considered expression
688    associations at the gene, exon and junction-level to the RPS scores from the first 5 clinical SNP
689    sets, corresponding to GWAS p-value thresholds of $p < 5e-8$ (s1), $p < 1e-6$ (s2), $p < 1e-4$ (s3), $p$
690    $< 0.001$ (s4), and $p < 0.01$ (s5) – subsequent SNP sets were ignored due to clinical risk
691    plateauing at s5. We also focused only on Caucasian individuals (96 cases, 113 controls), as
692    the s5 RPS was increased in patients relative to controls in this sample ($p=3.2x10^{-5}$), but did not
693    differ among African Americans ($p=0.9$). Within each expression feature type, we modeled
694    expression levels as a function of each RPS set (s1-s5), adjusting for 3 MDS components of the
695    genotype data, sex, and the first $K$ principal components (PCs) of the normalized expression
696    features, where $K$ was calculated using the Buja and Eyuboglu permutation-based algorithm[33] in
697    the "sva" Bioconductor package[34]. The resulting p-values of RPS on expression, adjusting for
698    the above factors, were subject to false discovery rate (FDR) control to account for multiple
699    testing.

700

701    Public data processing

702    *GTEx*: Raw RNA-seq reads from all brain samples with corresponding genotype data were
703    downloaded from SRA and aligned to the genome using TopHat2 [24] (version 2.0.14) using the
704    iGenomes transcriptome and genome annotations based on hg19. As above, featureCounts [25]
705    was used to quantify expression of genes and exons relative to Ensembl v75, and junctions
706    were quantified with regtools[26] as above. We used StringTie with the assembled merged GTF
707    from the LIBD DLPFC samples on the GTEx BAM files to quantify the same transcripts, and
708    used bwtool to quantify the coverage of the same expressed regions from the GTEx brain
709    samples. Genotype data from the two platforms (Illumina Omni 5M and 2.5M) were imputed
710    separately as described above and merged into a single plink[32] set.

711    *GEUVADIS*: Raw RNA-seq reads from all LCL samples were downloaded from SRA and
712    aligned to the genome using TopHat2 [24] (version 2.0.9) using the iGenomes transcriptome and
713    genome annotations based on hg19. As above, featureCounts[25] was used to quantify
714    expression of genes and exons relative to Ensembl v75, and junctions were quantified with
715    regtools[26] as above. We used StringTie with the assembled merged GTF from the LIBD DLPFC
716    samples on the GEUVADIS BAM files to quantify the same transcripts, and used bwtool to
717    quantify the coverage of the same expressed regions from the GEUVADIS LCL samples.

718    *CommonMind Consortium (CMC)*: 547 BAM files were downloaded from Synapse, which were
719    aligned with TopHat2 (version 2.0.9) using Ensembl v70 transcriptome annotation and the hg19
720    genome. As above, featureCounts [25] was used to quantify expression of genes and exons
721    relative to Ensembl v75, and junctions were quantified with regtools [26] as above. We used
722    StringTie with the assembled merged GTF from the LIBD DLPFC samples on the CMC BAM
723    files to quantify the same transcripts, and used bwtool to quantify the coverage of the same

724     expressed regions from the CMC brain samples. Genotypes were converted to plink file sets
725     from GEN files obtained from Synapse using posterior probabilities > 90%, resulting in genotype
726     data across 9,506,038 SNPs and 547 samples.

727

728     <u>Differential expression across brain development</u>

729     We modeled differential expression across age at each of the five feature summarizations
730     (gene, exon, junction, transcript, and ER) in the 320 control subjects across the lifespan. We
731     modeled expression, after transforming with log2 with an offset of 1, as a function of age after
732     creating using linear splines with breakpoints at ages: birth (0), 1, 10, 20, and 50, further
733     adjusting for sex and ancestry/ethnicity (first 3 MDS components). F-statistics were computed
734     comparing the model containing age (including the linear splines), sex, and ethnicity, to a
735     statistical model with just sex and ethnicity, with corresponding p-values calculated based on an
736     F-distribution with 11 and 308 degrees of freedom, and Bonferroni adjustment within each
737     feature type was performed using the number of features with non-zero expression (gene
738     RPKM > 0.01, exon RPKM > 0,1, and junction RP80M > 0.2 with non-novel annotation) across
739     all samples as the number of tests (which varied by feature type). We also computed post-hoc
740     statistics on the data, including the Pearson correlation between "cleaned" expression (after
741     regressing out the effects of sex and ethnicity, holding the age effects constant), and age to
742     determine if the expression of the fetal rose or fell across the lifespan, and also measured the
743     fetal versus postnatal $log_2$ fold changes.

744     Preferential isoform usage across aging was determined by identifying the subset of genes (by
745     Ensembl ID) that contained at least one Bonferroni-significant feature that had positive
746     correlation with age and another Bonferroni-significant feature that had negative correlation with
747     age. We also computed the difference in positive and negative correlations as a measure of the
748     magnitude of the preferential isoform use. Gene set analyses using pre-defined gene ontology
749     (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) sets were performed using the
750     clusterProfiler R/Bioconductor package[35], here using the genes (mapping from Ensembl to
751     Entrez ID) that had such preferential isoform use to those that were developmentally regulated
752     (having at least one feature that was associate with age at Bonferroni significance).
753     Enrichments with the PGC2 schizophrenia risk loci – defined by the chr:start-end roughly
754     corresponding to linkage disequilibrium blocks in the published manuscript - were performed
755     both parametrically, by overlapping the genomic coordinates of the 108 risk regions with those
756     genes that had preferential isoform usage, compared to a background of all genes with each set
757     of expressed features, as well as by permuting the locations of the 108 regions across the
758     genome 10,000 times and each time, re-computing the overlap within these null regions – see
759     additional details in Jaffe et al 2015[7]. Empirical p-values were calculated by counting the
760     number of the odds ratios across the 10,000 null permutations to each observed odds ratio.

761

762     <u>Schizophrenia differential expression analyses</u>

763 *Discovery dataset analysis*: we first filtered the subjects with RNA-seq to retain a more stringent
764 set of 155 SCZD cases and 196 controls (criteria: ages between 17-80, gene assignment rate >
765 0.5, mapping rate > 0.7, RIN > 6, not outlying on 2nd ancestry PC, only self-reported
766 Caucasians and African Americans). We fit three statistical models across each of the
767 expression summarizations, modeling $\log_2$ transformed expression (with an offset of 1) as a
768 function of:

769 (1) Adjusted ("_adj" suffix in supplementary tables): SCZD diagnosis, adjusting for age, sex,
770 ancestry (SNP PCs 1, 5, 6, 9, 10, which were at least marginally associated with diagnosis), and
771 then observed measures related to RNA quality: RIN, mitochondrial mapping rate, and gene
772 assignment rate.

773 (2) Adjusted + Quality Surrogate Variables ("_qsva" suffix in supplementary tables): SCZD
774 diagnosis adjusting for "Adjusted" model as well as the first 12 PCs from the degradation matrix
775 (see below) based on polyA+ libraries (selected using to using the BE algorithm [33] in the sva
776 Bioconductor package[34] while providing the adjusted model as input).

777 (3) Adjusted + Principal Components ("_pca" suffix in supplementary tables): SCZD diagnosis
778 adjusting for "Adjusted" model as well as the first *k* PCs from the expressed features (using the
779 50000 most variable features) depending on the feature type (gene: 23 PCs, exon: 20 PCs,
780 transcript: 26 PCs, junction: 26 PCs, ERs: 23 PCs).

781 We used the `lmTest` and `ebayes` functions in the limma Bioconductor package [36] to fit all of
782 the statistical models to estimate $\log_2$ fold changes, moderated T-statistics, and corresponding
783 p-values. Multiple testing correction via the false discovery rate (FDR) was applied using the set
784 of expressed features in this sample set for each summarization type: 24,122 genes (mean
785 RPKM > 0.1), 420,022 exons (mean RPKM > 0.2), 61,950 transcripts (mean FPKM > 0.2),
786 229,846 junctions (mean RP80M > 1), and the 275,885 ERs.

787

788 *RNA quality correction*: We summarize the RNA quality correction approach here – for more
789 detail, see the companion paper by Jaffe et al 2016 [15]. Briefly, the quality surrogate variable
790 analysis (qSVA) uses RNA sequencing data generated from five DLPFC tissue samples left
791 unfrozen for 0, 15, 30 and 60 minutes, resulting in 20 RNA samples. These samples were
792 sequenced with both polyA+ and RiboZero library preparations, and gene, exon and junction
793 counts were derived as above. We utilized the gene-level effects of degradation in these data in
794 Figure S5 to demonstrate residual confounding by RNA quality, which we call the "DEQual Plot".

795 For a given preparation type, we identified the genomic regions most susceptible to degradation
796 by correlating coverage at expressed regions [6] to degradation time, adjusting for donor. This
797 statistical modeling identified 515 regions significantly susceptible to degradation (at Bonferroni
798 significance) in the RiboZero libraries and the top 1000 regions most susceptible to degradation
799 (among the 35,287 at Bonferroni significance) in the polyA+ libraries – the BED files for these
800 degradation-susceptible regions are available at: https://github.com/nellore/region_matrix.

801 The algorithm then involves selecting the set of regions for a particular library type and
802 calculating total coverage within each region in the new user-provided samples (e.g. the 495
803 DLPFC RNA-seq polyA+ samples) to form the degradation matrix (which is either 515 or 1000
804 rows by N samples). Then PCA is performed on the log2 transformed degradation matrix (with
805 an offset of 1) and the top $K$ PCs are selected, for example using the BE algorithm [33], and
806 extracted – the set of these PCs are referred to as quality surrogate variables (qSVs), and are
807 included as adjustment variables in subsequent differential expression analyses.

808 *Replication dataset analysis:* we performed analogous sample selection procedures as in the
809 discovery dataset to select 159 patients and 172 controls (total gene assignment rate > 0.3,
810 alignment rate > 0.8, RIN > 6, ages between 18-80, non-outlying on genetic ancestry PCs 3 and
811 5 and keeping only reported Caucasians and African Americans). We similarly fit the three sets
812 of statistical models to all five feature summarizations, with the following differences compared
813 to the discovery analysis:

814 (1) Adjusted model: the model here was diagnosis adjusting for age, sex, race, brain bank, RIN,
815 gene assignment rate, alignment rate.

816 (2) qSVA model: the degradation matrix was constructed using the 515 regions based on the
817 RiboZero libraries in the degradation experiment.

818 (3) PC adjustment: for each feature summarization type, we included: 27 PCs for genes, 29 PCs
819 for exons, 39 PCs for transcripts, 39 PCs for junctions, and 33 PCs for ERs.

820 In these replication data we did not perform FDR correction. We were using the study for
821 replication, not discovery, and therefore only used the features that were expressed in our data
822 regardless of the expression levels in CMC. We considered features independently replicated if
823 they had the same directionality for the SCZD versus control $\log_2$ fold change and were
824 marginally significant (at p < 0.05) in the CMC dataset.

825 Gene set analyses on replicated differentially expressed features and genes were performed
826 with clusterProfiler[35] as described above. Set-level analyses on features in the GWAS risk
827 regions were conducted by assigning each expressed feature a binary variable for whether it
828 was in the risk regions or not. Then we fit a linear regression model of the t-statistics for
829 diagnosis, adjusted by the qSVA approach, as a function as whether the feature was in the risk
830 region, adjusting for its average expression level. This analysis was conducted across and then
831 within each of the five feature summarization types.

832

833 <u>eQTL discovery analyses</u>

834 We performed eQTL analyses separately by feature type (gene, exon, junction, transcript, and
835 ER) allowing for a 500kb window around each of the 7,421,423 common SNPs in the 412 age >
836 13 samples, adjusting for ancestry (first three MDS components from the genotype data), sex,
837 diagnosis, and the first $K$ principal components (PCs) of the normalized expression features,
838 where $K$ was calculated separately by feature type using the Buja and Eyuboglu permutation-

839     based algorithm[33] in the "sva" Bioconductor package[34] (gene: 22 PCs, exon: 19 PCs, junction:
840     26 PCs, transcript: 25 PCs, expressed regions: 20 PCs). The eQTL analyses were run using the
841     MatrixEQTL R package[37], which returned the $\log_2$ fold change per allele copy, and
842     corresponding T-statistic, p-value, and FDR for each SNP-feature pair. We further used the LD-
843     independent SNPs to estimate the effective number of tests (by counting the number of features
844     within a 500kb window around each LD independent SNP) for a more conservative Bonferroni
845     adjustment. For all five feature types, we retained all eQTLs with FDR < 1%. We retained the
846     best SNP for each expression feature for more detailed annotation and we calculated post-hoc
847     statistics separately by reported race (Caucasian or African American) using corresponding
848     subsets of the same design matrix.

849     <u>eQTL replication analyses</u>

850     We sought to replicate the best SNP-Feature pair for each eQTL in two independent datasets:
851     CommonMind Consortium and the GTEx project. As not every single SNP imputed in the
852     discovery dataset was present in each of CMC and GTEx, we retained the most highly ranked
853     (by eQTL p-value) SNP that was present in each dataset, and tested these SNP-Feature pairs
854     for replication.

855     *CommonMind Consortium*: for the ranked eQTL lists for each feature type (gene, exon, junction,
856     transcript, and ER), we retained the best discovery SNP present in the replication dataset. The
857     majority of these tested SNP-Feature pairs used the SNP with the best discovery p-value. For
858     example, among 1,812,008 FDR-significant SNP-Gene pairs, 1,539,218 pairs were present in
859     CMC (84.9%), and among the 18,394 genes with at least one FDR-significant SNP, 14,928
860     (81.1%) of the top ranked SNPs were present, and therefore used for replication. Another 2,279
861     genes had FDR-significant SNPs more lowly ranked in the discovery set that were tested for
862     replication, and 1,187 genes did not have any FDR-significant SNPs present in CMC. The
863     proportions of present top-ranked SNPs retained for replication were similar across the different
864     feature summarizations: gene=81.1%, exon=82.1%, junction=81.5%, ER=80.6%,
865     transcript=79.4%. The number of tested SNP-feature pairs across all expression
866     summarizations is provided in Table S9.

867     We calculated the corresponding principal components (PCs) across all expression features
868     separately in the 547 samples, as in the discovery data, and in our tests for replication, modeled
869     the log2-transformed expression data as a function of additive genotype effect, controlling for
870     the top 10 expression PCs per summary type and the first 3 ancestry MDS components from the
871     genetic data. We lastly used the allele information to ensure additive effects across study
872     were relative to the same allele, flipping the directionality of such effects where the major allele
873     in one study corresponded to the minor allele in the other study. The eQTLs that were
874     directionality consistent (same allelic directionality of the additive effect on expression) and had
875     marginal significance (p<0.01) were considered replicated and retained for summarization.

876     *GTEx*: we performed analogous replication analyses as described above across the 1184 GTEx
877     brain samples that had genotype and RNA-seq data. While we calculated eQTL replication
878     statistics separately by brain region, we adjusted for the first 10 expression PCs from the entire
879     sample (across all brain regions), as well as MDS ancestry components from across the

880     genome. We only considered eQTLs from the frontal cortex brain region for replication,
881     enforcing directional consistency as the data set was much smaller than both our discovery and
882     first replication dataset (N=99).

883     After computing replication statistics in the CMC and GTEx frontal cortex datasets, the eQTLs
884     that were directionality consistent (same allelic directionality of the additive effect on expression)
885     in both datasets and had marginal significance ($p<0.01$) in CommonMind were considered
886     replicated and retained for summarization (see Table S9). Factors positively influencing
887     replication were, as expected, the statistical magnitude of the eQTL effect (absolute T-statistic),
888     the minor allele frequency, and the mean expression level.

889

890

891     Clinical enrichment analyses:

892     We first calculated the "risk" effects of the schizophrenia GWAS variants using the set of all
893     GWAS results from: https://www.med.unc.edu/pgc/files/resultfiles/scz2.snp.results.txt.gz and
894     identified which allele at each SNP corresponded to increased risk, e.g. a positive odds ratio.
895     We then used these risk alleles to convert the directionality of our eQTL effects to reflect the
896     reference allele being the risk allele (rather than the minor allele). Lastly, we incorporated the
897     qSV-adjusted case-control fold changes to identify the genetic and clinical directional consistent
898     features.  Enrichment analyses for fetal effects, isoform shift genes, and directional consistency
899     was performed using Chi-squared tests with 1 df using the `chisq.test()` in R.

900

901     **References**

902     2      Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci.
903            *Nature* **511**, 421-+, doi:10.1038/nature13595 (2014).
904     4      Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from
905            RNA-seq reads. *Nature biotechnology* **33**, 290-295, doi:10.1038/nbt.3122 (2015).
906     6      Collado Torres, L. *et al.* Flexible expressed region analysis for RNA-seq with derfinder.
907            *Nucleic Acids Research* **In Press.**, doi:10.1101/015370 (2016).
908     7      Jaffe, A. E. *et al.* Developmental regulation of human cortex transcription and its clinical
909            relevance at single base resolution. *Nat Neurosci* **18**, 154-161, doi:10.1038/nn.3898
910            (2015).
911     15     Jaffe, A. E. *et al.* A framework for RNA quality correction in differential expression
912            analysis. *bioRxiv*, doi:10.1101/074245 (2016).
913     23     Lipska, B. K. *et al.* Critical factors in gene expression in postmortem human brain: Focus
914            on studies in schizophrenia. *Biol Psychiatry* **60**, 650-658,
915            doi:10.1016/j.biopsych.2006.06.019 (2006).
916     24     Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of
917            insertions, deletions and gene fusions. *Genome Biol* **14**, R36, doi:10.1186/gb-2013-14-4-
918            r36 (2013).

919 25    Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for
920       assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930,
921       doi:10.1093/bioinformatics/btt656 (2014).
922 26    regtools v. 0.1.0 (2016).
923 27    Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals
924       unannotated transcripts and isoform switching during cell differentiation. *Nature*
925       *biotechnology* **28**, 511-515, doi:10.1038/nbt.1621 (2010).
926 28    Frazee, A. C. *et al.* Ballgown bridges the gap between transcriptome assembly and
927       expression analysis. *Nature biotechnology* **33**, 243-246, doi:10.1038/nbt.3172 (2015).
928 29    Scharpf, R. B., Irizarry, R. A., Ritchie, M. E., Carvalho, B. & Ruczinski, I. Using the R
929       Package crlmm for Genotyping and Copy Number Estimation. *J Stat Softw* **40**, 1-32
930       (2011).
931 30    Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation
932       method for the next generation of genome-wide association studies. *PLoS genetics* **5**,
933       e1000529, doi:10.1371/journal.pgen.1000529 (2009).
934 31    Delaneau, O., Coulonges, C. & Zagury, J. F. Shape-IT: new rapid and accurate
935       algorithm for haplotype inference. *BMC bioinformatics* **9**, 540, doi:10.1186/1471-2105-9-
936       540 (2008).
937 32    Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based
938       linkage analyses. *American journal of human genetics* **81**, 559-575, doi:10.1086/519795
939       (2007).
940 33    Buja, A. & Eyuboglu, N. Remarks on Parallel Analysis. *Multivariate Behavioral Research*
941       **27**, 509-540, doi:10.1207/s15327906mbr2704_2 (1992).
942 34    Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package
943       for removing batch effects and other unwanted variation in high-throughput experiments.
944       *Bioinformatics* **28**, 882-883, doi:10.1093/bioinformatics/bts034 (2012).
945 35    Yu, G., Wang, L. G., Han, Y. & He, Q. Y. clusterProfiler: an R package for comparing
946       biological themes among gene clusters. *Omics : a journal of integrative biology* **16**, 284-
947       287, doi:10.1089/omi.2011.0118 (2012).
948 36    Smyth, G. K. Linear models and empirical Bayes methods for assessing differential
949       expression in microarray experiments. *Statistical applications in genetics and molecular*
950       *biology* **3**, Article 3 (2004).
951 37    Shabalin, A. A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations.
952       *Bioinformatics* **28**, 1353-1358, doi:10.1093/bioinformatics/bts163 (2012).

953

963

**Supplementary Information**

*Supplementary Figure Legends:*

**Figure S1**: Study overview and cartoon describing quantifying the five different expression summarizations.

**Figure S2**: Cartoon describing the four different splice junction annotation classes, relative to annotated exons (dark blue rectangles). (A) Annotated splice junctions map between two exons in a known transcript. (B) Exon-skipping splice junctions map to two annotated exons in different transcripts. (C) Alternative start/exon junctions map to only one annotated exon on either the 5' or 3' end. (D) Completely novel junction do not map to any known exon.

**Figure S3**: Venn diagram of developmentally regulated features mapped back to Ensembl Gene IDs by the five feature summarization methods. DER: differentially expressed region.

**Figure S4**: Venn diagram of Ensembl Gene IDs that contain significant isoform shifts by the four feature summarization methods that allow for multiple features per gene. DER: differentially expressed region.

**Figure S5**: Degradation quality (DEqual) plots for diagnosis-associated differential expression analyses at the gene level. Y-axis: T-statistic for the effect of RNA degradation on each gene, X-axis: diagnosis T-statistic, adjusting for (left) nothing/univariate effect, (middle) observed clinical and technical confounders and (right) the inclusion of the quality surrogate variables.

**Figure S6**: Replication rates for three differential expression statistical models. X-axis: discovery p-value, Y-axis: replication rate, e.g. the proportion of features that have directional consistency and are marginally significant (at p<0.05) in independent data. Point size is proportion to the number of features called significant at the indicated p-value threshold.

**Figure S7**: Scatter plot of effect sizes (fold changes) in discovery and replication datasets for those features significant and replicated. Colors have the same legend as Figure 3A.

**Figure S8**: GWAS loci set-level analysis for (A) all features together and then stratified by only (B) genes, (C) exons, (D) junctions, (E) transcripts and (F) expressed regions. P-values were based on the Wilcoxon rank sign test.

**Figure S9**: Boxplots of single gene eQTL effects for GWAS associated SNPs for the associations not featured in Figure 4.

*Supplementary Table Legends:*

**Table S1**: Demographic information for subjects in the present study, stratified by age and diagnosis group. Dx: diagnosis, N: sample size, F: Female, Cauc: Caucasian, SD: standard deviation, PCW: post-conception weeks. Antipsychotic use was measured using toxicology at

998      time of death. P-values for diagnosis differences in continuous variables are based on linear
999      regression and P-values for categorical variables are based on chi-squared tests.

1000     **Table S2**: Splice junction annotation and characterization in GTEx and GEUVADIS for any
1001     junction or highly expressed junctions (mean reads per 80M mapped reads, RP80M > 0 and 5).
1002     Each column represents a 2x2 table for presence of identified junctions in 495 DLPFC samples
1003     in two independent polyA+ datasets. This table is analogous to Table 1.

1004     **Table S3**: Summary statistics for those features significantly developmentally regulated in the
1005     control-only analyses across the lifespan.

1006     **Table S4**: Significant developmentally regulated features collapsed to Ensembl Gene ID, used
1007     to make Figure S3

1008     **Table S5**: Isoform shifts by Ensembl Gene ID and feature summarization type.

1009     **Table S6**: Gene set analyses for those genes with significant isoform shift, stratified by feature
1010     summarization type. Q-values, which control the false discovery rate, FDR, are shown.

1011     **Table S7**: Genes within the PGC schizophrenia GWAS risk regions that contain isoform shifts
1012     by feature summarization type. 21.8% of PGC2 genes had developmental isoform shifts using
1013     exon counts (N=96/440) and 31.9% showed this isoform shift association based on junction
1014     counts (N=137/430)

1015     **Table S8**: eQTL replication metrics for CommonMind Consortium, GTEx frontal cortex, and both
1016     among all significant eQTLs in the discovery adult sample.

1017     **Table S9**: eQTL statistics for the best SNP-feature pair in the discovery sample that replicated
1018     in CMC and GTEx frontal cortex, by feature type.

1019     **Table S10**: Differential expression statistics for those features that were significant and
1020     replicated in case-control comparisons.

1021     **Table S11**: Genes consistently differentially expressed by case-control analysis for the different
1022     feature summarizations.

1023     **Table S12**: Associations between diagnosis, RPS and expression at gene and exon levels. First
1024     two columns for each feature: p-values for gene set tests for the significant case-control
1025     features among statistics capturing the effect of RPS on expression. Second two columns for
1026     each feature: directionality between RPS on expression associations and diagnosis on
1027     expression associations.

1028     **Table S13**: Gene set analysis for genes with features differentially expressed by case-control
1029     status, stratified by directionality and feature summarization type.

1030     **Table S14**: GWAS region set-level analyses for diagnosis-associated differentially expressed
1031     features, testing whether features in the PGC risk loci were more or less expressed as a set in
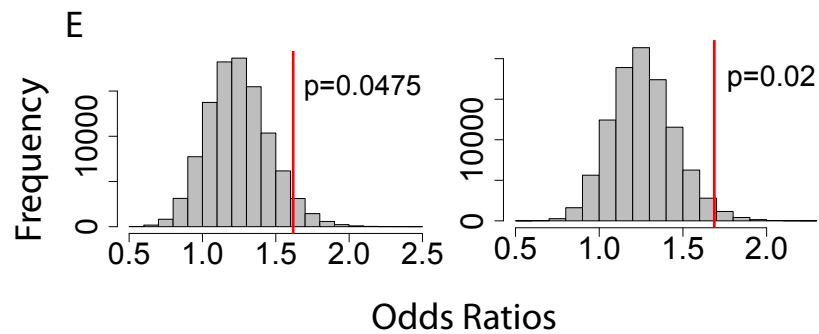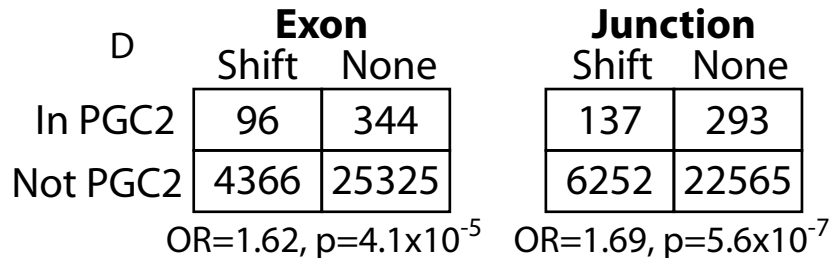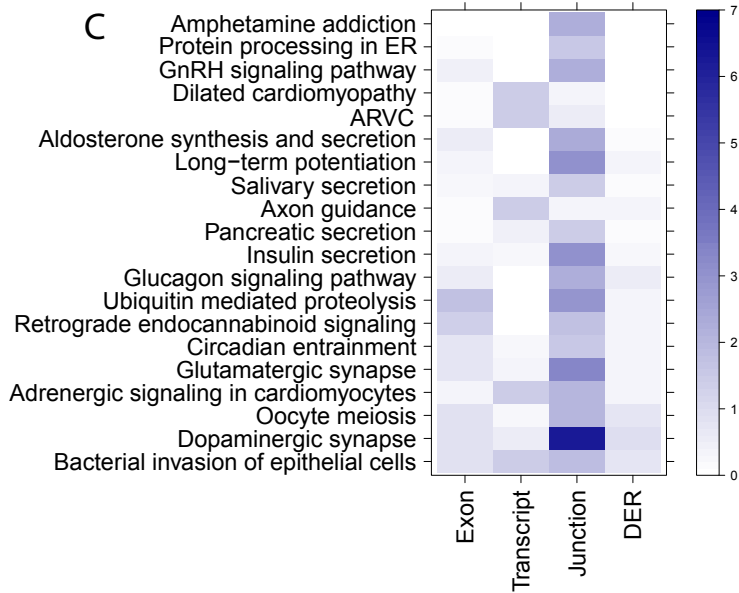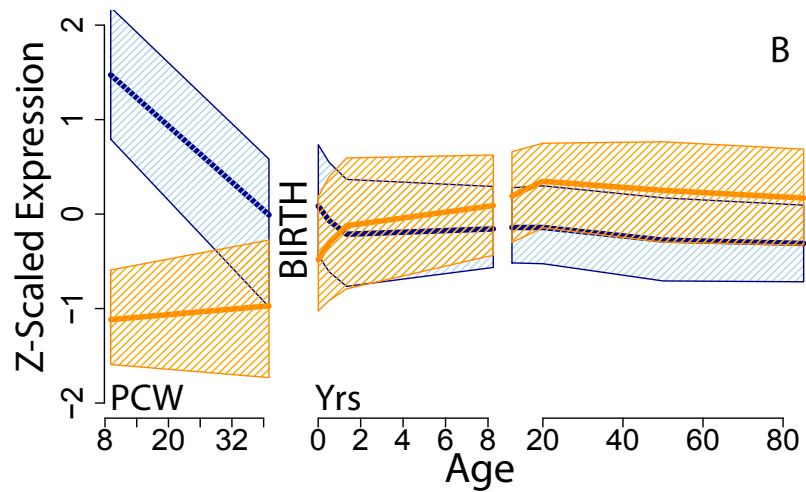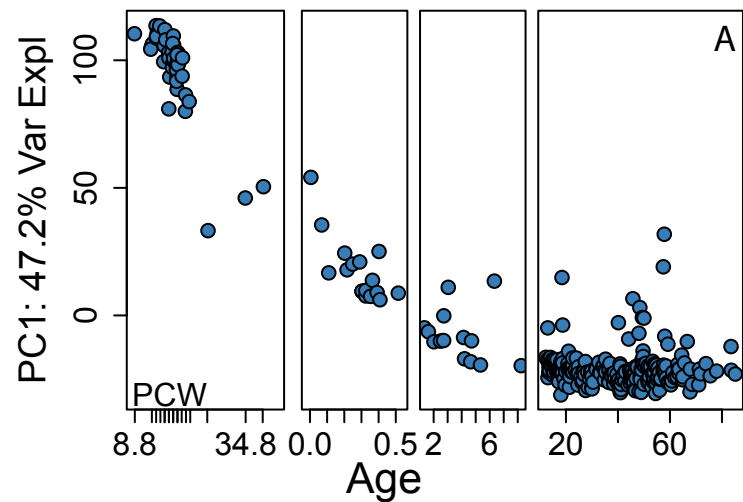
1032    cases compared to controls. Qual: qSVA adjusted analysis, Adj: observed covariate adjusted
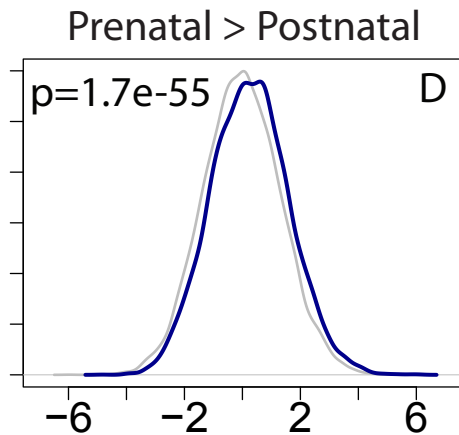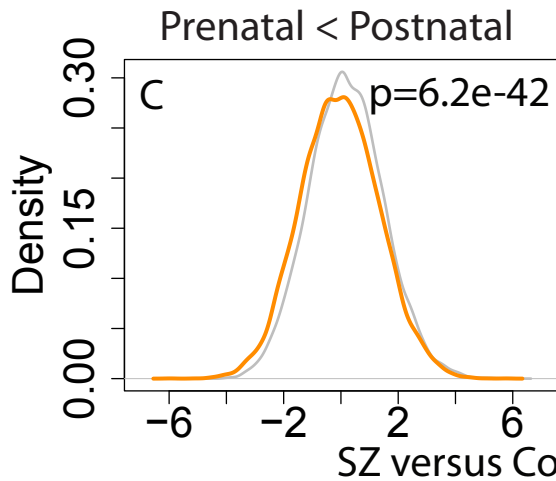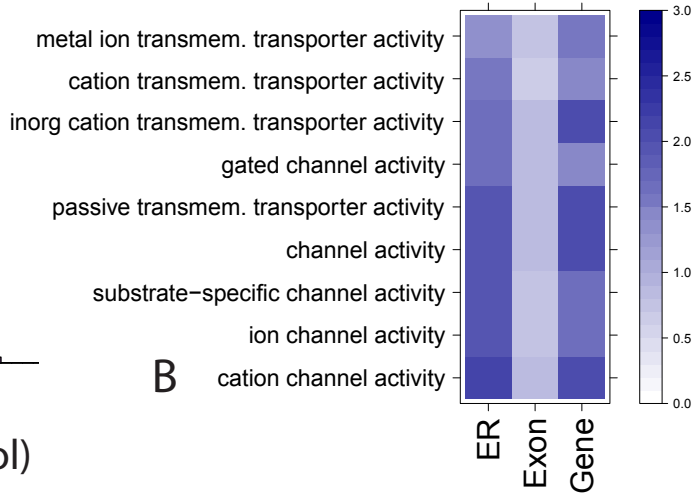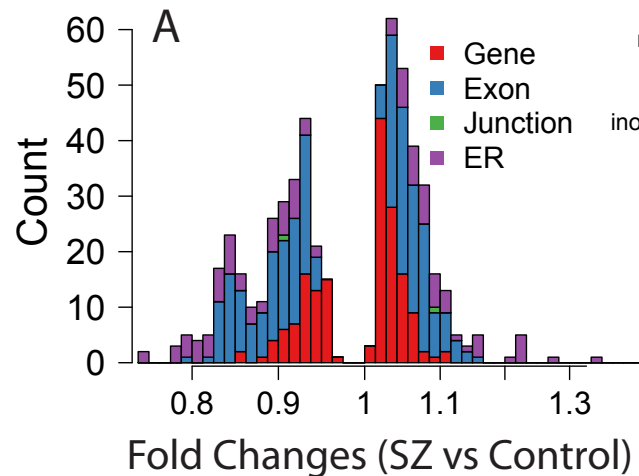
1033    analysis.

1034    **Table S15**: eQTL statistics for all PGC SCZD GWAS variants, including the risk direction (eQTL

1035    and GWAS) and clinical direction (case-control fold change)

1036    **Table S16**: eQTL statistics filtered to those that were directionally consistent between risk and
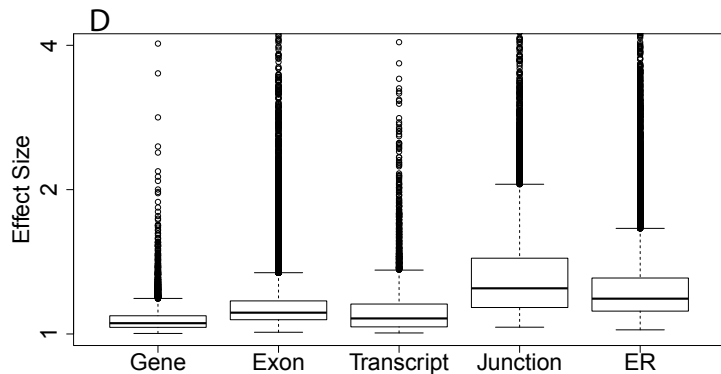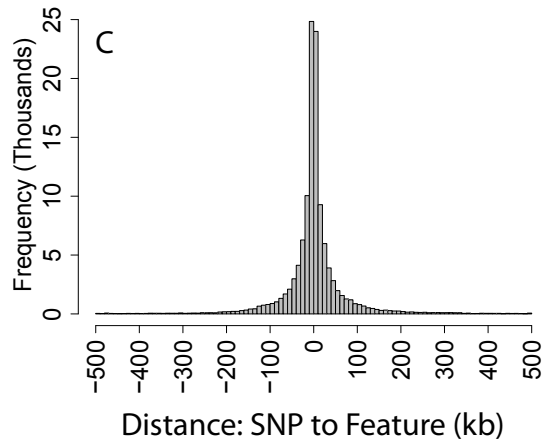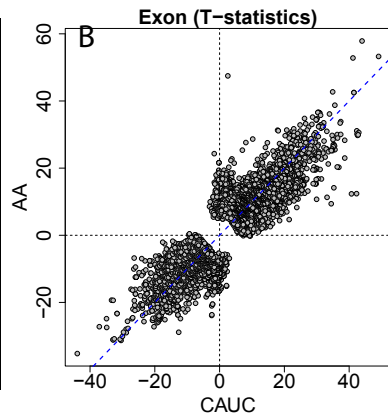
1037    clinical directionalities

1038    **Table S17**: Analogous to Table 2 with Ensembl transcript annotation data to highlight transcript
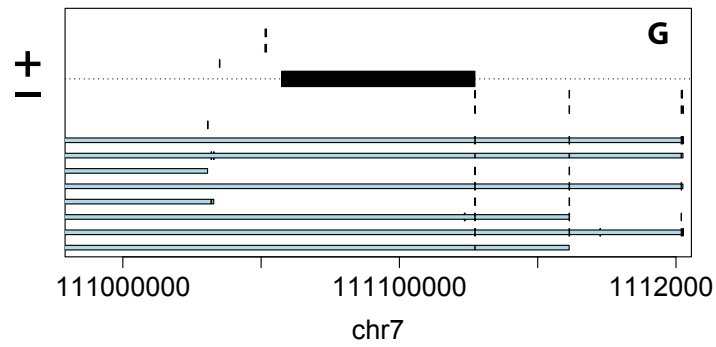
1039    specificity.

1040

| A | Gene | Exon | Transcript | Junction | DERs |
|---|---|---|---|---|---|
| # Bonf < 0.05 | 4934 | 36299 | 2993 | 5683 | 20295 |
| p-value cutoff | 8.39E-09 | 7.64E-10 | 1.71E-09 | 1.09E-09 | 1.28E-09 |
| # Ens Genes | 4934 | 6698 | 1708 | 2588 | 5381 |
| # Sym Genes | 3675 | 5195 | 1701 | 2376 | 4503 |
| Effect Size Median (IQR) | 1.07 (1.04-1.11) | 1.15 (1.1-1.23) | 1.1 (1.05-1.2) | 1.36 (1.19-1.63) | 1.25 (1.16-1.42) |
| Distance Median (IQR) | -40bp (-9.4kb-9.4kb) | -27bp (-14kb-13kb) | -130bp (-7.8kb-5.9kb) | -23bp (-10kb-8.8kb) | -23bp (-13kb-12kb) |
| % Novel | NA | NA | 37.5% | 11.3% | 52.7% |
| % Isoform-specific | NA | 67.0% | 100% | 31.1% | NA |

**A** — Schematic flow diagram:

GWAS: Allele ↑ Risk + eQTL: Allele ↑ Expression → Expression: SZ > Control

GWAS: Allele ↓ Risk + eQTL: Allele ↓ Expression → Expression: SZ < Control

**B** — PGC Rank: 47 (rs4523957), Exon − SRR − InEns, p=2.1e−12. Adjusted log2 Expression vs GG, GT, TT.

**C** — PGC Rank: 103 (rs324017), Exon − STAT6 − InEns, p=1.71e−07. Adjusted log2 Expression vs CC, CA, AA.

**D** — PGC Rank: 33 (rs6434928), Region − HSPD1 − Novel, p=3.34e−14. AA, AG, GG.

**F** — PGC Rank: 16 (rs13240464), Junction − IMMP2L − AltStartEnd, p=6.04e−10. TT, TC, CC.

**E** — chr2, HSPE1, HSPD1, 198350000 – 198370000.

**G** — chr7, 111000000 – 1112000.

Supplementary Results for: Developmental and genetic regulation of the human frontal cortex transcriptome in schizophrenia

Andrew E Jaffe, Richard E Straub, Joo Heon Shin, Ran Tao, Yuan Gao, Leonardo Collado Torres, Tony Kam-Thong, Hualin S Xi, Jie Quan, Carlo Colantuoni, Qiang Chen, Bill Ulrich, Brady J. Maher, Amy Deep-Soboslay,  The BrainSeq Consortium, Alan Cross, Nicholas J. Brandon, Jeffrey T Leek, Thomas M. Hyde, Joel E. Kleinman, Daniel R Weinberger

1.  Deep characterization of the human cortex transcriptome

We used Ensembl v75 annotation to quantify the normalized expression (via reads per kilobase per million reads mapped, RPKM) of 615,410 exons across 63,677 genes, and guided transcript assembly within and across samples to identify and quantify 188,578 transcripts with at least modest expression across all of the samples[1].  Only approximately half of the assembled transcripts were found in Ensembl annotation (N=100,932, 53.5%), while 41.8% (N=78,751) consisted of transcripts with novel splicing patterns. We also more directly interrogated the spliced alignments, which supported 3,582,199 unique exon-exon splice junctions across the 495 samples. We defined four classes of splice junctions based on the Ensembl gene/transcript annotation – present/existing, exon skipping, alternative exonic boundary, and completely novel (Figure S2, Table S2), and found moderate expression of many potentially unannotated transcripts in brain. To assess the replication and brain specificity of our findings, we identified junctions in both the brain samples from the GTEx project (N=1,393 samples from 201 individuals) and lymphoblastoid cell lines (LCLs) from the GEUVADIS project (N=666, see Methods). Overall, we find extensive replication in human brain GTEx samples of annotated junctions (95.0%), and high replication of potentially novel transcripts that capture exon skipping (75.6%) or alternative exonic boundaries (65.8%), with extremely high replication (>95%, Table 1). Much of this novel transcriptional activity appeared relatively-brain specific, as few junctions identified in DLPFC were present only in GEUVADIS LCLs and not GTEx brain samples further suggesting that this degree of junction discovery was not largely driven by alignment software. We lastly defined "expressed regions" (ERs) of contiguously expressed sequence – these ERs and junctions together can "tag" elements of transcripts in the data that are not constrained by limitations or incompleteness of existing annotation.

2.  Developmental regulation of human brain transcription

We sought in this study to more fully characterize developmental regulation of transcription across human brain development and aging, and modeled dynamic expression in 320 control samples using flexible linear splines in each of the five expression summarizations (see Methods). We found widespread developmental regulation across all five features(Table S3) corresponding to 30,490 unique Ensembl genes (across 20,596 gene symbols), including a core set of 15,603 genes with all five features showing convergent expression association with

development/aging (Figure S3). There were 89,865 previously unannotated expressed features corresponding to the majority of genes (17,374 IDs and 15,261 symbols) with genome-wide significant dynamic expression across brain development in these unaffected subjects (Table S4). The majority of these unannotated features were identified using the splice junction counts (N=43,512, 48.4%) of which 31,183 junctions (71.7%) tagged alternative exonic boundaries and 12,329 junctions (28.3%) corresponding to exonic skipping, while there were 26,010 differentially expressed regions (DERs, 28.9%) that all correspond to alternative exonic boundaries and 20,343 transcripts of which 19,138 (94.1%) correspond to exonic skipping.

3.  Expression associations with chronic schizophrenia illness

Given the apparent lower RNA quality of patient samples in virtually all schizophrenia brain case control studies including our own (see Table S1), we filtered the 384 samples above age 17 (209 controls and 175 patients) to a subset of 351 higher quality samples (196 controls and 155 cases) using metrics of age, ancestry, RNA quality, and cryptic ancestry (see Methods). Even after filtering, using gene-level expression, univariate comparisons between SCZD patients and controls identified 12,686 genes differentially expressed (DE) at a false discovery rate (FDR) < 5%, suggesting a large degree of bias. Regression modeling, adjusting for age, sex, ancestry, and observed RNA quality (see Methods) reduced the extent of confounding, resulting in 1,988 genes DE between patients and controls at FDR < 5%.   We identified the largest number of significant and replicated DE features using exon counts (N=274) followed by gene counts (N=170) and expressed region counts (N=110).  Very few DE junctions (N=2) and no DE transcripts were identified and independently replicated, which likely highlight the decreased statistical power using these approaches when annotated features are differentially expressed (see Discussion). Approximately one half of the genes annotated by these DE features were identified using gene-level summarization only (86/170, 50.6%), which utilize the largest number of reads collapsed across all possible transcript isoforms (Table S11). The majority of genes with DE signal were only supported by a single summarization type across both Ensembl IDs (142/237, 59.9%) and gene symbols (119/209 56.9%, Table S11).

Interestingly, analogous analyses for developmental regulation of schizophrenia-associated features without adjusting for the RNA quality qSVs were significant in the opposite directions, namely that schizophrenia-associated changes were further from, rather than closer to, fetal expression levels, which would be predicted as an effect of residual RNA quality confounding (as the quality of the samples were ranked fetal > adult control > adult SZ, see Table S1).  This directional difference in RNA quality between fetal and postnatal samples is also typical of earlier studies (e.g. Colantuoni, Geschwind, etc).

4.  GWAS implications of illness-associated expression differences

Only two genes in the 108 significant schizophrenia GWAS loci (*KLC1* and *PPP2R3A*) had expression features significantly differentially expressed and replicated. To potentially identify more subtle associations within these GWAS risk regions, we performed an exploratory "feature set" analysis (analogous to traditional gene set analyses) comparing the schizophrenia differential expression statistics of all expressed features within the loci to those outside the loci. We found decreased expression of the 15,050 expressed features within the risk loci in patients compared to controls, relative to the 996,775 expressed features outside of the loci, adjusting for mean expression level (Figure S8A, p =$3.32 \times 10^{-36}$, Table S13). While the absolute set-level effect sizes were small, these GWAS region-level associations were largely consistent when stratified by summarization type, with the exception of transcript counts which showed no association. Of particular importance is the observation that these findings also were only significant after quality adjustment in the differential expression analysis (Figures S8B-F) – analyses adjusting for only the usual observed confounders (e.g. clinical covariates and observed quality variable) showed no enrichment among PGC risk regions (p=0.2, Table S13). Therefore, while many of the case-control expression differences that meet both genome-wide statistical significance and replication criteria may be related to schizophrenia treatment and epiphenomena and depleted for GWAS regions, some differences in expression in some subsets of patient populations might be related to genetic risk for the disorder, though the biological interpretation of the feature distribution differences is not clear. Larger studies can likely improve power to detect expression changes within the GWAS risk regions amidst the clinical and molecular heterogeneity of schizophrenia.

5. Genetic control of RNA expression levels in the cortex

In order to elucidate the RNA features associated with schizophrenia risk variants themselves, rather than LD regions, we first performed a *cis* (<500kb) genome-wide expression quantitative trait loci (eQTL) analysis across the five expression feature summarizations. These analyses were performed within the post-adolescent subjects, including both cases and controls (age > 13, N=412) across densely observed and imputed common and high-quality SNP genotypes (MAF > 5%, N=7,421,423), adjusting for components related to ancestry and expression heterogeneity in each feature type (see Methods). Given the large number of tests, but also the need for comparability across the five expression summarizations, we adjusted for multiple testing using both the false discovery rate (FDR) with a cutoff of 1%, in line with the GTEX project [2] and also a more stringent Bonferroni-based cutoff adjusting for the number of LD-independent SNP-feature pairs. These multiple testing thresholds corresponded to $p<1.0 \times 10^{-4}$ and $p<2.6 \times 10^{-9}$, respectively, across the 5 feature summarizations.

We examined only eQTL associations of 682,828 unique genetic variants to a total of 70,204 nearby expression features, mapping back to 9,271 Ensembl IDs (of which 7,057 had HGNC symbols) that were at Bonferroni significance in our dataset and that replicated in the two independent datasets (Figure 3A, Table S10). We highlight the best SNP for each feature in Table S11, and have a complete database, including statistics and visualizations, of all significant and replicated eQTLs online (as well as significant control-only eQTLs) via a user

friendly browser at http://eqtl.brainseq.org.The majority of eQTLs (N=417,276, 61.1%) associated with features in a single Ensembl gene, and another large subset associated with features in two nearby genes (N=133,286, 19.5%). Conversely, the majority of genes were associated with more than 1 LD-independent SNP (76.3%) across the feature summarizations, including approximately half (49.0%) of genes having features associated with more than 3 LD-independent SNPs, suggesting that multiple variants can associate with expression independently of each other, and potentially result in more normally-distributed expression levels at the population-level. Furthermore, the eQTL effects (discovered in multiple/mixed ancestries) were largely consistent between Caucasian and African American subjects (Figure 2B). The most significant SNP for each eQTL feature was extremely proximal to the feature (based on nearest coordinate), with a median < 50bp (IQR less than -15kb to 15kb, Figure 2C). Lastly, we found the largest effect sizes on average for the junction and ER eQTLs (see Supplementary Results, Figure 2D).

While the majority of genes with features as eQTLs have multiple possible Ensembl isoforms, 67.0% of exon eQTLs and 31.1% of junction eQTLs were specific to a single Ensembl transcript when multiple transcripts existed for a given gene. Furthermore, 640 (11.2%), 1,123 (37.5%) and 10,698 (52.7%) of replicated junction, transcript, and expressed region eQTLs respectively, corresponded to novel transcriptional activity, including 387 alternative exonic boundary and 253 exon-skipping junction eQTLs and 4,036 alternative exonic boundary (crossing an intron and exon, and extended UTRs) and 6,662 intronic and intergenic ERs. These results indicate that performing both exon- and junction-level eQTL analyses together provide an encompassing view of genetic regulation of expression, as exon-level analyses identified eQTLs for the most genes as well as transcript-specific events, and junction-level analysis provided the largest biological effect sizes and the ability to identify potentially novel transcript isoforms. Together these eQTL results, based on analyses across multiple feature summarizations from RNA sequencing data, demonstrate more widespread genetic regulation of expression than previously reported, including extensive transcript specificity and regulation of unannotated sequence in the human brain.

There was moderate concordance across the different feature summarizations. Exon-level summarizations resulted in the largest number of unique genes as eQTLs (N=6,698) and transcript-level summarizations resulted in the fewest (N=1,708). While the junction eQTLs were less likely to appear in the GTEX project likely because of the lower coverage, and therefore replicate, those that did replicate tended to have larger effect sizes than other summarizations, with a median 36% (IQR: 19-63%) change in expression per allele copy (Figure 2D). Expressed region-level analysis identified the second most eQTLs to unique genes (N=5,381) and effect sizes (median: 25%, IQR: 16-42%). Gene-level eQTL analysis, which is the de-facto approach in RNA-seq analysis, resulted in the smallest eQTL effects, with a 7% (IQR: 4-11%) change in expression per allele copy, and identified fewer unique genes than exon and expressed region analysis (N=4,934). These data suggest that most eQTLs involve specific transcript isoforms whose signals may be diluted when examining only gene-level expression data. If this dilution when using gene counts results from averaging over potentially transcript-specific eQTLs, then mapping all expressed features (sans expressed regions) to all known Ensembl transcripts, and tabulating the proportion of eQTLs that appear relatively transcript specific should confirm this.

In fact, there was extensive evidence of transcript-specificity of replicated eQTLs across exons, junctions, and transcripts.

6. Characterization of schizophrenia risk-associated expressed features

We found significant enrichment of this directional consistency between the GWAS, case-control, and FDR-significant (34 loci across 444 eQTL features and 91 Ensembl genes, OR=2.00, p=4.2x10$^{-6}$) as well as with the more conservative Bonferroni-significant (19 loci across 161 features and 44 Ensembl genes, OR=4.67, p=7.97x10$^{-8}$) features

The majority of our eQTLs (SNP-feature pairs) had directionally consistent allelic effects in the 50 fetal samples (N=628, 82.2%, p < 10$^{-20}$) and a subset were marginally significant (at p < 0.05) even in this small sample size (N=162, 21.2%). As further evidence for similar genetic regulation, every eQTL that was marginally significant in the fetal samples was directionally consistent with the eQTL effect in the adult samples. For example, we find the most significant and replicated association with clinical directional consistency between rs4523957 (PGC rank 47) and exons (as well as other feature types, see Table S15) in Serine Racemase (*SRR*, Figure 4B), an NMDAR modulator previously associated with schizophrenia[3]. While the statistical significance of exon associations were only slightly better than at the gene level (2.1x10$^{-12}$ versus 5.44x10$^{-11}$), the effect size was almost three times larger (-0.16 versus -0.06). Similarly we identified significant eQTL association between rs324017 (PGC rank 103), and exon- and junction-level expression of *STAT6* (Figure 4C), a transcription factor involved in interleukin signaling[4] whereas there were no significant and replicated eQTLs at the gene-level.

The approach that was unconstrained by annotation also identified novel associations to schizophrenia risk, for example, the eQTL of rs6434928 (PGC rank 33) with intronic sequence in the *HSPD1* and *HSPE1* heat shock protein genes (Figure 4D,E), which are co-chaperonins linked by a common bidirectional promoter[5]. While both genes were recently highlighted in the subset of 20 PGC GWAS genes that are targets of currently approved drugs[6], the strictly intronic annotation of the eQTL association may suggest a different mechanism of risk from this locus in schizophrenia, and thus a different requirement for a therapeutic agent. Finally, we highlight the eQTL association between rs13240464 (PGC rank 16) and the expression level of a novel splice junction in *IMMP2L* (Figure 4F,G), a gene important in directing mitochondrial proteins to the mitochondria[7] and that has been previously associated with autism and Tourette Syndrome[8].

Supplementary Discussion

This proportion of genetic risk variants that are brain eQTLs – 42.5% – far exceeds previous reports using microarray and gene level (count-based) RNA-seq approaches. For example, of the 18 index SNPs labeled as eQTLs in Supplementary Table 4 In the PGC publication [9], only 8 risk loci had variants moderately linked (at $R^2 > 0.6$) to significant (at p < 10$^{-4}$) brain eQTLs. Our results highlight the enhanced detection of eQTL signals from GWAS-

positive risk variants using RNA-seq to quantify convergent measures of expression beyond the usual gene-level quantification in large collections of postmortem human brains.

In the original GWAS report, while there were only 18 reported index SNPs with brain eQTLs to a total of 30 genes, many of the reported eQTLs were themselves not GWAS-significant variants, being only weakly linked (at $R^2 < 0.6$) to the best GWAS SNP. There were only three of these 18 GWAS- and eQTL-significant index SNPs (putatively associating with *C10orf32 (BORCS7), AS3MT*, *OGFOD2*, *NAGA*, and *CENPM*) that were present in our larger RNA-seq data from the three independent studies that we surveyed. Similarly, a recent GWAS and eQTL integration analysis based on peripheral blood eQTLs identified very few overlapping genes[10]. These genes included *PCCB*, *SNX19*, *GATAD2A*, and *IRF3* as directionally-consistent eQTLs (with regard to the GWAS, eQTL, and case-control differences described above), and *SF3B1* and *PSMA4* as non-directionally consistent eQTLs, suggesting the need to query the disease-relevant tissue, i.e. human brain. Interrogating annotation-agnostic expression features added a considerable number of loci with eQTL signal further highlighting the value of RNA sequencing-based datasets. For example, we have recently dissected the GWAS signal in the 10q24.32 region in a more focused analysis based on the expression features highlighted in the current report. There we identified based on junction-level analysis a novel risk-associated isoform of AS3MT that is missing exons 2 and 3 [11]. This analysis showed that the previously referenced genetic association with this locus to AS3MT was inaccurate, as no association was found to the full length AS3MT [11], and importantly, the novel isoform does not have arsenic methyltransferase activity. Furthermore, eight additional loci, with clinical risk variants associating with expression of *SNX19*, *NRGN/VSIG2*, *HSPD1*, *GATAD2A/NDUGA13*, *VPS45*, *GPM6A*, *SOX2-O*T, and *BRINP2*, were identified here (and subsequently replicated) only via junction- and expressed region-level analysis.

## References

1    Frazee, A. C. *et al.* Ballgown bridges the gap between transcriptome assembly and expression analysis. *Nature biotechnology* **33**, 243-246, doi:10.1038/nbt.3172 (2015).
2    Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648-660, doi:10.1126/science.1262110 (2015).
3    Labrie, V. *et al.* Serine racemase is associated with schizophrenia susceptibility in humans and in a mouse model. *Human molecular genetics* **18**, 3227-3243, doi:10.1093/hmg/ddp261 (2009).
4    Shimoda, K. *et al.* Lack of IL-4-induced Th2 response and IgE class switching in mice with disrupted Stat6 gene. *Nature* **380**, 630-633, doi:10.1038/380630a0 (1996).
5    Ryan, M. T. *et al.* The genes encoding mammalian chaperonin 60 and chaperonin 10 are linked head-to-head and share a bidirectional promoter. *Gene* **196**, 9-17 (1997).
6    Lencz, T. & Malhotra, A. K. Targeting the schizophrenia genome: a fast track strategy from GWAS to clinic. *Molecular psychiatry* **20**, 820-826, doi:10.1038/mp.2015.28 (2015).
7    Nunnari, J., Fox, T. D. & Walter, P. A mitochondrial protease with two catalytic subunits of nonoverlapping specificities. *Science* **262**, 1997-2004 (1993).
8    Petek, E. *et al.* Molecular and genomic studies of IMMP2L and mutation screening in autism and Tourette syndrome. *Molecular genetics and genomics : MGG* **277**, 71-81, doi:10.1007/s00438-006-0173-1 (2007).

9       Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-+, doi:10.1038/nature13595 (2014).

10      Zhu, Z. *et al.* Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nature genetics* **48**, 481-487, doi:10.1038/ng.3538 (2016).

11      Li, M. *et al.* A human-specific AS3MT isoform and BORCS7 are molecular risk factors in the 10q24.32 schizophrenia-associated locus. *Nature medicine*, doi:10.1038/nm.4096 (2016).