

1 **Schizophrenia risk from locus-specific human endogenous retroviruses**

2

3 Rodrigo R.R. Duarte^{a,c}, Matthew L. Bendall^{c,d}, Miguel de Mulder^c, Christopher E. Ormsby^e, Greta A.
4 Beckerle^c, Sashika Selvackadunco^f, Claire Troakes^f, Gustavo Reyes-Terán^e, Keith A. Crandall^d,
5 Deepak P. Srivastava^{†b,g}, Douglas F. Nixon^{†c} and Timothy R. Powell^{*†a,c}

6 † Shared senior authors.

7

8 ^a Social, Genetic & Developmental Psychiatry Centre, Institute of Psychiatry, Psychology &
9 Neuroscience, King's College London, London, UK.

10 ^b Department of Basic & Clinical Neuroscience, Institute of Psychiatry, Psychology & Neuroscience,
11 King's College London, London, UK.

12 ^c Division of Infectious Diseases, Weill Cornell Medicine, Cornell University, New York, NY, USA.

13 ^d Computational Biology Institute, Milken Institute School of Public Health, The George Washington
14 University, Washington, DC, USA.

15 ^e Center for Research in Infectious Diseases, National Institute of Respiratory Diseases, Mexico City,
16 Mexico.

17 ^f MRC London Neurodegenerative Diseases Brain Bank, Institute of Psychiatry, Psychology &
18 Neuroscience, King's College London, London, United Kingdom.

19 ^g MRC Centre for Neurodevelopmental Disorders, King's College London, London, United Kingdom.

20

21 * Correspondence to: Dr. Timothy R. Powell, Social, Genetic & Developmental Psychiatry Centre,
22 Institute of Psychiatry, Psychology & Neuroscience, King's College London, 16 De Crespigny Park,
23 London SE5 8AF, United Kingdom. E-mail: timothy.1.powell@kcl.ac.uk. Tel: +44 (0)20 7848 5361.

24 **Abstract**

25 Schizophrenia genome-wide association studies highlight the substantial contribution of risk
26 attributed to the non-coding genome where human endogenous retroviruses (HERVs) are encoded.
27 These ancient viral elements have previously been overlooked in genetic and transcriptomic studies
28 due to their poor annotation and repetitive nature. Using a new, comprehensive HERV annotation,
29 we found that the fraction of the genome where HERVs are located (the 'retrogenome') is enriched
30 for schizophrenia risk variants, and that there are 148 disparate HERVs involved in susceptibility.
31 Analysis of RNA-sequencing data from the dorsolateral prefrontal cortex of 259 schizophrenia cases
32 and 279 controls from the CommonMind Consortium showed that HERVs are actively expressed in
33 the brain (n = 3,979), regulated in *cis* by common genetic variants (n = 1,759), and differentially
34 expressed in patients (n = 81). Convergent analyses implicate LTR25_6q21 and ERVLE_8q24.3h
35 as HERVs of etiological relevance to schizophrenia, which are co-regulated with genes involved in
36 neuronal and mitochondrial function, respectively. Our findings provide a strong rationale for
37 exploring the retrogenome and the expression of these locus-specific HERVs as novel risk factors
38 for schizophrenia and potential diagnostic biomarkers and treatment targets.

39

40 **Introduction**

41 Human endogenous retroviruses (HERVs) are remnants of genetic material acquired through our
42 evolutionary past which originated from the infection of germline cells with ancient retroviruses.
43 These viruses multiplied through a copy-and-paste mechanism and were eventually endogenized
44 (i.e., vertically transmitted), and now constitute approximately 8% of the genome^{1,2}. These repetitive
45 sequences were generally assumed to be transcriptionally inactive in the modern genome, having a
46 purely regulatory function due to the retainment of the viral promoters (long-terminal repeats, LTRs).
47 However, certain HERVs were co-opted to serve novel specialized roles, including in the regulation
48 of embryonic development^{3,4} and neural progenitor cells^{5,6}. They have also been implicated in
49 neuropsychiatric conditions such as amyotrophic lateral sclerosis⁷⁻⁹, major depressive disorder,
50 bipolar disorder, and schizophrenia¹⁰⁻¹². Despite their abundance in the genome and relevance to
51 disease and fundamental aspects of human biology, the location and function of most HERVs remain
52 elusive to-date.

53

54 Recent large genome-wide association studies (GWAS) comparing schizophrenia cases and non-
55 affected individuals enabled the identification of polymorphisms mediating risk for this disorder¹³⁻¹⁵.
56 These studies highlight the substantial contribution of risk attributed to the non-coding genome,
57 where HERVs are encoded, which are often overlooked in both genomic and transcriptomic
58 studies¹⁶. Until recently, there was no comprehensively annotated map of HERVs in the genome,
59 and there were no computationally efficient tools to analyze the expression of these repetitive
60 sequences with single-locus resolution. Consequently, previous studies were unable to test whether
61 locus-specific HERVs were genetically associated with traits of interest, or to assess their expression
62 in conditions of interest, or to distinguish expressed from dormant HERVs from within the same
63 family, which contributed to the generation of inconclusive findings. For example, Karlsson and
64 colleagues¹⁷ reported decreased ERV9 expression in the brain of schizophrenia patients, whereas
65 Diem and collaborators¹⁸ found the opposite. While the heterogeneity of schizophrenia is also likely
66 to be a contributing factor for these contradictory findings, there are approximately 120 copies of
67 ERV9 in the genome^{19,20} that likely also confounded these reports.

68

69 Recently, there have been significant advances in the annotation of HERVs in the genome²⁰⁻²⁴, and
70 in the development of tools that are able to map repetitive sequences in RNA-sequencing data to
71 their most likely source of origin in the genome^{20,25,26}. These improvements allow for the identification
72 and quantification of specific HERV copies associated with traits of interest. These advances,
73 combined with modern population genetic methods, enabled us to identify the contribution of the
74 ‘retrogenome’ (the full genomic complement of HERVs) and expression of locus-specific HERVs to
75 schizophrenia etiology and neurobiology.

76 **Results**

77 **The contribution of polymorphisms in the retrogenome to schizophrenia risk**

78 A recently developed annotation of putatively functional HERV elements describes 14,968 elements
79 dispersed throughout the genome, which were defined based on the presence of LTRs and remnants
80 of genetic elements which code for viral proteins like *env*, *gag* or *pol*²⁰. We assessed the contribution
81 of common genetic variants within HERVs to schizophrenia susceptibility using GARFIELD²⁷. This
82 pipeline quantifies the co-localization of GWAS results with variants in annotation categories (i.e.
83 variants in the retrogenome), assessing significance using linear models that control for minor allele
84 frequency and linkage disequilibrium. For comparison, enrichment was also calculated for heritable
85 neuropsychiatric traits and non-neuropsychiatric phenotypes: height, body mass index, coronary
86 artery disease, Crohn's Disease, type 2 diabetes, neuroticism, eczema, major depressive disorder,
87 bipolar disorder, Alzheimer's disease, attention deficit hyperactivity disorder, amyotrophic lateral
88 sclerosis, and autism spectrum disorder (**Figure 1**; full results on **Supplemental Table 1**).

89 Polymorphisms within HERVs were significantly enriched for genome-wide significant variants ($P <$
90 5.00×10^{-8}) associated with schizophrenia ($P_{\text{enrichment}} = 1.88 \times 10^{-5}$, $P_{\text{corrected}}$ [for 14 traits and 2 GWAS
91 thresholds tested] = 5.27×10^{-4} , $\beta = 0.90$, 95% CI [0.49, 1.31]). Analysis of variants associated with
92 these traits under a more relaxed P cut-off ($P < 5.00 \times 10^{-5}$) also showed an enrichment for
93 schizophrenia only ($P_{\text{enrichment}} = 4.56 \times 10^{-7}$, $P_{\text{corrected}} = 1.28 \times 10^{-5}$, $\beta = 0.57$, 95% CI [0.35, 0.80]).
94 These findings suggest a role for the retrogenome in the etiology of this neurodevelopmental
95 disorder.

96

97

<<< **Figure 1** >>>

98

99 **The contribution of locus-specific HERVs and HERV families to schizophrenia susceptibility**

100 We co-localized polymorphisms within individual HERVs and those associated with schizophrenia
101 by GWAS using MAGMA²⁸, which calculates gene-level statistics and weighted p-values based on
102 summary statistics whilst adjusting for gene size, single nucleotide polymorphism (SNP) density and
103 linkage disequilibrium. This gene-level enrichment analysis revealed that 148 HERVs from multiple
104 families were significantly enriched for risk variants implicated in schizophrenia, after correcting for

105 the number of HERVs tested ($P_{\text{corrected}} < 0.05 / 12,389$ HERVs in chromosomes 1-22, excluding those
106 at the major histocompatibility locus; **Supplemental Table 2**). The quantile-quantile plot highlights
107 the contribution of many locus-specific HERVs associated with risk, compared to an expected normal
108 distribution (**Figure 2A**), and the Manhattan plot shows their diverse genomic location (**Figure 2B**).

109

110 To explore the contribution of HERV families towards schizophrenia risk, we investigated whether
111 any of the 60 families defined in the HERV annotation were overrepresented in the list of
112 schizophrenia-associated HERVs using a gene-set enrichment analysis in MAGMA. Each HERV
113 was assigned to a family based on the RepBase model that most closely matched the internal region
114 sequences^{20,29}. We observed a nominal association between schizophrenia and the HERVL40 family
115 ($P = 0.02$, $\beta = 0.14 \pm 0.02$, $SE = 0.07$), but this did not survive multiple testing correction ($P_{\text{corrected}}$
116 [for 60 families] > 0.05 ; **Supplemental Figure 1**). These findings suggest that risk for schizophrenia
117 attributed to the retrogenome may occur via locus-specific sequences, as opposed to entire HERV
118 families.

119

120 <<< **Figure 2** >>>

121

122 **HERV expression in the dorsolateral prefrontal cortex**

123 To advance our understanding of HERV expression in the adult brain, we assessed global HERV
124 expression (the retrotranscriptome) in the dorsolateral prefrontal cortex (DLPFC) of 593 post-mortem
125 individuals ($N = 593$), including 279 unaffected controls, 259 schizophrenia patients, 47 bipolar
126 disorder patients and 8 cases broadly diagnosed with an affective disorder. This was achieved by
127 applying the Telescope pipeline²⁰ to the RNA-seq data from the CommonMind Consortium (CMC)
128 dataset. Analysis of these samples revealed that 3,979 HERVs were consistently expressed in the
129 DLPFC according to DESeq2 independent filtering criteria (mean normalized counts = 40.79 [37.8,
130 43.78]; **Figure 3A**). We performed RT-qPCR to confirm the expression of five arbitrarily selected
131 HERVs in an independent post-mortem cohort of control individuals from the London
132 Neurodegenerative Diseases Brain Bank ($N = 10$; **Supplemental Material; Supplemental Figure**
133 **2**).

134

135

<<< **Figure 3** >>>

136

137

138 **Cis-acting HERV eQTLs in the DLPFC**

139 To understand how HERVs are regulated in the brain, we performed an expression quantitative trait
140 loci (eQTL) analysis using all samples (N = 593). Such analysis aims to identify SNPs that explain
141 variation in the expression of HERVs residing in close proximity, to inform about the basic processes
142 responsible for HERV regulation, and to complement our genetic enrichment analyses. The genomic
143 coordinates from the expressed HERVs were remapped to hg19 positions to match genotype
144 information, and only HERVs from chromosomes 1-22 were analyzed (total = 5,349 HERVs, which
145 includes lowly expressed HERVs, according to DESeq2's internal filtering criteria). This analysis
146 revealed that 1,759 HERVs were regulated in *cis* by 1,622 SNPs located within a 1 Mb window
147 upstream or downstream the start site of the HERV annotation, under the false discovery rate of 5%
148 ($q < 0.05$). The majority of eQTLs were located within a 10 kb window upstream or downstream of
149 the annotation start site from the HERV they regulate (**Figure 3B**). Of the 148 HERVs that were
150 enriched for schizophrenia variants, we observed that 19 were significantly regulated by eQTLs in
151 the DLPFC (**Supplemental Tables 2 and 3**, highlighted in green; **Figure 3C**).

152

153 **Case-control differences in HERV expression**

154 We observed 81 HERVs as differentially expressed between cases (N = 259) and unaffected
155 individuals (N = 279) under the false discovery rate of 5%, independent of a genetic association with
156 schizophrenia ($q < 0.05$ [corrected for the 3,979 expressed HERVs]; **Supplemental Figure 4**,
157 **Supplemental Table 4**). Our analysis showed that 36 HERVs were downregulated and 45
158 upregulated in cases, and that expression differences were subtle, with log2 fold-changes of $0.18 \pm$
159 0.10 on average.

160

161 **Convergent analyses implicate ERVLE_8q24.3h and LTR25_6q21 as robust risk factors for**
162 **schizophrenia**

163 To make a more informative assessment of HERV expression differences associated with disease
164 status, we first explored whether HERVs identified as genetic risk factors for this disorder were
165 expressed in the DLPFC. Of the 148 HERVs identified via gene-level enrichment analysis, only 65
166 were consistently expressed across samples according to DESeq2 internal filtering criteria, with
167 mean normalized counts of 57.5, CI 95% [12.59, 102.3]. Of these, we observed that two were
168 significantly upregulated (ERV316A3_12q24.11, log₂ fold-change = 0.16, standard error = 0.06;
169 LTR25_6q21, log₂ fold-change = 0.09, standard error = 0.04), and one downregulated
170 (ERVLE_8q24.3h; log₂ fold-change = -0.09, standard error = 0.04) in schizophrenia cases relative
171 to unaffected controls, under the false discovery rate of 5% ($q < 0.05$ [corrected for 65 HERVs];
172 **Figures 3C and D; Supplemental Table 5**).

173
174 Of the three HERVs enriched for schizophrenia variants and differentially expressed in cases, we
175 observed that two were modulated by eQTLs (**Figures 3C and D**, highlighted in blue), which we
176 initially hypothesized could explain the case-control differences observed. These included
177 **ERVLE_8q24.3h** and **LTR25_6q21**, which were regulated by the top eQTLs, **rs4875048** and
178 **rs174399**, respectively. Importantly, **rs4875048** is associated with schizophrenia and is in linkage
179 disequilibrium with the top association signal at this locus, rs10552126; **rs174399** is in linkage
180 disequilibrium with a risk variant at the locus, rs11153302 (**Table 1**)¹⁴. Strikingly, no gene within a 1
181 Mb window upstream or downstream of either of the two top eQTLs was associated with
182 schizophrenia according to a recent analysis by PsychENCODE
183 (<http://resource.psychencode.org/>)³⁰.

184
185 We observed complex regulatory mechanisms governing expression of these two HERVs.
186 **ERVLE_8q24.3h** was significantly downregulated in cases, but, contrary to what was expected, the
187 risk (A-) allele of **rs4875048** was associated with increased expression of this HERV ($\beta = -0.35$, P_{β}
188 $_{\text{dist.}} = 0.002$, $q = 0.007$, **Figure 4A**). Similarly, **LTR25_6q21** was upregulated in patients, but the risk
189 (G-) allele of **rs174399** was associated with reduced expression of this HERV ($\beta = 0.29$,
190 $P_{\beta \text{ dist.}} = 0.004$, $q = 0.02$; **Figure 4B**). We investigated the effect of these eQTLs in cases and control
191 individuals separately, which revealed that they influenced HERV expression exclusively in

192 unaffected individuals, suggesting there is a compensatory mechanism in patients counteracting the
193 effects of the eQTLs (**Figures 4A and B; Supplemental Table 6**).

194

195 To understand about the regulation of these HERVs during neurodevelopment, we tested their
196 expression in an *in vitro* model of cortical development, which consisted of neural stem cells from
197 the CTX0E16 cell line and cells differentiated for 28 days³¹⁻³³ (**Supplemental Material,**
198 **Supplemental Figure 5**). We observed that both HERVs were differentially regulated during
199 differentiation, suggesting that risk to schizophrenia pertaining these HERVs starts during
200 neurodevelopment. Interestingly, we also observed that both HERVs are located in the antisense
201 strand of *SLC16A10* and *IQANK1* introns, respectively. Moreover, LTR25_6q21 is exclusively
202 present in humans, whereas ERVLE_8q24.3h is shared with primates (**Supplemental Figure 6**),
203 according to the UCSC Genome Browser³⁴.

204

205

<<< **Figure 4** >>>

206

207 **ERVLE_8q24.3h and LTR25_6q21 are co-regulated with genes implicated in mitochondrial**
208 **and synaptic function in the adult brain, respectively**

209

210 To explore the potential function of ERVLE_8q24.3h and LTR25_6q21, we determined the genes
211 that are co-expressed with these HERVs by applying Weighted Correlation Network Analysis
212 (WGCNA)³⁵ to the RNA-seq data from schizophrenia patients and controls combined. This systems
213 biology approach is based on the hypothesis that genes within the same co-regulated network share
214 the same function³⁶. We observed 19 modules of co-expression in these data (**Supplemental**
215 **Tables 7 and 8, Supplemental Figure 7**). LTR25_6q21 was assigned to the green module (module
216 membership statistic (MM) = 0.41, $P = 3.55 \times 10^{-23}$; gene significance (GS) statistic in relation to
217 case-control status: 0.12, GS $P = 0.005$), whereas ERVLE_8q24.3h was assigned to the turquoise
218 module (MM = 0.46, $P = 4.81 \times 10^{-29}$; GS = -0.15, GS $P = 5.14 \times 10^{-4}$). Gene Ontology (GO) analysis
219 of the green module indicates that this gene set is associated with neuronal function, with significant
220 terms including “synapse organization”, “presynapse” and “vesicle-mediated transport in synapse”

221 (q < 0.05, **Figure 5A**, upper panel, **Supplemental Table 9**). The turquoise module, in turn, was
222 significantly associated with mitochondrial function, with terms including the “respiratory chain” and
223 “mitochondrial matrix”, and “NADH dehydrogenase complex assembly” (q < 0.05, **Figure 5A**, lower
224 panel, **Supplemental Table 9**).

225

226 A correlation between the first principal component capturing the variability within each module
227 (module eigengene) and case-control status (‘Profile’) was performed in WGCNA, which revealed
228 that the green module is positively associated with case-control status (r = 0.13, P = 0.003), whereas
229 the turquoise module was negatively associated (r = -0.15, P = 4 x 10⁻⁴; **Figure 5B**). The co-
230 expression modules were detected assuming a signed network, which means that a positive module-
231 trait correlation entails higher expression of genes in the module in association with disease status,
232 and vice-versa, corroborating the case-control differences observed for each HERV in the previous
233 analysis (**Figures 4A and B**, left panels).

234

235 To understand the function of other HERVs associated with schizophrenia from the gene-level
236 enrichment analysis performed using MAGMA (**Supplemental Table 2**), we identified the co-
237 expression modules associated with each enriched HERV that was expressed in the brain
238 (**Supplemental Table 10**). We observed that several of these HERVs were assigned to modules
239 implicated in synaptic organization, including MER41_1p33 and MER41_5p12 alongside
240 LTR25_6q21 in the green module, as well as synaptic function, including HERVL_14q24.2d,
241 HERVL40_8q21.3b, HERVW_6q21c, HML3_5p12a and MER41_2q31.2 in the red module.
242 Interestingly, we observed that these two modules were clustered together in an unsupervised
243 hierarchical clustering analysis performed in WGCNA, corroborating a related or shared function
244 (**Supplemental Figure 8**).

245

246

<<< Figure 5 >>>

247

248 **LTR25_6q21 belongs to a module associated with increased neuronal counts, whereas**
249 **ERVLE_8q24.3h does not robustly correlate with major neural cell types**

250 We investigated the major cell types associated with each modules eigengene (the first principal
251 component of the module) to infer the cell types associated with ERVLE_8q24 and LTR25_6q21.
252 We estimated the proportion of neural cell types in the RNA-sequencing data by applying the R
253 package BRETIGEA³⁷ to the normalized gene counts of the RNA-sequencing data. BRETIGEA
254 estimates the cell type proportions that constitute the sequenced material based on a database of
255 single-cell RNA-sequencing data, ultimately generating coefficients that represent the proportion of
256 astrocytes, microglia, endothelial cells, oligodendrocytes, oligodendrocyte progenitor cells and
257 neurons in the sequenced samples³⁷. To infer the cell types associated with each module, we
258 performed correlations between each module and the cell-type proportion coefficients (**Figure 5B**).
259 The module assigned to LTR25_6q21 (green) was significantly correlated with the expression of
260 neuronal markers ($r = 0.24$, $P = 3 \times 10^{-8}$), and negatively associated with the expression of all other
261 cell types ($q < 0.05$, corrected for the six cell types tested and the two modules of interest), consistent
262 with a role for this module in the regulation of neuronal function, as indicated by the GO analysis.
263 The module assigned to ERVLE_8q24.3h (turquoise) was negatively correlated with the expression
264 of oligodendrocyte progenitor cell markers ($r = -0.21$, $P = 6 \times 10^{-7}$, $q < 0.05$), but showed no
265 correlation with the other cell types, which may suggest a specific role for this HERV in non-dividing
266 cells. Ultimately, these data indicate that LTR25_6q21 and ERVLE_8q24.3h are part of co-regulated
267 networks implicated in neuronal and mitochondrial function, respectively, suggesting their
268 involvement in these cell functions.

269
270
271

272 **Discussion**

273 HERVs are ancient viral genetic elements scattered throughout the genome, with previously
274 hypothesized influences on neurodevelopment and risk for schizophrenia. We took a comprehensive
275 approach to reconsider the role of HERVs at the omics level, leveraging on the recent advances in
276 the genomic annotation of HERVs, single-locus resolution quantification, systems biology methods,
277 and population genetic tools.

278

279 We investigated the combined contribution of HERVs (the retrogenome) to risk for schizophrenia
280 and other complex polygenic traits, by assessing the overlap between SNPs implicated in these traits
281 and those in genomic locations encompassing HERVs. We were surprised to find that schizophrenia
282 was the only tested trait significantly enriched for common variants within the retrogenome,
283 especially since there is evidence linking HERVs to conditions like amyotrophic lateral sclerosis⁷⁻⁹,
284 Crohn's disease³⁸, major depressive disorder and bipolar disorder^{10,11}. The fact that the retrogenome
285 is significantly enriched for polymorphisms implicated in schizophrenia suggests that HERVs
286 comprise an important set of risk factors for this disorder, within the 'non-coding' genome.

287

288 We investigated which HERV families and locus-specific HERVs could be particularly important in
289 moderating risk for schizophrenia at the genetic level by co-localizing risk variants with individual
290 HERV loci from across 60 families. We identified 148 specific HERVs enriched for schizophrenia-
291 associated SNPs. A subsequent gene-set analysis did not find a particular HERV family enriched for
292 schizophrenia variants, further suggesting that disparate HERVs scattered throughout the genome
293 moderate susceptibility rather than whole families. These findings have several implications for the
294 interpretation of previous studies in the literature, and may explain how our results differ from
295 previous data implicating other HERVs or HERV families as risk factors, which may not have been
296 identified here. Most HERV expression research to-date has been performed using microarrays,
297 antibodies or RT-qPCR probes, which do not provide sufficient specificity to assess single HERVs<sup>7-
298 11</sup>. Therefore, previous studies may have captured the expression of multiple HERVs concomitantly
299 due to their repetitive sequences. Based on our findings, it is unlikely that individual HERVs (even
300 from within the same family) contribute equally to risk. The HERV annotation developed by Bendall

301 and colleagues²⁰, as well as the Telescope pipeline, and the application of modern population
302 genetic methods, enabled us to revisit the role HERVs play in relation to schizophrenia risk on a
303 genome-wide scale, and allow for more robust inferences regarding their etiological relevance.

304
305 Next, to better understand HERV expression regulation in the brain, we analyzed RNA-sequencing
306 data from the CommonMind Consortium, which confirmed widespread HERV expression in the adult
307 brain and cis-regulatory mechanisms. We found 3,979 HERVs actively expressed in a key brain area
308 linked to schizophrenia pathophysiology, the DLPFC³⁹. We further identified 1,759 HERVs that were
309 regulated by 1,622 short-range (cis-) eQTLs in the DLPFC. The identification of eQTLs that impact
310 HERV expression informs us about the basic processes responsible for HERV regulation, and can
311 be useful for the interpretation of GWAS findings in the context of the retrogenome^{40,41}. In addition,
312 these results suggest that SNPs within HERVs are not simply affecting the expression of neighboring
313 protein-coding genes via their LTRs, rather it demonstrates that common genetic variation impacts
314 locally on HERV expression.

315
316 We used a complementary set of analyses to identify the most robust HERVs implicated in
317 schizophrenia. Two HERVs, **ERVLE_8q24.3h** and **LTR25_6q21**, identified from the gene-level
318 enrichment analysis, were found to be regulated by schizophrenia-associated eQTLs, and were
319 found to be differentially expressed in the DLPFC of patients and in an *in vitro* model of cortical
320 neurodevelopment. These findings suggest a potential risk mechanism for schizophrenia that starts
321 during neurodevelopment and persists through to adulthood, as observed for schizophrenia risk
322 genes such as *NT5C2*, *AS3MT*, and *BORCS7*^{33,42}. Importantly, we observed a complex regulation
323 of both HERVs, whereby the lead eQTLs only exerted their regulatory effects in unaffected
324 individuals. This suggests that compensatory mechanisms (e.g. epigenetic alterations) may be
325 acting to correct for the effects of the risk variants on the expression of these HERVs in the DLPFC
326 of schizophrenia patients.

327
328 We also describe here, for the first time, the co-regulation of several HERVs with known genes in
329 the adult brain, and the GO terms associated with each co-expression module. Our findings suggest

330 that LTR25_6q21 is implicated in neuronal function, whereas ERVLE_8q24.3h is involved in
331 mitochondrial function. WGCNA has been successfully used to predict the biological function of
332 unknown genes or non-coding RNAs in different organisms^{43,44}, and to identify clinically relevant cell
333 types when in combination with cell-type deconvolution analysis⁴⁵, and thus represents a powerful
334 approach to functionally characterize the HERVs expressed in the brain. For a long time HERVs
335 were assumed to be mere regulatory DNA sequences, but the discovery of their expression and co-
336 regulation with several other genes implicated in multiple biological processes in the brain, ranging
337 from neuronal, glial and mitochondrial regulation to splicing and cell motility (**Supplemental Table**
338 **9**), is a landmark for HERV research, and adds an extra layer of complexity to our understanding of
339 human neurobiology.

340

341 There are limitations to this study which should be acknowledged. Schizophrenia is a highly
342 polygenic, heterogeneous disorder, and as such large sample sizes are required for appropriate
343 comparisons. The CommonMind Consortium provides the largest and best characterized cohort of
344 schizophrenia cases and unaffected individuals with RNA-sequencing data to-date, but it might be
345 underpowered for case-control comparisons considering the heterogeneity of schizophrenia.
346 Nevertheless, we complemented case-control comparisons with genomic and eQTL analyses to
347 provide additional insights. Another limitation to our study is that it investigated RNA-sequencing
348 data from bulk DLPFC tissue only, which is composed of a heterogeneous mixture of several types
349 of neurons and glial cells, and it is possible that HERVs expressed in particular cell types, or in other
350 brain regions, are more relevant to risk⁴⁶. To address this, we performed cell type deconvolution
351 using BRETIGEA to determine cell-type specific effects, but ultimately the analysis of data from other
352 brain areas, developmental time points and single-cell datasets has the potential to reveal important
353 insights about the etiology of schizophrenia in relation to HERV expression. In addition, our post-
354 mortem and *in vitro* work suggest HERV expression is important in the DLPFC and during its
355 development, but we still do not know the function of these HERVs. To infer function we performed
356 WGCNA which provides insight into which processes risk HERVs moderate, but future functional
357 studies are required to definitively characterize how HERVs, particularly ERVLE_8q24.3h and
358 LTR25_6q21, influence the transcriptome, neural stem cell proliferation or neuronal differentiation in

359 the context of schizophrenia risk, as is currently being investigated in relation to protein-coding risk
360 genes^{32,33,47,48}.

361

362 The development of a retrogenome annotation, and advances in modern population genetic methods
363 and transcriptomic tools, now allows us to investigate HERVs at the omics level, in the context of
364 risk for many biological traits. Our work studying the role of HERVs in the brain, and their relationship
365 to schizophrenia ignites a new, provocative line of thought implicating HERVs as biological risk
366 factors for schizophrenia and confirms that these previously assumed 'dormant' sequences in the
367 brain may not be dormant after all.

368 **Online Methods**

369 We used a combination of gene expression, genetic and *in vitro* analyses to identify the most robust
370 HERVs implicated in schizophrenia risk (**Figure 6**). Further details are provided in the **Supplemental**
371 **Material**.

372 <<< **Figure 6** >>>

373

374 **Genetic enrichment analyses**

375 We estimated the contribution of genetic polymorphisms within the retrogenome towards risk of
376 developing multiple traits using GARFIELD²⁷. We downloaded summary statistics from well-powered
377 genome-wide association studies, including of schizophrenia (N = 105,318 individuals)¹⁴, height (N
378 = 693,529)⁴⁹, body mass index (N = 681,275)⁴⁹, coronary artery disease (N = 547,261)⁵⁰, Crohn's
379 Disease (N = 59,957)⁵¹, type 2 diabetes (N = 659,316)⁵², neuroticism (N = 2,370,390)⁵³, eczema (N
380 = 103,066)⁵⁴, major depressive disorder (N = 480,359)⁵⁵, bipolar disorder (N = 51,710)⁵⁶, Alzheimer's
381 disease (N = 74,046)⁵⁷, attention deficit hyperactivity disorder (N = 53,293)⁵⁸, amyotrophic lateral
382 sclerosis (N = 36,052)⁵⁹, and autism spectrum disorder (N = 46,350)⁶⁰, which analyzed European
383 cohorts only. GARFIELD performs greedy pruning of SNPs in GWAS summary statistics (those in
384 linkage disequilibrium, with $R^2 > 0.1$), and quantifies enrichments using odds ratios, assessing their
385 significance by employing generalized linear model testing, controlling for minor allele frequency,
386 and number of linkage disequilibrium proxies ($R^2 > 0.8$). Linkage disequilibrium and allele frequency
387 information were calculated based on the UK10K study. Enrichments were calculated based on
388 summary statistics from each trait using two $P_{\text{association}}$ thresholds: $P < 5 \times 10^{-8}$, to test the enrichment
389 of HERVs within genome-wide significant variants; and a more relaxed threshold, $P < 5 \times 10^{-5}$, to
390 allow signal capture in less powered GWAS. The enrichment significance was corrected for the
391 number of tests performed [14 traits and two P-value thresholds tested]. For consistency with the
392 GWAS data, the HERV annotation used in the expression analysis (hg38) was remapped to hg19
393 coordinates using liftOver³⁴.

394

395 To identify locus-specific HERVs and potential HERV families associated with schizophrenia, we
396 used MAGMA 1.07b²⁸. Briefly, MAGMA calculates gene-level enrichment by generating a gene-wide

397 statistic from summary statistics, adjusting for gene size, variant density, and linkage disequilibrium
398 using the 1000 Genomes Phase 3 European reference panel. SNPs from the summary statistics
399 were assigned to HERVs using an annotation window of 10 kb upstream and downstream of each
400 HERV (as suggested by the authors)²⁸. A Bonferroni correction was applied to identify significantly
401 enriched HERVs ($P_{\text{cut-off}} < 4.03 \times 10^{-6}$ [$0.05 / 12,393$ HERVs in chromosomes 1-22, excluding the
402 major histocompatibility locus]). Q-Q and Manhattan plots were generated using qqman 0.1.4⁶¹.
403 Gene-set enrichment analysis was additionally performed using MAGMA, to test whether HERVs
404 associated with schizophrenia in the previous step were enriched for any of the 60 HERV families
405 (excluding HERVs located in sex chromosomes or at the MHC locus – chromosome 6: 26 – 34 Mb).

406

407 **The CommonMind Consortium dataset**

408 To identify HERV expression differences in schizophrenia patients, or HERV expression quantitative
409 trait loci (eQTL) in the dorsolateral prefrontal cortex, we analyzed RNA-sequencing data from the
410 CommonMind Consortium (release 1.0, N = 593 individuals, <https://doi.org/10.7303/syn2759792>)³⁹.
411 This dataset consisted of dorsolateral prefrontal cortex (DLPFC) samples from 279 unaffected
412 individuals, 259 schizophrenia cases, 47 bipolar disorder patients and 8 cases broadly diagnosed
413 with an affective disorder. Access to this dataset, which includes expression, genotype and clinical
414 data, was granted under a Material Transfer Agreement with the NIMH Repository and Genomics
415 Resources. Briefly, autopsy samples from the Mount Sinai NIH Brain Bank and Tissue Repository,
416 the University of Pennsylvania Brain Bank of Psychiatric illnesses and Alzheimer's Disease Core
417 Center, and The University of Pittsburgh Brain Tissue Donation Program, were sent to the Icahn
418 School of Medicine at Mount Sinai for nucleic acid isolation and sequencing. Individuals were
419 diagnosed according to the Diagnostic and Statistical Manual of Mental Disorders, 4th Edition, as
420 determined in consensus conferences after review of medical records and interviews of family
421 members and care providers. Total RNA was extracted from autopsy tissue using the RNeasy kit
422 (QIAGEN, Hilden, Germany). Ribosomal RNA was depleted using the Ribo-Zero Magnetic Gold kit
423 (Illumina, San Diego, California, United States), libraries were constructed using the TruSeq RNA
424 Sample Preparation Kit v2 (Illumina), and samples were sequenced on an Illumina HiSeq 2500. For
425 whole-genome genotyping, DNA was extracted using the DNeasy Blood and Tissue Kit (QIAGEN)

426 according to the manufacturer's protocol, and samples were genotyped using Illumina Infinium
427 HumanOmniExpressExome 8 1.1b chips. Further details on quality control and sample processing
428 are described in Fromer and colleagues³⁹.

429

430 **RNA-sequencing data processing and HERV expression quantification**

431 Bam files containing mapped and unmapped RNA-sequencing reads aligned to the human reference
432 genome (hg19) using TopHat 2.0.9 and Bowtie 2.1.0, were downloaded to the King's College London
433 High Performance Computer Cluster Rosalind, using the synapse client (1.7.5). Bam files were
434 merged, and fastq files were extracted using samtools 1.5⁶² and the flag '-F 0x100'. Trimmomatic
435 0.38⁶³ was used to prune Illumina adaptors, low quality bases (leading/trailing sequences with phred
436 score < 3, or those with average score < 15 every four bases), or reads below 36 bases in length.
437 Trimmed reads were mapped to the human genome hg38 using bowtie2⁶⁴ and the parameters --
438 very-sensitive-local, -k 100, and --score-min L,0,1.6. Subsequently, Telescope 1.0.2 was used to
439 quantify expression of 14,968 HERVs (annotation version hg38), which we defined as the
440 retrogenome²⁰. We analyzed HERVs with counts > 10 across 4 samples at least, to avoid inflation
441 driven by lowly expressed elements. Case-control expression differences (N = 538 individuals in
442 total) were calculated in R⁶⁵ using Wald tests in DESeq2⁶⁶, where data was normalized (median of
443 ratios) and controlled for the main confounders of gene expression estimated by Fromer and
444 colleagues³⁹, which included institution of sample origin, RNA integrity number, gender, post-mortem
445 interval, age (determined in five bins: #1 = 13-29 years, #2 = 30-49 years, #3 = 50-69 years, #4 =
446 70-89 years, #5 = 90+ years), as well as the first five population covariates estimated using
447 multidimensional scaling in PLINK 1.9⁶⁷, and the first ten hidden HERV expression confounders
448 estimated using sva⁶⁸, considering schizophrenia and unaffected individuals only.

449

450 **Whole-genome genotype data processing**

451 Markers with zero alternate alleles, genotyping call rate < 0.98, Hardy-Weinberg $P < 5 \times 10^{-5}$, or
452 individuals with genotyping call rate < 0.90, were removed from the analysis, as described by Fromer
453 and colleagues³⁹. PLINK files were generated containing genotype information for 958,178 variants
454 for the 593 subjects. Marker alleles were phased to the forward strand, and ambiguously stranded

455 markers were removed. Additional genotype information was imputed from the 1000 Genomes
456 Phase 1 reference panel using minimac3 and Eagle v2.3 phasing with the Michigan Imputation
457 Server (<https://imputationserver.sph.umich.edu/index.html>). Genotype information from the 22
458 autosomes was concatenated using bcftools 1.9 (<https://samtools.github.io/bcftools/bcftools.html>),
459 non-single nucleotide polymorphisms (SNPs) were excluded, as well as sites with an imputation R^2
460 < 0.8 , minor allele frequency < 0.05 , or Hardy-Weinberg $P < 5 \times 10^{-5}$.

461

462 **eQTL analysis**

463 Normalized HERV counts per sample were obtained using DESeq2 (N = 593), and HERV expression
464 was tested for the effect of genotype at all variants located within a 1 Mb window upstream or
465 downstream from the annotation start site of each HERV using QTLtools⁶⁹, according to the authors'
466 manual. We covaried for the effect of case-control status, institution of sample origin, RNA integrity
467 number, gender, post-mortem interval, age (five bins, as described previously), the first five
468 population covariates, and ten hidden expression confounders estimated using sva⁶⁸
469 (**Supplemental Figure 3**). eQTL P-values were corrected through estimation of a beta distribution
470 using a minimum of 1,000 permutations and maximum of 10,000, and were further corrected for the
471 number of HERVs tested using the false discovery rate method ($q < 0.05$).

472

473 **Weighted Correlation Network Analysis (WGCNA)**

474 WGCNA is a systems biology approach that enables the identification of co-expressed genes in
475 transcriptomic data, which we used here to identify the genes co-expressed with schizophrenia
476 HERVs in order to infer their biological function³⁵. We used this tool to construct a signed network
477 consisting of HERVs and genes, which was created based on an adjacency matrix that informs about
478 the co-expression similarity observed between all pairs of genes and HERVs in the expression data
479 (i.e. genes and genes, genes and HERVs, HERVs and HERVs). Normalized HERV and gene counts
480 were variance-stabilized in DESeq2, and were further adjusted for all confounders previously
481 described, using the *removebatcheffect* function in limma⁷⁰. To achieve this, we combined gene and
482 HERV counts obtained from the brain of 538 individuals (279 unaffected individuals and 259
483 schizophrenia cases), and filtered out genes and HERVs that were lowly expressed, i.e. those with

484 < 10 counts in < 80% of samples), as these can drive spurious correlations³⁵. The normalized counts
485 were variance stabilized transformed using DESeq2 and adjusted for institution of sample origin,
486 gender, case-control status, age bins, post-mortem interval, the first five population dimensions
487 (estimated in plink), and RIN, using limma⁷⁰. WGCNA identifies modules by applying hierarchical
488 clustering to the adjacency matrix, further filtering spurious relationships through the application of a
489 topological overlap approach. We used an R^2 cut-off of 0.8, which corresponds to a $\beta = 12$, to
490 construct the network. Each module was assigned a color, and genes or HERVs not belonging to
491 any module were assigned to the gray module. The relationship between modules and specific cell
492 types was tested based on the correlation between the module eigengenes (ME), defined as the first
493 principal component of the module, and cell count estimates, as described below. We applied the
494 false discovery rate method to correct for the module-cell type associations ($q < 0.05$). Plots were
495 generated by WGCNA.

496

497 **Cell type estimates and module correlations**

498 We performed a BRain cEll Type specific Gene Expression Analysis (BRETIGEA)³⁷ to estimate the
499 abundance of major neural cell types in the 538 samples analyzed by WGCNA. Briefly, BRETIGEA
500 uses expression data from single cell RNA-sequencing data sets to identify the proportion of
501 astrocytes, microglia, endothelial cells, oligodendrocytes, oligodendrocyte progenitor cells and
502 neurons, in bulk brain gene expression data. More specifically, this tool uses a panel of 50 well-
503 established cell type-specific markers to generate coefficients that represent the proportion of each
504 cell type per sample, which were tested for association with each module.

505

506 **Gene Ontology (GO) analyses**

507 We performed GO analyses using the WEB-based GENE SeT AnaLysis Toolkit (Webgestalt)⁷¹ to
508 identify the function of the genes co-regulated with the schizophrenia HERVs in the brain, and thus
509 infer the potential function of these HERVs. All genes inputted to WGCNA were used as background
510 (reference) gene set. We used the false discovery rate method to correct for the GO enrichment
511 analyses within Webgestalt ($q < 0.05$) and report up to 10 significant GO terms per module. Volcano
512 plots were generated in Webgestalt.

513

514

515 **Statistical analysis and data visualization**

516 The co-localization of GWAS-supported variants with the retrogenome was calculated using linear
517 regressions in GARFIELD²⁷, and the gene-level and gene-set enrichment analyses were calculated
518 in MAGMA²⁸. Findings were corrected for multiple testing using the Bonferroni method. The case-
519 control HERV expression differences (N = 538 individuals) and effects of eQTLs on HERV
520 expression (N = 593 individuals) were calculated, respectively, using Wald tests in DESeq2, and
521 stepwise linear regressions in QTLtools, respectively. These were corrected using the false
522 discovery rate ($q < 0.05$), a more permissive multiple testing correction method, to increase our
523 detection power. The effect of genotype on specific HERVs within cases and control groups,
524 separately or combined, was calculated using linear regressions in IBM Statistics SPSS 25 (IBM
525 Corp., Armonk, NY, United States). Other analyses were performed in R⁶⁵. Graphs were generated
526 in R or Graph Pad Prism 7 (GraphPad Software, San Diego, CA, United States).

527

528 **Acknowledgements**

529 The work was supported in part by US National Institutes of Health grants: CA206488 (DFN),
530 AI076059 (DFN) and UL1TR001876 (KAC). This work was also supported by a grant from the
531 Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES, Brazil, Science without
532 Borders award no. BEX1279/13-0) and an NIHR Maudsley Biomedical Research Centre Career
533 Development Award to RRRD. TRP is funded by a Medical Research Council (MRC) Skills
534 Development Fellowship (MR/N014863/1). This study represents independent research part funded
535 by the NIHR-Wellcome Trust King's Clinical Research Facility and the National Institute for Health
536 Research (NIHR) Biomedical Research Centre at South London and Maudsley NHS Foundation
537 Trust and King's College London. The views expressed are those of the authors and not necessarily
538 those of the NHS, the NIHR or the Department of Health and Social Care. Data for this publication
539 were obtained from NIMH Repository & Genomics Resource, a centralized national biorepository for
540 genetic studies of psychiatric disorders. Data were generated as part of the CommonMind
541 Consortium supported by funding from Takeda Pharmaceuticals Company Limited, F. Hoffman-La
542 Roche Ltd and NIH grants R01MH085542, R01MH093725, P50MH066392, P50MH080405,
543 R01MH097276, RO1-MH-075916, P50M096891, P50MH084053S1, R37MH057881, AG02219,
544 AG05138, MH06692, R01MH110921, R01MH109677, R01MH109897, U01MH103392, and
545 contract HHSN271201300031C through IRP NIMH. Brain tissue for the study was obtained from the
546 Mount Sinai NIH Brain and Tissue Repository, the University of Pennsylvania Alzheimer's Disease
547 Core Center, the University of Pittsburgh NeuroBioBank and Brain and Tissue Repositories, and the
548 NIMH Human Brain Collection Core. CMC Leadership: Panos Roussos, Joseph Buxbaum, Andrew
549 Chess, Schahram Akbarian, Vahram Haroutunian (Icahn School of Medicine at Mount Sinai), Bernie
550 Devlin, David Lewis (University of Pittsburgh), Raquel Gur, Chang-Gyu Hahn (University of
551 Pennsylvania), Enrico Domenici (University of Trento), Mette A. Peters, Solveig Sieberts (Sage
552 Bionetworks), Thomas Lehner, Stefano Marengo, Barbara K. Lipska (NIMH). Tissue samples used
553 in RT-qPCR experiments were supplied by The London Neurodegenerative Diseases Brain Bank,
554 which receives funding from the Medical Research Council and as part of the Brains for Dementia
555 Research programme, jointly funded by Alzheimer's Research UK and Alzheimer's Society. We
556 thank Professor Cathryn Lewis for her comments on the manuscript and Daniel Bean (King's College

557 London) for his assistance with coding. We also thank Richard “Brad” Jones, Mario Ostrowski, and
558 Nathaniel Bachtel for their comments on this manuscript.

559

560 **Author contributions**

561 Study design: TRP, RRRD. Performed analyses and experiments: RRRD, TRP. Contributed
562 reagents, biological material, revised the manuscript: MLB, MM, CEO, GAB, SS, CT, GRT, KAC,
563 DPS, DFN. Wrote the paper: RRRD, TRP.

564

565 **Conflict of Interest**

566 The authors declare no conflict of interest.

567

568 **Data availability**

569 Telescope and the HERV annotation are available at <http://github.com/mlbendall/telescope>. Access
570 to the CommonMind Consortium dataset can be requested to the NIMH Repository & Genomics
571 Resource via <https://www.nimhgenetics.org/resources/commonmind>.

572

573 References

- 574 1. Gifford, R. & Tristem, M. The evolution, distribution and diversity of endogenous
575 retroviruses. *Virus Genes* **26**, 291-315 (2003).
- 576 2. Lander, E.S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**, 860-
577 921 (2001).
- 578 3. Mi, S. *et al.* Syncytin is a captive retroviral envelope protein involved in human placental
579 morphogenesis. *Nature* **403**, 785-9 (2000).
- 580 4. Grow, E.J. *et al.* Intrinsic retroviral reactivation in human preimplantation embryos and
581 pluripotent cells. *Nature* **522**, 221-225 (2015).
- 582 5. Brattas, P.L. *et al.* TRIM28 Controls a Gene Regulatory Network Based on Endogenous
583 Retroviruses in Human Neural Progenitor Cells. *Cell Rep* **18**, 1-11 (2017).
- 584 6. Fasching, L. *et al.* TRIM28 Represses Transcription of Endogenous Retroviruses in Neural
585 Progenitor Cells. *Cell reports* **10**, 20-28 (2015).
- 586 7. Li, W. *et al.* Human endogenous retrovirus-K contributes to motor neuron disease. *Sci*
587 *Transl Med* **7**, 307ra153 (2015).
- 588 8. Douville, R., Liu, J., Rothstein, J. & Nath, A. Identification of active loci of a human
589 endogenous retrovirus in neurons of patients with amyotrophic lateral sclerosis. *Annals of*
590 *neurology* **69**, 141-151 (2011).
- 591 9. Garson, J.A. *et al.* Quantitative analysis of human endogenous retrovirus-K transcripts in
592 postmortem premotor cortex fails to confirm elevated expression of HERV-K RNA in
593 amyotrophic lateral sclerosis. *Acta Neuropathol Commun* **7**, 45 (2019).
- 594 10. Perron, H. *et al.* Molecular characteristics of Human Endogenous Retrovirus type-W in
595 schizophrenia and bipolar disorder. *Transl Psychiatry* **2**, e201 (2012).
- 596 11. Weis, S. *et al.* Reduced expression of human endogenous retrovirus (HERV)-W GAG
597 protein in the cingulate gyrus and hippocampus in schizophrenia, bipolar disorder, and
598 depression. *J Neural Transm (Vienna)* **114**, 645-55 (2007).
- 599 12. Yolken, R.H., Karlsson, H., Yee, F., Johnston-Wilson, N.L. & Torrey, E.F. Endogenous
600 retroviruses and schizophrenia. *Brain Res Brain Res Rev* **31**, 193-9 (2000).
- 601 13. Roussos, P. *et al.* A role for noncoding variation in schizophrenia. *Cell reports* **9**, 1417-1429
602 (2014).
- 603 14. Pardiñas, A.F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant
604 genes and in regions under strong background selection. *Nature Genetics* **50**, 381-389
605 (2018).
- 606 15. Schizophrenia Working Group of the Psychiatric Genomics Consortium *et al.* Biological
607 insights from 108 schizophrenia-associated genetic loci. *Nature* **511**, 421 (2014).
- 608 16. Goke, J. & Ng, H.H. CTRL+INSERT: retrotransposons and their contribution to regulation
609 and innovation of the transcriptome. *EMBO Rep* **17**, 1131-44 (2016).
- 610 17. Karlsson, H. *et al.* Retroviral RNA identified in the cerebrospinal fluids and brains of
611 individuals with schizophrenia. *Proc Natl Acad Sci U S A* **98**, 4634-9 (2001).
- 612 18. Diem, O., Schäffner, M., Seifarth, W. & Leib-Mösch, C. Influence of Antipsychotic Drugs
613 on Human Endogenous Retrovirus (HERV) Transcription in Brain Cells. *PLoS ONE* **7**,
614 e30054 (2012).
- 615 19. Svensson, A.C. *et al.* Chromosomal distribution, localization and expression of the human
616 endogenous retrovirus ERV9. *Cytogenet Cell Genet* **92**, 89-96 (2001).
- 617 20. Bendall, M.L. *et al.* Telescope: Characterization of the retrotranscriptome by accurate
618 estimation of transposable element expression. *PLoS Computational Biology* **In press**,
619 398172 (2018).
- 620 21. Subramanian, R.P., Wildschutte, J.H., Russo, C. & Coffin, J.M. Identification,
621 characterization, and comparative genomic distribution of the HERV-K (HML-2) group of
622 human endogenous retroviruses. *Retrovirology* **8**, 90-90 (2011).

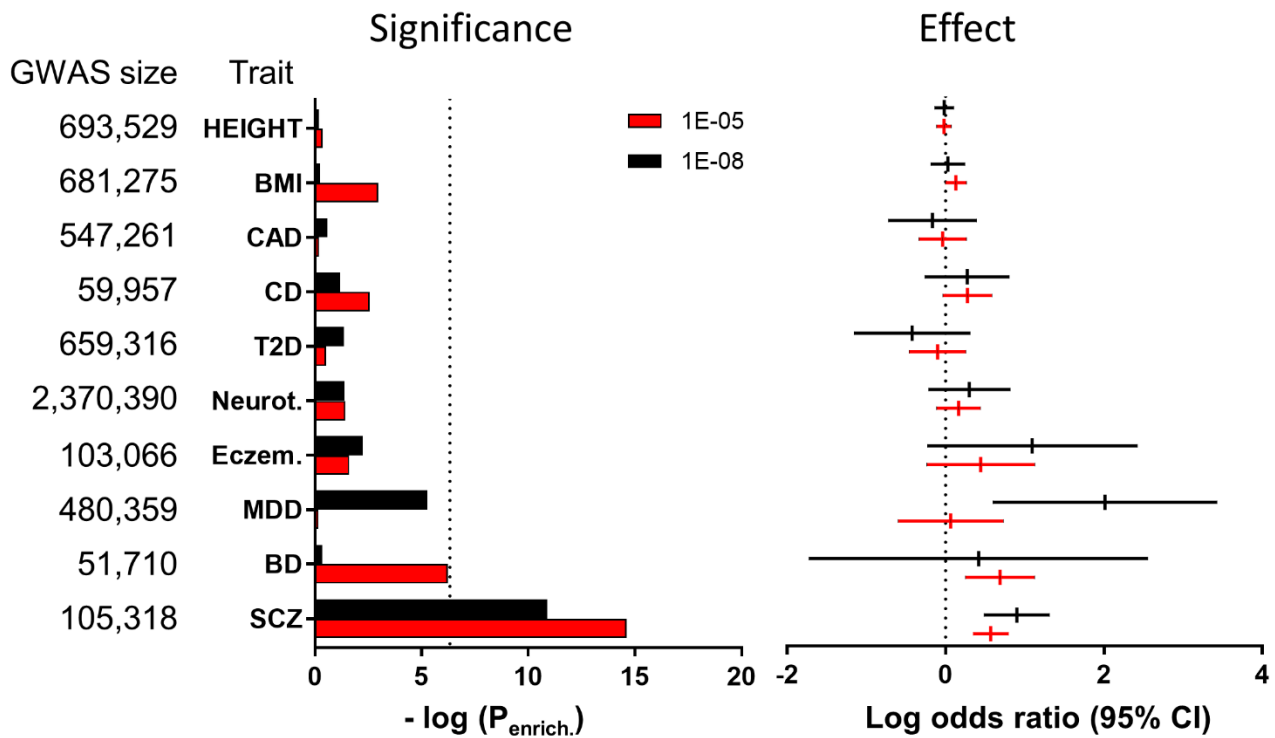
- 623 22. Grandi, N., Cadeddu, M., Blomberg, J. & Tramontano, E. Contribution of type W human
624 endogenous retroviruses to the human genome: characterization of HERV-W proviral
625 insertions and processed pseudogenes. *Retrovirology* **13**, 67-67 (2016).
- 626 23. Grandi, N. *et al.* Identification of a novel HERV-K(HML10): comprehensive
627 characterization and comparative analysis in non-human primates provide insights about
628 HML10 proviruses structure and diffusion. *Mobile DNA* **8**, 15 (2017).
- 629 24. Vargiu, L. *et al.* Classification and characterization of human endogenous retroviruses;
630 mosaic forms are common. *Retrovirology* **13**, 7 (2016).
- 631 25. Tokuyama, M. *et al.* ERVmap analysis reveals genome-wide transcription of human
632 endogenous retroviruses. *Proceedings of the National Academy of Sciences* **115**, 12565-
633 12572 (2018).
- 634 26. Jeong, H.-H., Yalamanchili, H.K., Guo, C., Shulman, J.M. & Liu, Z. An ultra-fast and
635 scalable quantification pipeline for transposable elements from next generation sequencing
636 data. in *Biocomputing 2018* 168-179.
- 637 27. Iotchkova, V. *et al.* GARFIELD classifies disease-relevant genomic features through
638 integration of functional annotations with association signals. *Nature Genetics* **51**, 343-353
639 (2019).
- 640 28. de Leeuw, C.A., Mooij, J.M., Heskes, T. & Posthuma, D. MAGMA: Generalized Gene-Set
641 Analysis of GWAS Data. *PLOS Computational Biology* **11**, e1004219 (2015).
- 642 29. Bao, W., Kojima, K.K. & Kohany, O. Repbase Update, a database of repetitive elements in
643 eukaryotic genomes. *Mobile DNA* **6**, 11 (2015).
- 644 30. Wang, D. *et al.* Comprehensive functional genomic resource and integrative model for the
645 human brain. *Science* **362**(2018).
- 646 31. Anderson, G.W. *et al.* Characterisation of neurons derived from a cortical human neural
647 stem cell line CTX0E16. *Stem Cell Res Ther* **6**, 149 (2015).
- 648 32. Deans, P.J.M. *et al.* Psychosis Risk Candidate ZNF804A Localizes to Synapses and
649 Regulates Neurite Formation and Dendritic Spine Structure. *Biol Psychiatry* **82**, 49-61
650 (2017).
- 651 33. Duarte, R.R.R. *et al.* The Psychiatric Risk Gene NT5C2 Regulates Adenosine
652 Monophosphate-Activated Protein Kinase Signaling and Protein Translation in Human
653 Neural Progenitor Cells. *Biol Psychiatry* **86**, 120-130 (2019).
- 654 34. Hinrichs, A.S. *et al.* The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res*
655 **34**, D590-8 (2006).
- 656 35. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network
657 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 658 36. van Dam, S., Vosa, U., van der Graaf, A., Franke, L. & de Magalhaes, J.P. Gene co-
659 expression analysis for functional classification and gene-disease predictions. *Brief*
660 *Bioinform* **19**, 575-592 (2018).
- 661 37. McKenzie, A.T. *et al.* Brain Cell Type Specific Gene Expression and Co-expression
662 Network Architectures. *Scientific Reports* **8**, 8868 (2018).
- 663 38. Klag, T. *et al.* Human Endogenous Retroviruses: Residues of Ancient Times Are
664 Differentially Expressed in Crohn's Disease. *Inflammatory Intestinal Diseases* **3**, 125-137
665 (2018).
- 666 39. Fromer, M. *et al.* Gene expression elucidates functional impact of polygenic risk for
667 schizophrenia. *Nature neuroscience* **19**, 1442-1453 (2016).
- 668 40. O'Brien, H.E. *et al.* Expression quantitative trait loci in the developing human brain and
669 their enrichment in neuropsychiatric disorders. *Genome biology* **19**, 194-194 (2018).
- 670 41. Bray, N.J. & O'Donovan, M.C. The genetics of neuropsychiatric disorders. *Brain and*
671 *Neuroscience Advances* **2**, 2398212818799271 (2018).
- 672 42. Duarte, R.R.R. *et al.* Genome-wide significant schizophrenia risk variation on chromosome
673 10q24 is associated with altered cis-regulation of BORCS7, AS3MT, and NT5C2 in the
674 human brain. *Am J Med Genet B Neuropsychiatr Genet* **171**, 806-14 (2016).

- 675 43. Liu, W. *et al.* Construction and Analysis of Gene Co-Expression Networks in *Escherichia*
676 *coli*. *Cells* **7**, 19 (2018).
- 677 44. van Dam, S., Vösa, U., van der Graaf, A., Franke, L. & de Magalhães, J.P. Gene co-
678 expression analysis for functional classification and gene–disease predictions. *Briefings in*
679 *Bioinformatics* **19**, 575-592 (2017).
- 680 45. Inkeles, M.S. *et al.* Cell-type deconvolution with immune pathways identifies gene networks
681 of host defense and immunopathology in leprosy. *JCI Insight* **1**(2016).
- 682 46. Tansey, K.E. & Hill, M.J. Enrichment of schizophrenia heritability in both neuronal and glia
683 cell regulatory elements. *Translational Psychiatry* **8**, 7-7 (2018).
- 684 47. Hill, M.J. *et al.* Knockdown of the schizophrenia susceptibility gene TCF4 alters gene
685 expression and proliferation of progenitor cells from the developing human neocortex.
686 *Journal of psychiatry & neuroscience : JPN* **42**, 181-188 (2017).
- 687 48. Hill, M.J., Jeffries, A.R., Dobson, R.J., Price, J. & Bray, N.J. Knockdown of the psychosis
688 susceptibility gene ZNF804A alters expression of genes involved in cell adhesion. *Hum Mol*
689 *Genet* **21**, 1018-24 (2012).
- 690 49. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass
691 index in approximately 700000 individuals of European ancestry. *Hum Mol Genet* **27**, 3641-
692 3649 (2018).
- 693 50. Harst, P.v.d. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an Expanded
694 View on the Genetic Architecture of Coronary Artery Disease. *Circulation Research* **122**,
695 433-443 (2018).
- 696 51. de Lange, K.M. *et al.* Genome-wide association study implicates immune activation of
697 multiple integrin genes in inflammatory bowel disease. *Nat Genet* **49**, 256-261 (2017).
- 698 52. Xue, A. *et al.* Genome-wide association analyses identify 143 risk variants and putative
699 regulatory mechanisms for type 2 diabetes. *Nature Communications* **9**, 2941 (2018).
- 700 53. Baselmans, B.M.L. *et al.* Multivariate genome-wide analyses of the well-being spectrum.
701 *Nature Genetics* **51**, 445-451 (2019).
- 702 54. The Early Genetics Lifecourse Epidemiology Eczema Consortium *et al.* Multi-ancestry
703 genome-wide association study of 21,000 cases and 95,000 controls identifies new risk loci
704 for atopic dermatitis. *Nature Genetics* **47**, 1449 (2015).
- 705 55. Wray, N.R. *et al.* Genome-wide association analyses identify 44 risk variants and refine the
706 genetic architecture of major depression. *Nature Genetics* **50**, 668-681 (2018).
- 707 56. Stahl, E.A. *et al.* Genome-wide association study identifies 30 Loci Associated with Bipolar
708 Disorder. *bioRxiv*, 173062 (2018).
- 709 57. Lambert, J.C. *et al.* Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci
710 for Alzheimer's disease. *Nat Genet* **45**, 1452-8 (2013).
- 711 58. Demontis, D. *et al.* Discovery of the first genome-wide significant risk loci for attention
712 deficit/hyperactivity disorder. *Nature Genetics* **51**, 63-75 (2019).
- 713 59. van Rheenen, W. *et al.* Genome-wide association analyses identify new risk variants and the
714 genetic architecture of amyotrophic lateral sclerosis. *Nat Genet* **48**, 1043-8 (2016).
- 715 60. Grove, J. *et al.* Identification of common genetic risk variants for autism spectrum disorder.
716 *Nature Genetics* **51**, 431-444 (2019).
- 717 61. Turner, S.D. qqman: an R package for visualizing GWAS results using Q-Q and manhattan
718 plots. *bioRxiv*, 005165 (2014).
- 719 62. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-
720 2079 (2009).
- 721 63. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
722 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 723 64. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nature Methods*
724 **9**, 357 (2012).
- 725 65. R Core Team. R: A language and environment for statistical computing. R Foundation for
726 Statistical Computing, Vienna, Austria. (2018).

- 727 66. Love, M.I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for
728 RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).
- 729 67. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer
730 datasets. *Gigascience* **4**, 7 (2015).
- 731 68. Jeffrey T. Leek *et al.* sva: Surrogate Variable Analysis. *R package Version 3.30.1*,
732 <https://bioconductor.org/packages/release/bioc/html/sva.html>(2019).
- 733 69. Delaneau, O. *et al.* A complete tool set for molecular QTL discovery and analysis. *Nature*
734 *Communications* **8**, 15452 (2017).
- 735 70. Ritchie, M.E. *et al.* limma powers differential expression analyses for RNA-sequencing and
736 microarray studies. *Nucleic Acids Research* **43**, e47-e47 (2015).
- 737 71. Wang, J., Vasaiakar, S., Shi, Z., Greer, M. & Zhang, B. WebGestalt 2017: a more
738 comprehensive, powerful, flexible and interactive gene set enrichment analysis toolkit.
739 *Nucleic acids research* **45**, W130-W137 (2017).
- 740

741 **Figures and Table**

Contribution of polymorphisms in HERVs to selected traits

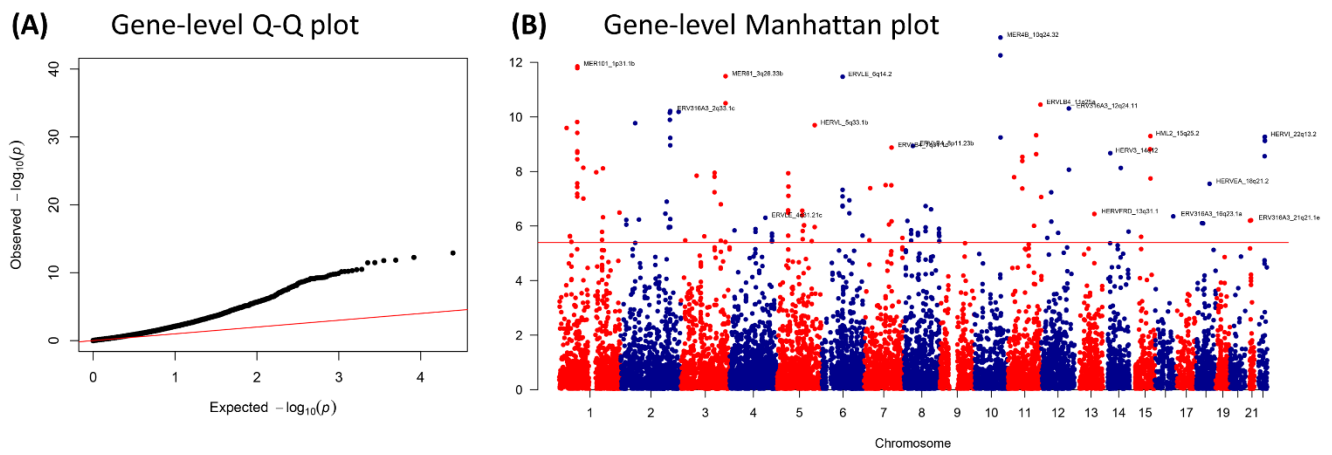


742

743 **Figure 1.** Variants within the retrogenome are enriched with schizophrenia-associated
744 polymorphisms. No association with any other trait was observed (see **Supplemental Table 1** for
745 full results). Calculated using GARFIELD²⁷.

746

747

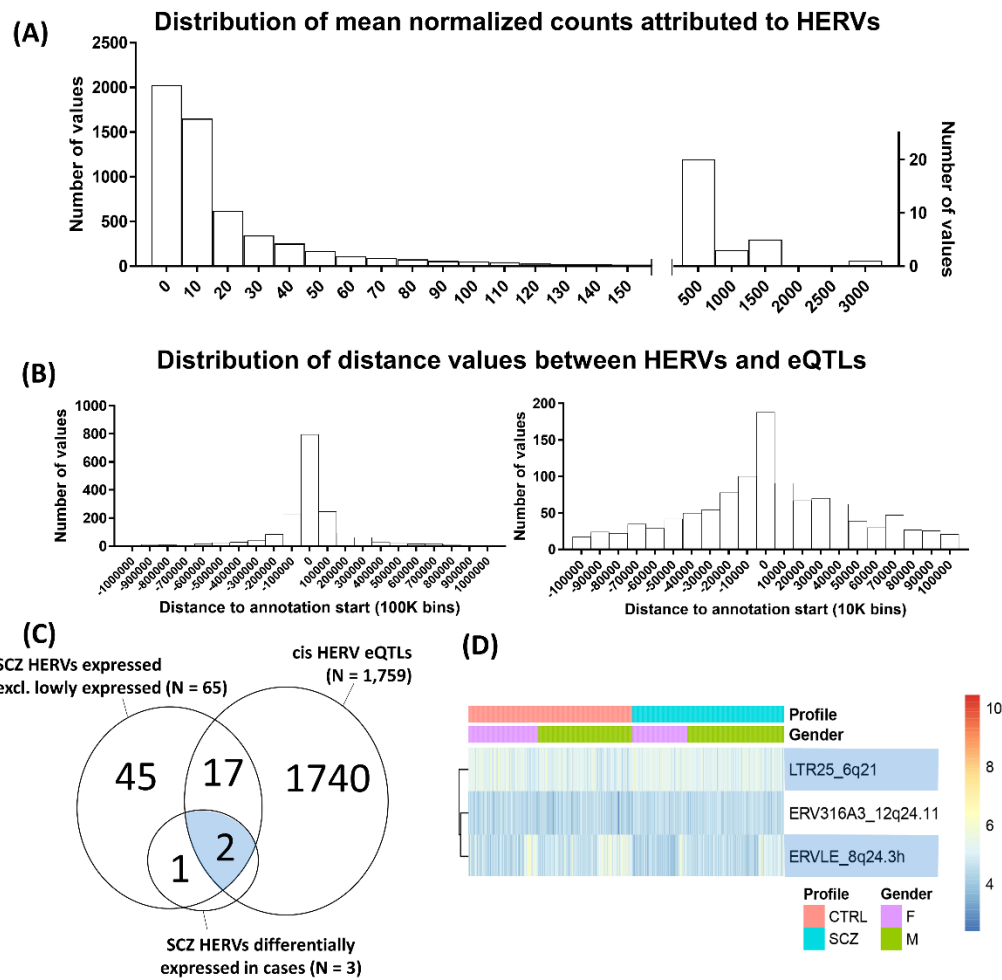


748

749 **Figure 2.** Gene-level enrichment analysis of the schizophrenia GWAS summary statistics using the
750 HERV annotation developed by Bendall and colleagues²⁰, calculated using MAGMA²⁸.
751 Chromosomes 1-22 only, extended MHC region excluded (chromosome 6, from 26-34 Mb). **(A)**
752 Quantile-quantile plot showing the contribution of several HERVs to schizophrenia genetics
753 compared to an expected normal distribution (red line). **(B)** Manhattan plot showing the location of
754 the HERVs associated with schizophrenia. Plots created using qqman⁶¹. All enriched HERVs are
755 shown on **Supplemental Table 2**.

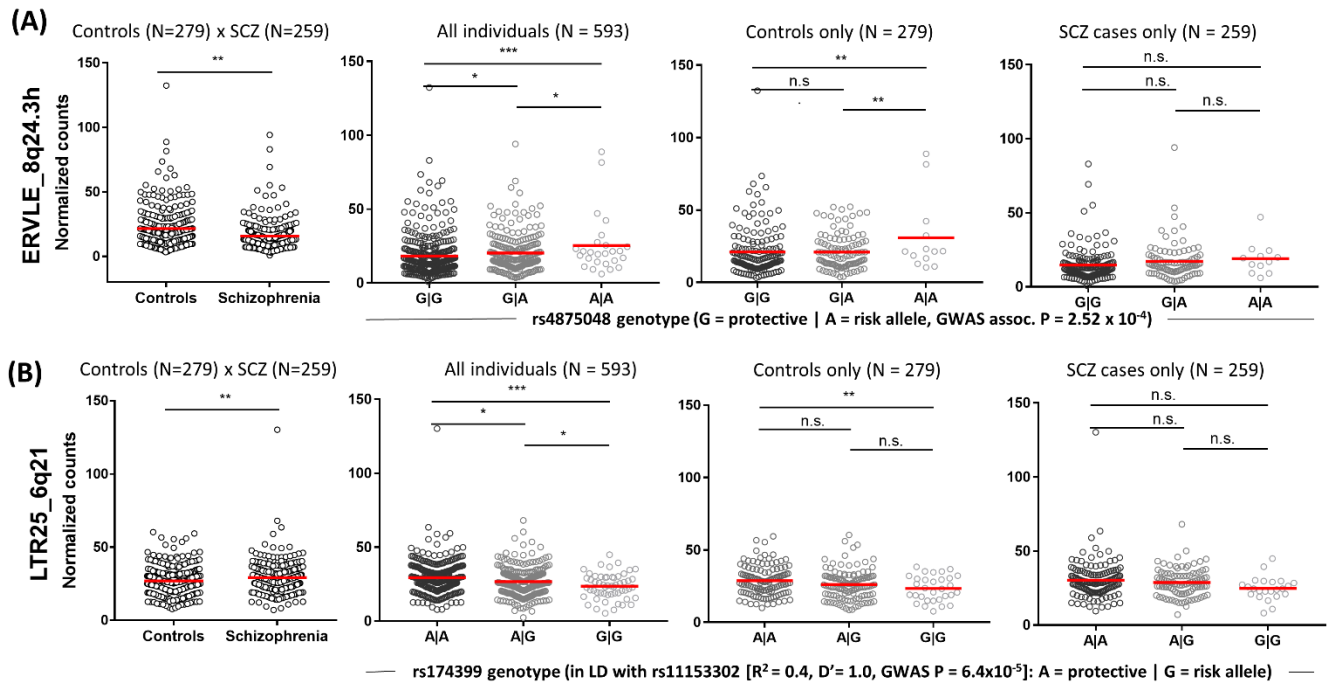
756

757



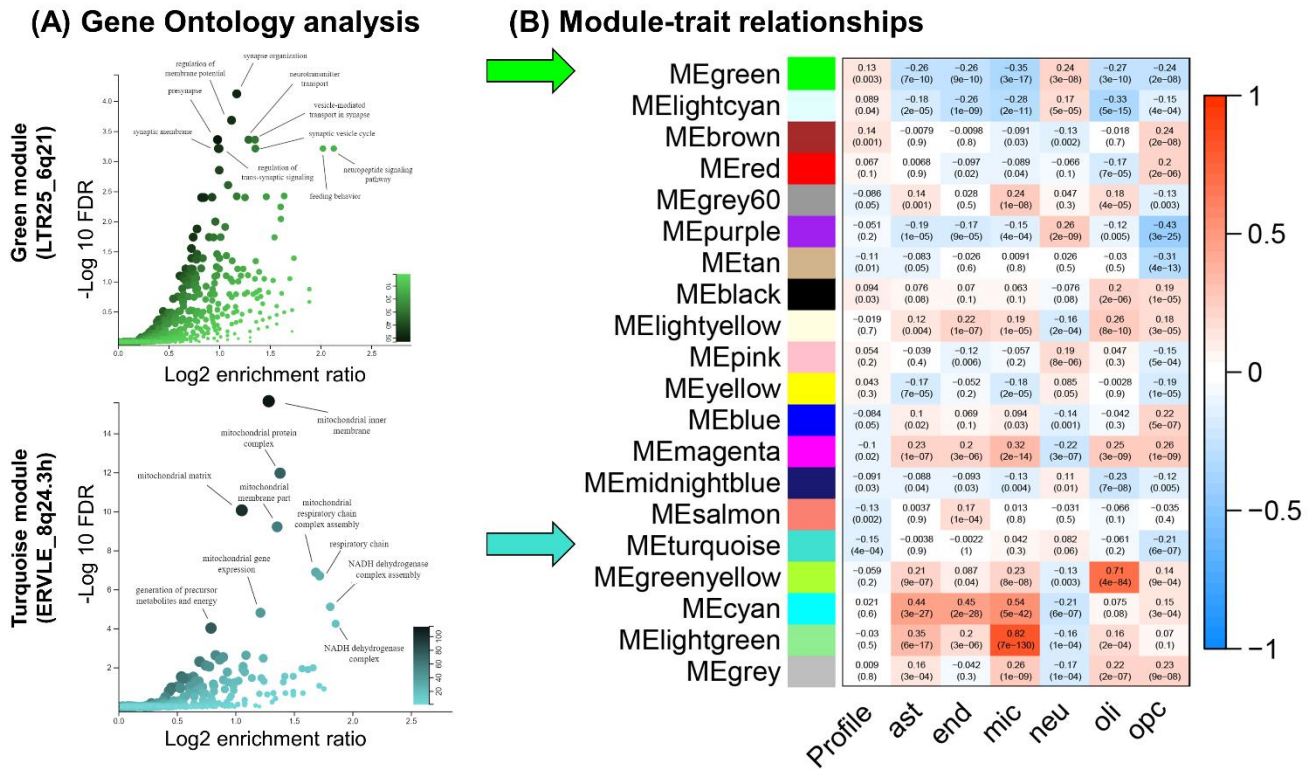
758

759 **Figure 3.** HERV expression in the dorsolateral prefrontal cortex, based on an analysis of the
 760 CommonMind Consortium dataset with Telescope. **(A)** Frequency distribution of the mean
 761 normalized counts per HERV across samples showing that the majority of HERVs are lowly
 762 expressed, according to an analysis of 593 post-mortem brains. Normalized counts do not include
 763 adjustments utilized in the analyses. **(B)** Distribution of values representing the distance between
 764 the eQTL for a HERV and the start site for that HERV. Left panel shows all data in bins of 100,000,
 765 and the right panel only data nearer the start site of the HERVs, in bins of 10,000. A large proportion
 766 of eQTLs was located within a 10kb window upstream or downstream the start coordinates of the
 767 HERV they regulate. **(C)** Overlap between the 65 schizophrenia-associated HERVs expressed in
 768 the brain, the 1,759 HERVs modulated by eQTLs, and the three HERVs enriched for schizophrenia
 769 variants and additionally differentially expressed between schizophrenia cases (N = 259) and
 770 unaffected individuals (N = 279). **(D)** Heatmap of the six HERVs enriched for schizophrenia variants
 771 and further differentially expressed in patients, separated by gender. In blue, HERVs that are
 772 modulated by eQTLs, as demonstrated in **Figure 4**.



773

774 **Figure 4.** HERVs enriched for schizophrenia variants and differentially expressed between cases
775 and controls, and according to genotype. eQTLs modulated HERV expression exclusively in control
776 individuals. Graphs show normalized counts associated with case-control status (left; Wald tests,
777 $**P < 0.01$) or per genotype considering all individuals, control individuals or schizophrenia patients
778 only (last three graphs, respectively; ANOVAs followed by pairwise comparisons corrected using the
779 Bonferroni method ($*P < 0.05$, $**P < 0.01$, $***P < 0.001$; n.s.: not significant; **Supplemental Table**
780 **6**). Data shown are for **(A)** ERVLE_8q24.3h and its top eQTL in the DLPFC, rs4875048, and **(B)**
781 LTR25_6q21 and its top eQTL, rs174399. Values are uncorrected for the factors and covariates
782 included in the eQTL analysis.

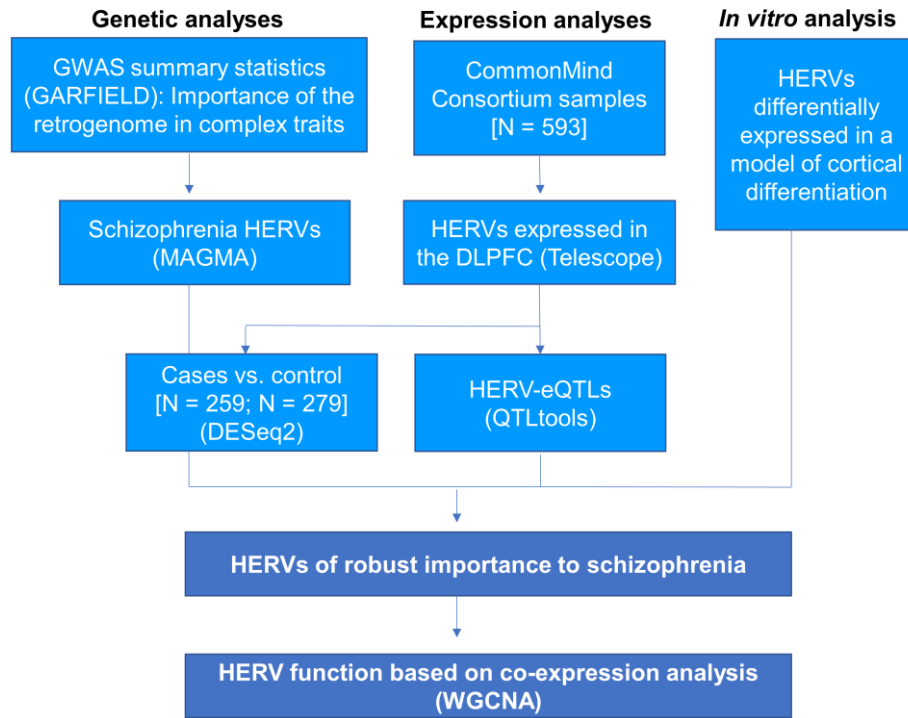


783

784 **Figure 5.** Gene ontology enrichment analysis of the co-expression modules associated with the
 785 green (LTR25_6q21) and turquoise modules (ERVLE_8q24.3h), and correlations between modules
 786 eigengenes, case-control status (Profile) and coefficients associated with cell counts for major neural
 787 cell types. **(A)** LTR25_6q21 belongs to the green module, which is significantly enriched for GO
 788 terms associated with neuronal function, as shown in the Volcano plot. ERVLE_8q24.3h, in turn,
 789 belongs to the turquoise module, which is significantly enriched for GO terms associated with
 790 mitochondrial function. **(B)** The green module is positively associated with neuronal counts, whereas
 791 the turquoise module does not correlate strongly with a specific cell type, apart from a negative
 792 correlation with oligodendrocyte progenitor cell counts. Each cell of the heatmap contains the
 793 Pearson's r coefficient followed by significance of the correlation (P). Ast: astrocytes, end: endothelial
 794 cells, mic: microglia, neu: neurons, oli: oligodendrocytes, opc: oligodendrocyte progenitor cells.

795

796



797

798 **Figure 6.** Analysis strategy. We performed a series of analyses using post-mortem brain RNA-
799 sequencing data, genetic enrichment analyses and *in vitro* cortical differentiation data to identify
800 HERVs of robust importance in schizophrenia. Bioinformatic tools used are indicated in
801 parentheses. DLPFC: dorsolateral prefrontal cortex.

Table 1. Association of the eQTLs (and variants in linkage disequilibrium) with schizophrenia.

HERV	SNP	Chromosome	Position	Risk allele	Other allele	Frq cases	OR	SE	P	R ² ‡	D' ‡
ERVLE_8q24.3h	rs4875048	8	144826671	A	G	0.1869	1.0553	0.011855	5.6E-06	-	-
	rs10552126	8	144844056	CAT	C	0.2018	1.05919	0.0116	7.1E-07	0.8954	1
LTR25_6q21	rs174399	6	111919619	G*	A	0.2565	1.0025	0.010488	0.81	-	-
	rs11153302	6	111918869	A	G	0.506	1.0397	0.00973	6.4E-05	0.38	1

* This SNP is not associated with schizophrenia, but is in linkage disequilibrium with another variant that is associated with this disorder.

‡ Linkage disequilibrium statistics in relation to the top association signal at the locus which is in LD with the HERV-eQTL.