

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

2

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

- 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](mailto:andrew.jaffe@libd.org) (co-corresponding)

28 & [drweinberger@libd.org](mailto:drweinberger@libd.org) (co-corresponding)

29

30

31 **Summary:**

32 GWAS have identified 108 loci that confer risk for schizophrenia, but risk mechanisms for  
33 individual loci are largely unknown. Using developmental, genetic, and illness-based RNA  
34 sequencing expression analysis, we characterized the human brain transcriptome around these  
35 loci and found enrichment for developmentally regulated genes with novel examples of shifting  
36 isoform usage across pre- and post-natal life. We found widespread expression quantitative trait  
37 loci (eQTLs), including many with transcript specificity and previously unannotated sequence  
38 that were independently replicated. We leveraged this eQTL database to show that 48.1% of  
39 risk variants for schizophrenia associated with nearby expression. Within patients and controls,  
40 we implemented a novel algorithm for RNA quality adjustment, and identified 237 genes  
41 significantly associated with diagnosis that replicated in an independent case-control dataset.  
42 These genes implicated synaptic processes and were strongly regulated in early development  
43 ( $p < 10^{-20}$ ). These data offer new targets for modeling schizophrenia risk in cellular systems.

44

45

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

48

## 49 **Introduction:**

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

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

84

## 85 **Results**

86 We performed deep polyA+ RNA-sequencing of 495 individuals, ranging in age from the  
87 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<sup>4</sup> as well as  
90 alignment-based 4) exon-exon splice junctions<sup>5</sup> and 5) expressed regions (ERs)<sup>6</sup>. These last  
91 two expression features were selected to reduce reliance on the potentially incomplete  
92 annotation of the brain transcriptome<sup>7</sup> (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 S2). 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

### 100 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<sup>8-12</sup>. Specifically, we and others<sup>7,13,14</sup> 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<sup>7</sup> (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 previously unannotated sequence (Table S4). Motivated in part by previous reports of  
112 preferential fetal isoform use among schizophrenia candidate risk genes<sup>10,11</sup> (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 We highlight a representative gene with isoform shifts in Figure S4 involving *CRTC2*, a  
118 transcription co-activator. There were 6672 Ensembl genes (23.7% of the set of  
119 developmentally regulated genes) with both positive and negative expression features having  
120 genome-wide significant correlations with age (each with  $p_{\text{bonf}} < 0.05$ , Figure 1B, Table S5, Figure  
121 S5). In other words, these represent alternate transcript isoforms of the same gene that show  
122 opposite patterns of expression across the prenatal-postnatal transition. In principle, this  
123 interaction would obscure developmental expression variation measured at the gene level.

124 We performed gene set analyses of those genes with shifting isoform usage compared  
125 to the larger set of genes with at least one developmentally regulated feature but without shifting  
126 isoform usage to identify more specific biological functions of this unique form of developmental  
127 regulation (Table S6). The set of developmentally shifting isoforms was relatively enriched for  
128 localization, catalytic activity, signaling-related processes, including synaptic transmission and

129 cell communication, and neuronal development, among many others. Interestingly, genes  
130 identified with shifting isoforms across development based exclusively on junction counts were  
131 enriched for both dopaminergic ( $FDR=1.67 \times 10^{-4}$ ) and glutamatergic ( $FDR=2.04 \times 10^{-4}$ ) synapse  
132 KEGG pathways (Figure 1C), the two neurotransmitter systems most prominently implicated in  
133 schizophrenia pathogenesis and treatment.

#### 134 Schizophrenia risk is associated with novel shifting isoform usage across brain development

135 Based on the KEGG analysis, we hypothesized that the genes with developmentally  
136 regulated isoform shifts may relate to risk for schizophrenia. Indeed, genes within the SCZD  
137 GWAS risk loci were more likely to harbor these novel isoform shifts occurring in the fetal-  
138 postnatal developmental transition compared with the rest of the expressed transcriptome  
139 (Figure 1D). For example, genes with developmental isoform shifts identified by exon, junction  
140 and expressed region counts were 75% ( $p=9.51 \times 10^{-6}$ ), 84% ( $p=1.63 \times 10^{-7}$ ) and 71% ( $p=2.0 \times 10^{-4}$ )  
141 more likely to lie within the PGC2 risk regions (with permutation-based  $p=0.02$ ,  $p=0.01$ , and  
142  $p=0.03$  respectively, see Methods) than developmentally regulated genes without isoforms  
143 shifts (Table S7). These results further underscore the role of changes in the regulation of  
144 transcription and splicing in the early brain developmental components of schizophrenia risk.

145

#### 146 Large-scale genetic regulation of transcript-specific and previously unannotated sequences

147 In order to elucidate the RNA features associated with schizophrenia risk variants  
148 themselves, rather than positional LD regions, we first performed a genome-wide *cis* (<500kb)  
149 expression quantitative trait loci (eQTL) analysis within the 412 post-adolescent subjects (see  
150 Methods) across the five convergent transcript features (Table 1). We hypothesized that, in  
151 general, analyzing transcript features like exons and junctions would increase statistical power  
152 for eQTL discovery if genetic variation regulated the expression levels of specific mRNA  
153 transcripts. At the gene-level, which collapses data from all transcripts into a single measure,  
154 which is the most common feature summarization for eQTL discovery, the vast majority of  
155 expressed genes were associated with the expression of at least one nearby genetic variant.  
156 There were eQTLs to 6748 Ensembl Gene IDs (of which 4955 genes had HGC symbols) at  
157 stringent Bonferroni-adjusted significance ( $p < 8.41 \times 10^{-9}$ , see Methods), and eQTLs to 18,416  
158 Ensembl Gene IDs at more liberal  $FDR < 1\%$  significance ( $p < 1.84 \times 10^{-4}$ ). However, we found a  
159 larger number of genes with eQTLs using exon-level analysis – 48,031 exons mapping to 8386  
160 Ensembl IDs - at Bonferroni significance (“eExons”,  $p < 7.64 \times 10^{-10}$ ). Exon-level analysis showed  
161 widespread transcript-specificity of eQTL associations. Almost all eExons mapped to genes with  
162 more than one annotated transcript ( $N=45,239$ , 94.2%), and the majority of these showed eQTL  
163 associations to exons belonging to a single transcript isoform ( $N=30,283$ , 66.9%). This  
164 transcript-specificity was also evident in the eQTL effect sizes, as the median additive effect  
165 size was approximately two-fold higher for exon- than gene-level analysis (15.6% versus 7.0%  
166 expression change per allele copy). Interestingly, while transcript-specific by nature, we actually  
167 found the fewest eQTLs to assembled-and-quantified transcripts (3,263 eTxns at  $p < 1.73 \times 10^{-9}$ ),  
168 in line with previous reports highlighting the difficulties in merging assemblies across  
169 replicates<sup>15</sup>. Lastly, there were an additional 3,022 eGenes identified with exon-level analysis

170 compared to the 5364 eGenes identified with both summarization levels. These results  
171 demonstrate extensive transcript specificity of many eQTL signals that are missed by gene-level  
172 analyses.

173 We next explored the extent of eQTLs to previously unannotated transcribed sequence  
174 using junction- and expressed-region feature summarizations which do not rely on existing gene  
175 annotation for quantification. Among the 18908 junctions with eQTL signal at Bonferroni  
176 significance (“eJxns”,  $p < 1.1 \times 10^{-9}$ ), 21.6% (N=4089) were previously unannotated, including  
177 1312 eJxns to exon-skipping splicing events and 2777 eJxns to shifted exonic boundaries  
178 (acceptor or donor splice sites). The eJxns also highlight a large degree of potential transcript  
179 specificity, both in the 4089 unannotated junctions as well as 3388 additional annotated eJxns  
180 that delineate individual transcript isoforms (when multiple isoforms are present for the gene). At  
181 the expressed region-level, among the 27,643 ERs with eQTL signal at Bonferroni significance  
182 (“eERs”,  $p < 1.28 \times 10^{-9}$ ), 14,890 were either fully or partially unannotated, with partial events  
183 including 4521 exon extensions into neighboring intronic sequence and 769 extended  
184 untranslated regions (UTRs) and fully unannotated events being strictly intronic (N=6,255) and  
185 intergenic (N=3,345) sequences. These two feature classes also had the largest eQTL effect  
186 sizes of the tested features, with 41.4% and 29.2% change in expression per allele copy for  
187 eJxns and eERs. Lastly, we found that 1,042 Ensembl genes had eQTLs exclusively to  
188 unannotated sequence with no corresponding eQTL signal to annotated features in the genes.  
189 Genetic regulation of previously unannotated sequence provides further evidence for biological  
190 relevance in the human brain.

191 Given the large degree of genetic regulation of transcript-specificity and unannotated  
192 sequences, we sought to assess the replication of the identified eQTLs (“LIBD”) in independent  
193 human brain RNA-seq data. We downloaded alignment-level data from the CommonMind  
194 Consortium (“CMC”) project, and quantified expression across the same five feature  
195 summarizations (in Table 1). Among those significant eQTL SNP-feature pairs that were well-  
196 imputed, polymorphic and expressed in the replication dataset (~84% of pairs, ~95% of  
197 eFeatures, see Methods, Figure S6), >94% had consistent directionality in the two datasets,  
198 between 75.7% (eTxns) and 81.5% (eJxns) were directionally consistent and marginally  
199 significant (at  $p < 0.01$ ), and just over half (52.1%-57.0%) were directionally consistent and  
200 FDR-corrected significant (published set,  $p < 10^{-5}$ ) in the DLPFC replication dataset. Meta-  
201 analysis between datasets demonstrated extensive significance and replication of the 9.3M  
202 SNP-feature Bonferroni-significant eQTL pairs including 97.6% at  $p < 1 \times 10^{-5}$  and 82.0% at  $p <$   
203  $10^{-9}$ . We further reprocessed and quantified GTEx v6 RNA-seq brain data (“GTEx”) from raw  
204 reads using the same pipeline, and assessed replication and regional specificity in these data  
205 using meta-analysis across 13 brain regions compared to frontal cortex alone. Here we found  
206 that many of the DLPFC-identified eQTLs showed strong concordant signal across all brain  
207 regions, suggesting an overall lack of regional specificity for the majority of our identified eQTLs  
208 (Figure S7). All significant eQTLs are searchable on our publicly available database:  
209 [eqtl.brainseq.org/phase1/eqtl/](http://eqtl.brainseq.org/phase1/eqtl/) which provides visualizations and eQTL statistics across three  
210 independent datasets.



## 212 Clinical enrichment of eQTL associations for schizophrenia and other traits

213 We sought to better determine the clinical relevance of our significant eQTLs particularly in the  
214 context of transcript feature-level and previously unannotated sequence associations. We cross-  
215 referenced our identified eQTLs with genome-wide association study (GWAS) risk variants.  
216 Here we used 3 significance levels to associate eQTLs with GWAS variants: a) more liberal  
217 FDR-significant eQTLs in the discovery dataset, b) these FDR-significant eQTLs with additional  
218 replication data support (meta-analysis p-values with CMC <  $10^{-8}$ ), and c) Bonferroni-significant  
219 eQTLs in the discovery dataset, e.g. Table 1, First we considered the proportion of common  
220 (MAF > 5%) and well-measured risk variants from the 128 index variants (N=106, see Methods)  
221 published in the latest PGC2 GWAS for schizophrenia<sup>2</sup> and their highly correlated proxies (see  
222 Methods). We identified FDR-significant eQTL associations to 51 risk SNP signals (of 106  
223 tested, 48.1%, Table S8), a substantially higher proportion of risk variants classified as brain  
224 eQTLs than previously reported<sup>16</sup> (Table 2). In total, there were 1,244 unique SNP-feature pairs  
225 that were genome-wide FDR-significant eQTLs (83 genes, 553 exons, 49 transcripts, 192  
226 junctions and 367 ERs) mapping to 194 unique Ensembl Gene IDs (of which 162 have HUGO  
227 gene symbols). Among these 51 risk SNPs, 17 were eQTLs only to exons, junctions or  
228 expressed regions, and 7 were eQTLs to only unannotated transcribed sequence. There were  
229 17 loci with annotated eQTLs to only a single gene and another 10 loci with eQTLs to two  
230 genes. More stringent meta-analysis significance ( $p < 10^{-8}$ ) retained eQTL evidence for 37  
231 variants including 17 to exons, junctions, and ERs, of which 6 were unannotated.

232 We also assessed enrichment of 23704 GWAS risk SNPs from the NHGRI GWAS catalog  
233 present and common in our genetic data (of 44,738 available), and found eQTL evidence for  
234 8988 variants (37.9%) at FDR < 0.01. These GWAS variants that were identified as eQTLs were  
235 from GWAS for the majority of all tested traits in the literature (68.1%, 1415 of 2078 present)  
236 across all sites in the body, suggesting that many of the identified eQTLs in brain are likely  
237 shared with other tissue sites as previously described<sup>17</sup>. Of the 8988 GWAS eQTL variants,  
238 2982 were eQTLs only to exons, junctions or expressed regions, of which 995 were only to  
239 unannotated sequence (Table 2). More stringent meta-analysis significance ( $p < 10^{-8}$ ) retained  
240 eQTL evidence for 5490 variants including 1824 to exons, junctions, and ERs, of which 671  
241 were unannotated. These results highlight the ability to identify more eQTL signal for clinical  
242 risk variants by casting a wider net of RNA-seq feature summarization, including previously  
243 unannotated transcribed sequences.

## 244 Refining risk transcripts through conditional analyses

245 We further sought to filter the eQTL hits to schizophrenia GWAS regions using  
246 conditional analysis in order to identify perhaps the most immediate downstream features of  
247 genetic risk. For each of the 51 eQTL-positive GWAS variants noted above, we conditioned on  
248 the most significant eQTL feature for each variant and then performed eQTL reanalysis of all  
249 other features. We then retained those eQTL features that remained at least marginally  
250 significant (at  $p < 0.05$ ) and repeated the conditional analysis now based on the two most  
251 independently associated expression features. We iteratively performed these conditional  
252 analyses until no other features were conditionally significantly associated eQTLs. These

253 analyses resulted in only 220 conditionally-independent SNP-feature eQTLs (35 genes, 66  
254 exons, 8 transcripts, 50 junctions and 61 ERs) to the 51 schizophrenia GWAS variants (Table  
255 S8) which mapped to 131 unique Ensembl Gene IDs (of which 106 have HUGO gene symbols).  
256 Conditional analysis resulted in an additional locus with eQTLs to a single gene (totaling 18 loci)  
257 and an additional four loci with eQTLs to features in two genes (totaling 14 loci, Table 3).  
258 Interestingly, these conditional analyses further highlighted the potential importance of  
259 transcript-specific and previously-unannotated eQTLs, as more loci were associated only with  
260 exons, junctions and ERs (27 versus 17), more were strictly unannotated (11 versus 7), and  
261 more showed eQTL associations to a single transcript isoform (18 versus 11).

262 We highlighted several representative eQTLs in Figure 2 for different classes of  
263 associations. The top GWAS risk variant rs1233578 associated with strictly intergenic sequence  
264 downstream of *ZSCAN23* (Figure 2A,2B,  $p=2.7\times 10^{-8}$ ) with replication in both CMC ( $p=0.01$ ) and  
265 GTEx ( $T=3.1$ ), suggesting potential novel transcribed sequence linked to schizophrenia risk. We  
266 also found significant eQTL signal to specific 5' junction and exon sequences of *CTNNA1* to  
267 rs3849046 (Figure 2C,2D; discovery  $p=6.2\times 10^{-8}$ , CMC replication  $p=1.4\times 10^{-8}$ ). Another example  
268 of eQTL associations of partially annotated sequence was rs9841616 exclusively associating  
269 with the 3' sequence of the most proximal short transcript isoform of *SOX2-OT* (Figure 2E,2F;  
270 discovery  $p=8.2\times 10^{-12}$ , replication  $p=2.9\times 10^{-8}$ ). We also found novel eQTL associations to  
271 annotated exons in *CD46* (Figure 2G,  $p=9.2\times 10^{-38}$ , replication  $p=2.9\times 10^{-14}$ ), *SRR* (Figure 2H,  
272  $p=2.0\times 10^{-12}$ , replication  $p=4.7\times 10^{-6}$ ) and *GPM6A* (Figure 2I,  $p=2.8\times 10^{-6}$ , replication  $p=0.02$ ).

273 We also found significant enrichment of these conditionally independent schizophrenia  
274 risk-associated eQTLs among genes with developmental isoform shifts identified above – 44.0%  
275 of genes with eQTLs compared to 23.6% without eQTLs ( $OR=2.54$ ,  $p=5.38\times 10^{-8}$ ). These  
276 conditional analyses could suggest potential regulatory roles of these unannotated transcribed  
277 sequences on annotated transcripts that play a putative role in the manifestation of  
278 schizophrenia risk in the brain. More generally, these eQTL results highlight significant and  
279 independently-replicated risk-associated schizophrenia candidate genes and their specific  
280 transcripts that comprise links in the causative chain of schizophrenia in the human brain.

### 281 Expression associations with chronic schizophrenia illness

282 We lastly explored the expression landscape of the prefrontal cortex of the  
283 schizophrenia illness state and its potential link with developmental regulation and genetic risk.  
284 We performed differential expression modeling using 351 high quality adult samples (age >16,  
285 196 controls, 155 cases), and found extensive bias by RNA degradation within both univariate  
286 analysis (where 12,686 genes were differentially expressed at  $FDR<5\%$ ) and even after  
287 adjusting for standard measured levels of RNA quality typical of all prior studies (Figure S6). We  
288 therefore implemented a novel statistical framework based on an independent molecular  
289 degradation experiment (see Methods, Results S3), called “quality surrogate variable analysis”  
290 (qSVA, see Methods)<sup>18</sup>. We further utilized potential replication RNA-seq data from the  
291 CommonMind Consortium (CMC) dataset, using a subset of age range-matched 159  
292 schizophrenia patients and 172 controls. Interestingly, adjusting for observed factors related to  
293 RNA quality that characterize all earlier studies of gene expression in schizophrenic brain,



294 including an earlier report using the CMC data <sup>16</sup>, the proportion of genes with differentially  
295 expressed features at genome wide significant FDR < 5% that replicate (with directionality and  
296 marginal significance at  $p < 0.05$ ) in the CMC dataset was small (only 11.0%, 244/2,215). In  
297 contrast, using our new statistical qSVA approach, 40.1% of differentially expressed genes at  
298 FDR < 5% (N=75/183) replicate in the CMC dataset. At genome-wide significant FDR<10%  
299 (see Methods), we identified 237 genes with 556 DE features that replicated in the CMC dataset  
300 (33.6% gene-level replication rate, Table S9, Table S10).

301 The differences in expression levels between cases and controls of these DE features  
302 were generally small in both our discovery and the replication datasets (Figure 3A, Figure S8),  
303 perhaps a direct result of the clinical and molecular heterogeneity of this disorder <sup>13,19</sup>. Gene  
304 ontology analysis implicated transporter- and channel-related signaling as significantly  
305 consistently downregulated in patients compared to controls across genes annotated in all three  
306 expression summarizations (Figure 3B, Table S11). These results suggested decreased  
307 signaling in patients with schizophrenia, but could raise the possibility that these replicated  
308 expression differences between patients and controls relate to epiphenomena of illness, such as  
309 treatment with antipsychotics which affect signaling in the brain<sup>14</sup>, as the majority of patients  
310 were on anti-psychotics at the time of death (64%, Table S1). Only two genes (*KLC1* and  
311 *PPP2R3A*) in the significant 108 schizophrenia GWAS loci were significantly differentially  
312 expressed. However, in an exploratory analysis, we found that overall the differential expression  
313 statistics within the loci were significantly different than those features outside the loci (Results  
314 S4, Figure S9, Table S12). We also investigated the relationships between transcription and  
315 genomic risk for schizophrenia using genome wide Polygene Risk Scores (PRS) from each  
316 subject calculated as previously described<sup>2</sup> (see Methods). Using the subset of 209 Caucasian  
317 samples, we largely found a lack of association between PRS and expression of individual  
318 expression features. We further found a lack of enrichment of PRS on expression comparing the  
319 differentially expressed and replicated case-control features to the rest of the transcriptome, as  
320 well as lack of directionally consistency between PRS- and diagnosis-associated statistics  
321 among expressed features (Table S13). These results further suggest that the significant case-  
322 control expression differences show little overlap with genetic risk for the disorder.

323 In an earlier study of the epigenetic landscape of frontal cortex of patients with  
324 schizophrenia, we showed that DNA methylation levels in patients were closer to fetal  
325 methylation levels than to those of adult control samples<sup>20</sup>. Here we tested for analogous effects  
326 in the RNA-seq data related to the illness state. Every significant gene with differentially  
327 expressed features in the adult case-control analysis and replicated in the independent dataset  
328 showed evidence for developmental regulation across at least two expression feature types. We  
329 further found that expression features more highly expressed in postnatal life tended to be more  
330 lowly expressed in patients compared to controls (max:  $p = 3.24 \times 10^{-11}$ , min:  $p = 1.05 \times 10^{-70}$ , Figure  
331 3C) and features more highly expressed in fetal life tended to be more highly expressed in  
332 patients with schizophrenia compared to controls (max:  $p = 6.86 \times 10^{-33}$ , min:  $p < 10^{-100}$ , Figure  
333 3D). Analogous analyses for developmental regulation of schizophrenia-associated features  
334 without adjusting for the RNA quality qSVs were significant in the opposite directions, namely  
335 that schizophrenia-associated changes were further from, rather than closer to, fetal expression  
336 levels, as might be predicted as a confounding artifact of residual RNA quality differences (as

337 the quality of the samples rank as fetal > adult control > adult SZ, see Table S1). These results  
338 further converge on a role for genes changing during brain development and maturation in the  
339 pathogenesis of schizophrenia, specifically that both DNA methylation and expression levels in  
340 adult patients appear to reflect levels in the developing brain more strongly than do those of  
341 unaffected individuals. These results also underscore the risk of spurious findings based on  
342 uncorrected RNA quality confounding.

343

## 344 Discussion

345 We have explored the diverse landscape of expression correlates of schizophrenia risk  
346 and illness state in the postmortem human frontal cortex across the lifespan. Using deep RNA  
347 sequencing to define convergent measures of gene expression and early brain development,  
348 we identified widespread developmental regulation of transcription, including novel discoveries  
349 related to preferential isoform usage across brain development. These unexpected isoform  
350 “shifts” were associated with genetic risk for schizophrenia, and the directionality of  
351 dysregulation of developmentally regulated features suggest a more fetal-like expression profile  
352 in patients with schizophrenia compared with controls. Our approach to transcript  
353 characterization, which included extensive characterization of unannotated sequence, revealed  
354 that many more schizophrenia risk associated SNPs are brain eQTLs than previously reported -  
355 many risk SNPs only associate with a single gene, or even a single transcript, and many of  
356 these adult-identified eQTLs show overlap with genes with dynamic isoform regulation across  
357 human brain development. Lastly, we identified significant and replicated genes differentially  
358 expressed in patients with schizophrenia compared to unaffected controls using a new  
359 experiment-based statistical framework to estimate and reduce the effects of latent RNA  
360 degradation bias which had not been accounted for in earlier studies. Without this new  
361 approach to RNA quality adjustment, replication across datasets is markedly limited if not  
362 negligible, and the directionality of the association with developmental isoform shifts is  
363 anomalous. These data suggest a convergence of developmental regulation and genetic risk for  
364 schizophrenia that appears relatively stable in patients ascertained at death, following decades  
365 of illness after diagnosis. We previously observed analogous stability of epigenetic marks  
366 highlighting prenatal life in adult patients with schizophrenia<sup>20</sup>, suggesting that both genetic and  
367 environmental risk factors implicated in schizophrenia illness involve early developmental  
368 events that are still observable in the brain tissue of adult individuals despite many years of  
369 illness.

370 While our approach utilizing convergent expression features – genes, exons, transcripts,  
371 junctions, and expressed regions – results in more complicated data processing and analysis, it  
372 can potentially cast a wider net in the search for valid biological signals in RNA sequencing  
373 datasets. Using all convergent features overcomes the limitations related to any given feature  
374 summarization, including the inability to measure and interrogate unannotated or novel  
375 transcribed sequences using gene and exon counts, and the difficulties in full transcript  
376 assembly from short sequencing reads<sup>21</sup>. We note that both quantifying and analyzing splice  
377 junctions, and also transcript-level data, rely on junction-spanning reads for statistical power. In

378 our data, there were approximately 3 times (IQR: 2.86-3.24) more reads available by gene/exon  
379 counting approaches than those that contain splice junctions, likely explaining why gene counts  
380 discovered more differentially expressed genes in the schizophrenia diagnosis analyses. Two  
381 relatively new approaches utilized here – direct quantification and statistical analyses of splice  
382 junction counts and expressed regions – can identify differential expression signal when it is  
383 outside of the annotated transcriptome. The junction-level approach can also identify previously  
384 uncharacterized novel transcribed sequences, which we replicated in other large publicly  
385 available datasets, as well as delineate individual transcripts or classes of transcripts that share  
386 a particular splice junction. As read lengths increase, the proportion of reads containing splice  
387 junctions will increase, making junction- and transcript-based approaches even more powerful,  
388 including those recently developed to identify splicing QTLs<sup>22</sup>.

389 Our analysis of RNA-seq data identified widespread shifts in preferential isoform use  
390 across brain development, which would have been impossible to identify using only gene-level  
391 data and incomplete with only exon-level data (Figure 2). The genes with these isoform shifts  
392 were significantly enriched for neurodevelopmental and cellular signaling processes, and as well  
393 as for genes in regions of genetic risk for schizophrenia. A prevalent hypothesis suggests that  
394 schizophrenia is a neurodevelopmental disorder that arises because of altered connectivity and  
395 plasticity in the early assembly of relevant neural circuits<sup>23</sup>, and the potential convergence of  
396 genetic risk with developing signaling processes across human brain development should point  
397 to specific candidate molecular disruptions occurring during the wiring of the fetal brain. Indeed,  
398 inefficient or disrupted signaling and tuning is thought to underlie the expression of illness in the  
399 adult brain<sup>23</sup>, and the most successful therapeutics work through improving these processes<sup>14</sup>.  
400 Consistent with this hypothesis, we find evidence for differences in the expression of genes  
401 coding for subunits of ion channels in the cortices of patients with schizophrenia compared to  
402 controls. We observed significant differential expression of both voltage-gated (*KCNA1*, *KCNC3*,  
403 *KCNK1*, *KCNN1*, *SCN9A*) and ligand gated ion channels (*GRIN3A*, *GABRA5*, *GABRB3*),  
404 transporters (*SLC16A2*, *ALC25A33*, *SLC26A11*, *SLC35F2*, *SLC7A3*), and ion channel auxiliary  
405 subunits (*KCNIP3*, *SCN1B*), supporting other evidence that the clinical phenomenology of  
406 schizophrenia is associated with altered neuronal excitability<sup>24</sup>. While these findings implicating  
407 basic mechanisms of cortical circuit dynamics may underlie fundamental aspects of the clinical  
408 disorder, the possibility that they are driven by the effects of pharmacological treatment and are  
409 thus state dependent epiphenomena cannot be excluded. Indeed, our failure to find association  
410 of genomic risk scores and differential gene expression in the illness state adds weight to the  
411 latter interpretation.

412 Our eQTL analyses are among the largest and most comprehensive to date in human  
413 brain tissue, based on stringent genome-wide significance and independent replication, and  
414 offer additional insights into the genetic regulation of RNA expression levels. Our data also  
415 suggest more widespread regulation of specific transcript isoforms, which we were able to  
416 identify using exon- and junction-level analyses. This transcript-specific genetic regulation was  
417 particularly prevalent among schizophrenia risk variants, where 66.9% of loci containing multiple  
418 transcripts showed clinically- and molecularly-consistent eQTL signal to a single Ensembl  
419 transcript isoform. Overall, we have identified many more eQTLs to genome-wide significant  
420 schizophrenia risk variants – 48.1% - than previously reported, experimentally implicating far

421 more potential “risk” genes within these loci than previously characterized. Our database of  
422 eQTLs – available at [eqtl.brainseq.org/phase1/eqtl](http://eqtl.brainseq.org/phase1/eqtl) – is searchable for candidate genes or SNPs  
423 and provides publication-ready visualizations (e.g. boxplots in Figure 2) and statistics eQTL  
424 associations. The database can serve as a “one stop shop” for eQTL statistics across three  
425 independent studies (LIBD, CMC, and GTEx) for both annotated and unannotated transcribed  
426 sequence in the human cortex, and can export results to the UCSC Genome Browser <sup>25</sup> for  
427 additional interrogation.

428         These eQTL associations within the genome-wide significant schizophrenia loci identify  
429 novel putative biological mechanisms underlying risk for the disorder. We have highlighted  
430 GWAS loci that contain significant and statistically independent eQTLs, as they often point to  
431 individual “risk” genes or even more specific “risk” transcripts. These “risk” genes and transcripts  
432 are targetable entry points for more focused cellular assays and model organism work to better  
433 characterize schizophrenia risk mechanisms. Moreover, these eQTLs of specific transcript  
434 features identifies a compelling strategy and directionality for target rescue, specifically to  
435 increase or decrease the function of the target transcript(s) and downstream effectors. Focusing  
436 solely on increased or decreased expression in brains of patients compared to controls, without  
437 considering genetic risk variants and their regulation of local gene expression, will likely  
438 predominantly highlight molecular changes resulting from the schizophrenia illness state, as we  
439 suggest with consistent down-regulation of ion channels. We stress the priority of identifying the  
440 most relevant cellular consequences of genetic risk, which we view as production of particular  
441 isoforms with predicted directionality, rather than trying to identify “causal” mutations tagged by  
442 “marker” risk SNPs from the GWAS. We suggest that identifying convergence between genetic  
443 risk and potential molecular consequences of the disorder is likely to result in better, or at least  
444 more consistent support for, targets for drug discovery efforts.

445

446

#### 447 **Author Contributions**

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

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

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

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

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

456 W.S.U. – created user-friendly database of eQTLs

457 A.D.S. – consented and clinically characterized human brain donors

458 T.M.H.,J.E.K.- collected, consented, characterized, and dissected human brain tissue;  
459 contributed to the design of the study

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

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

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

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

465 The remaining authors declare no competing financial interests.

466

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

469

#### 470 **Acknowledgements:**

471 We thank Dr. Ronald Zielke, Robert D. Vigorito, and Robert M. Johnson of the National Institute  
472 of Child Health and Human Development Brain and Tissue Bank for Developmental Disorders  
473 at the University of Maryland for their provision of fetal, child, and adolescent brain specimens;  
474 This work was supported by the funding from Lieber Institute for Brain Development and the  
475 Maltz Research Laboratories. The work was partially supported by R21MH109956 to A.E.J.  
476 L.C.T was supported by Consejo Nacional de Ciencia y Tecnología México 351535.

477 The Genotype-Tissue Expression (GTEx) Project was supported by the [Common Fund](#) of the  
478 Office of the Director of the National Institutes of Health. Additional funds were provided by the  
479 NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. Donors were enrolled at Biospecimen Source  
480 Sites funded by NCI\SAIC-Frederick, Inc. (SAIC-F) subcontracts to the National Disease  
481 Research Interchange (10XS170), Roswell Park Cancer Institute (10XS171), and Science Care,  
482 Inc. (X10S172). The Laboratory, Data Analysis, and Coordinating Center (LDACC) was funded  
483 through a contract (HHSN268201000029C) to The Broad Institute, Inc. Biorepository operations  
484 were funded through an SAIC-F subcontract to Van Andel Institute (10ST1035). Additional data  
485 repository and project management were provided by SAIC-F (HHSN261200800001E). The  
486 Brain Bank was supported by a supplements to University of Miami grants DA006227 &  
487 DA033684 and to contract N01MH000028. Statistical Methods development grants were made  
488 to the University of Geneva (MH090941 & MH101814), the University of Chicago (MH090951,  
489 MH090937, MH101820, MH101825), the University of North Carolina - Chapel Hill (MH090936  
490 & MH101819), Harvard University (MH090948), Stanford University (MH101782), Washington  
491 University St Louis (MH101810), and the University of Pennsylvania (MH101822). The data  
492 used for the analyses described in this manuscript were obtained from [dbGaP](#) accession  
493 number [phs000424.v6.p1](#) on October 6, 2015.



494 Data were generated as part of the CommonMind Consortium supported by funding from  
495 Takeda Pharmaceuticals Company Limited, F. Hoffman-La Roche Ltd and NIH grants  
496 R01MH085542, R01MH093725, P50MH066392, P50MH080405, R01MH097276, RO1-MH-  
497 075916, P50M096891, P50MH084053S1, R37MH057881 and R37MH057881S1,  
498 HHSN271201300031C, AG02219, AG05138 and MH06692. Brain tissue for the study was  
499 obtained from the following brain bank collections: the Mount Sinai NIH Brain and Tissue  
500 Repository, the University of Pennsylvania Alzheimer's Disease Core Center, the University of  
501 Pittsburgh NeuroBioBank and Brain and Tissue Repositories and the NIMH Human Brain  
502 Collection Core. CMC Leadership: Pamela Sklar, Joseph Buxbaum (Icahn School of Medicine  
503 at Mount Sinai), Bernie Devlin, David Lewis (University of Pittsburgh), Raquel Gur, Chang-Gyu  
504 Hahn (University of Pennsylvania), Keisuke Hirai, Hiroyoshi Toyoshiba (Takeda  
505 Pharmaceuticals Company Limited), Enrico Domenici, Laurent Essioux (F. Hoffman-La Roche  
506 Ltd), Lara Mangravite, Mette Peters (Sage Bionetworks), Thomas Lehner, Barbara Lipska  
507 (NIMH).

508 Members of the BrainSeq consortium include: Christian R. Schubert, Patricio O'Donnell, Jie  
509 Quan, Jens R. Wendland, Hualin S. Xi, Ashley R. Winslow, Enrico Domenici, Laurent Essioux,  
510 Tony Kam-Thong, David C. Airey, John N. Calley, David A. Collier, Hong Wang, Brian  
511 Eastwood, Philip Ebert, Yushi Liu, Laura Nisenbaum, Cara Ruble, James E. Scherschel, Ryan  
512 Matthew Smith, Hui-Rong Qian, Kalpana Merchant, Michael Didriksen, Mitsuyuki Matsumoto,  
513 Takeshi Saito, Nicholas J. Brandon, Alan J. Cross, Qi Wang, Hussein Manji, Hartmuth Kolb,  
514 Maura Furey, Wayne C. Drevets, Joo Heon Shin, Andrew E. Jaffe, Yankai Jia, Richard E.  
515 Straub, Amy Deep-Soboslay, Thomas M. Hyde, Joel E. Kleinman, Daniel R. Weinberger

516

## 517 **Figure Legends**

518 **Figure 1:** Developmental regulation of expression. (A) Principal component #1 of the gene-level  
519 expression data versus age; PCW: post-conception weeks, remaining ages are in years. (B)  
520 Expression features fall into two main development regulation signatures, increasing in  
521 expression from fetal to postnatal life (orange) or decreasing from fetal to postnatal life (blue). Y-  
522 axis is Z-scaled expression (to standard normal), dark lines represent median expression levels,  
523 and confidence bands represent 25<sup>th</sup>-75<sup>th</sup> percentiles of expression levels for each class of  
524 features. (C) KEGG pathways enriched for genes with isoform shifts, stratified by which feature  
525 type identified the gene as having a switch. Coloring/scaling represents  $-\log_{10}(\text{FDR})$  for gene  
526 set enrichment. Analogous data for GO gene sets (biological processes, BP, and molecular  
527 function, MF) are available in Table S6. DER: differentially expressed region. Enrichment  
528 analyses for isoform shift genes among PGC2 schizophrenia GWAS risk loci with exon and  
529 junction counts using both (D) parametric p-values) and (E) permutation-based p-values. OR:  
530 odds ratio.

531 **Figure 2:** Clinical enrichment of schizophrenia risk using representative eQTLs. (A) Association  
532 between rs1233578 and intergenic sequence downstream (B) of ZSCAN23. (B) Association  
533 between rs3849046 and a splice junction (C) of a particular longer isoform (D) of *CNNTA1*. (E)  
534 Association between rs9841616 and very proximal extended UTR (F) of SOX2-OT.

535 Associations between risk SNPs and annotated sequences are shown for (G) CD46, (H) SRR  
 536 and (I) GPM6A. In panels B, D, and F: thicker/dark blue: exon, thinner/light blue: intron;  
 537 coordinates relative to hg19.

538 **Figure 3:** Differential expression comparing patients with schizophrenia to controls. (A)  
 539 Histogram of fold changes of the diagnosis effect of those features that were significant and  
 540 independently replicated, colored by feature type. (B) Gene set analyses of genes with  
 541 decreased expression in patients compared to controls by feature type. Coloring/scaling  
 542 represents  $-\log_{10}(\text{FDR})$  for gene set enrichment. Significant directional effects of developmental  
 543 regulation among diagnosis-associated features for those features that (C) increased and (D)  
 544 decreased across development (i.e. those features shown in Figure 1B). P-values provided for  
 545 Wilcoxon rank sign test for those features developmentally regulated among case-control  
 546 differences to those not developmentally regulated.

547

## 548 Tables

549 Table 1: eQTL summary statistics at FDR and Bonferroni significance thresholds across five  
 550 feature summarizations. “logFC” is the log<sub>2</sub> fold change in expression per minor allele copy and  
 551 “% Unann” is the percent of features that were not strictly annotated.

	Type	eQTLs	# SNPs	# Features	p-cutoff	Ensembl Genes	Symbol Genes	log <sub>2</sub> FC	% Unann
FDR < 1%	Gene	1815172	1055186	18416	1.84E-04	18416	12874	0.061	NA
	Exon	13255860	1390362	157923	1.00E-04	20696	15697	0.13	NA
	Transcript	1465179	616346	26870	3.07E-05	11272	11219	0.094	50.7%
	Junction	4813472	1092615	67358	6.39E-05	14792	13204	0.33	21.3%
	ER	8115891	1367619	94200	1.25E-04	16379	12914	0.22	47.4%
Bonf < 5%	Gene	648597	431704	6748	8.41E-09	6748	4955	0.097	NA
	Exon	4019197	529237	48031	7.64E-10	8386	6439	0.21	NA
	Transcript	514563	236633	6349	1.73E-09	3263	3249	0.15	46.9%
	Junction	1557370	439920	18908	1.10E-09	5827	5205	0.55	21.6%
	ER	2575655	533978	27643	1.28E-09	6822	5643	0.37	53.9%

552

553

554

555 Table 2: eQTL summary metrics for GWAS variants from the latest schizophrenia GWAS and  
556 the more general genome-wide suggestive loci from the NHGRI GWAS catalog. “# SNPs  
557 Tested” were those that were observed or imputed with high quality and that were relatively  
558 common in our samples (MAF > 5%). “Unann” = unannotated, “Tx” = transcript

	SCZD GWAS			NHGRI GWAS Catalog		
	FDR<1%	FDR+Meta	Bonf<5%	FDR<1%	FDR+Meta	Bonf<5%
# SNPs Tested	106	106	106	23704	23704	23704
# SNP eQTLs	51	37	26	8988	5490	4255
> # w/o Gene	21	17	9	3763	2370	1891
> # w/o Gene+Tx	17	15	8	2982	1824	1445
> # Unann	47	28	17	5858	3470	2579
> # Only unann	7	6	3	995	671	589
> # Single Tx	11	10	5	1933	1156	976

559

560

561

562 Table 3: GWAS-significant index variants and eQTL associations, for those GWAS loci  
 563 associating with only one or two genes following conditional analysis.

SZ GWAS Locus	SNP	Gene	SZ GWAS Locus	SNP	Gene
1	rs1233578	Intergenic	59	rs10520163	<i>CLCN3</i>
1	rs1233578	<i>ZSCAN26</i>	63	rs9420	Intergenic
5	rs4129585	<i>TSNARE1</i>	73	rs3849046	<i>CTNNA1</i>
7	rs10650434	<i>MAD1L1</i>	82	rs6704641	<i>SATB2</i>
7	rs10650434	<i>FTSJ2</i>	84	rs1106568	<i>GPM6A</i>
11	rs4702	<i>FES</i>	86	rs10043984	<i>FAM53C</i>
11	rs4702	<i>AC068831.1</i>	86	rs10043984	<i>NME5</i>
12	rs75968099	<i>LRRFIP2</i>	88	rs7819570	<i>AC090568.2</i>
12	rs75968099	<i>AC011816.1</i>	96	rs8082590	<i>ATPAF2</i>
16	rs13240464	<i>LRRN3</i>	96	rs8082590	<i>DRG2</i>
16	rs13240464	<i>IMMP2L</i>	98	rs12325245	<i>GOT2</i>
17	rs10791097	<i>SNX19</i>	98	rs12325245	<i>NDRG4</i>
20	rs7893279	<i>NSUN6</i>	103	rs324017	<i>STAT6</i>
23	rs6704768	<i>C2orf82</i>	106	rs9841616	<i>SOX2-OT</i>
23	rs6704768	<i>GIGYF2</i>	109	rs149009306	<i>DFNA5</i>
24	rs55661361	<i>NRGN</i>	114	rs12421382	<i>AP003049.1</i>
30	rs11682175	<i>FANCL</i>	114	rs12421382	Intergenic
42	rs7432375	<i>AC117382.2</i>	117	rs75575209	<i>FANCL</i>
42	rs7432375	<i>PCCB</i>	119	rs14403	<i>AKT3</i>
47	rs4523957	<i>SRR</i>	119	rs14403	<i>SDCCAG8</i>
47	rs4523957	<i>TSR1</i>	120	rs6670165	<i>BRINP2</i>
52	rs140505938	Intergenic	120	rs6670165	Intergenic
57	rs34269918	<i>RERE</i>	121	rs7523273	<i>CD46</i>
57	rs34269918	<i>SNORA77</i>			

564

565

## 566 References

- 567 1 Birnbaum, R. & Weinberger, D. R. Genetic insights into the neurodevelopmental origins  
568 of schizophrenia. *Nature reviews. Neuroscience*, doi:10.1038/nrn.2017.125 (2017).
- 569 2 Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci.  
570 *Nature* **511**, 421+, doi:10.1038/nature13595 (2014).
- 571 3 Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in  
572 regulatory DNA. *Science* **337**, 1190-1195, doi:10.1126/science.1222794 (2012).
- 573 4 Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from  
574 RNA-seq reads. *Nature biotechnology* **33**, 290-295, doi:10.1038/nbt.3122 (2015).
- 575 5 Nellore, A. *et al.* Human splicing diversity across the Sequence Read Archive. *bioRxiv*,  
576 doi:10.1101/038224 (2016).
- 577 6 Collado Torres, L. *et al.* Flexible expressed region analysis for RNA-seq with derfinder.  
578 *Nucleic Acids Research In Press.*, doi:10.1101/015370 (2016).
- 579 7 Jaffe, A. E. *et al.* Developmental regulation of human cortex transcription and its clinical  
580 relevance at single base resolution. *Nature neuroscience* **18**, 154-161,  
581 doi:10.1038/nn.3898 (2015).
- 582 8 Tan, W. *et al.* Molecular cloning of a brain-specific, developmentally regulated  
583 neuregulin 1 (NRG1) isoform and identification of a functional promoter variant  
584 associated with schizophrenia. *The Journal of biological chemistry* **282**, 24343-24351,  
585 doi:10.1074/jbc.M702953200 (2007).
- 586 9 Kao, W. T. *et al.* Common genetic variation in Neuregulin 3 (NRG3) influences risk for  
587 schizophrenia and impacts NRG3 expression in human brain. *Proceedings of the*  
588 *National Academy of Sciences of the United States of America* **107**, 15619-15624,  
589 doi:10.1073/pnas.1005410107 (2010).
- 590 10 Tao, R. *et al.* Expression of ZNF804A in human brain and alterations in schizophrenia,  
591 bipolar disorder, and major depressive disorder: a novel transcript fetally regulated by  
592 the psychosis risk variant rs1344706. *JAMA psychiatry* **71**, 1112-1120,  
593 doi:10.1001/jamapsychiatry.2014.1079 (2014).
- 594 11 Hyde, T. M. *et al.* Expression of GABA signaling molecules KCC2, NKCC1, and GAD1 in  
595 cortical development and schizophrenia. *The Journal of neuroscience : the official*  
596 *journal of the Society for Neuroscience* **31**, 11088-11095,  
597 doi:10.1523/JNEUROSCI.1234-11.2011 (2011).
- 598 12 Birnbaum, R., Jaffe, A. E., Hyde, T. M., Kleinman, J. E. & Weinberger, D. R. Prenatal  
599 expression patterns of genes associated with neuropsychiatric disorders. *The American*  
600 *journal of psychiatry* **171**, 758-767, doi:10.1176/appi.ajp.2014.13111452 (2014).
- 601 13 Buchanan, R. W. & Carpenter, W. T. Domains of psychopathology: an approach to the  
602 reduction of heterogeneity in schizophrenia. *The Journal of nervous and mental disease*  
603 **182**, 193-204 (1994).
- 604 14 Winterer, G. & Weinberger, D. R. Genes, dopamine and cortical signal-to-noise ratio in  
605 schizophrenia. *Trends in neurosciences* **27**, 683-690, doi:10.1016/j.tins.2004.08.002  
606 (2004).
- 607 15 Niknafs, Y. S., Pandian, B., Iyer, H. K., Chinnaiyan, A. M. & Iyer, M. K. TACO produces  
608 robust multisample transcriptome assemblies from RNA-seq. *Nature methods* **14**, 68-70,  
609 doi:10.1038/nmeth.4078 (2017).
- 610 16 Fromer, M. *et al.* Gene expression elucidates functional impact of polygenic risk for  
611 schizophrenia. *Nature neuroscience* **19**, 1442-1453, doi:10.1038/nn.4399 (2016).
- 612 17 Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue  
613 gene regulation in humans. *Science* **348**, 648-660, doi:10.1126/science.1262110 (2015).



- 614 18 Jaffe, A. E. *et al.* qSVA framework for RNA quality correction in differential expression  
615 analysis. *Proceedings of the National Academy of Sciences of the United States of*  
616 *America* **114**, 7130-7135, doi:10.1073/pnas.1617384114 (2017).
- 617 19 Schizophrenia Working Group of the Psychiatric Genomics, C. Biological insights from  
618 108 schizophrenia-associated genetic loci. *Nature* **511**, 421-427,  
619 doi:10.1038/nature13595 (2014).
- 620 20 Jaffe, A. E. *et al.* Mapping DNA methylation across development, genotype and  
621 schizophrenia in the human frontal cortex. *Nature neuroscience* **19**, 40-47,  
622 doi:10.1038/nn.4181 (2016).
- 623 21 Steijger, T. *et al.* Assessment of transcript reconstruction methods for RNA-seq. *Nature*  
624 *methods* **10**, 1177-1184, doi:10.1038/nmeth.2714 (2013).
- 625 22 Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease.  
626 *Science* **352**, 600-604, doi:10.1126/science.aad9417 (2016).
- 627 23 Weinberger, D. R. & Levitt, P. in *Schizophrenia* 393-412 (Wiley-Blackwell, 2011).
- 628 24 Uhlhaas, P. J. & Singer, W. Abnormal neural oscillations and synchrony in  
629 schizophrenia. *Nature reviews. Neuroscience* **11**, 100-113, doi:10.1038/nrn2774 (2010).
- 630 25 Tyner, C. *et al.* The UCSC Genome Browser database: 2017 update. *Nucleic acids*  
631 *research* **45**, D626-D634, doi:10.1093/nar/gkw1134 (2017).

632

633

## 634 **Methods**

635

### 636 Postmortem brain samples

637 Post-mortem human brain tissue was obtained by autopsy primarily from the Offices of the Chief  
638 Medical Examiner of the District of Columbia, and of the Commonwealth of Virginia, Northern  
639 District, all with informed consent from the legal next of kin (protocol 90-M-0142 approved by the  
640 NIMH/NIH Institutional Review Board). Additional post-mortem fetal, infant, child and adolescent  
641 brain tissue samples were provided by the National Institute of Child Health and Human  
642 Development Brain and Tissue Bank for Developmental Disorders (<http://www.BTBank.org>)  
643 under contracts NO1-HD-4-3368 and NO1-HD-4-3383. The Institutional Review Board of the  
644 University of Maryland at Baltimore and the State of Maryland approved the protocol, and the  
645 tissue was donated to the Lieber Institute for Brain Development under the terms of a Material  
646 Transfer Agreement. Clinical characterization, diagnoses, and macro- and microscopic  
647 neuropathological examinations were performed on all samples using a standardized paradigm,  
648 and subjects with evidence of macro- or microscopic neuropathology were excluded. Details of  
649 tissue acquisition, handling, processing, dissection, clinical characterization, diagnoses,  
650 neuropathological examinations, RNA extraction and quality control measures were described  
651 previously in Lipska, *et al.*<sup>26</sup>. The Brain and Tissue Bank cases were handled in a similar  
652 fashion (<http://medschool.umaryland.edu/BTBank/ProtocolMethods.html>). Antipsychotic use  
653 was measured using toxicology at time of death.

654

### 655 RNA extraction and sequencing

656 Post-mortem tissue homogenates of dorsolateral prefrontal cortex grey matter (DLPFC)  
657 approximating BA46/9 in postnatal samples and the corresponding region of PFC in fetal  
658 samples were obtained from all subjects. Total RNA was extracted from ~100 mg of tissue using  
659 the RNeasy kit (Qiagen) according to the manufacturer's protocol. The poly-A containing RNA  
660 molecules were purified from 1 µg DNase treated total RNA and sequencing libraries were  
661 constructed using the Illumina TruSeq® RNA Sample Preparation v2 kit. Sequencing  
662 indices/barcodes were inserted into Illumina adapters allowing samples to be multiplexed in  
663 across lanes in each flow cell. These products were then purified and enriched with PCR to  
664 create the final cDNA library for high throughput sequencing using an Illumina HiSeq 2000 with  
665 paired end 2x100bp reads.

666

### 667 RNA sequencing data processing

668

669 The Illumina Real Time Analysis (RTA) module performed image analysis, base calling, and the  
670 BCL Converter (CASAVA v1.8.2), generating FASTQ files containing the sequencing reads.  
671 These reads were aligned to the human genome (UCSC hg19 build) using the spliced-read  
672 mapper TopHat (v2.0.4) using the reference transcriptome to initially guide alignment, based on  
673 known transcripts of the previous Ensembl build GRCh37.67 (the “-G” argument in the  
674 software)<sup>27</sup>. We achieved a median of 85.3 million (IQR: 71.7M-111.2M) aligned reads per  
675 sample (see Table S1).

676

677 We characterized the transcriptomes of these 495 samples using five convergent  
678 measurements of expression (“feature summarizations”)– (1) gene and (2) exon counts, and (3)  
679 transcript-level quantifications that rely on existing gene annotation, and two annotation-  
680 agnostic approaches we have developed that are determined solely from the read alignments –  
681 (4) read coverage supporting exon-exon splice junctions (e.g. coordinates of potentially intronic  
682 sequence that are spliced out of mature transcripts captured by a single read) and (5) read  
683 coverage overlapping each base in each sample which we have summarized into contiguous  
684 “expressed regions” (ERs, see Methods, Figure S1). These last three measurements generate  
685 expression for features of interest that can “tag” elements of transcripts in the data that are not  
686 constrained by limitations or incompleteness of existing annotation, and the counts for these  
687 features can then be directly used for differential expression analysis.

- 688 1. Gene counts were generated using the featureCounts tool<sup>28</sup> (v1.4.3-p1) based on the more  
689 recent Ensembl v75, which was the last stable release for the hg19 genome build, using  
690 single end read counting [featureCounts -a \$GTF -o \$OUT \$BAM]. We converted counts to  
691 RPKM values using the total number of aligned reads across the autosomal and sex  
692 chromosomes (dropping reads mapping to the mitochondria chromosome).
- 693 2. Exon counts were also generated using the featureCounts tool<sup>28</sup> (v1.4.3-p1) based on the  
694 more recent Ensembl v75, using single end read counting, and allowing reads to be  
695 assigned to multiple exons (e.g. those with splice junctions) [featureCounts -O -f -a \$GTF -  
696 o \$OUT \$BAM]. We converted counts to RPKM values using the total number of aligned  
697 reads across the autosomal and sex chromosomes (dropping reads mapping to the  
698 mitochondria chromosome).

- 699 3. Junction counts were generated by first filtering the TopHat BAM file to primary alignments  
700 only [samtools view -bh -F 0x100 \$BAM > \$NEWBAM ] and regtools<sup>29</sup> (v 0.1.0) was used to  
701 extract analogous junction information (coordinates and number of reads supporting) as the  
702 TopHat output. We found that native TopHat output (junctions.bed) was based on both  
703 primary and secondary alignments, which could influence the degree of potentially novel  
704 splice junctions. We used a modified version of TopHat's "bed\_to\_juncs" program to retain  
705 the number of supporting reads (in addition to returning the coordinates of the spliced  
706 sequence, rather than the maximum fragment range), and used R code (see Supplementary  
707 Code) to combine and annotate these junctions across all samples. We identified splice  
708 junctions using Ensembl v75 – while the initial alignment was guided by Ensembl v67, novel  
709 junctions, by definition, are identified in the second genome alignment, rather than the initial  
710 guided transcriptome alignment step. We converted counts to "RP80M" values, or "reads per  
711 80 million mapped" using the total number of aligned reads across the autosomal and sex  
712 chromosomes (dropping reads mapping to the mitochondria chromosome), which can be  
713 interpreted as the number of reads supporting the junction in an average library size (we  
714 were targeting 80M reads in the sequencing). Most junctions were lowly expressed in our  
715 homogenate tissue, with fewer than 1 average normalized supporting read (N=3,330,642;  
716 92.98%) including approximately half unique to a single individual (N= 1,779,241, 49.67%).
- 717 4. Transcripts were assembled using StringTie<sup>4</sup> (version 1.1.2) guided by Ensembl v75  
718 annotation within each sample [stringtie \$BAM -o \$OUT -G \$GTF]. We then used  
719 "CuffMerge"<sup>30</sup> to merge all assembled transcriptomes across all samples, and then re-  
720 quantified the expression of each transcript isoform in each sample again using StringTie to  
721 this global set of transcripts [stringtie \$BAM -B -e -o \$OUT -G \$GTF\_ALL] to have  
722 expression measurements on the same transcripts across all samples. We then used the  
723 "ballgown" tool<sup>31</sup> to merge all assembled and quantified transcripts across all samples (N=  
724 733,339), and used liberal filtering to remove lowly or uniquely expressed transcripts (mean  
725 FPKM > 0.025), resulting in 188,578 transcripts across the 495 samples.
- 726 5. Expressed regions (ERs) were calculated using the "derfinder" R Bioconductor package<sup>6</sup>  
727 using a cutoff of 5 normalized (to 80M reads) read coverage, which identified 389,797 ERs.  
728 We retained the 275,885 ERs that were at least 12 basepairs, and annotated the ERs to  
729 Ensembl v75.

730

### 731 Genotype data processing

732 SNP genotyping with HumanHap650Y\_V3 (N=135), Human 1M-Duo\_V3 (N=357), and Omni5  
733 (N=3) BeadChips (Illumina, San Diego, CA) was carried out according to the manufacturer's  
734 instructions with DNA extracted from cerebellar tissue. Genotype data were processed and  
735 normalized with the crlmm R/Bioconductor package<sup>32</sup> separately by platform. Genotype  
736 imputation was performed on high-quality observed genotypes (removing low quality and rare  
737 variants) using the prephasing/imputation stepwise approach implemented in IMPUTE2<sup>33</sup> and  
738 Shape-IT<sup>34</sup>, with the imputation reference set from the full 1000 Human Genomes Project Phase  
739 3 data set, separately by platform. We retained common SNPs (MAF > 5%) that were present in  
740 the majority of samples (missingness < 10%) that were in Hardy Weinberg equilibrium (at  $p >$   
741  $1 \times 10^{-6}$ ) using the Plink<sup>35</sup> version 1.9 tool kit [ plink --bfile \$BFILE --geno 0.1 --maf 0.05 --hwe  
742 0.000001 ]. We then identified linkage disequilibrium (LD)-independent SNPs to use in genome-

743 wide clustering of samples and in the number of independent eQTL tests performed [ `plink –  
744 bfile \$BFILE --indep 100 10 1.25`]. Multidimensional scaling (MDS) was performed on the  
745 autosomal LD-independent construct genomic ancestry components on each sample, which can  
746 be interpreted as quantitative levels of ethnicity – the first component separated the Caucasian  
747 and African American samples. This processing and quality control steps resulted in 7,421,423  
748 common variants in this dataset of 495 subjects.

749

750 Polygene risk score (PRS) analysis: Using the allelic dosage files following imputation described  
751 above and the SNPs from provided by the PGC to the Lieber Institute that did not contain  
752 completely different clinical subjects used in the GWAS<sup>2</sup>. We considered expression  
753 associations at the gene, exon and junction-level to the PRS scores from the first 5 clinical SNP  
754 sets, corresponding to GWAS p-value thresholds of  $p < 5e-8$  (s1),  $p < 1e-6$  (s2),  $p < 1e-4$  (s3),  $p$   
755  $< 0.001$  (s4), and  $p < 0.01$  (s5) – subsequent SNP sets were ignored due to clinical risk  
756 plateauing at s5. We also focused only on Caucasian individuals (96 cases, 113 controls), as  
757 the s5 PRS was increased in patients relative to controls in this sample ( $p=3.2 \times 10^{-5}$ ), but did not  
758 differ among African Americans ( $p=0.9$ ). Within each expression feature type, we modeled  
759 expression levels as a function of each PRS set (s1-s5), adjusting for 3 MDS components of the  
760 genotype data, sex, and the first  $K$  principal components (PCs) of the normalized expression  
761 features, where  $K$  was calculated using the Buja and Eyuboglu permutation-based algorithm<sup>36</sup> in  
762 the “sva” Bioconductor package<sup>37</sup>. The resulting p-values of PRS on expression, adjusting for  
763 the above factors, were subject to false discovery rate (FDR) control to account for multiple  
764 testing.

765

#### 766 Public data processing

767 *GTEX:* Raw RNA-seq reads from all brain samples with corresponding genotype data were  
768 downloaded from SRA and aligned to the genome using TopHat2<sup>27</sup> (version 2.0.14) using the  
769 iGenomes transcriptome and genome annotations based on hg19. As above, featureCounts<sup>28</sup>  
770 was used to quantify expression of genes and exons relative to Ensembl v75, and junctions  
771 were quantified with regtools<sup>29</sup> as above. We used StringTie with the assembled merged GTF  
772 from the LIBD DLPFC samples on the GTEX BAM files to quantify the same transcripts, and  
773 used bwtool<sup>38</sup> to quantify the coverage of the same expressed regions from the GTEX brain  
774 samples. Genotype data from the two platforms (Illumina Omni 5M and 2.5M) were imputed  
775 separately as described above and merged into a single plink<sup>35</sup> set.

776 *GEUVADIS:* Raw RNA-seq reads from all LCL samples were downloaded from SRA and  
777 aligned to the genome using TopHat2<sup>27</sup> (version 2.0.9) using the iGenomes transcriptome and  
778 genome annotations based on hg19. As above, featureCounts<sup>28</sup> was used to quantify  
779 expression of genes and exons relative to Ensembl v75, and junctions were quantified with  
780 regtools<sup>29</sup> as above. We used StringTie with the assembled merged GTF from the LIBD DLPFC  
781 samples on the GEUVADIS BAM files to quantify the same transcripts, and used bwtool  
782 to quantify the coverage of the same expressed regions from the GEUVADIS LCL samples.

783 *CommonMind Consortium (CMC)*: 547 BAM files were downloaded from Synapse, which were  
784 aligned with TopHat2 (version 2.0.9) using Ensembl v70 transcriptome annotation and the hg19  
785 genome. As above, featureCounts<sup>28</sup> was used to quantify expression of genes and exons  
786 relative to Ensembl v75, and junctions were quantified with regtools<sup>29</sup> as above. We used  
787 StringTie with the assembled merged GTF from the LIBD DLPFC samples on the CMC BAM  
788 files to quantify the same transcripts, and used bwtool to quantify the coverage of the same  
789 expressed regions from the CMC brain samples. Genotypes were converted to plink file sets  
790 from GEN files obtained from Synapse using posterior probabilities > 90%, resulting in genotype  
791 data across 9,506,038 SNPs and 547 samples.

792

### 793 Differential expression across brain development

794 We modeled differential expression across age at each of the five feature summarizations  
795 (gene, exon, junction, transcript, and ER) in the 320 control subjects across the lifespan. We  
796 modeled expression, after transforming with log<sub>2</sub> with an offset of 1, as a function of age after  
797 creating using linear splines with breakpoints at ages: birth (0), 1, 10, 20, and 50, further  
798 adjusting for sex and ancestry/ethnicity (first 3 MDS components). F-statistics were computed  
799 comparing the model containing age (including the linear splines), sex, and ethnicity, to a  
800 statistical model with just sex and ethnicity, with corresponding p-values calculated based on an  
801 F-distribution with 11 and 308 degrees of freedom, and Bonferroni adjustment within each  
802 feature type was performed using the number of features with non-zero expression (gene  
803 RPKM > 0.01, exon RPKM > 0.1, and junction RPKM > 0.2 with non-novel annotation) across  
804 all samples as the number of tests (which varied by feature type). We also computed post-hoc  
805 statistics on the data, including the Pearson correlation between “cleaned” expression (after  
806 regressing out the effects of sex and ethnicity, holding the age effects constant), and age to  
807 determine if the expression of the fetal rose or fell across the lifespan, and also measured the  
808 fetal versus postnatal log<sub>2</sub> fold changes.

809 Preferential isoform usage across aging was determined by identifying the subset of genes (by  
810 Ensembl ID) that contained at least one Bonferroni-significant feature that had positive  
811 correlation with age and another Bonferroni-significant feature that had negative correlation with  
812 age. We also computed the difference in positive and negative correlations as a measure of the  
813 magnitude of the preferential isoform use. Gene set analyses using pre-defined gene ontology  
814 (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) sets were performed using the  
815 clusterProfiler R/Bioconductor package<sup>39</sup>, here using the genes (mapping from Ensembl to  
816 Entrez ID) that had such preferential isoform use to those that were developmentally regulated  
817 (having at least one feature that was associate with age at Bonferroni significance).

818 Enrichments with the PGC2 schizophrenia risk loci – defined by the chr:start-end roughly  
819 corresponding to linkage disequilibrium blocks in the published manuscript - were performed  
820 both parametrically, by overlapping the genomic coordinates of the 108 risk regions with those  
821 genes that had preferential isoform usage, compared to a background of all genes with each set  
822 of expressed features, as well as by permuting the locations of the 108 regions across the  
823 genome 10,000 times and each time, re-computing the overlap within these null regions – see



824 additional details in Jaffe et al 2015<sup>7</sup>. Empirical p-values were calculated by counting the  
825 number of the odds ratios across the 10,000 null permutations to each observed odds ratio.

826

827

### 828 eQTL discovery analyses

829 We performed eQTL analyses separately by feature type (gene, exon, junction, transcript, and  
830 ER) allowing for a 500kb window around each of the 7,421,423 common SNPs in the 412 age >  
831 13 samples, adjusting for ancestry (first three MDS components from the genotype data), sex,  
832 diagnosis, and the first  $K$  principal components (PCs) of the normalized expression features,  
833 where  $K$  was calculated separately by feature type using the Buja and Eyuboglu permutation-  
834 based algorithm<sup>36</sup> in the “sva” Bioconductor package<sup>37</sup> (gene: 22 PCs, exon: 19 PCs, junction:  
835 26 PCs, transcript: 25 PCs, expressed regions: 20 PCs). The eQTL analyses were run using the  
836 MatrixEQTL R package<sup>40</sup>, which returned the  $\log_2$  fold change per allele copy, and  
837 corresponding T-statistic, p-value, and FDR for each SNP-feature pair. We further used the LD-  
838 independent SNPs to estimate the effective number of tests (by counting the number of features  
839 within a 500kb window around each LD independent SNP) for a more conservative Bonferroni  
840 adjustment. For all five feature types, we retained all eQTLs with FDR < 1%.

841

### 842 eQTL replication analyses

843 We sought to replicate all significant SNP-Feature pairs for each eQTL in two independent  
844 datasets across all five feature summarizations: CommonMind Consortium and the GTEx  
845 project. We used chromosome and position of variants to attempt to match across dataset –  
846 almost all SNPs in the discovery sample were present in each replication samples. Within each  
847 dataset, we tested all polymorphic SNPs (e.g. not monomorphic) and corresponding expressed  
848 features, adjusting for the first 10 PCs of each feature summarization type and the first 5 MDS  
849 components of the corresponding common genotype data. Analyses within CMC were  
850 performed on the 285 controls and analyses in GTEx were performed within each brain region  
851 separately. After identifying and matching back on SNP-feature eQTL pairs, we checked  
852 whether the counted alleles were the same within the discovery and replication datasets and  
853 flipped the directionality of eQTL associations where the alleles were discordant. Note that in  
854 GTEx, some residual discordancy was still present across dataset (e.g. off-diagonal points in  
855 Figures S6, S7A and S7B) but not within a dataset (Figure S7C). Meta-analysis between  
856 discovery (LIBD) and CMC was performed using Stouffer’s Methods<sup>41</sup>, by summing the T-  
857 statistics and dividing by the square-root of the number of datasets (N=2). Meta-analysis within  
858 GTEx brain regions was performed using the same approach, here dividing by the square root  
859 of number of datasets/brain regions (N=13). When replication statistics were not present in  
860 replication datasets due no/low expression or being monomorphic, the discovery eQTL was  
861 “penalized” by setting the replication statistic to 0 prior to meta-analysis.

862

### 863 eQTL clinical enrichment analyses

864 We downloaded the 128 linkage-disequilibrium-independent variants that reached genome-wide  
865 significance in combined analysis from the latest schizophrenia GWAS (their Supplementary  
866 Table 2) and matched those variants to our data by chromosome and position relative to hg19.  
867 Of the 128 variants, only 106 were present in our final QC'd and common (MAF>5%) genotype  
868 data. Most were excluded due to MAFs less than 5% although several variants were dropped  
869 for other reasons (not present in 1000 Genomes, failed Hardy Weinberg equilibrium, poorly  
870 imputed, etc). We therefore interrogated only those 106 schizophrenia-associated variants  
871 among our eQTL associations. We utilized a similar strategy for the latest NHGRI GWAS  
872 catalog (downloaded 7/24/2017) with an additional step of lifting over our variants to hg38 and  
873 again matching by variant coordinates. Here, only approximately half of the variants were well-  
874 measured in our samples (see Table 2).

875

### 876 eQTL conditional analyses

877 We performed conditional analyses within the eQTLs for each schizophrenia risk variant to  
878 remove highly correlated signal and improve resolution of associations. We used the residuals  
879 of the statistical model described above within each feature type (regressing out PCs, MDS  
880 components and diagnosis) to allow for analyses across feature types. We iteratively  
881 conditioned on the expression level of the most significant eQTL feature and recomputed the  
882 eQTL p-values for all other features to the risk SNP. Those features that were still marginally  
883 significantly (at  $p < 0.05$ ) were retained, and then next-best expression feature (following  
884 conditioning) was additionally adjusted for in the statistical model. This procedure of iteratively  
885 testing for conditional independence among remaining features and subsequently adjusting for  
886 the most significant feature continued until no additional features were independently associated  
887 with the genetic risk variant at  $p < 0.05$ . This procedure was performed separately within each of  
888 the 51 loci with eQTL signal.

889

### 890 Schizophrenia differential expression analyses

891 *Discovery dataset analysis:* we first filtered the subjects with RNA-seq to retain a more stringent  
892 set of 155 SCZD cases and 196 controls (criteria: ages between 17-80, gene assignment rate >  
893 0.5, mapping rate > 0.7, RIN > 6, not outlying on 2nd ancestry PC, only self-reported  
894 Caucasians and African Americans). We fit three statistical models across each of the  
895 expression summarizations, modeling  $\log_2$  transformed expression (with an offset of 1) as a  
896 function of:

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

901 (2) Adjusted + Quality Surrogate Variables ("\_qsva" suffix in supplementary tables): SCZD  
902 diagnosis adjusting for "Adjusted" model as well as the first 12 PCs from the degradation matrix  
903 (see below) based on polyA+ libraries (selected using to using the BE algorithm <sup>36</sup> in the sva  
904 Bioconductor package<sup>37</sup> while providing the adjusted model as input).

905 (3) Adjusted + Principal Components ("\_pca" suffix in supplementary tables): SCZD diagnosis  
906 adjusting for "Adjusted" model as well as the first  $k$  PCs from the expressed features (using the  
907 50000 most variable features) depending on the feature type (gene: 23 PCs, exon: 20 PCs,  
908 transcript: 26 PCs, junction: 26 PCs, ERs: 23 PCs).

909 We used the `lmTest` and `ebayes` functions in the limma Bioconductor package <sup>42</sup> to fit all of  
910 the statistical models to estimate  $\log_2$  fold changes, moderated T-statistics, and corresponding  
911 p-values. Multiple testing correction via the false discovery rate (FDR) was applied using the set  
912 of expressed features in this sample set for each summarization type: 24,122 genes (mean  
913 RPKM > 0.1), 420,022 exons (mean RPKM > 0.2), 61,950 transcripts (mean FPKM > 0.2),  
914 229,846 junctions (mean RP80M > 1), and the 275,885 ERs.

915

916 *RNA quality correction:* We summarize the RNA quality correction approach here – for more  
917 detail, see the companion paper by Jaffe et al 2017. Briefly, the quality surrogate variable  
918 analysis (qSVA) uses RNA sequencing data generated from five DLPFC tissue samples left  
919 unfrozen for 0, 15, 30 and 60 minutes, resulting in 20 RNA samples. These samples were  
920 sequenced with both polyA+ and RiboZero library preparations, and gene, exon and junction  
921 counts were derived as above. We utilized the gene-level effects of degradation in these data in  
922 Figure S5 to demonstrate residual confounding by RNA quality, which we call the “DEQual Plot”.

923 For a given preparation type, we identified the genomic regions most susceptible to degradation  
924 by correlating coverage at expressed regions <sup>6</sup> to degradation time, adjusting for donor. This  
925 statistical modeling identified 515 regions significantly susceptible to degradation (at Bonferroni  
926 significance) in the RiboZero libraries and the top 1000 regions most susceptible to degradation  
927 (among the 35,287 at Bonferroni significance) in the polyA+ libraries – the BED files for these  
928 degradation-susceptible regions are available in Jaffe et al 2017<sup>18</sup>

929 The algorithm then involves selecting the set of regions for a particular library type and  
930 calculating total coverage within each region in the new user-provided samples (e.g. the 495  
931 DLPFC RNA-seq polyA+ samples) to form the degradation matrix (which is either 515 or 1000  
932 rows by  $N$  samples). Then PCA is performed on the  $\log_2$  transformed degradation matrix (with  
933 an offset of 1) and the top  $K$  PCs are selected, for example using the BE algorithm <sup>36</sup>, and  
934 extracted – the set of these PCs are referred to as quality surrogate variables (qSVs), and are  
935 included as adjustment variables in subsequent differential expression analyses.

936 *Replication dataset analysis:* we performed analogous sample selection procedures as in the  
937 discovery dataset to select 159 patients and 172 controls (total gene assignment rate > 0.3,  
938 alignment rate > 0.8, RIN > 6, ages between 18-80, non-outlying on genetic ancestry PCs 3 and  
939 5 and keeping only reported Caucasians and African Americans). We similarly fit the three sets

940 of statistical models to all five feature summarizations, with the following differences compared  
941 to the discovery analysis:

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

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

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

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

953 Gene set analyses on replicated differentially expressed features and genes were performed  
954 with clusterProfiler<sup>39</sup> as described above. Set-level analyses on features in the GWAS risk  
955 regions were conducted by assigning each expressed feature a binary variable for whether it  
956 was in the risk regions or not. Then we fit a linear regression model of the t-statistics for  
957 diagnosis, adjusted by the qSVA approach, as a function as whether the feature was in the risk  
958 region, adjusting for its average expression level. This analysis was conducted across and then  
959 within each of the five feature summarization types.

960

## 961 **References**

- 962 2 Ripke, S. *et al.* Biological insights from 108 schizophrenia-associated genetic loci.  
963 *Nature* **511**, 421+, doi:10.1038/nature13595 (2014).
- 964 4 Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from  
965 RNA-seq reads. *Nature biotechnology* **33**, 290-295, doi:10.1038/nbt.3122 (2015).
- 966 6 Collado Torres, L. *et al.* Flexible expressed region analysis for RNA-seq with derfinder.  
967 *Nucleic Acids Research In Press.*, doi:10.1101/015370 (2016).
- 968 7 Jaffe, A. E. *et al.* Developmental regulation of human cortex transcription and its clinical  
969 relevance at single base resolution. *Nature neuroscience* **18**, 154-161,  
970 doi:10.1038/nn.3898 (2015).
- 971 18 Jaffe, A. E. *et al.* qSVA framework for RNA quality correction in differential expression  
972 analysis. *Proceedings of the National Academy of Sciences of the United States of*  
973 *America* **114**, 7130-7135, doi:10.1073/pnas.1617384114 (2017).
- 974 26 Lipska, B. K. *et al.* Critical factors in gene expression in postmortem human brain: Focus  
975 on studies in schizophrenia. *Biological psychiatry* **60**, 650-658,  
976 doi:10.1016/j.biopsych.2006.06.019 (2006).
- 977 27 Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of  
978 insertions, deletions and gene fusions. *Genome biology* **14**, R36, doi:10.1186/gb-2013-  
979 14-4-r36 (2013).

- 980 28 Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for  
981 assigning sequence reads to genomic features. *Bioinformatics* **30**, 923-930,  
982 doi:10.1093/bioinformatics/btt656 (2014).  
983 29 regtools v. 0.1.0 (2016).  
984 30 Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals  
985 unannotated transcripts and isoform switching during cell differentiation. *Nature*  
986 *biotechnology* **28**, 511-515, doi:10.1038/nbt.1621 (2010).  
987 31 Frazee, A. C. *et al.* Ballgown bridges the gap between transcriptome assembly and  
988 expression analysis. *Nature biotechnology* **33**, 243-246, doi:10.1038/nbt.3172 (2015).  
989 32 Scharpf, R. B., Irizarry, R. A., Ritchie, M. E., Carvalho, B. & Ruczinski, I. Using the R  
990 Package crlmm for Genotyping and Copy Number Estimation. *J Stat Softw* **40**, 1-32  
991 (2011).  
992 33 Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation  
993 method for the next generation of genome-wide association studies. *PLoS genetics* **5**,  
994 e1000529, doi:10.1371/journal.pgen.1000529 (2009).  
995 34 Delaneau, O., Coulonges, C. & Zagury, J. F. Shape-IT: new rapid and accurate  
996 algorithm for haplotype inference. *BMC bioinformatics* **9**, 540, doi:10.1186/1471-2105-9-  
997 540 (2008).  
998 35 Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based  
999 linkage analyses. *American journal of human genetics* **81**, 559-575, doi:10.1086/519795  
1000 (2007).  
1001 36 Buja, A. & Eyuboglu, N. Remarks on Parallel Analysis. *Multivariate Behavioral Research*  
1002 **27**, 509-540, doi:10.1207/s15327906mbr2704\_2 (1992).  
1003 37 Leek, J. T., Johnson, W. E., Parker, H. S., Jaffe, A. E. & Storey, J. D. The sva package  
1004 for removing batch effects and other unwanted variation in high-throughput experiments.  
1005 *Bioinformatics* **28**, 882-883, doi:10.1093/bioinformatics/bts034 (2012).  
1006 38 Pohl, A. & Beato, M. bwtool: a tool for bigWig files. *Bioinformatics* **30**, 1618-1619,  
1007 doi:10.1093/bioinformatics/btu056 (2014).  
1008 39 Yu, G., Wang, L. G., Han, Y. & He, Q. Y. clusterProfiler: an R package for comparing  
1009 biological themes among gene clusters. *OmicS : a journal of integrative biology* **16**, 284-  
1010 287, doi:10.1089/omi.2011.0118 (2012).  
1011 40 Shabalin, A. A. Matrix eQTL: ultra fast eQTL analysis via large matrix operations.  
1012 *Bioinformatics* **28**, 1353-1358, doi:10.1093/bioinformatics/bts163 (2012).  
1013 41 Stouffer, S. A. S., E.A.; DeVinney, L.C.; Star, S.A.; Williams, R.M. Jr. The American  
1014 Soldier, Vol.1: Adjustment during Army Life. . *Princeton University Press, Princeton*,  
1015 (1949).  
1016 42 Smyth, G. K. Linear models and empirical Bayes methods for assessing differential  
1017 expression in microarray experiments. *Statistical applications in genetics and molecular*  
1018 *biology* **3**, Article 3 (2004).

1019

1020 **Author information:** sequencing reads and genotype data are available through SRA and  
dbGaP at accession numbers: [TBD].

1022 [Correspondence and requests for materials should be addressed to](#)

1023 Andrew Jaffe ([andrew.jaffe@libd.org](mailto:andrew.jaffe@libd.org)). The following authors have competing interests:

1024 Tony Kam-Thong is employed by F. Hoffmann-La Roche

1025 Hualin S Xi and Jie Quan are employees of Pfizer Inc.



1026 Alan Cross, and Nicholas J. Brandon were full time employees and shareholders in AstraZeneca  
1027 at the time these studies were conducted.

1028 The remaining authors declare no competing financial interests.

## 1029 **Supplementary Information**

1030 *Supplementary Figure Legends:*

1031 **Figure S1:** Study overview and cartoon describing quantifying the five different expression  
1032 summarizations.

1033 **Figure S2:** Cartoon describing the four different splice junction annotation classes, relative to  
1034 annotated exons (dark blue rectangles). (A) Annotated splice junctions map between two exons  
1035 in a known transcript. (B) Exon-skipping splice junctions map to two annotated exons in different  
1036 transcripts. (C) Alternative start/exon junctions map to only one annotated exon on either the 5'  
1037 or 3' end. (D) Completely novel junction do not map to any known exon.

1038 **Figure S3:** Venn diagram of developmentally regulated features mapped back to Ensembl  
1039 Gene IDs by the five feature summarization methods. DER: differentially expressed region.

1040 **Figure S4:** Example of *CRTC2* (A) containing a developmental isoform shift. (B) Gene-level  
1041 analysis shows no developmental regulation but at the junction-level (C) one splice junction  
1042 significantly decreases in expression and (D) another splice junction significantly increases in  
1043 expression over the lifespan. Exons in panels (E), (F), and (H) show some marginal increases in  
1044 expression across the lifespan, but only the exon in (G) is unique to a single isoform and shows  
1045 significant decreases in expression.

1046 **Figure S5:** Venn diagram of Ensembl Gene IDs that contain significant isoform shifts by the four  
1047 feature summarization methods that allow for multiple features per gene. DER: differentially  
1048 expressed region.

1049 **Figure S6:** Discovery (LIBD) and replication (CMC) T-statistics for eQTLs identified in the  
1050 DLPFC for the best SNP-feature pair for each feature across 5 feature summarization types.

1051 **Figure S7:** Assessing regional specificity of eQTLs in GTEx for the best SNP-feature pair for  
1052 each feature across 5 feature summarization types. (A) Significant replication of many eQTLs  
1053 within discovery (LIBD) and Frontal Cortex samples. (B) These DLPFC-identified eQTLs  
1054 showed very significant meta-analysis T-statistics across the 13 brain regions in GTEx. (C)  
1055 These DLPFC-identified eQTLs showed lack of regional specificity even within GTEx.

1056 **Figure S8:** Scatter plot of effect sizes (fold changes) in discovery and replication datasets for  
1057 those features significant and replicated. Colors have the same legend as Figure 3A.

1058 **Figure S9:** GWAS loci set-level analysis for (A) all features together and then stratified by only  
1059 (B) genes, (C) exons, (D) junctions, (E) transcripts and (F) expressed regions. P-values were  
1060 based on the Wilcoxon rank sign test.

1061 *Supplementary Table Legends:*

- 1062 **Table S1:** Demographic information for subjects in the present study, stratified by age and  
1063 diagnosis group. Dx: diagnosis, N: sample size, F: Female, Cauc: Caucasian, SD: standard  
1064 deviation, PCW: post-conception weeks. Antipsychotic use was measured using toxicology at  
1065 time of death. P-values for diagnosis differences in continuous variables are based on linear  
1066 regression and P-values for categorical variables are based on chi-squared tests.
- 1067 **Table S2:** Splice junction annotation and characterization in GTEx and GEUVADIS for any  
1068 junction or highly expressed junctions (mean reads per 80M mapped reads, RP80M > 0, > 1  
1069 and > 5). Each column represents a 2x2 table for presence of identified junctions in 495 DLPFC  
1070 samples in two independent polyA+ datasets.
- 1071 **Table S3:** Summary statistics for those features significantly developmentally regulated in the  
1072 control-only analyses across the lifespan.
- 1073 **Table S4:** Significant developmentally regulated features collapsed to Ensembl Gene ID, used  
1074 to make Figure S3
- 1075 **Table S5:** Isoform shifts by Ensembl Gene ID and feature summarization type.
- 1076 **Table S6:** Gene set analyses for those genes with significant isoform shift, stratified by feature  
1077 summarization type. Q-values, which control the false discovery rate, FDR, are shown.
- 1078 **Table S7:** Genes within the PGC schizophrenia GWAS risk regions that contain isoform shifts  
1079 by feature summarization type. 21.8% of PGC2 genes had developmental isoform shifts using  
1080 exon counts (N=96/440) and 31.9% showed this isoform shift association based on junction  
1081 counts (N=137/430)
- 1082 **Table S8:** Significant eQTLs to schizophrenia GWAS index variants, including replication  
1083 statistics and additional annotation metrics for variants and expressed features. "condIndep"  
1084 column refers to those associations that were conditionally independent.
- 1085 **Table S9:** Differential expression statistics for those features that were significant and replicated  
1086 in case-control comparisons.
- 1087 **Table S10:** Genes consistently differentially expressed by case-control analysis for the different  
1088 feature summarizations.
- 1089 **Table S11:** Gene set analysis for genes with features differentially expressed by case-control  
1090 status, stratified by directionality and feature summarization type.
- 1091 **Table S12:** GWAS region set-level analyses for diagnosis-associated differentially expressed  
1092 features, testing whether features in the PGC risk loci were more or less expressed as a set in  
1093 cases compared to controls. Qual: qSVA adjusted analysis, Adj: observed covariate adjusted  
1094 analysis.
- 1095 **Table S13:** Associations between diagnosis, RPS and expression at gene and exon levels. First  
1096 two columns for each feature: p-values for gene set tests for the significant case-control

1097 features among statistics capturing the effect of RPS on expression. Second two columns for  
1098 each feature: directionality between RPS on expression associations and diagnosis on  
1099







