

## Elevated L1 expression in ataxia telangiectasia likely explained by an RNA-seq batch effect

Geoffrey J. Faulkner<sup>1,2,\*</sup>

<sup>1</sup>Mater Research Institute - University of Queensland, Woolloongabba, QLD, 4102, Australia.

<sup>2</sup>Queensland Brain Institute, University of Queensland, St. Lucia, QLD, 4067, Australia.

\*Correspondence: [faulknergj@gmail.com](mailto:faulknergj@gmail.com)

### 1 **Abstract**

2 A recent study (Takahashi et al., *Neuron*, 2022) concluded LINE-1 (L1)  
3 retrotransposon activation drives cerebellar ataxia and neurodegeneration. This  
4 position was based on L1 upregulation in ataxia telangiectasia (AT) patient  
5 cerebellum samples, as measured by RNA-seq, and observation of ataxia and  
6 neurodegeneration in mice where cerebellar L1 expression was induced via dCas9-  
7 CRISPR. Here, a re-analysis of the RNA-seq data, which were obtained by rRNA  
8 depletion rather than polyA+ selection, revealed a high fraction (38.4%) of intronic  
9 reads. Significantly ( $p=0.034$ ) more intronic reads were present in the AT data than  
10 the matched controls. This finding provides an alternative and robust explanation for  
11 a key result reported by Takahashi et al.: intronic L1 sequences are abundant in pre-  
12 mRNAs, and more pre-mRNAs were retained in the AT libraries. This apparent batch  
13 effect deserves further examination, as claims of L1-mediated pathogenesis could  
14 shape future efforts to treat AT by trying to attenuate L1 activity.

15

### 16 **Main**

17 The retrotransposon LINE-1 (L1) is an autonomous mobile genetic element whose  
18 ~500,000 copies compose more than 17% of the human genome. A full-length,  
19 retrotransposition-competent L1 is 6kb long and incorporates a 5' internal sense  
20 promoter that drives transcription of a bicistronic mRNA, which in turn encodes two  
21 proteins (ORF1p and ORF2p). These L1 proteins strongly prefer to mobilize their  
22 encoding mRNA, and there are fewer than 100 retrotransposition-competent human-  
23 specific L1 (L1HS) copies per individual. The vast majority of L1s scattered

24 throughout the genome thus do not have intact ORFs and are immobile (Kazazian  
25 and Moran, 2017), even if they reside in the introns of protein-coding genes.

26 Ataxia telangiectasia (AT) is a severe neurodegenerative disorder caused by defects  
27 in the AT mutated (ATM) DNA damage repair gene. Prior works have found ATM  
28 either facilitates L1 mobilization (Gasior et al., 2006; Wallace et al., 2013) or, if  
29 mutated, changes the character and, moderately, increases the frequency of L1  
30 insertions found in neuronal cells (Coufal et al., 2011). ATM is not however thought  
31 to repress L1 promoter activity or limit L1 protein expression (Coufal et al., 2011).

32 In recent work, Takahashi et al. concluded that L1 can drive neurodegeneration in  
33 AT, based on increased cerebellar L1 RNA abundance in AT patients, as well as  
34 ataxia and neurodegeneration in mice where L1 transcription was activated in the  
35 cerebellum via a dCas9-CRISPR approach (Takahashi et al., 2022). This intriguing  
36 study could shape strategies to treat or prevent AT in children and, for this reason,  
37 its findings warrant careful consideration.

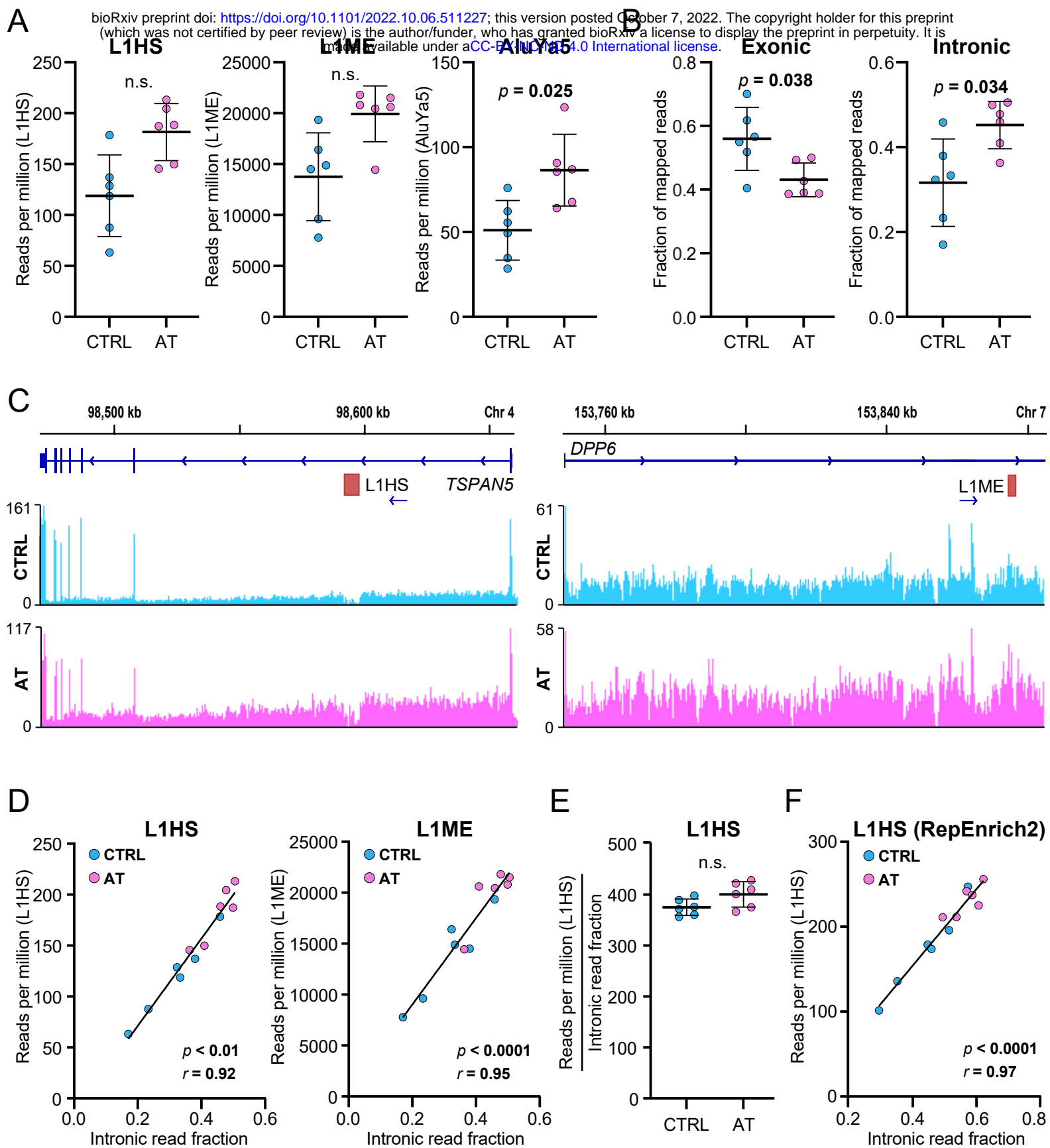
38 Takahashi et al. used bulk RNA-seq, where libraries were prepared via rRNA  
39 depletion, to compare the cerebellar transcriptomes of AT patients ( $n=6$ ) and control  
40 individuals ( $n=6$ ). They reported significant upregulation of the L1HS subfamily in AT.  
41 Accurate quantification of L1HS transcription, and that of other recently emerged, or  
42 young, L1 subfamilies with RNA-seq presents many challenges, principally due to  
43 their copy number and limited sequence divergence (Faulkner et al., 2009; Iñiguez et  
44 al., 2019; Jin et al., 2015; Lanciano and Cristofari, 2020). When evaluating  
45 differential expression amongst numerous L1 subfamilies, Takahashi et al. did not  
46 incorporate a multiple testing correction despite performing t tests on several  
47 retrotransposon families simultaneously, which would have rendered the noted L1HS  
48 upregulation in AT patients ( $p=0.042$ , two-sided t test) no longer significant. Multiple  
49 testing correction and the incorporation of reads that align to multiple genomic  
50 locations, or “multi-mapping” reads are essential to the robust and reproducible  
51 RNA-seq quantification of young retrotransposon subfamilies, such as L1HS (Jin et  
52 al., 2015; Lanciano and Cristofari, 2020).

53 Re-analyzing the RNA-seq data with TETranscripts (Jin et al., 2015) and the same  
54 STAR (Dobin et al., 2013) alignment parameters used by Takahashi et al., a

55 statistically significant difference was not apparent for L1HS in AT patients compared  
56 to controls, and incorporating multi-mapping reads did not change this outcome  
57 (**Figure S1A**). Older L1 subfamilies such as L1ME, which are immobile and tend to  
58 be 5' truncated and lack a canonical promoter and intact ORFs, accounted for far  
59 more RNA-seq reads than L1HS (**Figure S1A**). These older L1s were also typically  
60 more highly expressed in the AT samples than in controls, as were young PolIII-  
61 transcribed *Alu* retrotransposon families, such as AluYa5 (**Figure S1A**). Visual  
62 inspection of a selection of the most highly expressed L1 loci, such as an L1HS  
63 intronic to the *TSPAN5* gene (**Figure S1C**) and an L1ME intronic to the *DPP6* gene  
64 (**Figure S1C**), indicated they did not possess intact ORFs and were typically flanked  
65 by pronounced and pervasive intronic RNA-seq signals (**Figure S1C**).

66 RNA-seq libraries prepared via rRNA depletion, as opposed to polyA+ selection, can  
67 contain considerably more nascent, unspliced pre-mRNAs (Zhao et al., 2018). While  
68 the average percentage of intragenic reads did not appreciably differ in control  
69 (87.6%) and AT samples (88.3%), the average percentage of intronic reads was both  
70 unusually high overall (38.4%) and significantly ( $p=0.034$ , two-tailed t test with  
71 Bonferroni correction) higher in AT samples (45.2%) than in control samples (31.6%)  
72 (**Figure S1B**). Normalized L1HS, L1ME and AluYa5 read counts were each strongly  
73 correlated (Pearson  $r > 0.92$ , two-tailed  $p < 0.01$ ) with intronic read fraction in both AT  
74 and control samples (**Figure S1D**). Further normalization by dividing the L1HS read  
75 count by the fraction of intronic reads detected in each sample brought L1HS  
76 expression to virtual parity in AT and control samples (**Figure S1E**). The expression  
77 of L1HS repressors downregulated in AT, as noted by Takahashi et al., was also  
78 strongly anticorrelated with intronic read fraction, such as for *TRIM28* (Pearson  $r = -$   
79  $0.90$ , two-tailed  $p < 0.0001$ ), which has a relatively high (1:1) exon:intron sequence  
80 ratio and, notably, regulates older L1 subfamilies and not L1HS transcription (Castro-  
81 Diaz et al., 2014). Despite the availability of robust human L1 ORF1p antibodies, an  
82 immunoblot was not shown to support differential L1 protein expression in AT patient  
83 samples.

84 To corroborate the RNA-seq results, qPCR was performed. However, this approach  
85 depends on high quality input RNA and, as the target L1HS RNA is not spliced, is  
86 prone to gDNA contamination. These factors, and off-target effects, can influence L1



**Figure S1: Intronic enrichment from retained pre-mRNAs in AT patient RNA-seq data explains apparent retrotransposon upregulation.** (A) Normalized RNA-seq read counts in retrotransposon subfamilies obtained from Tetranscripts and allowing multi-mapping reads. Statistical significance was calculated using DESeq2 and BH multiple testing correction. (B) Fractions of exonic and intronic reads per sample, using only uniquely mapping reads, determined with featureCounts. Significance testing was via two-tailed t test with Bonferroni correction. (C) Integrative Genomics Viewer (IGV) view of aggregate RNA-seq coverage plots at the *TSPAN5* locus (left) and *DPP6* locus (right). Each L1 has broken ORFs in the annotated reference genome sequence and is oriented on the same strand as the host gene. (D) Correlation of L1 subfamily normalized read count, obtained with Tetranscripts, and intronic read fraction per sample, with a linear regression line and Pearson correlation ( $r$  and two-tailed  $p$ ) shown. (E) Normalized L1HS read count divided by intronic read fraction per sample. (F) As per (D), except showing normalized L1HS read counts obtained using RepEnrich2 and intronic read fractions determined using the bowtie2 alignment input files for RepEnrich2 via rnaseqc. Note: in (A), (B) and (E), individual data points, each representing a sample, are marked, and also represented as mean  $\pm$  S.D.

87 qPCR assays, particularly when primers are designed against the L1HS ORFs, as  
88 done here, and not the 5'UTR, which is retained only by full-length L1s. Finally, to  
89 demonstrate the robustness of the present claims, normalized L1HS read counts  
90 were obtained by aligning the RNA-seq data with bowtie2 (Langmead and Salzberg,  
91 2012) and then analyzed with RepEnrich2 (Criscione et al., 2014), as per the  
92 approach of Takahashi et al., confirming a strong correlation (Pearson  $r=0.97$ , two-  
93 tailed  $p<0.0001$ ) with intronic read fraction across AT and control samples (**Figure**  
94 **S1F**). This analysis confirmed the strong correlation of L1HS and intronic read count  
95 remained, regardless of the alignment algorithm, TE transcript analysis software, or  
96 how exonic/intronic regions were defined (see **Methods**).

97 Taken together, these results point to a technical explanation (Zhao et al., 2018) for  
98 increased L1 transcript abundance in AT patients: pre-mRNAs, whose introns are  
99 substantially composed of L1 and other retrotransposons, were more abundant in  
100 the AT sample RNA-seq libraries than in the matched control samples. The most  
101 likely explanation for this difference is a batch effect where the RNA quality was  
102 lower for the AT samples than for the control samples, or perhaps an effect arising  
103 during the rRNA depletion-based library processing. Ideally, the authors could  
104 provide further information distinguishing these possibilities. A much less plausible  
105 explanation is that pre-mRNAs are for unknown reasons more abundant in AT  
106 cerebellum than in controls, and even this scenario does not yield support for  
107 elevated protein-coding L1 transcription in AT.

108 A finding of L1 not being upregulated in AT would be concordant with earlier data  
109 suggesting the L1HS promoter is not more active, and L1 proteins not more  
110 abundant, in ATM-deficient human neuronal precursor cells (Coufal et al., 2011). It  
111 would subtract the human disease rationale for the remainder of the Takahashi et al.  
112 study, contradict the model of AT pathology being potentially driven by elevated and  
113 ectopic L1-mediated reverse transcription, and, importantly, argue against the  
114 therapeutic potential of L1 reverse transcriptase inhibitors in AT.

115 Even if the available evidence of L1 dysregulation in human AT is discounted, the  
116 animal experiments conducted by Takahashi et al. provide a remarkable proof-of-  
117 principle of L1 activation *in vivo* via the dCas9-CRISPR system. Previously, AT  
118 mouse models have not recapitulated the neurodegeneration seen in AT patients

119 (Lavin, 2013) or the L1 activator dCas9-CRISPR mice generated by Takahashi et al.  
120 As a caveat they note, these *in vivo* data were generated using one L1 sgRNA  
121 matching thousands of genomic loci, elevating the probability of an off-target effect.  
122 This issue is potentially more acute given inadvertent basal dCas9 expression in the  
123 L1 activation system, and compounded by the lack of cerebellar RNA-seq data from  
124 the dCas9-CRISPR mice, as L1 sequences can act as promoters for genes  
125 expressed during development and in neurons (Gerdes et al., 2022; Jönsson et al.,  
126 2019). The loci containing the integrated L1 targeting and scrambled sgRNA  
127 transgenes were not identified, and hence the effects of their integration, and their  
128 genomic context, in combination with the genetic background of the dCas9-CRISPR  
129 line, on phenotype was not indicated. While treatment with the L1 reverse  
130 transcriptase inhibitor lamivudine appeared to attenuate disease progression in  
131 dCas9-CRISPR animals, a relatively low concentration of lamivudine, which crosses  
132 the blood-brain barrier very poorly (Osborne et al., 2020), was used. The cerebellar  
133 concentration of lamivudine, and whether it was sufficient to alter L1 reverse  
134 transcriptase activity, was therefore not determined.

135 In light of these open questions, more data is needed to ascertain whether L1  
136 expression is dysregulated in AT. Independent replication experiments using RNA-  
137 seq protocols that avoid pre-mRNAs, which neither rRNA depletion-based or single-  
138 nucleus RNA-seq (snRNA-seq) approaches easily achieve, could be useful in this  
139 regard. These data would help clarify the causal role, if any, played by L1 in human  
140 neurodegeneration, and indicate its value as a therapeutic target.

141

## 142 **Methods**

143 Human AT patient RNA-seq fastq files (accession GSE175776) were downloaded  
144 from the Gene Expression Omnibus and aligned to the hg38 reference genome  
145 using STAR (Dobin et al., 2013). Reads were aligned in two ways. Firstly, to retain  
146 only uniquely aligned reads, the same STAR alignment parameters as Takahashi et  
147 al. were used. Secondly, to incorporate multi-mapping reads, the STAR alignment  
148 parameter `outFilterMultimapNmax` was changed from 1 to 10000. Each set of  
149 alignments were then processed with TETranscripts (Jin et al., 2015) with default  
150 parameters, using “`uniq`” for the first set of alignments and “`multi`” for the second set



151 of alignments. RefSeq gene and retrotransposon annotation files were obtained from  
152 the UCSC Table Browser and the Tetrascripts GitHub, respectively. Differential  
153 gene and retrotransposon expression analyses were performed with DESeq2 (Love  
154 et al., 2014) and default (BH) multiple testing correction, as implemented as part of  
155 Tetrascripts. Exonic and intronic read fractions were obtained using the RefSeq  
156 gene annotation file and featureCounts (Liao et al., 2014) with parameters -p -T 8 -B  
157 -O. Tetrascript quantification of retrotransposon families included multi-mapping  
158 reads; all other analyses used uniquely mapping reads only. For the analysis shown  
159 in Figure S1F, normalized L1HS read counts were obtained using RepEnrich2  
160 (Criscione et al., 2014) and bowtie2 (Langmead and Salzberg, 2012), whilst intronic  
161 read fractions were determined using the bowtie2 alignment output files and rnaseqc  
162 (DeLuca et al., 2012) and Gencode transcript annotations.

163

## 164 **References**

- 165 Castro-Diaz, N., Ecco, G., Coluccio, A., Kapopoulou, A., Yazdanpanah, B., Friedli,  
166 M., Duc, J., Jang, S.M., Turelli, P., and Trono, D. (2014). Evolutionally dynamic L1  
167 regulation in embryonic stem cells. *Genes Dev.* *28*, 1397–1409.
- 168 Coufal, N.G., Garcia-Perez, J.L., Peng, G.E., Marchetto, M.C.N., Muotri, A.R., Mu,  
169 Y., Carson, C.T., Macia, A., Moran, J.V., and Gage, F.H. (2011). Ataxia  
170 telangiectasia mutated (ATM) modulates long interspersed element-1 (L1)  
171 retrotransposition in human neural stem cells. *Proc. Natl. Acad. Sci. U. S. A.* *108*,  
172 20382–20387.
- 173 Criscione, S.W., Zhang, Y., Thompson, W., Sedivy, J.M., and Neretti, N. (2014).  
174 Transcriptional landscape of repetitive elements in normal and cancer human cells.  
175 *BMC Genomics* *15*, 583.
- 176 DeLuca, D.S., Levin, J.Z., Sivachenko, A., Fennell, T., Nazaire, M.-D., Williams, C.,  
177 Reich, M., Winckler, W., and Getz, G. (2012). RNA-SeQC: RNA-seq metrics for  
178 quality control and process optimization. *Bioinformatics* *28*, 1530–1532.
- 179 Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P.,  
180 Chaisson, M., and Gingeras, T.R. (2013). STAR: ultrafast universal RNA-seq aligner.  
181 *Bioinformatics* *29*, 15–21.
- 182 Faulkner, G.J., Kimura, Y., Daub, C.O., Wani, S., Plessy, C., Irvine, K.M., Schroder,  
183 K., Cloonan, N., Steptoe, A.L., Lassmann, T., et al. (2009). The regulated  
184 retrotransposon transcriptome of mammalian cells. *Nat. Genet.* *41*, 563–571.
- 185 Gasiior, S.L., Wakeman, T.P., Xu, B., and Deininger, P.L. (2006). The human LINE-1  
186 retrotransposon creates DNA double-strand breaks. *J. Mol. Biol.* *357*, 1383–1393.
- 187 Gerdes, P., Lim, S.M., Ewing, A.D., Larcombe, M.R., Chan, D., Sanchez-Luque, F.J.,

- 188 Walker, L., James, C., Knaupp, A.S., Carreira, P.E., et al. (2022). Retrotransposon  
189 instability dominates the acquired mutation landscape of mouse induced pluripotent  
190 stem cells. *bioRxiv* <https://doi.org/10.1101/2022.02.16.480772>.
- 191 Iñiguez, L.P., de Mulder Rougvié, M., Stearrett, N., Jones, R.B., Ormsby, C.E.,  
192 Reyes-Terán, G., Crandall, K.A., Nixon, D.F., and Bendall, M.L. (2019).  
193 Transcriptomic analysis of human endogenous retroviruses in systemic lupus  
194 erythematosus. *Proc. Natl. Acad. Sci. U. S. A.* *116*, 21350–21351.
- 195 Jin, Y., Tam, O.H., Paniagua, E., and Hammell, M. (2015). TETranscripts: a package  
196 for including transposable elements in differential expression analysis of RNA-seq  
197 datasets. *Bioinformatics* *31*, 3593–3599.
- 198 Jönsson, M.E., Ludvik Brattås, P., Gustafsson, C., Petri, R., Yudovich, D., Piracs, K.,  
199 Verschuere, S., Madsen, S., Hansson, J., Larsson, J., et al. (2019). Activation of  
200 neuronal genes via LINE-1 elements upon global DNA demethylation in human  
201 neural progenitors. *Nat. Commun.* *10*, 3182.
- 202 Kazazian, H.H., Jr, and Moran, J.V. (2017). Mobile DNA in Health and Disease. *N.*  
203 *Engl. J. Med.* *377*, 361–370.
- 204 Lanciano, S., and Cristofari, G. (2020). Measuring and interpreting transposable  
205 element expression. *Nat. Rev. Genet.* *21*, 721–736.
- 206 Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie  
207 2. *Nat. Methods* *9*, 357–359.
- 208 Lavin, M.F. (2013). The appropriateness of the mouse model for ataxia-  
209 telangiectasia: neurological defects but no neurodegeneration. *DNA Repair* *12*, 612–  
210 619.
- 211 Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general  
212 purpose program for assigning sequence reads to genomic features. *Bioinformatics*  
213 *30*, 923–930.
- 214 Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change  
215 and dispersion for RNA-seq data with DESeq2. *Genome Biol.* *15*, 550.
- 216 Osborne, O., Peyravian, N., Nair, M., Daunert, S., and Toborek, M. (2020). The  
217 Paradox of HIV Blood–Brain Barrier Penetrance and Antiretroviral Drug Delivery  
218 Deficiencies. *Trends Neurosci.* *43*, 695–708.
- 219 Takahashi, T., Stoiljkovic, M., Song, E., Gao, X.-B., Yasumoto, Y., Kudo, E.,  
220 Carvalho, F., Kong, Y., Park, A., Shanabrough, M., et al. (2022). LINE-1 activation in  
221 the cerebellum drives ataxia. *Neuron* <https://doi.org/10.1016/j.neuron.2022.08.011>.
- 222 Wallace, N.A., Gasior, S.L., Faber, Z.J., Howie, H.L., Deininger, P.L., and Galloway,  
223 D.A. (2013). HPV 5 and 8 E6 expression reduces ATM protein levels and attenuates  
224 LINE-1 retrotransposition. *Virology* *443*, 69–79.
- 225 Zhao, S., Zhang, Y., Gamini, R., Zhang, B., and von Schack, D. (2018). Evaluation of  
226 two main RNA-seq approaches for gene quantification in clinical RNA sequencing:



227 polyA+ selection versus rRNA depletion. *Sci. Rep.* 8, 4781.

228

## 229 **Figure Legends**

230 **Figure S1: Intronic enrichment from retained pre-mRNAs in AT patient RNA-**

231 **seq data explains apparent retrotransposon upregulation. (A)** Normalized RNA-

232 seq read counts in retrotransposon subfamilies obtained from TETranscripts and

233 allowing multi-mapping reads. Statistical significance was calculated using DESeq2

234 and BH multiple testing correction. **(B)** Fractions of exonic and intronic reads per

235 sample, using only uniquely mapping reads, determined with featureCounts.

236 Significance testing was via two-tailed t test with Bonferroni correction. **(C)**

237 Integrative Genomics Viewer (IGV) view of aggregate RNA-seq coverage plots at the

238 *TSPAN5* locus (left) and *DPP6* locus (right). Each L1 has broken ORFs in the

239 annotated reference genome sequence and is oriented on the same strand as the

240 host gene. **(D)** Correlation of L1 subfamily normalized read count, obtained with

241 TETranscripts, and intronic read fraction per sample, with a linear regression line and

242 Pearson correlation ( $r$  and two-tailed  $p$ ) shown. **(E)** Normalized L1HS read count

243 divided by intronic read fraction per sample. **(F)** As per (D), except showing

244 normalized L1HS read counts obtained using RepEnrich2 and intronic read fractions

245 determined using the bowtie2 alignment input files for RepEnrich2 via rnaseqc. Note:

246 in (A), (B) and (E), individual data points, each representing a sample, are marked,

247 and also represented as mean  $\pm$  S.D.

248

## 249 **Acknowledgements**

250 The author thanks Dr Sandra Richardson, Dr Adam Ewing, and members of the

251 Faulkner laboratory for helpful discussions. The author receives funding from the

252 Australian NHMRC (GNT1173711) and ARC (DP200102919), Cancer Australia

253 (2003170), and the Mater Foundation.

254

## 255 **Declaration of Interests**

256 The author declares no competing interests.