1 **Transposable elements and their KZFP controllers are drivers of transcriptional innovation in the**

2 **developing human brain**

3 Christopher J. Playfoot[1], Julien Duc[1], Shaoline Sheppard[1], Sagane Dind[1], Alexandre Coudray[1], Evarist

4 Planet[1] and Didier Trono[1, 2]

5 [1]School of Life Sciences, Ecole Polytechnique Fédérale de Lausanne (EPFL), Lausanne, Switzerland

6 [2]Corresponding author. Didier.Trono@epfl.ch

7 **Abstract:**

8 Transposable elements (TEs) constitute 50% of the human genome and many have been co-opted

9 throughout human evolution due to gain of advantageous regulatory functions controlling gene

10 expression networks. Several lines of evidence suggest these networks can be fine-tuned by the largest

11 family of TE controllers, the KRAB-containing zinc finger proteins (KZFPs). One tissue permissive for TE

12 transcriptional activation (termed 'transposcription') is the adult human brain, however

13 comprehensive studies on the extent of this process and its potential contribution to human brain

14 development are lacking.

15 In order to elucidate the spatiotemporal transposcriptome of the developing human brain, we have

16 analysed two independent RNA-seq datasets encompassing 16 distinct brain regions from eight weeks

17 post-conception into adulthood. We reveal an anti-correlated, KZFP:TE transcriptional profile defining

18 the late prenatal to early postnatal transition, and the spatiotemporal and cell type specific activation

19 of TE-derived alternative promoters driving the expression of neurogenesis-associated genes. We also

20 demonstrate experimentally that a co-opted antisense L2 element drives temporal protein re-

21 localisation away from the endoplasmic reticulum, suggestive of novel TE dependent protein function

22 in primate evolution. This work highlights the widespread dynamic nature of the spatiotemporal

23 KZFP:TE transcriptome and its potential importance throughout neurotypical human brain

24 development.

**Introduction**

KZFPs constitute the largest family of transcription factors encoded by mammalian genomes. These proteins harbor an N-terminal Krüppel-associated box (KRAB) domain and a C-terminal zinc finger array, which, for many, mediates sequence-specific DNA recognition. The KRAB domain of a majority of KZFPs recruits the transcriptional co-repressor KAP1 (KRAB-associated protein 1, also known as Tripartite motif protein 28, TRIM28), which acts as a scaffold for heterochromatin inducers such as the histone methyl-transferase SETDB1, the histone deacetylating NuRD complex, heterochromatin protein 1 (HP1) and DNA methyltransferases (Ecco et al. 2017). Many KZFPs bind to and repress TEs, a finding that led to the 'arms race' hypothesis, which states that waves of genomic invasion by TEs throughout evolution drove the selection of KZFP genes after they first emerged in the last common ancestor of tetrapods, lung fish and coelacanth some 420 million years ago (Jacobs et al. 2014; Imbeault et al. 2017). While partly supportive of this proposal, functional and phylogenetic studies point to a more complex model, strongly suggesting that KZFPs have facilitated the co-option of TE-embedded regulatory sequences (TEeRS) into transcriptional networks throughout tetrapod evolution (Najafabadi et al. 2015; Imbeault et al. 2017; Helleboid et al. 2019). TEeRS indeed host an abundance of transcription factor (TF) binding sites (Bourque et al. 2008; Sundaram et al. 2014), and KZFPs and their TE targets influence a broad array of biological processes from early embryogenesis to adult life, conferring a high degree of species specificity (Trono 2015; Pontis et al. 2019; Chuong et al. 2013, 2016; Turelli et al. 2020). TEeRS can act as enhancers, repressors, promoters, terminators, insulators or via post-transcriptional mechanism (Garcia-Perez et al. 2016; Chuong et al. 2017). While these co-opted TE functions are key to human biology, their deregulation can also contribute to pathologies such as cancer and neurodegenerative diseases (Jang et al. 2019; Attig et al. 2019; Chuong et al. 2016; Li et al. 2015; Ito et al. 2020; Jönsson et al. 2020).

KZFPs and TEs are broadly expressed during human early development, playing key roles in embryonic genome activation and controlling transcription in pluripotent stem cells (Theunissen et al. 2016;

50    Pontis et al. 2019; Turelli et al. 2020). However, how much TEeRS and their polydactyl controllers

51    influence later developmental stages and the physiology of adult tissues is still poorly defined.

52    Intriguingly, KZFPs are collectively more highly expressed in the human brain than in other adult

53    tissues, suggesting a prominent impact for these epigenetic regulators and their TEeRS targets in the

54    function of this organ (Nowick et al. 2009; Imbeault et al. 2017; Farmiloe et al. 2020; Turelli et al. 2020).

55    In line with this hypothesis, we recently described how ZNF417 and ZNF587, two primate specific KZFPs

56    repressing HERVK (human endogenous retrovirus K) and SVA (SINE-VNTR-Alu) integrants in human

57    embryonic stem cells (hESC), are expressed in specific regions of the human developing and adult brain

58    (Turelli et al. 2020). Through the control of TEeRS, these KZFPs influence the differentiation and

59    neurotransmission profile of neurons and prevent the induction of neurotoxic retroviral proteins and

60    an interferon-like response (Turelli et al. 2020). Furthermore, expression of LINE1, another class of TEs,

61    has been noted in human neural progenitor cells (hNPCs) and in the adult human brain, occasionally

62    leading to *de novo* retrotransposition events (Muotri et al. 2005; Coufal et al. 2009; Muotri et al. 2010;

63    Upton et al. 2015; Erwin et al. 2016; Guffanti et al. 2018). Finally, various patterns of TE de-repression

64    have been reported in several neurodevelopmental and neurodegenerative disorders, indicating that

65    a de-regulated 'transposcriptome' may be detrimental to brain development or homeostasis (Tam et

66    al. 2019; Jönsson et al. 2020).

67    A growing number of genomic studies relying on bulk RNA sequencing (RNA-seq), single cell RNA

68    sequencing (scRNA-seq), assay for transposase accessible chromatin using sequencing (ATAC-seq) and

69    other types of epigenomic analyses are teasing apart the transcriptional landscape of the developing

70    human brain, revealing its dynamism and the complexity of the underlying cellular make-up (Kang et

71    al. 2011; Miller et al. 2014; Fullard et al. 2018; Li et al. 2018; Keil et al. 2018; Zhong et al. 2018; Cardoso-

72    Moreira et al. 2019). The present work was undertaken to explore the contribution of TEs and their

73    KZFP controllers to this process. Our results identify KZFPs and TEeRS as important spatiotemporal

74    contributors to gene expression in both the developing and adult brain, and reveal how neurological

75    proteins with modified characteristics can arise from TE-mediated transcriptional innovations.
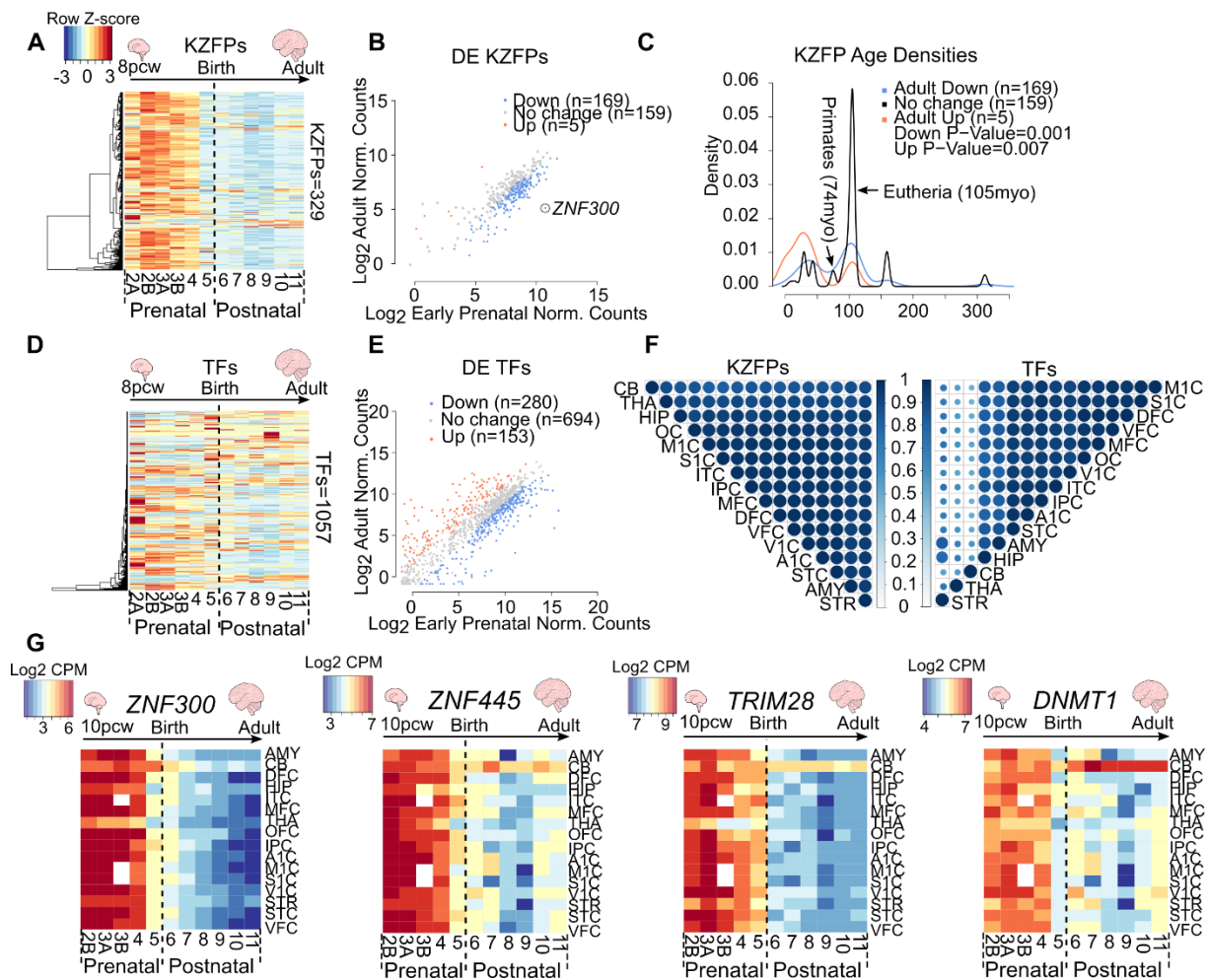
3

76 **Results**

77 **Spatiotemporal patterns of KZFP gene expression during brain development**

78 In order to determine the spatiotemporal patterns of KZFPs and TE expression in human neurogenesis,

79 we analysed RNA-seq data from 507 samples corresponding to 16 different brain regions and 12

80 developmental stages (from 4 weeks post-conception to adulthood) available through the Brainspan

81 Atlas of the Human Brain (Miller et al. 2014) and through Cardoso-Moreira et al. 2019 (Supplemental

82 Fig. S1A & B). While the latter dataset comprises 114 samples exclusively from dorsolateral frontal

83 cortex (DFC) and cerebellum (CB), transcriptomes for these regions were largely concordant with those

84 documented in Brainspan, justifying the two resources as suitable for reciprocal validation

85 (Supplemental Fig. S1C & D; Supplemental Table 1 & 2). We first examined KZFP gene expression in

86 these two brain regions, which are representative of the forebrain and the hindbrain, respectively

87 (Supplemental Fig. S1B). The large majority of KZFPs expressed in the DFC exhibited higher levels at

88 early prenatal stages to drop shortly before birth and remain low onwards (Fig. 1A). When comparing

89 early prenatal (2A-3B; 8-18 post-conception weeks) and adult (11; age 20-60+ years) stages, about half

90 (169/333) of KZFPs were more expressed in the former and only 1.5% (5/333) in the latter, the rest

91 being stable (Fig. 1B). This temporal pattern was less striking in the cerebellum (Supplemental Fig. S2A),

92 with only 15.9% (53/333) and 2.1% (7/333) of KZFPs more strongly expressed in early prenatal and in

93 adult respectively (Supplemental Fig. S2B). Thus, KZFP gene expression patterns are characterized by

94 both temporal and regional specificity.

95 KZFP genes have emerged continuously during higher vertebrate evolution, collectively undergoing a

96 high turnover in individual lineages. Amongst some 360 human KZFPs, about half are primate-

97 restricted, whereas a few are highly conserved, with orthologous sequences present in species that

98 diverged more than 300 million years ago (Imbeault et al. 2017; Huntley et al. 2006). To determine if

99 the differentially

**Figure 1. KZFP genes exhibit a global pre to postnatal decrease in expression.** (A) Heatmaps of KZFP expression across human neurogenesis in the DFC. Scale represents the row Z-score. See also Supplemental Table 2 (B) Dot plot of differential expression analysis of KZFP genes in the DFC comparing adult (stage 11) to early prenatal stages (stage 2A to 3B) of neurogenesis. Only KZFPs differentially expressed in both datasets are shown. Up (orange) represents KZFPs significantly upregulated in adult versus early prenatal (Fold change ≥ 2, FDR ≤ 0.05). Down (blue) represents KZFPs significantly downregulated in adult (Fold change ≤ -2, FDR ≤ 0.05). See also Supplemental Table 3. (C) Density plot depicting estimated age of KZFPs of each category in (B) (*P*≤0.05, Wilcoxon test). (D) Heatmaps of TF expression across human neurogenesis in the DFC. Scale same as in A. (E) Dot plot of differential expression analysis of TFs (as defined in Lambert et al., 2018) in the DFC, excluding KZFP genes, comparing adult (stage 11) to early prenatal stages (stage 2A to 3B) of neurogenesis. Only TFs differentially expressed in both datasets are shown. Up (orange) represents TFs significantly upregulated in adult versus early prenatal (Fold change ≥ 2, FDR ≤ 0.05). Down (blue) represents KZFPs significantly downregulated in adult (Fold change ≤ -2, FDR ≤ 0.05). See also Supplemental Table 3. (F) Correlation plots representing the Pearson correlation coefficient of temporal KZFP expression (left) and TF expression (right) between all 16 regions. Size of spot and colour both represent the correlation coefficient. 0=no correlation, 1=strong correlation. (G) Heatmaps depicting the log2 counts per million (CPM) for selected KZFPs and TFs over the 16 regions included. See also Supplemental Table 1 & 2. All plots show expression data from Brainspan.

120    expressed KZFPs arose at particular times in evolution, we determined their ages. We found KZFPs

121    either significantly downregulated or upregulated from early prenatal to adult stages to be significantly

122    younger than those displaying no differences between these developmental periods (Fig. 1C, Wilcox

123    test $p<=0.01$). This delineates two subsets amongst KZFPs participating in brain development, one

124    evolutionarily recent and more transcriptionally dynamic, the other more conserved and

125    transcriptionally static.

126    Of note, KZFPs appeared distinct amongst TFs (as defined in Lambert et al. 2018), as other members

127    of this functional family exhibited far more diverse patterns of expression throughout development,

128    whether in the forebrain or in the cerebellum (Fig. 1D & E; Supplemental Fig. S2C & D). Only about a

129    quarter of TFs were indeed more highly expressed in early prenatal stages in either region, against

130    around 10% in the adult brain (Fig. 1E; Supplemental Fig. S2C & D). Furthermore, temporal expression

131    patterns of KZFP genes were highly correlated across all 16 brain regions, albeit to a lesser extent in

132    the cerebellum (Fig. 1F). In contrast, other TFs displayed far more diverse behaviours, with the CB,

133    mediodorsal nucleus of the thalamus (THA) and striatum (STR) exhibiting reduced correlation values

134    compared to other regions (Fig. 1F). Thus, KZFPs are collectively subjected to a remarkable degree of

135    spatiotemporal coordination in spite of the diversity of their genomic targets and of cell types present

136    in the various regions of the brain. The KZFP gene most differentially expressed in prenatal versus

137    postnatal DFC was the hematopoietic differentiation associated *ZNF300* (Xu et al. 2010) (Fig. 1B;

138    Supplemental Table 3). This was true in all brain regions, although its transcripts persisted longer in

139    the cerebellum compared to other areas (Fig. 1G; Supplemental Table 1 & 2). *ZNF445*, which binds and

140    controls imprinted loci in humans (Takahashi et al. 2019), similarly exhibited comparable patterns

141    across all brain regions but its expression was largely maintained in the cerebellum all the way to
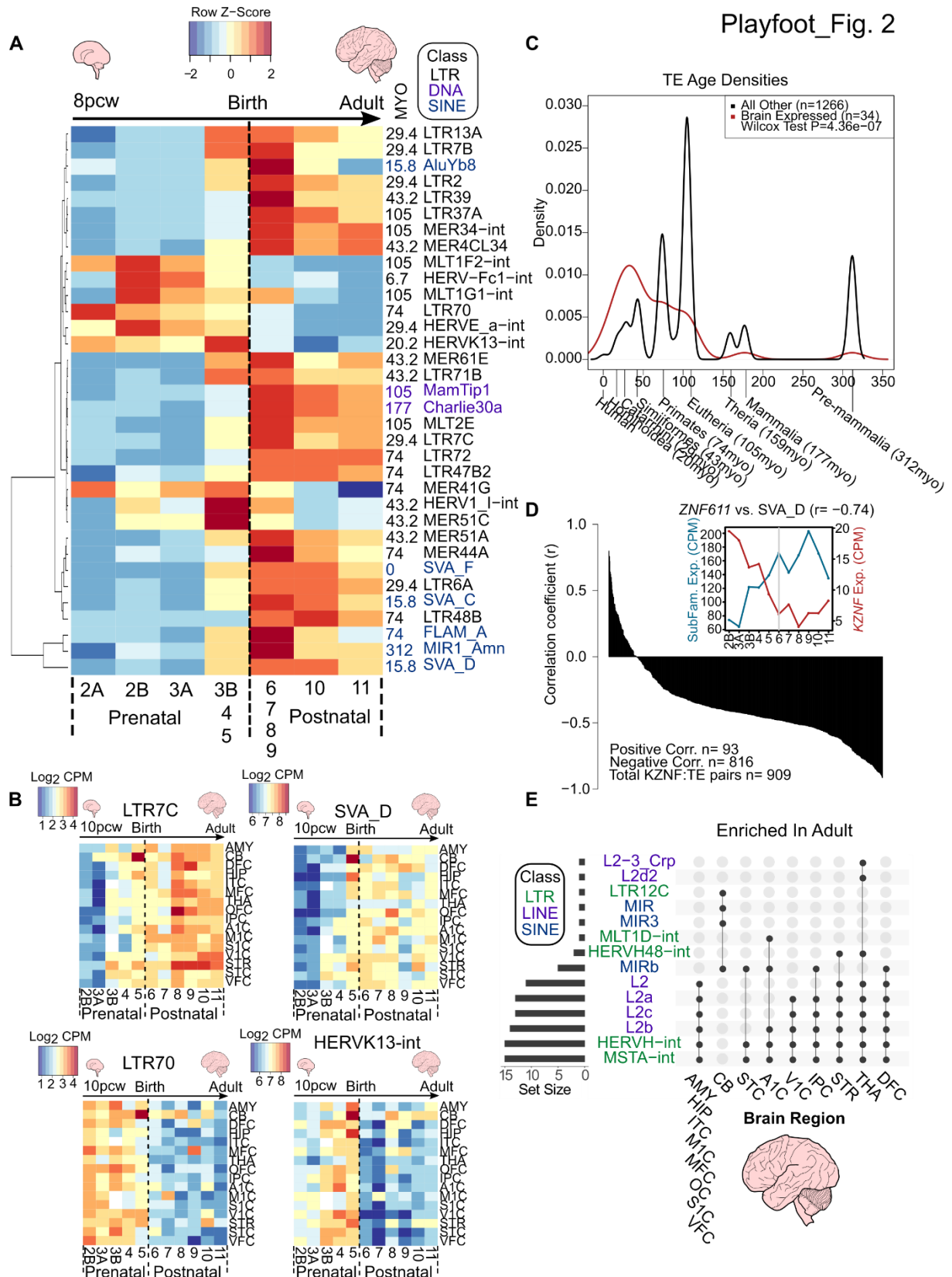
142    adulthood (Fig. 1G; Supplemental Table 1 & 2).

143    We next examined *KAP1,* which encodes a protein that serves as corepressor for many KZFP (Ecco et

144    al. 2017). Its expression levels were globally higher than those of any KZFP, albeit also with a drop from

145   prenatal to postnatal stages except in the cerebellum (Fig. 1G; Supplemental Table 1 & 2). We also

146   probed *DNMT1*, which encodes the maintenance DNA methyltransferase important for TE repression

147   in neural progenitor cells and other somatic tissues beyond the early embryonic period (Jönsson et al.

148   2019). Although displaying overall patterns comparable to those seen for *KZFPs* and *KAP1*, *DNMT1*

149   expression progressively increased in the cerebellum to reach its highest level in the adult (Fig. 1G;

150   Supplemental Table 1 & 2).  In sum, KZFPs and their main epigenetic cofactors exhibit a largely

151   homogenous, dynamic spatiotemporal reduction in expression during human brain development.


152   **TE subfamilies are dynamically expressed throughout development**


153   Having determined that the expression of most KZFPs drops at late stages of prenatal brain

154   development, we examined the behaviour of their TE targets. Young TEs are highly repetitive, which

155   complicates the mapping of TE-derived RNA-seq reads to unique genomic loci, thus biasing against the

156   scoring of their expression. We therefore first analysed RNA-seq reads mapping to multiple TE loci

157   within the same subfamily, regardless of positional information. In the DFC, discrete subfamilies,

158   predominantly from the LTR class and to a lesser extent the SINE class, exhibited temporally distinct

159   dynamics, concordant between datasets (Pearson correlation coefficient ≥ 0.7) (Fig. 2A; Supplemental

160   Table 4). The same was true for the cerebellum, but with moderately different subfamilies passing our

161   threshold for concordance between datasets (Supplemental Fig. S3A; Supplemental Table 4). In the

162   DFC, for example, the LTR7C and SVA-D subfamilies exhibited higher postnatal expression, whereas

163   LTR70 and HERVK13-int behaved inversely, albeit without marked differences between brain regions

164   (Fig. 2B; Supplemental Table 4 & 5). Similarly to KZFP genes, TEs have emerged continuously

165   throughout evolution, with both young integrants and relics of ancient TEs reflective of different waves

166   of genomic invasion. Using TE subfamily age estimates from DFAM (Hubley et al. 2016), we found that

167   dynamically expressed TEs, concordant between both datasets, were significantly younger than non-

168   concordantly expressed subfamilies in the DFC and cerebellum (Fig. 2C; Supplemental Fig. S3B).

Playfoot_Fig. 2

169

170

171

**Figure 2. TE subfamilies and unique loci exhibit spatiotemporal expression patterns.** (A) Heatmap of TE subfamilies with concordant expression behaviours between both datasets (Pearson correlation coefficient ≥ 0.7) across human neurogenesis in the DFC. See also Supplemental Table 4. The mean expression values for stages 3B, 4 and 5, and also stages 6, 7, 8 and 9 were combined and averaged to reduce inherent variability due to low numbers of samples for some stages (see Supplemental Fig. S1B). Scale represents the row Z-score. TE subfamily age in million years old (MYO) and class is shown to the right of the plot. (B) Heatmaps of TE subfamily expression across human neurogenesis in all 16 regions. See also Supplemental Table 4 & 5. Scale represents log2 CPM. Stage 2A was omitted due to lack of samples for some brain regions (see Supplemental Fig. S1B). (C) Density plot depicting estimated age of TEs in A (P≤0.05, Wilcoxon test). Evolutionary stages and corresponding ages are shown beneath the plot. (D) Barplot showing the Pearson correlation coefficient of KZFP expression and their target TE subfamily expression. 1=highly correlated, -1=highly anti-correlated. (D Inset) Line plot showing expression in counts per million of *ZNF611* and its main TE target subfamily, SVA_D and their Pearson correlation coefficient (-0.74, p-value=0.006). Grey line indicates birth at stage 6. See also Supplemental Table 6. (E) UpSet plot showing the significantly enriched differentially expressed subfamilies between adult and early pre-natal stages per region from unique mapping analyses. Set size represents the number of regions the specific TE was significantly differentially enriched in. Joined points represent combinations of significantly differentially expressed TE subfamilies. See also Supplemental Table 7 and 8. All plots show expression data from Brainspan.

We next analysed the temporal dynamics of the expression of KZFPs and their TE targets in the DFC, for which samples were available in highest abundance. For this, we matched KZFP ligands to their significantly bound TE subfamilies using an in-house algorithm on a large collection of ChIP-exo data (Imbeault et al. 2017). The results revealed that an overwhelming majority of KZFP:TE subfamily pairs (816 vs. 93) were anti-correlated in their expression, consistent with the known role of KZFPs as TE repressors (Fig. 2D; Supplemental Table 6). For example, *ZNF611* is a previously characterised major regulator of SVA-D in early embryogenesis (Pontis et al. 2019), and the two exhibited strongly anti-correlated expression throughout human brain development (Fig. 2D inset).

We next expanded our study by examining the expression of individual TE integrants, assigning RNA-seq reads to their genomic source loci and comparing early prenatal (stages 2A to 3B) and adult (stage 11) samples for the 16 available brain regions (Supplemental Fig. S1A & B). We found between 5,000 and 7,000 significant differentially expressed TE loci in each region, with 4,000 loci common to both DFC and CB datasets (Supplemental Fig. S3C; Supplemental Table 7 & 8). Integrants belonging to fourteen TE subfamilies from the LTR, LINE and SINE classes were significantly more expressed in adult

206    samples, with HERVH-int, MSTA-int and L2 elements significantly enriched in most brain regions (Fig.

207    2E). The cerebellum again exhibited distinct patterns, with significant enrichment of LTR12C and MIR

208    elements instead (Fig. 2E). Conversely, integrants from 11 TE subfamilies were more expressed in the

209    early prenatal period, largely in specific brain regions (Supplemental Fig. S3D). Together, these results

210    highlight the spatiotemporal dynamic nature of the transposcriptome in the developing human brain.

211    **Transpochimeric gene transcripts during human brain development**
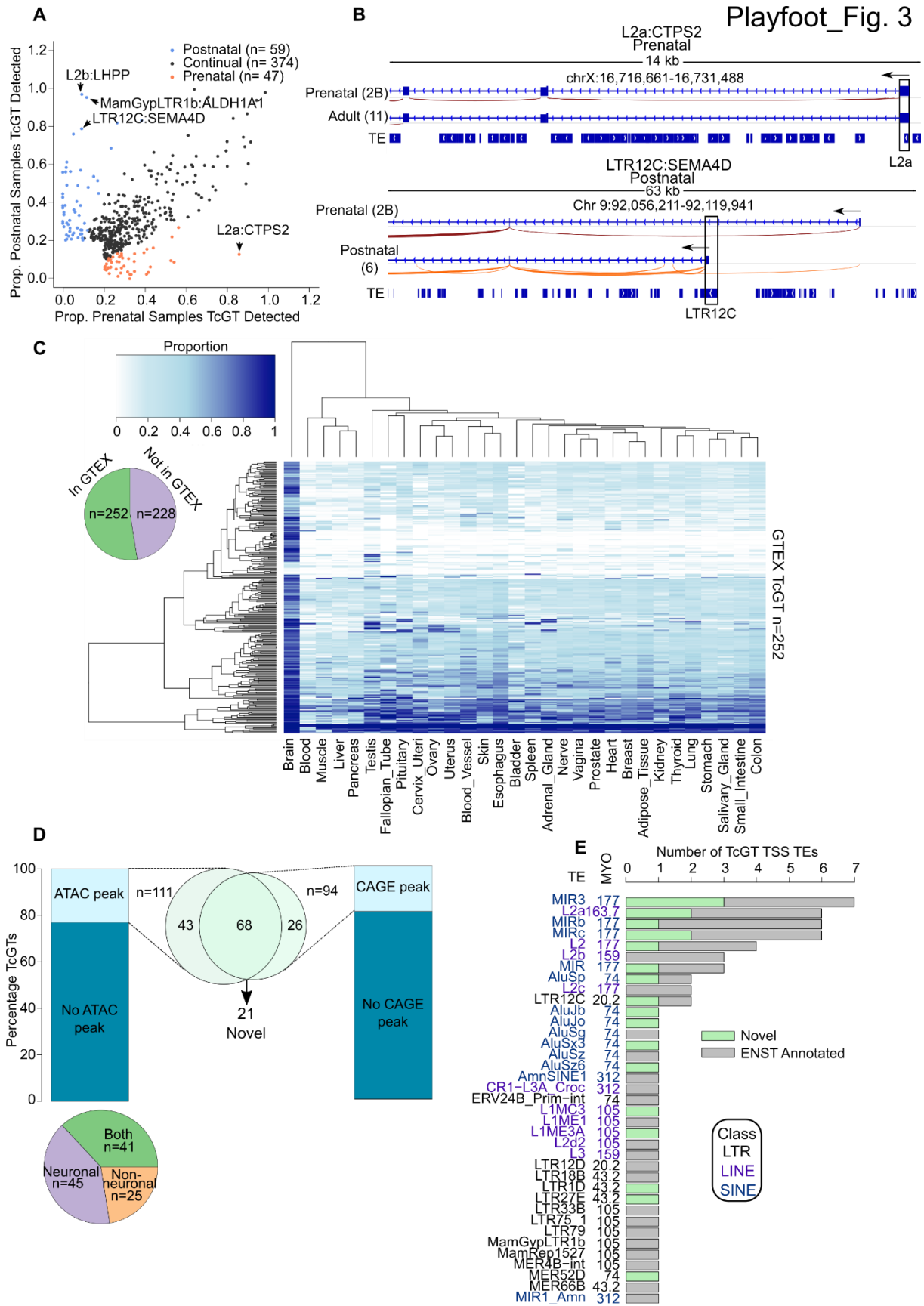
212    TE expression may be reflective of either 'passive' co-transcription from genic transcripts or *bona fide*

213    TE promoter activity (reviewed in Lanciano and Cristofari 2020). Transpochimeric gene transcripts

214    (TcGTs), that is, gene transcripts driven by TE-derived promoters, are the most easily interpretable and

215    direct manifestation of the influence of TEeRS on gene expression. Some evidence for a role of TcGTs

216    in the brain was provided by the recent observation that DNMT1 represses in hNPCs the expression of

217    hominoid-restricted LINE1 elements, which subsequently act as alternative promoters for genes

218    involved in neuronal functions (Jönsson et al. 2019). To explore more broadly the potential role of

219    TcGTs in human brain development and function, we performed *de novo* transcript assembly,

220    searching for mature transcripts with a TE-derived sequence at their 5' end and the coding sequence

221    of a cellular gene downstream. Due to the striking anti-correlation in KZFP and global TE expression

222    between prenatal (stage 2A to stage 5) and postnatal stages (stage 6 to stage 11), we concentrated on

223    these two periods, retaining only TcGTs present in greater than 20% of either prenatal, postnatal or

224    both categories of samples and behaving in the same temporal manner in the two independent

225    datasets. If there was a two-fold difference in the proportion of prenatal versus postnatal, the TcGT

226    was annotated as either pre- or postnatal, whereas those below this threshold were deemed continual.

227    Our search yielded 480 high confidence TcGTs, of which 9.8% (47/480) were prenatal, 12.3% (59/480)

228    postnatal and 72.3% (374/480) continual (Fig. 3A; Supplemental Table 9). Amongst pre- or postnatal

229    TcGTs, developmental trajectories differed substantially, with some detected exclusively at either

230    stage. For example, an L2a-driven isoform of *CTP synthase 2* (*CTPS2*), whose product catalyses CTP

10

231 formation from UTP (van Kuilenburg et al. 2000), was found in 86% of all prenatal samples but only

232 12% of postnatal samples (Fig. 3A & B), whereas the inverse was observed for a MamGypLTR1b-driven

233 isoform of the astrocyte associated *Aldehyde Dehydrogenase 1 Family Member A1* (*ALDH1A1*) (Adam

234 et al. 2012) (12% vs. 95%) and an L2b-driven isoform of *Phospholysine Phosphohistidine Inorganic*

235 *Pyrophosphate Phosphatase* (*LHPP*) (0.9% vs. 97%) (Fig. 3A), the host of intronic single nucleotide

236 polymorphisms (SNPs) associated with major depressive disorder (Neff et al. 2009; Cui et al. 2016). The

237 previously reported LTR12C-driven transcript of *Semaphorin 4D* (SEMA4D), the product of which

238 participates in axon guidance (Cohen et al. 2009; Kumanogoh and Kikutani 2004), was detected in 79%

239 of postnatal and only 0.9% of prenatal samples where it was instead expressed from a non-TE

240 promoter, indicating a promoter switch during neurogenesis (Fig. 3A & B).

241 We next examined the broader expression pattern of the 480 TcGTs detected during brain

242 development. By applying our pipeline to the Genotype Tissue Expression (GTEX) dataset (Melé et al.

243 2015), we detected around half of them in this collection of predominantly adult samples (Fig. 3C;

244 Supplemental Table 9). Some were present in all available tissues, but the vast majority were brain

245 restricted (Fig. 3C).

**TcGTs exhibit cell type-specific modes of expression**

247 We next analysed the state of the chromatin at the transcription start site (TSS) of the 480 TcGTs

248 expressed during brain development by intersecting their proximal, TE-residing TSS (+/-200bp) with

249 ATAC-seq consensus peaks from neuronal (NeuN+) and non-neuronal (NeuN-) cells across 14 distinct

250 adult brain regions from the Brain Open Chromatin Atlas (BOCA) (Fullard et al. 2018). About a quarter

251 (111/480) of these TcGTs TSS overlapped with ATAC-seq peaks in the adult brain, indicating that their

252 chromatin was opened in this setting (Fig. 3D). Of these, two-thirds exhibited cell type

11

Playfoot_Fig. 3

**Figure 3. TE co-option as genic promoters drives spatiotemporal gene expression in human neurogenesis.** (A) Dot plot showing the proportion of pre or postnatal samples TcGTs were detected in and behaving similarly in both datasets (prenatal, postnatal or continual). (B) Sashimi browser plots from IGV showing the splicing events in representative samples for prenatal enriched TcGT L2a:CTPS2 and the postnatal enriched LTR12C:SEMA4D. (C) Heatmap indicating the proportion of samples per GTEX tissue each TcGT from A was detected in. Each row represents an individual TcGT and each column a different tissue. (C inset) Pie chart indicating the proportion of neurodevelopmental TcGTs detected in GTEX. (D) Stacked barplot indicating the proportion of TcGT TE TSS loci overlapping an ATAC-seq peak from BOCA (left) and a pie chart indicating their cell type distribution (bottom left). Stacked barplot (right) indicating the proportion of TcGT TE TSS loci overlapping a FANTOM5 defined CAGE peak. Pie chart (centre) showing the overlap of ATAC-seq and CAGE peak associated TcGTs and highlighting 21 novel, non-ENSEMBL annotated transcripts. (E) Stacked barplots indicating the TE subfamily, TE class, TE age and the ENSEMBL overlap of each TcGT TE TSS loci. See also Supplemental Table 9 for all TcGT information.
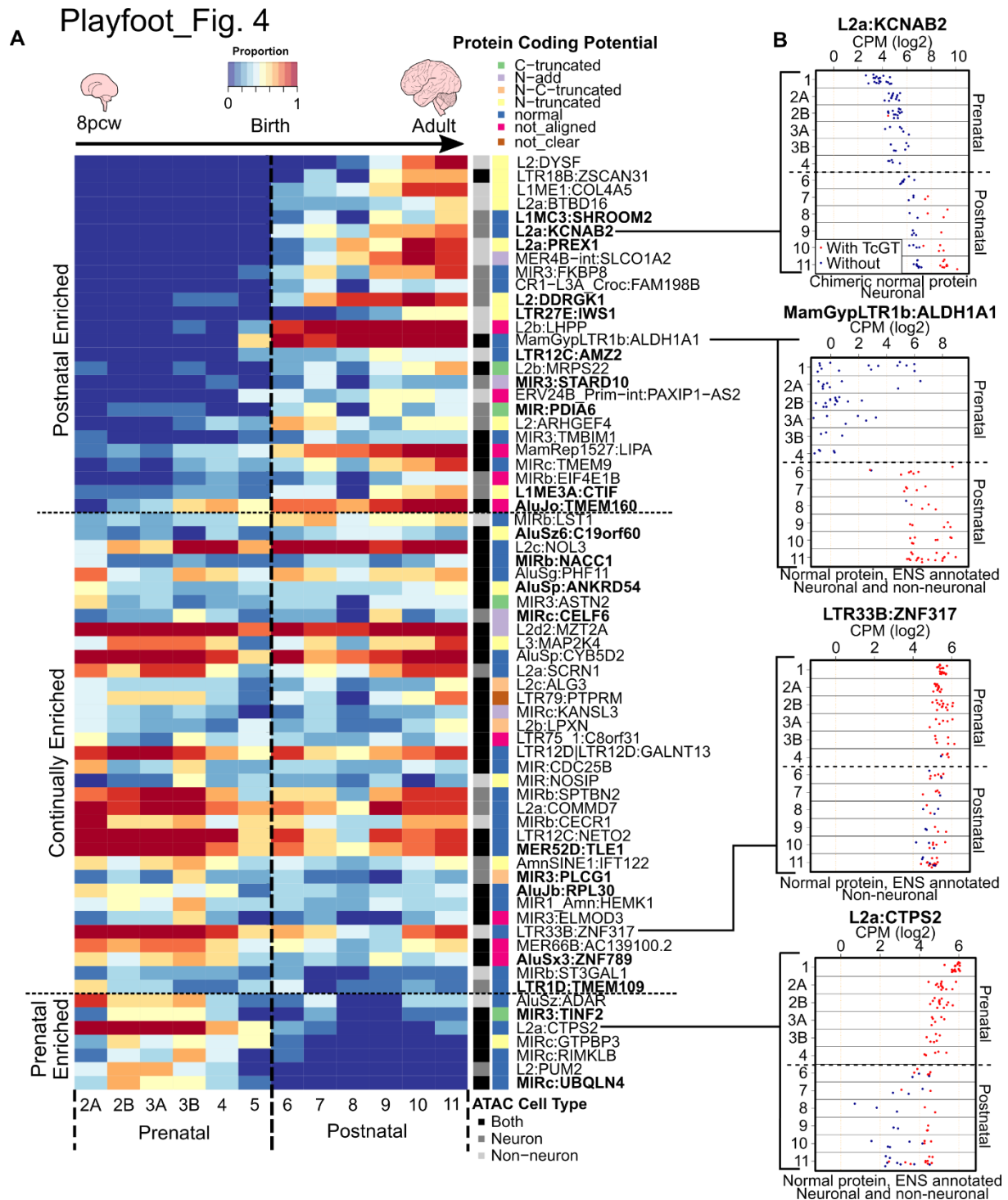
specificity, either to neurons (40.5%, 45/111) or to non-neuronal cells (22.5% 25/111), whereas a third (41/111) were present in both cell subsets (Fig. 3D; Supplemental Table 9). These cell-restricted patterns were generally independent of the brain region considered, as illustrated by two postnatal enriched TcGTs, the non-neuronal L2-driven *Dysferlin* (*DYSF*) (Supplemental Fig. S4A), a gene mutations of which are associated with limb girdle muscular dystrophy 2B (Bashir et al. 1998; Liu et al. 1998), and the neuronal L2a-driven *Potassium Voltage-Gated Channel Subfamily A Regulatory Beta Subunit 2* (*KCNAB2*) encoding a regulator of neuronal excitability (McCormack et al. 2002) (Supplemental Fig. S4B).

To confirm that transcription of the TcGTs detected in the developing human brain was starting at the identified TE, we intersected their TSS with CAGE (cap analysis of gene expression) peaks previously defined in around 1,000 human cell lines and tissues by the FANTOM5 consortium (Forrest et al. 2014; Lizio et al. 2015). About a fifth of the TcGTs TSS (19.5%, 94/480) overlapped with CAGE peaks, of which 68 also corresponded to ATAC-seq peaks, providing a subset of high confidence TE-derived TSS loci driving gene transcription in the developing brain (Fig. 3D; Supplemental table 9). Of these, 21 were not annotated in ENSEMBL (Fig. 3D; Supplemental table 9), indicating that co-opted TEs acting as promoter elements are contributing to a previously undetected TE-derived neurodevelopmental transcription network.

13

287    We concentrated deeper analyses on the 68 high confidence brain developmental TcGTs. Thirty-seven

288    different TE subfamilies accounted for their promoters but MIRs and L2s, belonging respectively to the

289    SINE and LINE families, contributed almost half, perhaps due in part to their high prevalence in the

290    genome (MiR3 and L2a: 87,870 and 166,340 integrants, respectively) (Fig. 3E), and LTRs about a fifth.

291    A large range of evolutionary ages were represented, from the ~20 myo (million year old) LTR12C to

292    the ~177 myo MIRs and L2s.

293    Of these 68 high-confidence TcGTs, 38.2% (26/68) were postnatal-specific, 51.5% (35/68) were

294    continually detected and 10.3% (7/68) were prenatal-restricted (Fig. 4A). Furthermore, the 5' end of

295    these TcGTs coincided with ATAC-seq peaks from neurons in 26.5% (18/68), from non-neuronal cells

296    in 22% (15/68), and from both in 51.5% (35/68) of cases (Fig. 4A). Some TcGTs were present in all brain

297    regions, whereas others exhibited regional specificity (Supplemental Fig. S5A). For example,

298    L2:DDRGK1 and L2a:KCNAB2, among others, were detected both postnatally and in a higher proportion

299    of neocortex regions compared to the cerebellum (Fig. 4A; Supplemental Fig. S5A). We next aimed to

300    determine if the detected TcGTs had the capacity to code for protein. Importantly, *in silico* prediction

301    of the protein coding potential of these TcGTs, found that about half (31/68) likely encoded the

302    canonical protein sequence and a fifth (15/68) an N-truncated isoform, while other configurations (N-

303    terminal addition, C- or N- and C-truncation) were less frequent (Fig. 4A; Supplemental Table 9).

304    To estimate the relative contribution of the TE and non-TE promoters to the expression of the 68 genes

305    involved in high confidence TcGTs, we compared their transcription levels in samples where the TcGT

306    was or was not detected (Fig. 4B). In some cases, the TcGT was associated with higher levels of gene

307    expression in a temporal manner such as the postnatally detected L2a:KCNAB2 (top) and most

308    strikingly MamGypLTR1b:ALDH1A1 (top mid), compared to their non-TE-driven counterparts (Fig. 4B).

309    The continually detected, non-neuronal LTR33B:ZNF317 (bottom mid) was associated with high

310    expression throughout brain development, suggestive of a constitutive TE derived promoter.

311    Conversely, some TcGTs were associated with higher prenatal expression, such as with L2a:CTPS2

**Figure 4. TcGTs are temporally expressed throughout neurogenesis in a cell type specific manner, exhibit protein coding potential and drive transcript expression.** (A) Heatmap showing the proportion of samples per developmental stage the 68 TcGTs (from Fig. 3D) were detected in the Brainspan dataset, alongside their ATAC-seq cell type overlaps and protein coding potential determined via *in silico* translation. Bold indicates novel transcripts not annotated in ENSEMBL. See also Supplemental Table 9. (B) Dot plots showing the gene expression level per stage for the specified gene for samples where the TcGT was detected (red) and where it was not (blue) from Cardoso-Moreira dataset as comparison to (A). Dashed line represents birth at stage 6.

321

322    (bottom), while  for other genes there were  more moderate expression differences in samples with

323    and without TcGT detection as seen for the postnatally detected neuronal TcGT L2:DDRGK1

324    (Supplemental Fig. S5B).

325    **Experimental validation of brain-detected TcGTs**

326    To verify that the TE and genic exon belonged to the same mRNA transcript, we next aimed to

327    experimentally confirm TcGT candidates in the SH-SY-5Y neuroblastoma cell line. Using qRT-PCR

328    primers within the TE TSS and subsequent genic exon, we detected appreciable expression of TcGTs in

329    this cell system (Supplemental Fig. S6A).  However, this did not formally demonstrate that transcription

330    was driven by the TE. To address this point, we targeted a CRISPR-based activation system (CRISPRa)

331    to the TSS region of TcGTs in 293T cells (Chavez et al. 2015) (Fig. 5A). We picked candidates based on

332    the ease of gRNA design and the potential mechanistic or biological relevance of their protein product.

333    We selected three anti-sense L2-driven, cell type-specific TcGTs predicted to encode for proteins

334    involved in brain development: KCNAB2, DYSF and DDRGK1, the first in its canonical protein isoform

335    and the other two as N-truncated isoforms. Activation of each of these three TcGTs could be induced

336    with the CRISPRa system, confirming that they were indeed driven by their respective TE promoters
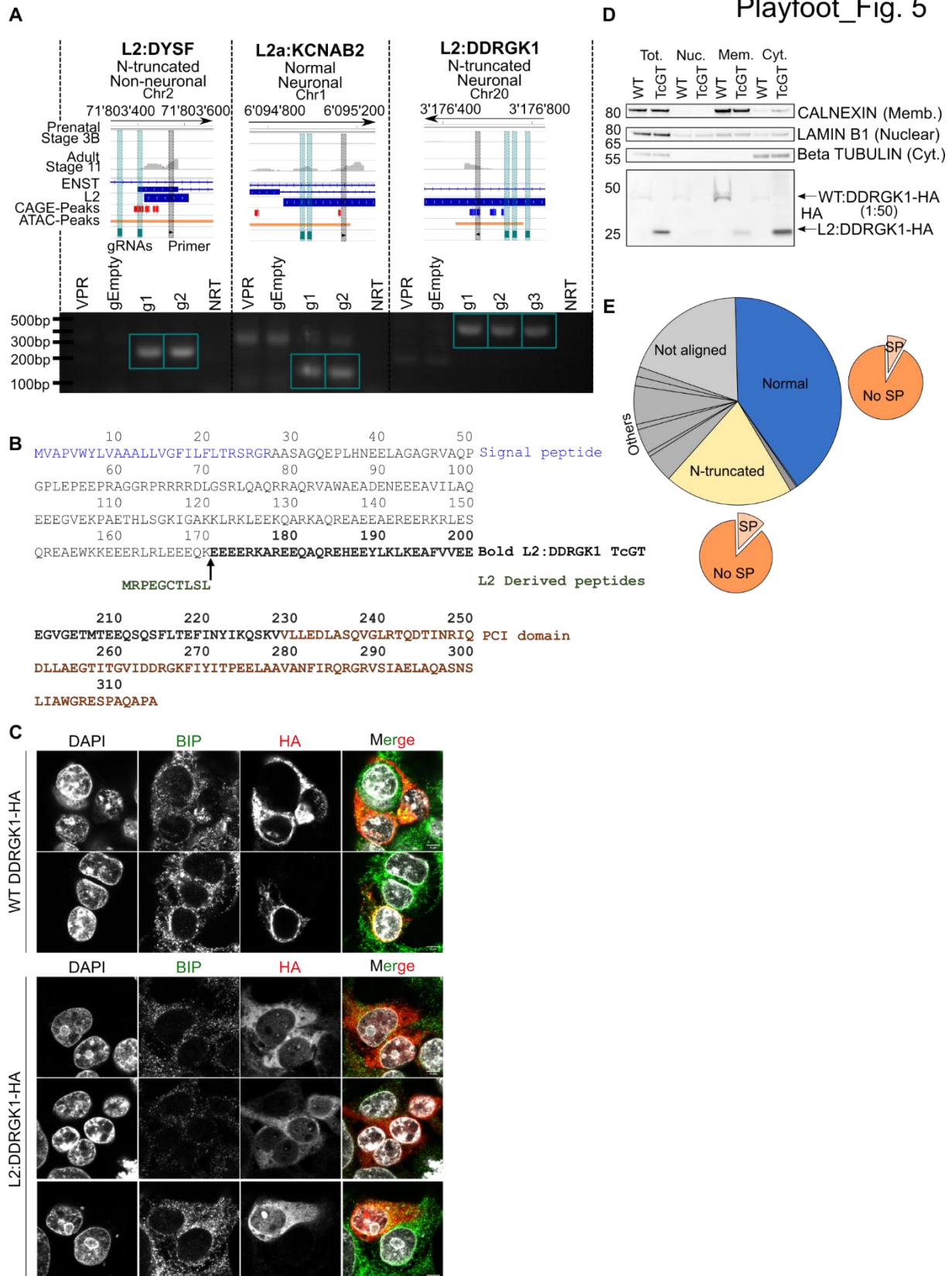
337    (Fig. 5A).

338    **TcGT-encoded protein isoforms can display differential subcellular localisation**

339    Having noted that 22% of high-confidence TcGTs were predicted to encode N-truncated proteins (Fig.

340    4A), we hypothesised that this could, in some cases, result in derivatives deprived of important

341    subcellular localization domains, such as the endoplasmic reticulum (ER)-targeting N-terminal signal

342    peptide. We focused on L2:DDRGK1 as it was enriched postnatally, neuron-specific, not annotated in

343    ENSEMBL and experimentally validated by our 293T-based CRISPRa experiment (Fig. 4A; Fig. 5A;

344    Supplemental Fig. S6B; Supplemental Table 9). GWAS studies have also identified a DDRGK1 associated

345    risk locus for Parkinson's disease (Nalls et al. 2014; Chang et al. 2017). The canonical DDRGK1 protein

346    product is anchored to the ER membrane by an N-terminal 27 amino acid signal peptide (Fig. 5B) and

347    plays a role in ER homeostasis and ER-phagy (Liang et al. 2020; Liu et al. 2017). In the predicted

348    translated product of the L2:DDRGK1 TcGT, the signal peptide is replaced by a 10 amino acid L2-

349    encoded sequence, conserved in new-world primates, but harboring non-synonymous substitutions in

350    old-world primates (Fig. 5B; Supplemental Fig. S7A). Of note, this L2 integrant is absent in mice

351    (Supplemental Fig. S7A). Furthermore, the L2:DDRGK1 TcGT is detected in the Rhesus Macaque

352    developing brain with the same prenatal to postnatal expression dynamics as in humans (Supplemental

353    Fig. S7B & C). We therefore transfected HEK293T cells with plasmids expressing HA-tagged versions of

354    either the canonical "wild-type" (WT) DDRGK1 transcript or its TcGT counterpart and examined the

355    subcellular localization of the resulting proteins by indirect immunofluorescence (Fig. 5C) and by

356    cellular fractionation followed by western blotting (Fig. 5D). Confocal microscopy revealed that

357    WT:DDRGK1-HA largely co-localized with BIP, an ER membrane marker, while L2:DDRGK1-HA displayed

358    a diffuse cytosolic pattern (Fig. 5C). Cellular fractionation further confirmed that the WT DDRGK1

359    isoform was sequestered in the membrane fraction, whereas the L2:DDRGK1 counterpart was

360    enriched in cytosol (Fig. 5D).

361    As N-truncated isoforms made up the largest category of *in-silico* predicted TcGT products besides full-

362    length proteins, we next asked how widespread this type of TE-induced protein re-localisation might

363    be. For this, we intersected a database of signal peptide-containing proteins with our initial list of 480

364    TcGT-encoded protein products (Fig. 5E; Supplemental Table 9). Of 94 TcGT products predicted to be

365    N-truncated, 12 contained a putative signal peptide in the canonical isoform. This prediction was

366    supported in 11 cases *in silico* by signalP 5.0 (Almagro Armenteros et al., 2019), which predicted that

367    in all of these instances the TcGT isoforms lacked this putative signal peptide (Supplemental Fig. S8).

368    Therefore, subcellular re-targeting may be a frequent consequence of TE-driven protein innovation.

369

**Playfoot_Fig. 5**

373 **Figure 5. Antisense L2 elements directly drive TcGTs and contribute to chimeric protein formation**
374 **and cytosolic re-localisation of the ER-membrane associated DDRGK1.** (A) Schematic of TcGT TE TSS
375 loci for indicated genes and representative prenatal (stage 3B) and adult (stage 11) RNA-seq tracks.
376 Their associated protein coding potential and cell type specificity are highlighted and CAGE peak loci
377 (red sense strand, blue anti-sense strand), CRISPRa gRNAs (green vertical bar) and TE associated PCR
378 primers are shown (black vertical bar) (top). RT-PCR on cDNA generated from HEK293T cells transiently
379 transfected with dCAS9-VPR plasmid and individual gRNA plasmids containing sequences targeting the
380 TcGT TE TSS loci denoted in the schematic. dCAS9-VPR (VPR) or empty gRNA plasmids (gEmpty) alone
381 were used as controls. Green box indicates bands of correct PCR product size absent in controls. NRT=
382 no reverse transcriptase. (B) Canonical DDRGK1 and TcGT L2:DDRGK1 derived protein sequence. (C)
383 Overexpression of canonical WT DDRGK1-HA and L2:DDRGK1-HA in HEK293T cells followed by
384 immunofluorescent staining for BIP (an ER-membrane associated protein) and HA tag, followed by
385 confocal imaging (scale bar = 5μm). (D) Overexpression of canonical DDRGK1-HA (WT) and L2:DDRGK1-
386 HA (TcGT) in HEK293T cells followed by cellular fractionation and western blot for the indicated marker
387 proteins (right of western blot) and HA tag. For WT DDRGK1 50x less protein lysate compared to
388 L2:DDRGK1 was loaded for the HA blot due to high levels of protein expressed. Image is representative
389 of two independent experiments. (E) Pie charts showing the *in silico* protein coding potential of the
390 480 TcGTs identified in Fig. 3A with the proportion containing a signal peptide shown with the orange
391 pie charts. See also Supplemental Table 9.

392

393 **Discussion**

394 An increasing number of studies are aimed at unravelling the transcriptional dynamics of human

395 neurogenesis (Li et al. 2018; Cardoso-Moreira et al. 2019; Keil et al. 2018), yet, so far, little attention

396 has been paid to the participation of TEeRS in this process. While retrotransposition of L1HS elements

397 has been suggested to contribute to neuronal plasticity, experimental support for this model is lacking,

398 and the vast majority of TEs hosted by the human genome have long lost the ability to spread (Muotri

399 et al. 2005; Brouha et al. 2003). This prompted us to hypothesise that TEs might exert far greater

400 influences on brain development through their ability to shape gene expression. As a first step towards

401 testing this model, we analysed two independent human neurogenesis RNA-seq datasets with a 'TE

402 centric' approach. This led us to uncover that the transposcriptome undergoes profound changes at

403 each stage of brain development, with the expression of individual TE subfamilies largely anti-

404 correlating to that of their cognate KZFP controllers. Strikingly, KZFP genes were globally

405 downregulated at postnatal versus prenatal stages, coincident with the upregulation of their TE

406 targets. Recent indications from an analysis of TEs resistant to loss of DNA methylation during the wave

19

407    of epigenetic reprogramming in human primordial germ cells (hPGCs) showed modest anti-correlations

408    of KZFPs and their target TE subfamilies in prenatal neurogenesis (Dietmann et al. 2020). The proposal

409    that KZFPs may mediate the exaptation of TEs as developmental enhancers marked in hPGCs is

410    intriguing and, combined with our analyses, suggests a multifaceted KZFP and TE mediated

411    spatiotemporal transcriptional network, not only in prenatal stages but also highly prevalent after

412    birth, with TEeRS playing important roles as alternative promoters, in addition to enhancers,

413    throughout.  Indeed, correlative expression studies on genic KZFP targets suggest that KZFPs may also

414    directly regulate gene promoters during human neurogenesis independently from their TE binding

415    ability (Farmiloe et al. 2020), and KZFPs were amongst genes previously found to be most differentially

416    expressed between the chimpanzee and human brain (Nowick et al. 2009). Increasing evidence also

417    supports a regulatory role for KZFP-targeted TEs in this and other developmental contexts (Ecco et al.

418    2016, 2017; Chen et al. 2019; Pontis et al. 2019; Turelli et al. 2020). For example, we recently

419    demonstrated that two primate-restricted KZFPs, ZNF417 and ZNF587, control the expression of

420    neuronal genes such as *PRODH* and *AADAT* via the regulation of HERVK-based TEeRS (Turelli et al.

421    2020). Furthermore, studies on the transcriptional co-repressors KAP1 and DNMT1 in hNPCs have

422    highlighted their roles in the regulation of TEs and secondarily of cellular genes (Brattås et al. 2017;

423    Jönsson et al. 2019). However, *in vitro* models do not recapitulate the global spatiotemporal

424    complexity of gene and TE expression in the brain, nor its diverse cell type milieu throughout

425    development, hence the interest of performing large scale 'TE centric' bioinformatics analyses on large

426    post-mortem brain RNA-seq datasets.

427    De-repression of TEs, specifically of the LTR class, has been associated with various neurological

428    disorders such as amyotrophic lateral sclerosis (ALS), Alzheimer's disease (AD) and multiple sclerosis

429    (MS) (Tam et al. 2019; Jönsson et al. 2020). The upregulation of LTR class elements in adult versus early

430    prenatal brain is intriguing, as it suggests that LTR transposcription *per se* is a developmentally

431    regulated feature of neurogenesis, which when deregulated is associated with a disease state. We

432    propose that increased postnatal TE expression may possibly be reflective of the development of cell

433   types not present in early prenatal stages, such as astrocytes, microglia and oligodendrocytes, the

434   developmental and transcriptional trajectories of which were identified by scRNA-seq analyses (Li et

435   al. 2018). To determine the transposcriptome in scRNA-seq data remains technically challenging

436   because many TE-derived transcripts are lowly abundant, a limitation that will hopefully be alleviated

437   by progress in sequencing techniques and computational approaches (Linker et al. 2020). Of note, TEs

438   heavily contribute to long non-coding RNAs (lncRNAs), which are abundant in the human brain (Derrien

439   et al. 2012; Kelley and Rinn 2012; Zimmer-Bensch 2019). It is plausible that upregulated TE transcripts

440   play a role in this context, thereby exerting not cis- but trans-acting influences, the identification of

441   which is far more challenging.

442   One increasingly well-characterised aspect of TE co-option is the engagement of TEeRS as alternative

443   promoters. A wide range of oncogene-encoding TE-driven TcGTs have been documented in recent

444   surveys of cancer databases (Jang et al. 2019; Attig et al. 2019), but the role of these transcript variants

445   in physiological conditions remains largely undefined. Tissue-specific TcGTs have also been detected

446   in the mouse developing intestine, liver, lung, stomach and kidney (Miao et al. 2020). Here, we

447   demonstrate not only the spatially and temporally orchestrated expression of TcGTs in the developing

448   human brain, but also that these TcGTs are largely organ- and cell type-specific. Some of them appear

449   to be solely responsible for the expression of the involved gene, whereas others were present

450   alongside canonical non-TE-driven transcripts, indicating sophisticated levels of regulation.

451   By experimental activation of a selected subset of antisense L2-driven TcGTs with CRISPRa and

452   functional analyses of the product of the L2:DDRGK1 transcript, we highlight the functional relevance

453   of this phenomenon for human neurogenesis. DDRGK1 is an ER membrane-associated protein with

454   critical roles in UFMylation, an ubiquitin-like modification, and is involved in the unfolded protein

455   response and ER-phagy (Liu et al. 2017; Liang et al. 2020). DDRGK1 is essential to target interactors like

456   UFL1, the UFMylation ligase, to the ER membrane. The novel cytosolic chimeric L2:DDRGK1 protein,

457   where a short N-terminal sequence derived from the L2 integrant replaces the signal peptide

21

458    characteristic of its canonical counterpart, may therefore exert novel functions in the cytosol of

459    postnatal to adult neurons. As signal peptide excision seems to affect a number of other TcGT products,

460    this example may illustrate a more general phenomenon, whereby TE-driven genome evolution

461    generates novel protein isoforms altering critical cell functions.

462    Our study indicates that the exaptation of TE-embedded regulatory sequences and its facilitation by

463    TE-targeting KZFP controllers have significantly contributed to the complexity of transcriptional

464    networks in the developing human brain. This warrants efforts aimed at delineating the evolutionary

465    and functional impact of this phenomenon, and at defining how its alterations, notably in the context

466    of inter-individual differences at these genomic loci, translates into variations in brain development,

467    function and disease susceptibility.

468    **Methods**

469    **Datasets**

470    Raw RNA-seq fastq files for human and Rhesus macaque brain development (Cardoso-Moreira et al.

471    2019) were downloaded from the European Nucleotide Archive (datasets PRJEB26969 and

472    PRJEB26956, respectively).

473    Raw RNA-seq fastq files for the GTEX and Brainspan (phs000424.v7.p2, phs000755.v2.p1), were

474    downloaded from the dbGaP authorized access platform. Processed bed files containing regional

475    neuronal or non-neuronal ATAC-seq peak loci from the Brain Open Chromatin Atlas (Fullard et al. 2018)

476    were downloaded for hg19. To generate consensus neuronal and non-neuronal ATAC-peak bed files,

477    bed coordinates from all regions were combined and overlapping peak coordinates merged using

478    bedtools merge. Processed bed files for CAGE-seq peak loci from FANTOM5 (Forrest et al. 2014) were

479    downloaded for hg19 (Lizio et al. 2015). Signal peptide containing proteins in human were downloaded

480    from http://signalpeptide.com/index.php. Processed bed files from KZFP ChIP-exo experiments were

481    used from our previous study (Imbeault et al. 2017).

**RNA-seq analysis**

Reads were mapped to the human (hg19), or macaque (rheMac8) genome using hisat2 (Kim et al. 2015) with parameters hisat2 -k 5 --seed 42. Counts on genes and TEs were generated using featureCounts (Liao et al. 2014). To avoid read assignation ambiguity between genes and TEs, a gtf file containing both was provided to featureCounts. For repetitive sequences, an in-house curated version of the Repbase database was used (fragmented LTR and internal segments belonging to a single integrant were merged), generated as previously described (Turelli et al. 2020). Minor modifications to the repeat merging pipeline described in Turelli et al., 2020 were made for Macaque (RepeatMasker 4.0.5 20160202) with the distance between two LTR elements of the same orientation to an ERV-int fragment being less than 400bp. For genes the ensemble release 75 annotation was used. Only uniquely mapped reads were used for counting on genes and TEs with the command 'featureCounts -t exon -g gene_id -Q 10'. For the Brainspan dataset, samples with less than 10 million unique mapped reads on genes were discarded from the analysis. TEs that did not have at least one sample with 50 reads or overlapped an exon were discarded from the mapping TE integrant analysis. For estimating TE subfamilies expression level, reads were summarized using the command featureCounts -M -- fraction -t exon -g gene_id -Q 0 then, for each subfamily, counts on all TE members were added up. As the Cardoso-Moreira et al., 2019 RNA-seq was stranded data, reads on both strands were combined for TEs to facilitate comparison to the non-stranded Brainspan dataset. Normalization for sequencing depth was done for both genes and TEs using the TMM method as implemented in the limma package of Bioconductor (Gentleman et al. 2004) and using the counts on genes as library size. Differential gene expression analysis was performed using voom (Law et al. 2014) as it has been implemented in the limma package of Bioconductor (Gentleman et al. 2004). A gene (or TE) was considered to be differentially expressed when the fold change between groups was greater than two and the p-value was smaller than 0.05. A moderated t-test (as implemented in the limma package of R) was used to test significance. P-values were corrected for multiple testing using the Benjamini-Hochberg's method (Benjamini and Hochberg 1995). Temporal expression correlation analyses of individual genes, TE

23

508    integrants or subfamilies were performed between Brainspan and Cardoso datasets using the

509    'Pearson' method. For inter-regional correlations within the Brainspan dataset, only expressed genes

510    or TEs common to all regions were considered. Bam files and sashimi plots were visualised using the

511    Integrative Genomics Viewer (Katz et al. 2015; Robinson et al. 2011).

512    **TcGT detection pipeline**

513    First, a per sample transcriptome was computed from the RNA-seq bam file using Stringtie (Kovaka et

514    al. 2019) with parameters –j 1 –c 1. Each transcriptome was then crossed using BEDTools (Quinlan and

515    Hall 2010), to ensembl hg19 (or rheMac8) coding exons and curated RepeatMasker to extract TcGTs

516    with one or more reads spliced between a TE and genic exon for each sample. Second, a custom python

517    program was used to annotate and aggregate the sample level TcGTs into counts per stages (defined

518    in Supplemental Fig. S1B). In brief, for each dataset, a GTF containing all annotated TcGTs was created

519    and TcGTs having their first exon overlapping an annotated gene, or TSS not overlapping a TE were

520    discarded. From this filtered file, TcGTs associated with the same gene and having a TSS within 100bp

521    of each other were aggregated. Finally, for each aggregate, its occurrence per group was computed

522    and a consensus transcript was generated for each TSS aggregate. For each exon of TcGT aggregate,

523    its percentage of occurrence across the different samples was computed and integrated in the

524    consensus if present in more than 30% of the samples the TcGT was detected in. All samples available

525    in both datasets were used regardless of mapped read count.

526    From the resulting master file, additional criteria were applied to determine prenatal, postnatal or

527    continually expressed TcGTs. 1. Only TcGTs that were present in at least 20% of prenatal, postnatal or

528    20% of both pre and postnatal samples (continual) were kept for each dataset. 2. To ensure TcGTs

529    were robustly detectable in the different datasets, TcGT files were merged based on the same TSS TE

530    and associated gene name. 3. TcGTs were required to exhibit the same temporal transcriptional

531    behaviour in both datasets. I.E a 2 fold change in TcGT detection pre vs postnatal and vice versa or a

532    lower fold change in both datasets (continual). This resulted in the 480 robustly detectable temporal

24

533    TcGTs in Fig. 3A and Supplemental Table 9. These TcGTs were further filtered for strong promoter

534    regions using a Bedtools intersect of the 200bp up and downstream of the TcGT TSS with FANTOM5

535    CAGE-seq (Forrest et al. 2014) and BOCA neuronal and non-neuronal consensus ATAC-seq peak bed

536    files (Fullard et al. 2018). TcGT TSS loci were also intersected with ENSEMBL (GRCh37.p13)

537    transcriptional start sites to determine non-annotated transcripts.

**Protein product prediction**

539     DNA sequences were retrieved for each TcGTs consensus and protein products were derived from the

540    longest ORF in the three reading frames using biopython (Cock et al. 2009). The resulting translation

541    products were aligned against the protein sequence of the most similar cognate gene isoforms (exons

542    intersect between TcGTs and each gene isoform) and classified into several categories. Proteins with

543    no alignment for any isoform were classified as out-of-frame, therefore not clear or not aligned. In-

544    frame peptides were further classified according to their N-terminal modifications: Normal, TcGT ORF

545    peptides align perfectly with cognate ORF peptides; N-add, TcGT ORF peptides encode novel in-frame

546    N-terminal amino acids followed by the full length cognate protein sequence; N-truncated, TcGT ORF

547    peptides lack parts of the cognate N-terminal protein sequence and might contain novel in-frame N-

548    terminal amino acids. TcGTs that we could not clearly classify were grouped in the 'other' category,

549    such as TcGTs including C-terminal modifications. If the classification was ambiguous for different

550    protein isoforms, the normal category was always privileged.

**TE and KZFP age estimation**

552    TE subfamily ages were downloaded from DFAM (Hubley et al. 2016). To compare KZFP ages we

553    developed a score we called Complete Alignment of Zinc Finger (CAZF) (as we described in Thorball et

554    al. 2020), which rely on the alignments of zinc finger domains, using only the four amino acid

555    presumably touching DNA. Briefly, alignment scores made with BLOSUM80 matrix were used,

556    normalised by the 'perfect' alignment score (alignment against itself) and by the length of the

557    alignment. To compute an age for KZFPs, we relied on inter-species clusters of KZFPs made with CAZF

558    score. KZFPs with CAZF>0.5 were clustered together, using a bottom-up approach. The divergence time

559    between human and the farthest species present in the cluster was used as the age of individual KZFPs

560    in the cluster. Multiz alignments for L2:DDRGK1 locus were extracted from the UCSC genome browser

561    **Cell culture**

562    Human embryonic kidney 293T (HEK293T) cells and SH-SY-5Y neuroblastoma cells were cultured in

563    DMEM supplemented with 10% fetal calf serum and 1% penicillin/streptomycin.

564    **Transfection**

565    Transient transfection of HEK293T cells was performed with FuGENE HD (Promega) as per the

566    manufacturer's recommendation. Cells were harvested 48 hours after transfection for either RNA

567    extraction or immunofluorescence.

568    **CRISPRa**

569    The SP-dCas9-VPR (Addgene 63798) (Chavez et al. 2015) and the gRNA cloning vector (Addgene 41824)

570    (Mali et al. 2013) were gifts from George Church. gRNAs were designed with CRISPOR (Concordet and

571    Haeussler 2018), using input DNA sequence 50 to 300bp upstream of the TE resident CAGE peak and

572    the most 5' location of RNA-seq reads mapping to the TcGT TE TSS loci. Multiple gRNAs were selected

573    for each TcGT to control for gRNA specific effects and increase experimental robustness. gRNA

574    oligonucleotides were synthesised (Microsynth) with the recommended overhangs (Supplemental

575    Table 10) for integration into the gRNA cloning vector (Mali et al. 2013). gRNA oligonucleotides were

576    annealed and extended using Phusion High Fidelity DNA polymerase master mix (NEB) with thermal

577    cycling conditions of 98°C two minutes (1x), 98°C 10 seconds + 72°C 20 seconds (3x) and 72°C for five

578    minutes. 10μg of SP-dCas9-VPR was digested with Af1II (NEB) in CutSmart buffer for two hours at 37°C,

579    followed by gel electrophoresis and purification of the correct sized band of linearised plasmid with

580    E.Z.N.A Gel Extraction Kit (Omega Bio-tek). The resulting linearised plasmid and double stranded

581  oligonucleotides were ligated using Gibson Assembly Master Mix (NEB) as per manufacturer's

582  recommendations. The resulting gRNA containing plasmid was transformed into HB101 chemically

583  competent *E.coli,* with colonies containing the transformed plasmid selected on agar plates containing

584  kanamycin, followed by colony picking for growth in kanamycin agar broth followed by GeneJET

585  Plasmid Miniprep (ThermoFisher). gRNA plasmids were Sanger sequenced to detect the correct

586  insertion of specific gRNA sequences. 300,000 HEK293T cells were seeded per well of a six well plate.

587  24 hours later, co-transfection was performed with 1μg each of SP-dCas9-VPR and TcGT targeting gRNA

588  containing gRNA cloning vector. SP-dCas9-VPR or empty gRNA cloning vector alone were transfected

589  as non-targeting controls. Cells were harvested for RNA 48 hours post-transfection.

590  **RT-PCR and qRT-PCR**

591  Primers to detect TcGTs were designed with Primer3 (Untergasser et al. 2012) by inputting DNA

592  sequences covering and flanking the splice junction between the TE and genic exon (Supplemental

593  Table 10). One primer was required to be present in the TE sequence where RNA-seq reads were

594  detected downstream of a CAGE-peak, whilst the other was present in the first or second genic exon.

595  BLAT (Kent 2002) of primer sequences against the human genome ensured only uniquely mapping

596  primers were used. RNA was extracted from cells using the NucleoSpin RNA mini kit (Macherey-Nagel)

597  with on-column deoxyribonuclease treatment. 1ug RNA was used in the cDNA synthesis reaction with

598  the Maxima H minus cDNA synthesis master mix (ThermoFisher) and RT-PCR was performed with

599  Phusion High Fidelity DNA polymerase master mix (NEB) each with the manufacturer recommended

600  PCR thermal cycles, on a 9800 Fast Thermal Cycler (Applied Bioscience). PCR products were visualised

601  by 1.5% agarose gel electrophoresis stained with SYBR Safe DNA gel stain (ThermoFisher) and imaged

602  with a BioDoc-It imaging system (UVP). Bands of the correct size were excised, gel purified with E.Z.N.A

603  Gel Extraction Kit (Omega Bio-tek) and Sanger sequenced using primers used for PCR. The correct PCR

604  product was confirmed using BLAT (Kent 2002) of the Sanger sequencing results against the human

605  genome. qRT-PCR was performed with PowerUp SYBR Green Master Mix on a QuantStudio 6 Flex Real-

606    Time PCR system. The standard curve method was used to quantify expression normalised to *BETA*

607    *ACTIN* with no amplification in the no reverse transcriptase control.

608    **Cloning of WT:DDRGK1 and L2:DDRGK1**

609    The WT:DDRGK1 cDNA Clone (Genbank accession:HQ448262 ImageID:100071664) was obtained from

610    the ORFeome Collaboration (http://www.orfeomecollaboration.org/) in the *pENTR223* vector without

611    a stop codon. The L2:DDRGK1 sequence was PCR amplified with Phusion High Fidelity DNA polymerase

612    master mix (NEB), using cDNA generated in the L2:DDRGK1 CRISPRa experiment with gRNA 1. This

613    ensured the *bona fide* L2 driven transcript was cloned. Cloning primers used are shown in

614    Supplemental Table 10, with the forward primer containing a CACC Kozak sequence and the reverse

615    primer omitting the stop codon. Thermal cycling conditions were 98°C 30 seconds (1x), 98°C 10 seconds

616    + 60°C 15 seconds + 72°C 15 seconds (35x) and 72°C for 10 minutes. A 466bp PCR fragment was

617    extracted after agarose gel electrophoresis, purified with E.Z.N.A Gel Extraction Kit (Omega Bio-tek),

618    transformed into chemically competent HB101 E.coli, colonies picked and mini-prepped. WT:DDRGK1

619    and L2:DDRGK1 in the *pENTR* vectors were then shuttled into pTRE-3HA (Imbeault et al. 2017) with the

620    Gateway LR Clonase II Enzyme mix (ThermoFisher) as per manufacturer's instructions. pTRE-3HA

621    produces proteins with three C-terminal HA tags in a doxycyclin-dependent manner.

622    **Cellular fractionation**

623    Approximately 400,000 HEK293T cells in different wells of a 6 well plate were transfected with either

624    pTRE-WT:DDRGK1-HA or pTRE-L2:DDRGK1-HA whose expression was induced for 48 hours by adding

625    1µg/ml doxycycline to the media. After 48 hours wells were washed with 1ml ice cold PBS and cells

626    were scraped and transferred to Eppendorf tubes on the second wash. After centrifugation at 300rcf

627    for five minutes at 4°C, PBS was aspirated, cells re-suspended in 400µl ice-cold cytoplasmic isolation

628    buffer (10mM KOAc, 2mM MgOAC, 20mM HEPES pH7.5, 0.5mM DTT, 0.015% digitonin) and

629    centrifuged at 900rcf for five minutes at 4°C. Supernatant was collected as the cytoplasmic fraction

630    and the remaining pellet was re-suspended in 400µl of membrane isolation buffer (10mM HEPES,

28

631    10mM KCl, 0.1mM EDTA pH8, 1mM DTT, 0.5% Triton X-100, 100mM NaF), then centrifuged for 10

632    minutes at 900rcf at 4°C to pellet nuclei with the supernatant collected as the membrane fraction.

633    Pelleted nuclei were resuspended in 400μl of lysis buffer (1% NP-40, 500mM Tris-HCL pH8, 0.05% SDS,

634    20mM EDTA, 10mM NaF, 20mM benzamidine) for 10 minutes on ice, centrifuged for 10 minutes at

635    900rcf at 4°C and the supernatant collected as the nuclear fraction. 100μl of 4x NuPAGE LDS sample

636    buffer (ThermoFisher) was added to the 400μl cellular fractions and samples boiled at 95°C for five

637    minutes.

**Western blot**

639    20μl of each cellular fraction was used for SDS-PAGE in a NuPAGE 4-12% Bis-TRIS gel and MOPs running

640    buffer (ThermoFisher). For subcellular fraction marker proteins, the same amount of lysate was added

641    from each sample but for the HA blot, pTRE-WT:DDRGK1-HA samples were diluted 1:50 due to high

642    over-expression levels compared to pTRE-L2:DDRGK1-HA. Proteins were transferred to a nitrocellulose

643    membrane using an iBLOT 2 dry blotting system (ThermoFisher) and analysed by immunoblotting using

644    CALNEXIN (Bethyl A303-696A, 1:2000), LAMIN B1 (Abcam ab16048, 1:1000), β TUBULIN (Sigma T4026,

645    1:1000), HA-HRP conjugated (Roche 12013819001, 1:2000). HRP-conjugated anti-mouse (GE

646    Healthcare NA931V, 1:10000) and HRP-conjugated anti-rabbit (Santa Cruz sc-2004 1:5000) antibodies

647    were used where appropriate and the blot was visualised using the Fusion SOLO S (Vilber).

**Immunofluorescence**

649    HEK293T cells were plated on glass coverslips and immunofluorescence was performed as previously

650    described (Helleboid et al. 2019) 48 hours post-transfection and expression induction with 1μg/ml

651    doxycycline for pTRE-WT:DDRGK1-HA or pTRE-L2:DDRGK1-HA. Once 70% confluent, cells were washed

652    three times with PBS, fixed in ice-cold methanol for 20 minutes at -20°C then washed three more times

653    with PBS. Cells were blocked with 1% BSA/PBS for 30 minutes and then incubated with antibodies for

654    HA.11 (BioLegend MMS-101P, 1:2000) and BIP (Abcam ab21685, 1:1000) in 1% BSA/PBS for one hour.

655    Three washes with PBS were performed, followed by incubation with anti-mouse and anti-rabbit Alexa

656     488 or 568 (ThermoFisher 1:800) for one hour. DAPI (1:10000) was added in the last 10 minutes of

657     incubation, samples washed three times with PBS and coverslips mounted on slides with ProLong Gold

658     Antifade Mountant (ThermoFisher). Images were acquired on a SP8 upright confocal microscope

659     (Leica) and processed in ImageJ.

660     **Data Access**

661     No additional high throughput data was generated in this study.

662     **Acknowledgements**

663     We thank all members of the Trono Lab for helpful and insightful discussions, along with Samuel

664     Corless and Nezha Benabdallah for critical reading of the manuscript.

665     **Funding**

669     **Author Contributions**

670     C.P. and D.T. conceived the study, interpreted the data, and wrote the manuscript. C.P. performed

671     bioinformatics analyses and all experiments. J.D. and S.S. developed key code and performed

672     bioinformatics analyses. S.D. performed the GTEX TcGT analysis and *in silico* translation of TcGTs. A.C.

673     performed the KZFP aging analysis and determined KZFP TE subfamily targets. E.P. contributed to

674     bioinformatics tools and code. All authors reviewed the manuscript.

675     **Disclosure Declaration**

676     The authors declare they have no competing interests.

677

**References**

Adam SA, Schnell O, Pöschl J, Eigenbrod S, Kretzschmar HA, Tonn J-C, Schüller U. 2012. ALDH1A1 is a Marker of Astrocytic Differentiation during Brain Development and Correlates with Better Survival in Glioblastoma Patients. *Brain Pathol* **22**: 788–797. doi:10.1111/j.1750-3639.2012.00592.x

Attig J, Young GR, Hosie L, Perkins D, Encheva-Yokoya V, Stoye JP, Snijders AP, Ternette N, Kassiotis G. 2019. LTR retroelement expansion of the human cancer transcriptome and immunopeptidome revealed by de novo transcript assembly. *Genome Res* **29**: 1578–1590. doi:10.1101/gr.248922.119

Bashir R, Britton S, Strachan T, Keers S, Vafiadaki E, Lako M, Richard I, Marchand S, Bourg N, Argov Z, et al. 1998. A gene related to Caenorhabditis elegans spermatogenesis factor fer-1 is mutated in limb-girdle muscular dystrophy type 2B. *Nat Genet* **20**: 37–42. doi: 10.1038/1689

Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple Testing. *J R Stat Soc* **57**: 289–300. doi:10.1111/j.2517-6161.1995.tb02031.x

Bourque G, Leong B, Vega VB, Chen X, Lee YL, Srinivasan KG, Chew J-L, Ruan Y, Wei C-L, Ng HH, et al. 2008. Evolution of the mammalian transcription factor binding repertoire via transposable elements. *Genome Res* **18**: 1752–1762. doi:10.1101/gr.080663.108

Brattås PL, Jönsson ME, Fasching L, Nelander Wahlestedt J, Shahsavani M, Falk R, Falk A, Jern P, Parmar M, Jakobsson J. 2017. TRIM28 Controls a Gene Regulatory Network Based on Endogenous Retroviruses in Human Neural Progenitor Cells. *Cell Rep* **18**: 1–11. doi:10.1016/j.celrep.2016.12.010

Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran J V, Kazazian HH. 2003. Hot L1s account for the bulk of retrotransposition in the human population. *Proc Natl Acad Sci U S A* **100**: 5280–5285. doi:10.1073/pnas.0831042100

Cardoso-Moreira M, Halbert J, Valloton D, Velten B, Chen C, Shao Y, Liechti A, Ascenção K, Rummel C, Ovchinnikova S, et al. 2019. Gene expression across mammalian organ development. *Nature* **571**: 505–509. doi:10.1038/s41586-019-1338-5

Chang D, Nalls MA, Hallgrímsdóttir IB, Hunkapiller J, van der Brug M, Cai F, Kerchner GA, Ayalon G, Bingol B, Sheng M, et al. 2017. A meta-analysis of genome-wide association studies identifies 17 new Parkinson's disease risk loci. *Nat Genet* **49**: 1511–1516. doi:10.1038/ng.3955

Chavez A, Scheiman J, Vora S, Pruitt BW, Tuttle M, P R Iyer E, Lin S, Kiani S, Guzman CD, Wiegand DJ, et al. 2015. Highly efficient Cas9-mediated transcriptional programming. *Nat Methods* **12**: 326–328. doi:10.1038/nmeth.3312

Chen W, Schwalie PC, Pankevich E V., Gubelmann C, Raghav SK, Dainese R, Cassano M, Imbeault M, Jang SM, Russeil J, et al. 2019. ZFP30 promotes adipogenesis through the KAP1-mediated activation of a retrotransposon-derived Pparg2 enhancer. *Nat Commun* **10**: 1809. doi:10.1038/s41467-019-09803-9

Chuong EB, Elde NC, Feschotte C. 2017. Regulatory activities of transposable elements: from conflicts to benefits. *Nat Rev Genet* **18**: 71–86. doi:10.1038/nrg.2016.139

Chuong EB, Elde NC, Feschotte C. 2016. Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science (80- )* **351**: 1083–1087. doi:10.1126/science.aad5497

720  Chuong EB, Rumi M a K, Soares MJ, Baker JC. 2013. Endogenous retroviruses function as species-
721      specific enhancer elements in the placenta. *Nat Genet* **45**: 325–329. doi:10.1038/ng.2553

722  Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, Friedberg I, Hamelryck T, Kauff F, Wilczynski
723      B, et al. 2009. Biopython: freely available Python tools for computational molecular biology and
724      bioinformatics. *Bioinformatics* **25**: 1422–1423. doi:10.1093/bioinformatics/btp163

725  Cohen CJ, Lock WM, Mager DL. 2009. Endogenous retroviral LTRs as promoters for human genes: A
726      critical assessment. *Gene* **448**: 105–114. doi:10.1016/j.gene.2009.06.020

727  Concordet J-P, Haeussler M. 2018. CRISPOR: intuitive guide selection for CRISPR/Cas9 genome editing
728      experiments and screens. *Nucleic Acids Res* **46**: W242–W245. doi:10.1093/nar/gky354

729  Coufal NG, Garcia-Perez JL, Peng GE, Yeo GW, Mu Y, Lovci MT, Morell M, O'Shea KS, Moran J V., Gage
730      FH. 2009. L1 retrotransposition in human neural progenitor cells. *Nature* **460**: 1127–1131.
731      doi:10.1038/nature08248

732  Cui L, Gong X, Tang Y, Kong L, Chang M, Geng H, Xu K, Wang F. 2016. Relationship between the LHPP
733      Gene Polymorphism and Resting-State Brain Activity in Major Depressive Disorder. *Neural Plast*
734      **2016**: 1–8. doi:10.1155/2016/9162590

735  Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D, Merkel A, Knowles
736      DG, et al. 2012. The GENCODE v7 catalog of human long noncoding RNAs: Analysis of their gene
737      structure, evolution, and expression. *Genome Res* **22**: 1775–1789. doi:10.1101/gr.132159.111

738  Dietmann S, Keogh MJ, Tang WW, Magnusdottir E, Kobayashi T, Chinnery P, Surani A. 2020.
739      Transposable elements resistant to epigenetic resetting in the human germline are epigenetic
740      hotspots for development and disease. *bioRxiv*. doi:10.1101/2020.03.19.998930

741  Ecco G, Cassano M, Kauzlaric A, Duc J, Coluccio A, Offner S, Imbeault M, Rowe HM, Turelli P, Trono D.
742      2016. Transposable Elements and Their KRAB-ZFP Controllers Regulate Gene Expression in Adult
743      Tissues. *Dev Cell* **36**: 611–623. doi:10.1016/j.devcel.2016.02.024

744  Ecco G, Imbeault M, Trono D. 2017. KRAB zinc finger proteins. *Development* **144**: 2719–2729.
745      doi:10.1242/dev.132605

746  Erwin JA, Paquola ACM, Singer T, Gallina I, Novotny M, Quayle C, Bedrosian TA, Alves FIA, Butcher CR,
747      Herdy JR, et al. 2016. L1-associated genomic regions are deleted in somatic cells of the healthy
748      human brain. *Nat Neurosci* **19**: 1583–1591. doi:10.1038/nn.4388

749  Farmiloe G, Lodewijk GA, Robben SF, van Bree EJ, Jacobs FMJ. 2020. Widespread correlation of KRAB
750      zinc finger protein binding with brain-developmental gene expression patterns. *Philos Trans R*
751      *Soc B Biol Sci* **375**: 20190333. doi:10.1098/rstb.2019.0333

752  Forrest ARR, Kawaji H, Rehli M, Baillie JK, De Hoon MJL, Haberle V, Lassmann T, Kulakovskiy I V., Lizio
753      M, Itoh M, et al. 2014. A promoter-level mammalian expression atlas. *Nature* **507**: 462–470.
754      doi:10.1038/nature13182

755  Fullard JF, Hauberg ME, Bendl J, Egervari G, Cirnaru M-D, Reach SM, Motl J, Ehrlich ME, Hurd YL,
756      Roussos P. 2018. An atlas of chromatin accessibility in the adult human brain. *Genome Res* **28**:
757      1243–1252. doi:10.1101/gr.232488.117

758  Garcia-Perez JL, Widmann TJ, Adams IR. 2016. The impact of transposable elements on mammalian
759      development. *Development* **143**: 4101–4114. doi:10.1242/dev.132639

760  Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge Y, Gentry J, et
761      al. 2004. Bioconductor: open software development for computational biology and
762      bioinformatics. *Genome Biol* **5**: R80. doi:10.1186/gb-2004-5-10-r80

763  Guffanti G, Bartlett A, Klengel T, Klengel C, Hunter R, Glinsky G, Macciardi F. 2018. Novel Bioinformatics
764      Approach Identifies Transcriptional Profiles of Lineage-Specific Transposable Elements at Distinct
765      Loci in the Human Dorsolateral Prefrontal Cortex ed. I. Arkhipova. *Mol Biol Evol* **35**: 2435–2453.
766      doi:10.1093/molbev/msy143/5056710

767  Helleboid P, Heusel M, Duc J, Piot C, Thorball CW, Coluccio A, Pontis J, Imbeault M, Turelli P, Aebersold
768      R, et al. 2019. The interactome of <scp>KRAB</scp> zinc finger proteins reveals the evolutionary
769      history of their functional diversification. *EMBO J* **38**: 1–16. doi:10.15252/embj.2018101220

770  Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, Smit AFA, Wheeler TJ. 2016. The Dfam
771      database of repetitive DNA families. *Nucleic Acids Res* **44**: D81–D89. doi:10.1093/nar/gkv1272

772  Huntley S, Baggott DM, Hamilton AT, Tran-Gyamfi M, Yang S, Kim J, Gordon L, Branscomb E, Stubbs L.
773      2006. A comprehensive catalog of human KRAB-associated zinc finger genes: Insights into the
774      evolutionary history of a large family of transcriptional repressors. *Genome Res* **16**: 669–677.
775      doi:10.1101/gr.4842106

776  Imbeault M, Helleboid P-Y, Trono D. 2017. KRAB zinc-finger proteins contribute to the evolution of
777      gene regulatory networks. *Nature* **543**: 550–554. doi:10.1038/nature21683

778  Ito J, Kimura I, Soper A, Coudray A, Koyanagi Y, Nakaoka H, Inoue I, Turelli P, Trono D, Sato K. 2020.
779      Endogenous retroviruses drive KRAB zinc-finger protein family expression for tumor suppression.
780      *Sci Adv* **6**: eabc3020. doi:10.1126/sciadv.abc3020

781  Jacobs FMJ, Greenberg D, Nguyen N, Haeussler M, Ewing AD, Katzman S, Paten B, Salama SR, Haussler
782      D. 2014. An evolutionary arms race between KRAB zinc-finger genes ZNF91/93 and SVA/L1
783      retrotransposons. *Nature* **516**: 242–245. doi:10.1038/nature13760

784  Jang HS, Shah NM, Du AY, Dailey ZZ, Pehrsson EC, Godoy PM, Zhang D, Li D, Xing X, Kim S, et al. 2019.
785      Transposable elements drive widespread expression of oncogenes in human cancers. *Nat Genet*
786      **51**: 611–617. doi:10.1038/s41588-019-0373-3

787  Jönsson ME, Garza R, Johansson PA, Jakobsson J. 2020. Transposable Elements: A Common Feature of
788      Neurodevelopmental and Neurodegenerative Disorders. *Trends Genet* **36**: 610–623.
789      doi:10.1016/j.tig.2020.05.004

790  Jönsson ME, Ludvik Brattås P, Gustafsson C, Petri R, Yudovich D, Pircs K, Verschuere S, Madsen S,
791      Hansson J, Larsson J, et al. 2019. Activation of neuronal genes via LINE-1 elements upon global
792      DNA demethylation in human neural progenitors. *Nat Commun* **10**: 3182. doi:10.1038/s41467-
793      019-11150-8

794  Kang HJ, Kawasawa YI, Cheng F, Zhu Y, Xu X, Li M, Sousa AMM, Pletikos M, Meyer KA, Sedmak G, et al.
795      2011. Spatio-temporal transcriptome of the human brain. *Nature* **478**: 483–489.
796      doi:10.1038/nature10523

797  Katz Y, Wang ET, Silterra J, Schwartz S, Wong B, Thorvaldsdóttir H, Robinson JT, Mesirov JP, Airoldi EM,
798      Burge CB. 2015. Quantitative visualization of alternative exon expression from RNA-seq data.
799      *Bioinformatics* **31**: 2400–2402. doi:10.1093/bioinformatics/btv034

800  Keil JM, Qalieh A, Kwan KY. 2018. Brain Transcriptome Databases: A User's Guide. *J Neurosci* **38**: 2399–
801      2412. doi:10.1523/JNEUROSCI.1930-17.2018

802  Kelley DR, Rinn JL. 2012. Transposable elements reveal a stem cell specific class of long noncoding
803      RNAs. *Genome Biol* **13**: R107. doi:10.1186/gb-2012-13-11-r107

804  Kent WJ. 2002. BLAT---The BLAST-Like Alignment Tool. *Genome Res* **12**: 656–664.
805      doi:10.1101/gr.229202

806  Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory requirements.
807       *Nat Methods* **12**: 357–360. doi:10.1038/nmeth.3317

808  Kovaka S, Zimin A V., Pertea GM, Razaghi R, Salzberg SL, Pertea M. 2019. Transcriptome assembly from
809       long-read RNA-seq alignments with StringTie2. *Genome Biol* **20**: 278. doi:10.1186/s13059-019-
810       1910-1

811  Kumanogoh A, Kikutani H. 2004. Biological functions and signaling of a transmembrane semaphorin,
812       CD100/Sema4D. *Cell Mol Life Sci* **61**: 292–300. doi:10.1007/s00018-003-3257-7

813  Lambert SA, Jolma A, Campitelli LF, Das PK, Yin Y, Albu M, Chen X, Taipale J, Hughes TR, Weirauch MT.
814       2018. The Human Transcription Factors. *Cell* **172**: 650–665. doi:10.1016/j.cell.2018.01.029

815  Lanciano S, Cristofari G. 2020. Measuring and interpreting transposable element expression. *Nat Rev
816       Genet* **21**: 721–736. doi:10.1038/s41576-020-0251-y

817  Law CW, Chen Y, Shi W, Smyth GK. 2014. voom: precision weights unlock linear model analysis tools
818       for RNA-seq read counts. *Genome Biol* **15**: R29. doi:10.1186/gb-2014-15-2-r29

819  Li M, Santpere G, Imamura Kawasawa Y, Evgrafov O V., Gulden FO, Pochareddy S, Sunkin SM, Li Z, Shin
820       Y, Zhu Y, et al. 2018. Integrative functional genomic analysis of human brain development and
821       neuropsychiatric risks. *Science (80- )* **362**: eaat7615. doi:10.1126/science.aat7615

822  Li W, Lee M-H, Henderson L, Tyagi R, Bachani M, Steiner J, Campanac E, Hoffman DA, von Geldern G,
823       Johnson K, et al. 2015. Human endogenous retrovirus-K contributes to motor neuron disease. *Sci
824       Transl Med* **7**: 307ra153. doi:10.1126/scitranslmed.aac8201

825  Liang JR, Lingeman E, Luong T, Ahmed S, Muhar M, Nguyen T, Olzmann JA, Corn JE. 2020. A Genome-
826       wide ER-phagy Screen Highlights Key Roles of Mitochondrial Metabolism and ER-Resident
827       UFMylation. *Cell* **180**: 1160-1177.e20. doi:10.1016/j.cell.2020.02.017

828  Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose program for assigning
829       sequence    reads    to    genomic    features. *Bioinformatics*    **30**:    923–930.
830       doi:10.1093/bioinformatics/btt656

831  Linker SB, Randolph-Moore L, Kottilil K, Qiu F, Jaeger BN, Barron J, Gage FH. 2020. Identification of
832       bona fide B2 SINE retrotransposon transcription through single-nucleus RNA-seq of the mouse
833       hippocampus. *Genome Res* **30**: 1643–1654. doi:10.1101/gr.262196.120

834  Liu J, Aoki M, Illa I, Wu C, Fardeau M, Angelini C, Serrano C, Urtizberea JA, Hentati F, Hamida M Ben, et
835       al. 1998. Dysferlin, a novel skeletal muscle gene, is mutated in Miyoshi myopathy and limb girdle
836       muscular dystrophy. *Nat Genet* **20**: 31–36. doi:10.1038/1682

837  Liu J, Wang Y, Song L, Zeng L, Yi W, Liu T, Chen H, Wang M, Ju Z, Cong Y-S. 2017. A critical role of
838       DDRGK1 in endoplasmic reticulum homoeostasis via regulation of IRE1α stability. *Nat Commun
839       8*: 14186. doi:10.1038/ncomms14186

840  Lizio M, Harshbarger J, Shimoji H, Severin J, Kasukawa T, Sahin S, Abugessaisa I, Fukuda S, Hori F,
841       Ishikawa-Kato S, et al. 2015. Gateways to the FANTOM5 promoter level mammalian expression
842       atlas. *Genome Biol* **16**: 22. doi:10.1186/s13059-014-0560-6

843  Mali P, Yang L, Esvelt KM, Aach J, Guell M, DiCarlo JE, Norville JE, Church GM. 2013. RNA-Guided Human
844       Genome Engineering via Cas9. *Science (80- )* **339**: 823–826. doi:10.1126/science.1232033

845  McCormack K, Connor JX, Zhou L, Ho LL, Ganetzky B, Chiu S-Y, Messing A. 2002. Genetic Analysis of the
846       Mammalian K + Channel β Subunit Kvβ2 ( Kcnab2 ). *J Biol Chem* **277**: 13219–13228.
847       doi:10.1074/jbc.M111465200

848    Melé M, Ferreira PG, Reverter F, DeLuca DS, Monlong J, Sammeth M, Young TR, Goldmann JM,
849        Pervouchine DD, Sullivan TJ, et al. 2015. The human transcriptome across tissues and individuals.
850        *Science (80- )* **348**: 660–665. doi:10.1126/science.aaa0355

851    Miao B, Fu S, Lyu C, Gontarz P, Wang T, Zhang B. 2020. Tissue-specific usage of transposable element-
852        derived promoters in mouse development. *Genome Biol* **21**: 255. doi:10.1186/s13059-020-
853        02164-3

854    Miller JA, Ding S-L, Sunkin SM, Smith KA, Ng L, Szafer A, Ebbert A, Riley ZL, Royall JJ, Aiona K, et al. 2014.
855        Transcriptional landscape of the prenatal human brain. *Nature* **508**: 199–206. doi:
856        10.1038/nature13185

857    Muotri AR, Chu VT, Marchetto MCN, Deng W, Moran J V., Gage FH. 2005. Somatic mosaicism in
858        neuronal precursor cells mediated by L1 retrotransposition. *Nature* **435**: 903–910. doi:
859        10.1038/nature03663

860    Muotri AR, Marchetto MCN, Coufal NG, Oefner R, Yeo G, Nakashima K, Gage FH. 2010. L1
861        retrotransposition in neurons is modulated by MeCP2. *Nature* **468**: 443–446.
862        doi:10.1038/nature09544

863    Najafabadi HS, Mnaimneh S, Schmitges FW, Garton M, Lam KN, Yang A, Albu M, Weirauch MT,
864        Radovani E, Kim PM, et al. 2015. C2H2 zinc finger proteins greatly expand the human regulatory
865        lexicon. *Nat Biotechnol* **33**: 555–562. doi:10.1038/nbt.3128

866    Nalls MA, Pankratz N, Lill CM, Do CB, Hernandez DG, Saad M, DeStefano AL, Kara E, Bras J, Sharma M,
867        et al. 2014. Large-scale meta-analysis of genome-wide association data identifies six new risk loci
868        for Parkinson's disease. *Nat Genet* **46**: 989–993. doi:10.1038/ng.3043

869    Neff CD, Abkevich V, Packer JCL, Chen Y, Potter J, Riley R, Davenport C, DeGrado Warren J, Jammulapati
870        S, Bhathena A, et al. 2009. Evidence for HTR1A and LHPP as interacting genetic risk factors in
871        major depression. *Mol Psychiatry* **14**: 621–630. doi:10.1038/mp.2008.8

872    Nowick K, Gernat T, Almaas E, Stubbs L. 2009. Differences in human and chimpanzee gene expression
873        patterns define an evolving network of transcription factors in brain. *Proc Natl Acad Sci* **106**:
874        22358–22363. doi:10.1073/pnas.0911376106

875    Pontis J, Planet E, Offner S, Turelli P, Duc J, Coudray A, Theunissen TW, Jaenisch R, Trono D. 2019.
876        Hominoid-Specific Transposable Elements and KZFPs Facilitate Human Embryonic Genome
877        Activation and Control Transcription in Naive Human ESCs. *Cell Stem Cell* **24**: 724-735.e5.
878        doi:10.1016/j.stem.2019.03.012

879    Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
880        *Bioinformatics* **26**: 841–842. doi:10.1093/bioinformatics/btq033

881    Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP. 2011.
882        Integrative genomics viewer. *Nat Biotechnol* **29**: 24–26. doi:10.1038/nbt.1754

883    Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, Snyder MP, Wang T. 2014. Widespread contribution
884        of transposable elements to the innovation of gene regulatory networks. *Genome Res* **24**: 1963–
885        1976. doi:doi/10.1101/gr.168872.113

886    Takahashi N, Coluccio A, Thorball CW, Planet E, Shi H, Offner S, Turelli P, Imbeault M, Ferguson-Smith
887        AC, Trono D. 2019. ZNF445 is a primary regulator of genomic imprinting. Genes Dev 33: 49–54.
888        doi:10.1101/gad.320069.118

889    Tam OH, Ostrow LW, Gale Hammell M. 2019. Diseases of the nERVous system: retrotransposon activity
890        in neurodegenerative disease. *Mob DNA* **10**: 32. doi:10.1186/s13100-019-0176-1

891  Theunissen TW, Friedli M, He Y, Planet E, O'Neil RC, Markoulaki S, Pontis J, Wang H, Iouranova A,
892      Imbeault M, et al. 2016. Molecular Criteria for Defining the Naive Human Pluripotent State. *Cell*
893      *Stem Cell* **19**: 502–515. doi:10.1016/j.stem.2016.06.011

894  Thorball CW, Planet E, de Tribolet-Hardy J, Coudray A, Fellay J, Turelli P, Trono D. 2020. Ongoing
895      evolution of KRAB zinc finger protein-coding genes in modern humans. *bioRxiv*
896      doi:10.1101/2020.09.01.277178

897  Trono D. 2015. Transposable Elements, Polydactyl Proteins, and the Genesis of Human-Specific
898      Transcription Networks. *Cold Spring Harb Symp Quant Biol* **80**: 281–288.
899      doi:10.1101/sqb.2015.80.027573

900  Turelli P, Playfoot C, Grun D, Raclot C, Pontis J, Coudray A, Thorball C, Duc J, Pankevich E V., Deplancke
901      B, et al. 2020. Primate-restricted KRAB zinc finger proteins and target retrotransposons control
902      gene expression in human neurons. *Sci Adv* **6**: eaba3200. doi:10.1126/sciadv.aba3200

903  Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG. 2012. Primer3—new
904      capabilities and interfaces. *Nucleic Acids Res* **40**: e115–e115. doi:10.1093/nar/gks596

905  Upton KR, Gerhardt DJ, Jesuadian JS, Richardson SR, Sánchez-Luque FJ, Bodea GO, Ewing AD, Salvador-
906      Palomeque C, van der Knaap MS, Brennan PM, et al. 2015. Ubiquitous L1 Mosaicism in
907      Hippocampal Neurons. *Cell* **161**: 228–239. doi:10.1016/j.cell.2015.03.026

908  van Kuilenburg AB., Meinsma R, Vreken P, Waterham HR, van Gennip AH. 2000. Identification of a
909      cDNA encoding an isoform of human CTP synthetase. *Biochim Biophys Acta - Gene Struct Expr*
910      **1492**: 548–552. doi:10.1016/S0167-4781(00)00141-X

911  Xu J-H, Wang T, Wang X-G, Wu X-P, Zhao Z-Z, Zhu C-G, Qiu H-L, Xue L, Shao H-J, Guo M-X, et al. 2010.
912      PU.1 can regulate the ZNF300 promoter in APL-derived promyelocytes HL-60. Leuk Res 34: 1636–
913      1646. doi:10.1016/j.leukres.2010.04.009

914  Zhong S, Zhang S, Fan X, Wu Q, Yan L, Dong J, Zhang H, Li L, Sun L, Pan N, et al. 2018. A single-cell RNA-
915      seq survey of the developmental landscape of the human prefrontal cortex. *Nature* **555**: 524–
916      528. doi:10.1038/nature25980

917  Zimmer-Bensch G. 2019. Emerging Roles of Long Non-Coding RNAs as Drivers of Brain Evolution. *Cells*
918      **8**: 1399. doi:10.3390/cells8111399

919

920