1    **Heteromeric RNP assembly at LINEs controls lineage-specific**

2    **RNA processing**

3    Jan Attig[1,2,*,#], Federico Agostini[1,#], Clare Gooding[3], Aarti Singh[2,4], Anob M

4    Chakrabarti[1,5], Nejc Haberman[1,2], Warren Emmett[1,2,5], Christopher WJ Smith[3],

5    Nicholas M Luscombe[1,5,6] and Jernej Ule [1,2]*


6

7    [1] The Francis Crick Institute, Midland Road 1, Kings Cross, London NW1 1AT

8    [2] Department of Molecular Neuroscience, UCL Institute of Neurology, Queen
9    Square, London, WC1N 3BG, UK

10   [3] Department of Biochemistry, University of Cambridge, Tennis Court Road,
11   Cambridge, CB2 1QW, UK

12   [4] Department of Comparative Biomedical Sciences, The Royal Veterinary
13   College, Royal College Street, London NW1 0TU, UK

14   [5] Department of Genetics, Environment and Evolution, UCL Genetics Institute,
15   Gower Street, London WC1E 6BT, UK

16   [6] Okinawa Institute of Science & Technology Graduate University, 1919-1
17   Tancha, Onna-son, Kunigami-gun, Okinawa 904-0495, Japan

18   [#] These authors contributed equally to this work.

19

20   *Correspondence: jan.attig@crick.ac.uk; jernej.ule@crick.ac.uk

21

22   *Lead author:* Jernej Ule

23

24   Keywords: splicing, pre-mRNA processing, LINE repeats, MATR3, PTBP1

27 **ABSTRACT**

28 It is challenging for RNA processing machineries to select exons within long

29 intronic regions. We find that intronic LINE repeat sequences (LINEs)

30 contribute to this selection by recruiting dozens of RNA-binding proteins

31 (RBPs). This includes MATR3, which promotes binding of PTBP1 to

32 multivalent binding sites in LINEs. Both RBPs repress splicing and 3' end

33 processing within and around LINEs, as demonstrated in cultured human

34 cells and mouse brain. Notably, repressive RBPs preferentially bind to

35 evolutionarily young LINEs, which are confined to deep intronic regions.

36 These RBPs insulate both LINEs and surrounding regions from RNA

37 processing. Upon evolutionary divergence, gradual loss of insulation

38 diversifies the roles of LINEs. Older LINEs are located closer to exons, are a

39 common source of tissue-specific exons, and increasingly bind to RBPs that

40 enhance RNA processing. Thus, LINEs are hubs for assembly of repressive

41 RBPs, and contribute to evolution of new, lineage-specific transcripts in

42 mammals.

43

44 **INTRODUCTION**

45 Human introns are replete with sequences that resemble splice sites and
46 polyA sites, creating a demand for mechanisms that can help processing
47 machineries distinguish true from cryptic RNA processing sites (Semlow and
48 Staley, 2012, Sibley et al., 2016). Inappropriate recognition of such sites
49 initiates inclusion of cryptic exons, which can disrupt gene expressing by
50 changing the reading frame, introducing premature stop codons, and
51 decreasing transcript stability. Mutations that activate splicing of cryptic exons
52 therefore cause a number of hereditary human diseases (Vorechovsky, 2010,
53 Sibley et al., 2016). It is well understood that exon definition mechanisms
54 maintain splicing fidelity by combinatorial recognition of 3' and 5' splice sites
55 and exonic enhancer sequences. Moreover, RNA-binding proteins (RBPs)
56 can contribute to splicing fidelity by directly repressing the cryptic splice sites
57 of RNA processing (Reed, 2000). Therefore, mutations disrupting their binding
58 sites can activate cryptic exons and cause disease (Eom et al., 2013,
59 Vorechovsky, 2010). However, the RBPs that promote splicing fidelity by
60 assembling over deep intronic regions, and the elements that coordinate
61 assembly of ribonucleoprotein complexes (RNPs) across introns are yet to be
62 fully identified.

63 The human genome contains over 1.5 million LINE repeats, retrotransposons,
64 many of which are located in introns. The two most common LINE repeat
65 families in mammals are L1 and L2, which first inserted before the split of
66 monotremata and therians approximately 200 million years ago (O'Leary et
67 al., 2013), and new subfamilies have amplified in bursts ever since (Huang et
68 al., 2010). The consensus L1 sequence contains strong cryptic splice sites in
69 both sense and antisense orientation (Belancio et al., 2006, Merkin et al.,
70 2015). The effect of intragenic LINEs on splicing has been studied mainly in
71 the context of pathological conditions, in which intronic LINE insertions disrupt
72 expression of their host gene. For instance, creation of new cryptic LINE-
73 derived exons disrupt expression of the *CYBB* gene in individuals with chronic
74 granulomatous disease (Meischl et al., 2000) and the *DMD* gene in X-linked
75 dilated cardiomyopathy (Yoshida et al., 1998), and a less well characterised

76    intronic LINE insertion disrupts expression of *XRP2* in X-linked retinitis
77    pigmentosa 2 (Schwahn et al., 1998). Yet, only 60-80 LINEs were found
78    capable of retrotransposition in the human genome, and account for all the *de*
79    *novo* LINE insertions observed in human populations or *in vitro* (Beck et al.,
80    2010, Brouha et al., 2003). In the remaining LINEs, mutations have disrupted
81    the ability to retrotranspose, and most are degenerated and truncated
82    compared to the consensus sequence. In many genes, such degenerated
83    LINEs form part of their introns. Several RBPs are known to bind active LINEs
84    and thereby interfere with their retrotransposition (Goodier et al., 2013, Taylor
85    et al., 2013, Goodier et al., 2012), but the RBPs binding intronic LINEs, and
86    the regulatory potential of these LINEs, are poorly understood.

87    Here, we surveyed iCLIP and eCLIP data to identify 28 RBPs that have
88    enriched binding to intragenic LINEs. Notably many of these RBPs, including
89    MATR3 and PTBP1, primarily bind to deep intronic regions i.e. more than
90    >500nt away from any known exon. MATR3 promotes binding of PTBP1 to
91    LINEs at clusters rich in CU-binding motifs and, together, the two RBPs create
92    a repressive environment by blocking the recognition of cryptic polyA-sites
93    and splice sites. Analysis of distinct evolutionary classes of intragenic LINEs
94    demonstrated that repressive RBPs are most enriched on younger, primate-
95    specific L1 elements, which are depleted in the vicinity of exons. In contrast,
96    the binding of repressive RBPs is decreased on evolutionary older LINEs,
97    especially those preserved across mammals, with a concomitant increase in
98    binding of RBPs that are involved in recognition of 3' and 5' splice sites, and
99    polyA sites. These older LINEs are located in closer proximity to exons, and
100   are a source of new mammalian exons, with higher inclusion levels and
101   differential regulation across human tissues. Thus, while most LINEs recruit
102   repressive RBPs to insulate the deep intronic regions from RNA processing,
103   many older LINEs have started to escape from this repression. Hence, LINEs
104   facilitate evolutionary innovations in the emergence of mammal-specific
105   transcripts.

106

107

108 **RESULTS**

109 **LINE-derived sequences recruit dozens of RBPs to deep intronic**
110 **regions.**

111 The primary goal of our study was to identify repressive RBPs that assemble
112 over deep intronic regions in a coordinated manner to distinguish genuine
113 exons from other exon-like sequences that are present within introns.
114 Because many cryptic exons arise from retrotransposable elements, which
115 are pervasive in introns (Smit et al., 1996-2010a, Deininger and Batzer,
116 2002), we hypothesised that retrotransposons might be RNP assembly points
117 in deep intronic regions. We focused on LINEs, because they constitute the
118 largest proportion of intronic sequence, and are generally excluded from
119 mRNAs but not pre-mRNA, as is evident from their presence in nuclear but
120 not cytoplasmic transcriptomes in HeLa, K562 and HepG2 cell lines, and from
121 their depletion in nuclear polyA+ compared to polyA- RNA (Figure 1A). To
122 identify RBPs that bind to L1-derived sequences, we examined iCLIP data for
123 17 RBPs published by our laboratory, and eCLIP data from K562 and HepG2
124 cells for 112 RBPs available from ENCODE (Sloan et al., 2016, Van Nostrand
125 et al., 2017). We ranked these RBPs by the proportion of crosslink events
126 mapping to sense or antisense L1 elements (Figure 1B), as well as by their
127 propensity to bind to deep intronic regions (Figure 1C). Overall, find that RBPs
128 with most enrichment on L1 elements also rank highly for deep intronic
129 binding.

130 While many RBPs show binding in deep intronic regions, MATR3 and PTBP1
131 ranked highest in iCLIP data (Figure 1C). Both RBPs have strong enrichment
132 on L1s: 19% of MATR3 and PTBP1 iCLIP crosslink events were in antisense
133 L1 repeats, which is a strong overrepresentation compared to the median of
134 6.4% across all examined iCLIP data (Figure 1B and SupplTable1). This
135 agrees with a previous study that also found enrichment of PTBP1 iCLIP
136 reads in LINEs compared to a genomic null model (Kelley et al., 2014). In
137 eCLIP data (for is lacking for MATR3), ~10% of PTBP1 crosslink events
138 mapped to antisense L1 elements, compared with ~4% across all RBPs
139 examined. The decreased enrichment in eCLIP compared to iCLIP likely

140 stems from differences in data processing (see Methods for comment).

141 Within our set of iCLIP data, we also found CELF2, ELAVL1 and TARDBP
142 overrepresented at antisense L1s. Nine additional RBPs showed enrichment
143 on L1 elements in eCLIP data: SUGP2, hnRNPM, KHSRP, hnRNPC, FUBP3
144 and SFPQ on antisense L1 repeats, and HLTF, KHDRBS1 and SAFB2 on L1
145 elements in sense. We also examined RBP binding to L2 elements, which are
146 about three times less common in human genome than L1s. Correspondingly,
147 we found that L2s accounted for a smaller proportion of RBP binding sites
148 than L1s. SUGP2, MATR3, PTBP1 and HNRNPK had strongest enrichment in
149 sense L2s, and HNRNPA1, TAF15, HNRNPU and SAFB2 in antisense L2s
150 (Suppl. Table 2).

151 In our analysis we used reads mapping uniquely to the genome, which could
152 miss the reads mapping to highly repetitive sequences. To account for them,
153 we also examined eCLIP RBP binding to different sub-families of LINEs by
154 using the TEtranscripts method (Jin et al., 2015). Median enrichment of LINE
155 subfamilies recapitulated our ranking, with equal enrichment for all of the
156 subfamilies for most of the RBPs and with strongest variation between
157 families seen for HNRNPM and SFPQ (Figure S1A). In total, TEtranscript
158 identified >2-fold enrichment on L1 or L2s for 25 RBPs in the eCLIP data. In
159 conclusion, we find strong enrichment of MATR3 and PTBP1 binding on L1
160 and L2 elements, which appears to be related to their deep intronic binding
161 profiles.

162

163 **MATR3 stabilises PTBP1-RNA interactions to promote L1 binding.**

164 MATR3 and PTBP1 directly interact (Coelho et al., 2015), but it is not known if
165 this affects their endogenous RNA binding. Given that both proteins are
166 enriched on antisense L1s, we wished to understand if their binding sites
167 overlap. Since MATR3 eCLIP data are not yet available, we focused solely on
168 iCLIP analysis. We analysed the five RBPs with most LINE binding in iCLIP
169 and performed unsupervised clustering on the 50 most strongly bound LINEs
170 of each RBP. The strongest correlation on individual LINEs was observed

171    between MATR3 and PTBP1 (pearson coefficient = 0.83), with less overlap
172    shared with CELF2, essentially no overlap with TARDBP-bound elements,
173    and a slight anti-correlation with ELAVL1 occupied loci (Figure S2A).
174    Conversely, MATR3 binding was enriched in proximity of PTBP1 binding
175    peaks and this enrichment is significantly stronger at peaks located within
176    rather than outside of LINEs (p-value < 2.2e-16, Figure S2B). Hence, LINEs
177    appear to act as a platform for simultaneous binding of MATR3 and PTBP1.

178    Next, we examined if MATR3 and PTBP1 affect each other in their binding to
179    LINEs. We performed iCLIP with PTBP1 in HEK293 cells depleted of MATR3,
180    and iCLIP with MATR3 in HEK293 cells depleted of PTBP1 and PTBP2, and
181    in both cases performed control iCLIP from cells transfected with control
182    siRNA (Figure 2A and Figure S2C-D). Efficient siRNA-depletion of MATR3
183    and PTBP1 was validated by Western blot (Figure 2A and Figure S3D).
184    Notably, the amount of RNA crosslinked to PTBP1 was visibly decreased
185    upon MATR3 depletion, as measured by $^{32}$P labelling of immunoprecipiated
186    RNA (Figure 2A; replicates in Figure S2C). This decrease was not a result of
187    decreased abundance of PTBP1 (Figure 2A). We did not observe any
188    decrease in MATR3 crosslinked RNA upon depletion of PTBP1/PTBP2
189    (Figure S2D). This indicates that MATR3 is required for efficient crosslinking
190    of PTBP1 to endogenous transcripts, but not vice versa.

191    To analyse changes in PTBP1 binding upon MATR3 depletion, we identified
192    peaks of PTBP1 crosslinking, focused on those with sufficient coverage, and
193    classified them based on the change in normalised counts between
194    MATR3-depleted    sample    and    control    into    MATR3-dependent,
195    MATR3-independent and remaining peaks (Figure 2B). MATR3-dependent
196    PTBP1 peaks were shorter than MATR3-independent ones, and had more
197    MATR3 binding in their vicinity (Figure 2C and D). As expected, PTBP1
198    binding peaks were highly enriched for CT-tetramers, which is most prominent
199    within the peak, but is also seen over a 200nt region around the peak (Figure
200    2E). Intriguingly, we found that the enrichment over this 200nt region is
201    stronger at the MATR3-dependent PTBP1 peaks compared to the remaining
202    peaks (Figure 2E).  Thus, the MATR3-dependent PTBP1 binding peaks are

8

203　shorter but are composed of the highest density of binding motifs over a 200nt

204　region. This indicates that MATR3 supports the interactions between PTBP1

205　and multivalent RNA binding sites, i.e., those that contain multiple repeats of

206　high-affinity PTBP1 binding motifs

207　Next, we examined enrichment of different categories of peaks within

208　repetitive elements. Importantly, MATR3-dependent PTBP1 binding peaks

209　were most enriched within an antisense L1 elements compared to the

210　remaining peaks (Figure 2F). PTBP1 also binds CT- and T-rich microsatellite

211　repeats (Ling et al., 2016), but this accounts for only ~0.2% of all binding

212　peaks in unperturbed HEK293 cells. While no L1 enrichment is seen for

213　MATR3-independent PTBP1 peaks, they more frequently overlap with the

214　microsatellite repeats. In the reciprocal analysis of MATR3 iCLIP after PTBP1

215　depletion (Figure S2D-F), we found that PTBP1 is not required for MATR3

216　binding to LINE repeats, indicating that MATR3 is recruited to LINEs either

217　through its own specificity, or through interactions with additional RBPs.

218　Our analysis suggested that *in vivo* binding of PTBP1 to L1 repeats is

219　stabilised by MATR3. To confirm that MATR3 directly aids RNA-binding of

220　PTBP1, we purified recombinant PTBP1 (rPTBP1) and MATR3 fragments

221　(rMATR3) that do or do not interact with PTBP1. We previously found a

222　PTBP1 RRM2 Interacting (PRI) motif N-terminal of the MATR3 RRMs is

223　essential for interaction with PTBP1 RRM2 (Coelho et al., 2015). We purified

224　a MATR3 fragment comprising its two RRMs but missing the PRI motif

225　('RRMs'), as well as the RRMs with the PRI motif ('PRI-RRMs') and a mutated

226　sequence with point mutations in the PRI disrupting PTBP1 binding ('mPRI-

227　RRMs'). We designed an *in vitro* synthesised RNA with two MATR3 RNA

228　compete motifs (ATCTT, Ray et al., 2013) as well as small CT-stretches,

229　which could allow multivalent binding of PTBP1 in vicinity. In agreement,

230　rPTBP1 and all rMATR3 fragments bound to this RNA. We found that the

231　non-interacting rMATR3 RRMs fragment competes with PTBP1 for RNA

232　binding at equimolar concentrations (Figure 2G). Unlike the RRM rMATR3

233　fragment, the PRI-RRM rMATR3 did not block crosslinking of PTBP1 to the

234　RNA even at excess molarity of rMATR3. The mPRI-RRMs rMATR3 blocked

9

235    PTBP1 binding, demonstrating the dependency on the interaction motif for

236    formation of a heteromeric PTBP1*MATR3 complex on the RNA. As a next

237    step, we added rMATR3 to HeLa nuclear extracts with endogenous PTBP1,

238    and assayed binding to two RNA probes. We used the probe with two ATCTT

239    motifs (as in Figure 2G), and in addition a probe with six CTCTT motifs (the

240    RNA compete motif for PTBP1), for which we expected stronger binding of

241    PTBP1. Addition of rMATR3 promoted binding of endogenous PTBP1 to the

242    exogenous $ATCTT_2$ RNA through the PRI motif (Figure S2G). On the $CTCTT_6$

243    probe, we observed increased binding of endogenous PTBP1 in absence of

244    recombinant MATR3 compared to the $ATCTT_2$ probe, and PTBP1 was

245    completely displaced by non-interacting rMATR3 RRMs but not by the

246    PRI-RRM rMATR3. Hence, the PRI motif in MATR3 allows formation of a

247    heteromeric complex on substrate RNAs.

248    Together, we show that PTBP1 and MATR3 *in vivo* overlap at antisense

249    L1-derived binding sites, which are rich in multiple repeats of high-affinity

250    PTBP1 binding motifs. We find that the PRI-mediated interaction between

251    MATR3 and PTBP1 is crucial to promote simultaneous binding of both

252    proteins to an RNA, and that MATR3 can recruit PTBP1 to RNA in a

253    sequence-dependent manner *in vitro* and *in vivo*. We conclude MATR3

254    promotes binding of PTBP1 to multivalent binding sites within antisense L1

255    repeats.

256

257    **MATR3 and PTBP1 co-repress exons and poly(A) sites close to LINE**

258    **repeats**

259    To resolve the role that coordinated LINE binding of MATR3 and PTBP1

260    might play in RNA-processing, we first re-analysed our previous splice

261    junction microarray data on repression of alternatively spliced exons by

262    MATR3 (Coelho et al., 2015). Of the 421 exons that were found to be

263    repressed by MATR3, 64 contained at least one of their splice sites within a

264    LINE repeat, and were therefore considered 'LINE-derived exons'; this

265    represents a 2.3-fold enrichment for LINE-derived exons compared to all

266    exons covered in the array design. For PTBP1, we found 50 significantly

267    repressed LINE-derived exons. We evaluated the frequency of L1 and L2

268    repeats in the introns flanking MATR3/PTBP1 repressed exons (Figure 3A,

269    and Figure S3A), and found ~2-fold enrichment for antisense L1 sequence in

270    a window of 2kb, even after removing all LINE-derived exons (Figure S3B).

271    Hence, we found both LINE-derived as well as LINE-proximal exons are

272    overrepresented among exons repressed by MATR3/PTBP1.

273    Next we generated total RNAseq data of cytoplasmic and nuclear RNA from

274    HeLa cells depleted of MATR3 with two independent siRNAs, or

275    PTBP1/PTBP2, or all three factors simultaneously (Figure S3C and Suppl.

276    Table 3). We detected 1,430 LINE-derived exons, each supported by at least

277    one splice-junction read mapping to a LINE element; 1,114 within

278    protein-coding genes and the remaining in long non-coding RNAs. Of the

279    1,430 exons, 858 (~77%) were cryptic, i.e. not annotated in UCSC (Suppl

280    Table 2). LINE sequences can donate either 5' or 3' splice sites, and in ~50%

281    of LINE-derived exons both splice sites were LINE-derived (Figure S3D).

282    Depletion of both MATR3 and PTBP1 led to increased use of 131  (9.1%) of

283    the LINE-derived exons (Figure 3B), with a median increase of more than five

284    fold. Repression of these exons by the two proteins is strongly synergistic,

285    since exon usage increased by about 1.6fold depleting MATR3 or PTBP1

286    individually. We tested changes in inclusion of 16 splicing events significantly

287    regulated by co-depletion of MATR3/PTBP1 by semi-quantitative RT-PCR,

288    including six cryptic and ten annotated exons (Figure S3E), and found

289    synergistic repression for two out of nine LINE-derived exons and two out of

290    five LINE-proximal exons.

291    Since antisense L1 elements are rich in cryptic polyA-signals (Han et al.,

292    2004, Lee et al., 2008), we also produced 3' end sequencing data to

293    investigate if MATR3 and PTBP1 repress poly(A) sites in a LINE-dependent

294    fashion. We used the expressRNA platform (Rot et al., 2017) to find

295    alternative poly(A) site usage. We thereby annotated poly(A) site pairs in

296    5,189 genes, in which two different polyA sites each account for at least 5% of

297    this gene's signal (referred to as pA1 and pA2). Of these, 240 pA-sites

298    originated from a LINE repeat. LINE repeats were enriched at proximal polyA
299    sites repressed by MATR3/PTBP1 for an extended region of ~2kb (Figure
300    3A), reminiscent of the pattern observed on repressed exons. Overall
301    changes in polyA site usage suggest a primarily repressive function of
302    MATR3/PTBP1 (Figure 3C). We split all significantly regulated proximal pA-
303    sites into those within 2kb of a LINE, and those further away from any LINE
304    repeat, and found LINE-proximal sites to be slightly more responsive to
305    MATR3 depletion than LINE-distant sites (Figure 3C). This is mirrored in
306    individual examples (Figure S3 E; for instance in MROH1, an annotated
307    alternative terminal exon with a LINE-derived 3' SS (indicated by red dashed
308    line) is used ~70% in control cells, but entirely replaces the canonical pA site
309    in MATR3 depleted cells accompanied by exonisation of additional sequences
310    and an additional pA site, all from the adjacent LINEs. In PIGN1 (Figure S3F),
311    a stretch of LINE sequences give rise to a cryptic LINE-derived terminal exon,
312    which is used partially upon MATR3 depletion, and fully upon combined
313    depletion of MATR3 and PTBP1, as indicated by loss of all signal on the
314    downstream exon.

315    Metaprofiles of iCLIP binding on regulated splice- and polyA-sites showed
316    increased binding of MATR3 and PTBP1, confirming direct targeting of these
317    loci. LINE-derived exons were enriched in MATR3 and PTBP1 binding
318    compared to non-repeat derived exons (Figure 3D) and those LINE-derived
319    exons most susceptible to depletion of MATR3/PTBP1 showed strongest
320    enrichment. MATR3 binding was extended for ~2kb upstream to ~1kb
321    downstream of the exons, covering both splice sites. At polyA sites, MATR3
322    and PTBP1 binding was enriched at repressed pA1 sites, with extended
323    binding on those pA sites that were proximal to a L1 repeat.

324    We conclude that MATR3/PTBP1 are potent repressors of RNA processing at
325    LINE repeats, thus preventing exonisation of LINEs. Similarly, polyA sites are
326    repressed in vicinity of LINE repeats. Together, this strongly suggests that
327    LINEs are the specificity element in directing MATR3 to alternative exons,
328    linking its function in alternative splicing to its binding on repeat elements, and
329    explaining the lack of a short binding motif of MATR3 *in vivo* we described in

330    the past (Coelho et al., 2015). The binding pattern of PTBP1 on LINE-derived
331    exons was consistent with co-targeting of the same elements by MATR3 and
332    PTBP1. Lastly, changes in abundance of LINE-derived exons suggest
333    functional synergy of MATR3 and PTBP1 on LINE-derived exons but not on
334    non-repeat derived alternative exons, suggesting co-ordinated assembly of
335    both proteins is necessary to ensure complete repression of cryptic exons
336    originating from LINE repeats.

337

338    **LINE-derived exons reduce transcript abundance through NMD**

339    The majority of LINE-derived exons that were detected in MATR3/PTBP1
340    depleted cells are cryptic exons, i.e. not annotated by UCSC or ENSEMBL
341    (Suppl Table 2). Retrotransposon-derived exons, and in particular Alu-exons,
342    are known to be prone to spurious inclusion which generally reduces
343    expression of the host gene through nonsense-mediated mRNA decay (NMD)
344    (Attig et al., 2016). Given the involvement of PTBP1 in repression of LINE-
345    derived exons, we used RNAseq data produced by Ge et al. (Ge et al., 2016)
346    to evaluate if LINE-derived exons detected in HEK293 cells trigger NMD.
347    Depletion of PTBP1 alone produced a marked change in abundance of
348    LINE-derived exons, while depletion of UPF1 alone drastically increased the
349    number of LINE-derived exons detected (Figure S4A); and this number almost
350    doubled after combined depletion of UPF1 and PTBP1. Importantly, genes
351    containing any of those LINE-derived exons showed increased expression in
352    UPF1-depleted cells (Figure S4B). We conclude most LINE-derived exons are
353    cryptic exons that, when included, render the resulting transcript susceptible
354    to NMD.

355

356    **Deletion of an intronic LINE disrupts MATR3-dependent repression of a**
357    **cryptic exon in *ACAD9***

358    To confirm MATR3 and PTPB1 directly repress exons flanked by a LINE
359    within 2kb of their splice site, even if they are not LINE-derived exons, we

360   made a splice reporter plasmid. Among 16 exons for which we validated the
361   role of MATR3 and PTBP1 by RT-PCR (Figure S3E), 6 exons were such
362   LINE-proximal exons, including an exon within intron1 of ACAD9.
363   Endogenous ACAD9 splices efficiently at intron1, with two known splice
364   isoforms of exon1 using different 5' splice sites (here referred to as *exon 1a*
365   and *exon 1b*). We observed a two-fold loss of ACAD9 expression in cells
366   depleted of MATR3, and a three-fold loss of expression in cells co-depleted of
367   MATR3 and PTBP1 (Figure S4 C and D). Intron1 of ACAD9 contains three L2
368   repeat elements in sense orientation, which all showed pronounced binding
369   by MATR3 and PTBP1 in cultured human cells as well as binding by MATR3
370   and PTBP2 in mouse brain (Figure S4C). We confirmed by RT-PCR that
371   individual depletion of MATR3, but not PTBP1, led to inclusion of an
372   alternative exon starting 323 bp upstream of the L2 repeats (Figure 4B), and
373   verified its splice sites by Sanger sequencing. Notably, inclusion of the exon is
374   markedly elevated after co-depletion of MATR3 and PTBP1 (Figure 4B).

375   To confirm that the LINE nucleotide sequence recruits MATR3 and PTBP1
376   and causes distant splicing repression, we created a splice reporter of ACAD9
377   comprising exon1 and exon2 and the complete intronic sequence including all
378   three L2 repeats (called wildtype), and a mutant splice reporter that lacked
379   two out of the three L2 repeats (called ΔLINE, see Figure 4A). The wildtype
380   reporter reproduced the splicing pattern of the endogenous sequence in
381   non-perturbed cells and in cells depleted of MATR3, PTBP1 or both (Figure
382   4C), except of overall more frequent usage of the 5' splice site of exon 1b.
383   Importantly, the ΔLINE splice reporter showed increased usage of the LINE-
384   proximal 3' splice site in unperturbed cells, with little to no further change in
385   incusion upon MATR3/PTBP1 depletion (Figure 4C). Hence, the L2 sequence
386   downstream of the exon was essential to confer responsiveness to
387   MATR3/PTBP1. This supports our transcriptome-wide finding that MATR3
388   and PTBP1 repress LINE-proximal exons, in addition to regulating
389   LINE-derived exons.

390

391   **PTBP2 prevents LINE-exon inclusion in mouse brain**

14

392  Having identified the role of MATR3 and PTBP1 in repressing the splicing of

393  cryptic LINE-derived exons, we sought to explore if they might play such a

394  role also in the brain, given the known role of the PTBP1 orthologue PTBP2 in

395  regulating splicing during neuronal development (Li et al., 2014). We first

396  generated iCLIP data of MATR3 from mouse brain, and compared the

397  enrichment on LINEs in the mouse brain for PTBP2, MATR3, CELF4, FUS

398  and TARDBP. Enrichment was most pronounced for MATR3 and highest on

399  antisense L1 sequences, to a similar extent as in HEK293 cells (Figure S5A).

400  Interestingly, we found MATR3 and PTBP1 show stronger enrichment on

401  rodent-specific L1 families than on evolutionary older L1 families. A MATR3

402  knockout mouse is not available (MGI:1298379); therefore we focused on

403  RNAseq data from PTBP2$^{-/-}$ mouse brain (Li et al., 2014, Vuong et al., 2016).

404  In nestin-Cre-PTBP2$^{-/-}$ E18 mouse brain, we found LINE-derived exons were

405  more likely to be significantly deregulated than SINE-derived exons or non-

406  repeat derived exons (Figure S5B; $\chi^2$-test, p-value < $10^{-5}$) and were repressed

407  by PTBP2 (Figure S5B; $\chi^2$-test, p-value < $10^{-5}$). Hence, we suggest repression

408  of LINE-derived exons is redundant between PTBP1 and PTBP2. Focusing on

409  exons with inclusion levels above 10% (measured as percent spliced index, or

410  PSI, see SupplTable3 and 5 of (Vuong et al., 2016)), PTBP2-regulated exons

411  in Emx-Cre-PTBP2$^{-/-}$ mouse brain include LINE-derived exon3 of *Fam124A*,

412  exon5 of *Osblp9* and exon 2 of Ub3g1, all of which modify the protein

413  sequence. These exons are absent or lowly included in E14 wildtype brains

414  but their inclusion increases in the adult brain at P10 (Figure S5C to E). Out of

415  13 LINE-derived exons selected for this pattern, 8 are rodent-specific

416  insertions that are not shared between mouse and human. This suggests that

417  PTBP2 preferentially represses exons derived from evolutionarily young

418  LINEs during brain development, and several of these exons become more

419  highly included in the wildtype adult brain, thus gaining the potential to alter

420  the species-specific neuronal transcriptome.

421

422  **Evolutionarily old LINEs are a major source of mammalian alternative**

423  **exons**

15

424  After observing exonisation of LINEs in mouse brain, we decided to survey

425  the inclusion of LINE-derived exons in human tissues by using the extensive

426  data available from the GTEx consortium (V6p data (Consortium, 2015)). We

427  tested the percent inclusion of a total of 45,940 exons in RNAseq data from

428  51 human tissue types, covering all exons of the 4,566 genes that contain a

429  known LINE- or Alu-derived exon. We detected 1,154 LINE-derived exons

430  with at least 5% inclusion in at least one tissue. The LINE-derived exons

431  showed higher degree of exonisation than Alu-derived exons, measured by

432  maximum inclusion level (PSI) across tissues (Figure S6A), but showed a

433  similarly high degree of tissue-specificity (Figure S6B). In contrast to well-

434  established alternative exons, Alu- and LINE-derived exons are virtually never

435  switch-like events (Figure S6B).

436  Since we observed enrichment of PTBP2-regulated exons derived from

437  evolutionarily younger LINE families in mouse brain (Figure S5C to E), we

438  also further explored the evolutionary history of human LINE-derived exons.

439  For this purpose, we determined the evolutionary age of all human L1 and L2

440  repeats by cross-species comparison with two primate genomes (gorilla and

441  rhesus macaque), two rodents (mouse and rat), and one each of the carnivore

442  and laurasiatherian genomes (dog and cow, respectively), which were chosen

443  due to the high quality of their genome assemblies. In this manner, we

444  annotated    LINEs    as    primate-specific,    euarchontoglires-specific    or

445  mammalian-wide insertions (Figure 5A). We were able to categorise

446  mammalian-wide insertions further by assigning if they were present in dog

447  and cow or only one out of the two, which might indicate differences in

448  selection pressure. We ignored elements for which the evolutionary history

449  remained unclear, or which were present but largely sequence truncated in

450  dog or cow. The number of substitutions of the elements from family

451  consensus validated the average age in our annotation (Figure S6D),

452  although we found it to be more robust on L1 than on L2 elements. Human L2

453  elements are much older than L1s (Deininger and Batzer, 2002), which

454  means their insertion age is frequently older than the divergence of the

455  genomes we used. Since it remains unclear if any L2 family has remained

456  active in early ancestors of the euarchontoglires lineage, we focused on L1

457    elements for analysis of young LINE insertions.

458    Next, we examined the proportion of L1 repeats from each phylogenetic LINE
459    group that is capable of seeding exons. We were surprised to find that
460    L1-derived exons are a rich source of exons in the regions of the genome that
461    encode the highly variable and species-specific immunoglobulin variable
462    chain region (the Ig-region on chromosomes 2, 14, 15, 16 and chr22, Figure
463    5B). The Ig-regions are densely packed with 1,845 LINEs, 1,152 of which
464    produce exons according to exon annotation by UCSC. The LINE-derived
465    exons in these regions are almost exclusively seeded by primate-specific L1s
466    (Figure 5B), and we consider them as cryptic exons, since we did not detect
467    them by our analysis of the GTEx data. However, even when they are
468    included, the exons are unlikely to map to the reference genome due to the
469    rearrangement of the variable chain region during B- and T-Cell maturation.
470    Detailed analysis of B and T cell receptor sequences will be needed to further
471    examine the contribution of these young L1-derived exons to the expression
472    of immunoglobulin genes. After excluding the Ig-regions, we find that less
473    than 0.4% of the evolutionary young L1s can seed exons compared to 0.8%
474    of the well-preserved older L1 elements, demonstrating that the older L1s
475    more frequently contribute to the established transcriptomes across tissues
476    (Figure S6C).

477    To quantify the differential regulation of LINE-derived exons across tissues,
478    we calculated the maximum difference in inclusion between any pair of
479    tissues. Interestingly, the inclusion of exons seeded by young L1 elements
480    was more tissue-specific (Figure 5C, Figure S6E). However, exons derived
481    from evolutionarily older LINEs generally showed higher maximum inclusion,
482    comparing young L1 elements to either old L1 elements, or to exons derived
483    from L2 and CR1 elements (Figure 5D). We found 594 L2- and 150 CR1-
484    derived exons, which had inclusion levels similar to the evolutionary old L1
485    insertions (Figure 5D). In fact, CR1-derived exons were the group with highest
486    inclusion levels of all groups examined, which agrees with them generally
487    being the evolutionarily oldest in human. Between tissues, we found most
488    LINE-derived exons in tissues of the reproductive system (Testis, Fallopian

17

489 tube and Cervix) and the brain (considering all 13 regions of the brain; Figure
490 S6G). Taken together, our analyses show that the evolutionarily older LINE
491 insertions are a major source of mammal-specific alternative exons, some of
492 which have reached high inclusion levels in different human tissues.

493

494 **Loss of repressor binding drives the exonisation of LINE-derived exons**

495 To explain the differences in the inclusion level of the different evolutionary
496 categories of LINE-derived exons, we compared their splice site strength, but
497 did not find any marked differences (Figure S6F). Therefore, we reasoned that
498 instead of changes in splice site strength, changes in the binding of different
499 RBPs might determine exonisation of LINE-derived exons. To test this
500 hypothesis, we exploited the available iCLIP and eCLIP data to analyse
501 trends in RBP binding across the phylogenetic groups of L1 insertions. To
502 ensure that we assessed elements that are part of expressed transcripts, we
503 selected the 10% of L1 elements with highest coverage by any of the 121
504 RBPs. All phylogenetic groups were represented in this selection in expected
505 proportions. Next, we averaged the binding of each RBP against the sum of
506 all RBPs, generating a relative binding metric among all RBPs (ranging from 0
507 to 1). We then visualised any preferences in binding to a phylogenetic group
508 as enrichment, considering all 49 RBPs that had above-average binding to
509 LINEs (see Figure 1A). Strikingly, MATR3 is the RBP with strongest
510 enrichment on primate-specific L1s among iCLIP experiments, and BCCIP
511 and hnRNPM among eCLIP (Figure 5E). Both iCLIP and eCLIP, in both cell
512 lines, also show PTBP1 enriched on primate-specific L1s. In general, known
513 splicing repressors are enriched on primate-specific L1s, with the exception of
514 hnRNPC. In contrast, RBPs that are well known to enhance splicing or
515 3' processing also bind to evolutionarily older L1s, which includes SR
516 proteins, and RBPs that recognise sequences close to 3' and 5' splice sites,
517 or the polyA sites (Figure 5E, lower panels). We conclude that the stronger
518 binding of repressive RBPs to the young L1s is the likely reason for their lower
519 inclusion. The loss of these repressive RBPs, accompanied with binding by
520 splice-promoting factors, can thus explain why the evolutionarily older L1s are

18

521 the most common source of exons, and why they tend to be more highly
522 included.

523

524 **Sequence divergence of the evolutionarily old L1 elements results in**
525 **loss of repressor binding**

526 To understand why the evolutionarily older LINEs do not bind repressive
527 RBPs as well as younger insertions, we analysed the density of sequence
528 motifs known to interact with these RBPs (Figure 6A). We found binding
529 motifs in the literature for ELAVL1, HNRNPK, HNRNPM, KHDRBS1, QKI,
530 PTBP1/2, RBFOX1 and TARDBP (see Suppl. Table 1 for details and
531 references). We tested the distribution of all 256 tetramer sequences and
532 found clear differences in line with the expected AT-richness of antisense L1
533 sequences (Suppl. Table 6), with TG- and TA-rich motifs being quite abundant
534 in antisense L1 and CG-rich motifs being most depleted. We found
535 evolutionarily older L1s contained fewer binding motifs of hnRNPM, TARDBP
536 and ELAV1 (at FDR < 0.1). We found they contained on average more
537 binding motifs of KHDRBS1, hnRNPC and QKI, though a large proportion of
538 evolutionarily old L1s did not contain any QKI motif and none of the motifs
539 associated with these three was among the most enriched motifs. PTBP1
540 motifs were highly abundant (a median of 1.26 motifs per 100nt) in all L1s,
541 irrespective of their genomic age. We conclude the unequal binding towards
542 L1s of RBPs, especially of splicing repressive RBPs such as hnRNPM,
543 ELAVL1 and TARDBP, is a consequence of the L1 sequence and its decline
544 through accumulation of sequence mutations.

545

546 **The evolutionarily young LINEs maintain the repression of deep intronic**
547 **regions**

548 Given that we show assembly of mostly splice-repressive RBPs at and across
549 evolutionary young LINE sequences, we hypothesised that these LINEs are in
550 a repressed state. Furthermore, at least MATR3 and PTBP1 inhibited splicing

19

551    also in nearby regions, which raised the intriguing question of whether

552    evolutionarily young LINEs are generally prohibitive for splicing. To test if

553    young L1s act as intronic splice silencers, we examined their distribution

554    around annotated exons as well as the inclusion of these exons across

555    human tissues. Strikingly, we found that evolutionarily young LINEs were

556    excluded from an approximately 3kb region around constitutive and

557    alternative exons (Figure 6B). However, they were not excluded around those

558    exons with very low inclusion across human tissues (maxPSI<15%),

559    indicating that they may contribute to the repression of these exons (Figure

560    6B). In total contrast to young antisense L1 sequences, the primate-specific

561    Alu repeats were only excluded from the immediate vicinity of exons, but not

562    from flanking intronic regions. Older L1 elements are well tolerated up to

563    250bp at all exons, and their incidence decreases only within ±200nt of

564    constitutive exons. As independent validation, we repeated the analysis on

565    mouse exons, and found mouse- and rodent-specific LINEs excluded in a

566    large window around their splice sites, a pattern recapitulating the

567    primate-specific insertions in human (Figure 6C). Thus, the evolutionarily

568    younger antisense L1 elements are more depleted from the vicinity of exons

569    both in primates and rodents. This could be a consequence of them being

570    particularly potent in recruiting repressive RBPs such as MATR3 and PTBP1,

571    which leads to repression of exons in their vicinity.

572

573

20

574     **DISCUSSION**

575     We find that tens of thousands of LINEs recruit a diverse set of 28 RBPs to

576     deep intronic loci. Insulating RNA from the splicing and polyadenylation

577     machinery is a known mechanism of repression used by a number of RBPs

578     (Witten and Ule, 2011). Of the RBPs binding of LINEs many are splicing

579     repressors, such as MATR3 and PTBP1, which repress the recognition of

580     LINE-derived exons and polyA sites, both in cultured human cells and in

581     PTBP2$^{-/-}$ mouse brain. Importantly, we demonstrate that MATR3 promotes

582     efficient PTBP1 binding to L1s by stabilising its interaction with multivalent

583     binding motifs.

584     The repetitive nature of LINEs and their evolutionary divergence allowed us to

585     demonstrate a dual role of young and old LINEs in RNA processing (Figure

586     7). Repressive RBPs preferentially bind to the young LINEs and insulate both

587     the LINEs and the surrounding regions from the processing machineries. As a

588     consequence, young LINEs are confined to deep intronic regions. We

589     propose their insulation allows the accumulation of cryptic RNA processing

590     sites and facilitate evolutionary innovation. Through the accumulation of

591     sequence mutations, the density of repressive binding motifs is gradually

592     decreased, and these processing sites become gradually de-repressed. This

593     is evident by the closer proximity of older LINEs to exons, by their binding to

594     RBPs that enhance RNA processing, and by their increased contribution to

595     tissue-specific transcript isoforms. Thus, we find that intronic LINEs play a

596     dual role: they recruit repressors to insulate the deep intronic regions from

597     processing machineries, but after long evolutionary periods also act as a

598     source of RNA processing sites that facilitate the formation of mammal-

599     specific transcripts.

600

601     **MATR3 and PTBP1 bind LINEs to synergistically repress RNA**

602     **processing**

603     We found that antisense L1 and sense L2 elements recruit MATR3 and

604     PTBP1 to deep intronic regions, where both RBPs repress RNA processing at

21

605 and around the bound LINEs. Binding of both proteins significantly overlaps at
606 these LINEs, and MATR3 is required for efficient PTBP1 binding to L1s.
607 Antisense L1s contain many PTBP1 binding motifs and MATR3-dependent
608 binding sites of PTBP1 are characterised by increased density of binding
609 motifs over a broad 200nt region. A model to explain our observations is that
610 LINEs provide multivalent binding sides, and complex formation with MATR3
611 promotes PTBP1 binding to those. PTBP1 is capable of multivalent RNA
612 binding through its four RRM domains (Oberstrass et al., 2005). Analysis of
613 liquid phase separation properties of PTBP1 *in vitro* recently demonstrated
614 that its binding is mediated in part by multivalent binding sites on the RNA,
615 and is further stabilised by fusing PTBP1 to intrinsically disordered regions
616 (IDRs) of different proteins, due to IDR-mediated protein-protein interactions
617 (Lin et al., 2015). It was therefore proposed that the PTB-RNA and IDR
618 interactions could act together to produce larger oligomeric assemblies with
619 increased affinity for RNA. PTBP1 RRM2 interacts with a short linear motif
620 within an IDR in MATR3 (the PRI motif, Coelho et al., 2015), and MATR3
621 interacts with itself (Zeitz et al., 2009). It seems likely that PTBP1 and MATR3
622 assemble across the antisense L1 sequences as a larger oligomeric assembly
623 through multivalent RNA binding. Notably, we find that one of the previously
624 studied MATR3-repressed exons is derived from a sense L2 insertion (exon
625 11 in ST7), and repression of this exon depends on its PRI motif (Coelho et
626 al., 2015), indicating that the repressive function of MATR3 involves formation
627 of a multiprotein complex with PTBP1 and possibly additional LINE-binding
628 RBPs. Indeed, MATR3 has been reported as part of several nuclear
629 multimeric complexes (Damianov et al., 2016, Kula et al., 2011, Zhang and
630 Carmichael, 2001). One of these is the LASR complex, which includes
631 hnRNPM (Damianov et al., 2016), an RBP we find preferentially binds to
632 young antisense L1 elements much like MATR3. Taken together, our results
633 suggest that L1 elements are sites of multivalent binding of PTBP1 and
634 possibly other RBPs. This can provide the high avidity for assembly of
635 repressive ribonucleoprotein complexes in order to insulate the antisense L1s
636 and nearby RNA from RNA processing machineries.

637

22

638 **LINE-derived exons are highly tissue-specific**

639 Evolutionary young LINE insertions are bound by a large number of splice
640 repressors. This is the likely reason why only small numbers of LINEs form
641 canonical exons, even though most LINEs contain strong 3' and 5' splice site
642 sequences. Based on the often ubiquitous expression of MATR3, hnRNPM
643 and other repressive RBPs across human tissues (Petryszak et al., 2016),
644 LINEs need to escape this repression before they can be spliced into exons.
645 In agreement with this, we found that LINE-derived exons are alternatively
646 spliced, with large differences in inclusion between tissues and often
647 completely absent from most of them. For instance, 398 of the 1,169
648 detectable LINE-derived exons are restricted to less than five tissues in
649 humans; and in wildtype adult mouse brain, where activity of PTBP1 and
650 PTBP2 is decreased, only a handful of LINE-derived exons become included
651 at 10% or higher. However, we found strong de-repression of LINE-derived
652 exons in PTBP2$^{-/-}$ mouse brain and in cultured human cells in the absence of
653 MATR3 and PTBP1. Since many different RBPs bind to intronic LINEs, it is
654 likely that the regulation of LINE elements is combinatorial, such that the
655 abundance of multiple RBPs determines inclusion of LINE-derived exons in
656 each tissue. In addition, most PTBP1-repressed LINE-derived exons trigger
657 NMD, which is likely to be common across evolutionarily young LINE-derived
658 exons. Hence, a *bona fide* LINE-derived exon has to overcome both the
659 splicing repressive mechanisms and NMD in order to form alternative, tissue-
660 specific exons.

661

662 **LINEs facilitate evolution of RNA processing**

663 To understand the forces that drive the evolutionary emergence of new LINE-
664 derived alternative exons and poly(A) sites, we asked how the evolutionary
665 age of LINEs affects their distribution in vicinity to exons, and the binding
666 patterns of RBPs as surveyed by iCLIP/eCLIP. We find that intronic LINEs
667 that recruit strong splicing repressors such as MATR3 can have repressive
668 effects on nearby exons, which agrees with the lower inclusion of exons that

23

669    are located nearby young LINEs. Conversely, young LINEs are depleted from

670    the vicinity of constitutive exons, most likely as a result of purifying selection

671    against the insertion of novel LINEs near existing exons, or selection of exons

672    outside the repressive environment created by LINEs. While LINEs were

673    known to be depleted in immediate vicinity of splice sites (Zhang et al., 2011),

674    we now find that the extent of this depletion is distinctly dependent on their

675    evolutionary age.

676    Our analysis of RBP binding patterns on LINEs demonstrates the difference in

677    RNP assembly at evolutionarily young and old LINEs, which mediates their

678    functional differences. The binding of repressive RBPs is most enriched in

679    young LINEs, whereas binding of known splicing enhancers belonging to

680    U2AF, TIA and SR families, and CPSF machinery (CSTF2 and CPSF6), is

681    either increased or unchanged in the older LINEs. Hence, repressive RBPs

682    prevent recognition of cryptic splice sites in thousands of young LINEs. The

683    further the sequence of a LINE diverges, the more likely binding of one or

684    more repressive RBPs is lost, thus allowing individual elements to seed

685    lineage-specific and highly tissue-specific exons with low or modest inclusion.

686    These exons become susceptible to evolutionary selection, which sets the

687    scene for the emergence of a few *bona fide* exons with higher inclusion,

688    seeded by evolutionarily old LINEs. The relationship between splicing

689    repressors and LINEs is in many ways similar to the evolutionary dynamic of

690    KAP1/KRAB transcription factors, which repress transcription at

691    retrotransposons and confer robustness to transcriptional networks while

692    facilitating evolutionary novelty (Thomas and Schneider, 2011, Imbeault et al.,

693    2017).

694

695    **Conclusion**

696    We propose that MATR3, PTBP1 and other repressive RBPs insulate the

697    intronic LINEs to allow accumulation of cryptic elements. This could explain

698    why strong RNA processing sites are prevalent in LINEs, and why LINEs

699    remain cryptic without any deleterious effects. Evolutionarily young LINEs

700      form a main hub for the recruitment of repressive RBPs, which in turn

701      demarcate introns by insulating cryptic elements both within and around the

702      LINE from processing machineries. These repressive RBPs are crucial to

703      protect gene expression from the cryptic exons derived from LINE insertions,

704      and it appears that a network of more than two dozen of LINE-binding RBPs

705      contribute to this repression, and some of them possibly in a cooperative

706      manner. We note this aligns with (1) previously proposed models of exon

707      emergence (Modrek and Lee, 2003, Xing and Lee, 2006, Attig et al., 2016), in

708      which lowly included alternative exons are the test bed for emergence of new

709      exons; as well as (2) the proposal that genomes accumulate cryptic variation

710      between lineages which is only apparent upon perturbation (Ward et al., 2013,

711      Payne and Wagner, 2015, Tirosh et al., 2010). The consequences LINEs

712      have on the transcriptome are apparent in the evolution of novel transcript

713      isoforms, and the frequency of hereditary diseases occurring if one of these

714      elements is unmasked.

715

716

717

718 **METHODS**

719 **Cell culture and siRNA transfection**

720 Hela and HEK293T cells were maintained in DMEM with 10% FBS at 37°C
721 with 5% CO2 injection, and routinely passaged twice a week. Cells were
722 regularly cultured for three days in antibiotic-free medium and tested for
723 mycoplasma using either the LookOut Mycoplasma PCR Detection Kit or the
724 MycoAlert mycoplasma detection kit (Lonza).

725 To deliver siRNAs, Lipofectamin RNAiMax (Life Technologies) was used
726 according to manufacturer's recommendations. HeLa cells were grown in
727 antibiotic-free medium and forward transfected with siRNA targeting PTBP1 at
728 10nM (AACUUCCAUCAUUCCAGAGAA) and PTBP2 at 5nM (AAGAGAGGAUCUGACGAACUA),
729 synthesized by Dharmacon, or siRNAs purchased from Invitrogen targeting
730 MATR3 mRNA at 5nM (HSS114732) or 20nM (HSS114730), as well as
731 control siRNA Negative Control with medium GC content (20nM, Invitrogen,
732 Cat. number 12935-300).

733

734 **Nucleo-cytoplasmic fractionation**

735 Cells were washed ice-cold PBS and lysed with NP40E-CSK (350µl per well
736 of a 6-well-plate or 600µl per 10cm dish). NP40E-CSK buffer is similar to the
737 cytoskeleton buffer used in (Reyes et al., 1997) and composed of 50 mM Tris-
738 HCl (pH 6.5), 100 mM NaCl, 300 mM sucrose, 3mM MgCl2, 0.15% NP40 and
739 40 mM EDTA). Lysis was allowed to proceed for 5 minutes on ice. HeLa cells
740 had to be scraped off due to their strong adhesion to the culture dish.
741 Cytoplasmic supernatant and pelleted nuclei were separated at 4°C, 5000 x g
742 for 3 minutes. The cytoplasmic supernatant was cleared with another spin at
743 4°C, 5000 x g for 3 minutes and another spin at 4°C, 10000 x g for 10
744 minutes. Nuclei were washed with 400µl NP40E-CSK and incubated for 5
745 minutes under rotation to ensure complete cell lysis. After repeat of the
746 centrifugation step, nuclei were lysed in 300µl CLIP lysis buffer and sonicated
747 at 5x 30sec pulses in a BioRuptor waterbath device. Subsequently, RNA was

26

748 isolated using Trizol LS (Invitrogen) and Zymo RNA isolation columns

749 (Zymogen) according to manufacturer's recommendations. For preparation of

750 RNA for RNAseq, an additional wash step with 180µl NP40E-CSK was done

751 before nuclei rupture.

752

753 **Semi-quantitative RT-PCRs**

754 Reverse transcription was done with 500ng of RNA using RevertAid enzyme

755 (Fermentas) according to manufacturer's recommendations. The reverse

756 transcription was primed with equal parts of random N6 and N15

757 oligonucleotides (Sigma) at 100µM concentration. For semi-quantitative PCR,

758 we run 35 cycles of amplification with the primer combinations as indicated in

759 each figure (primers are listed in Supplementary Table 1), and quantified the

760 abundance of each product using Qiaxcel™ (Qiagen) gel electrophoresis.

761

762 **UV crosslinking assay on recombinant proteins**

763 For in vitro assays, we made two artificial sequences. The first contained two

764 embedded AUCUU motifs (shown in bold) and CT-rich stretches in their

765 vicinity (underlined):

766 GAATACGAATTCCATATATGATCGATAAATATATGGTACCTTGCTATCTTACATCTTTTTACGGATCCCATATATG

767 ATCGATATATATAAGCT.

768 The second RNA probe contained six CTCC motifs (shown in bold):

769 GAATACGAATTCC**CTCTT**TGAATCGATAA**CTCTT**TGGTACCC**CTCTT**TGATCGATAA**CTCTT**TGGATCCC**CTCTT**T

770 GATCGAT**CTCTT**TAAGCTT.

771 The RNA probes were labelled with $^{32}$P-UTP using SP6 RNA polymerase. We

772 purified full-length N-terminal His-tagged recombinant PTBP1 (rPTBP1) and

773 three MATR3 fragments (rMATR3, amino acids 362-592 or 'RRMs', and

774 amino acids 341-592 or 'RRM-PRI' with or without mutations in the PRI motif),

775 using Blue Sepharose 6 and HisTrap HP columns. In UV crosslinking assays

27

776   with recombinant proteins, we used 10fmol of RNA, 0.5µM rPTBP and titrated

777   increasing amounts of rMATR3 fragments against it (0 to 2 µM). After

778   incubation at 30°C for 20 minutes, the sample was UV cross-linked on ice in a

779   Stratalinker with 1920 milliJoule. The binding reaction was then incubated for

780   10 minutes at 37°C together with 0.28 mg/ml RNase A1 and 0.8 U/ml RNase

781   T1. SDS loading buffer was added and the samples heated to 90°C for 5

782   minutes before loading on 15% denaturing polyacrylamide gel. To assay

783   binding in HeLa nuclear extract, we prepared standard nuclear extract

784   (Dignam et al., 1983), and combined 10fmol of RNA probe with 0.5 µM

785   rMATR3 and 20% extract.

786

787   **Generation of iCLIP data**

788   HEK293T cells were grown on 10 cm$^2$ dishes, incubated for 8 h with 100 µM

789   4SU and crosslinked with 2x 400mJ/cm$^2$ 365nm UV light. Protein A

790   Dynabeads were used for immunoprecipitations (IP). 80 µl of beads were

791   washed in iCLIP lysis buffer (50 mM Tris-HCl pH 7.4, 100 mM NaCl,

792   1% NP-40, 0.1% SDS, 0.5% sodium deoxycholate). For the preparation of the

793   cell lysate, 2 million cells were lysed in 1 ml of iCLIP lysis buffer (50 mM Tris-

794   HCl pH 7.4, 100 mM NaCl, 1% NP-40, 0.1% SDS, 0.5% sodium

795   deoxycholate) buffer, and the remaining cell pellet was dissolved in 50 µL

796   MSB lysis buffer (50mM Tris-HCl pH 7.4, 100mM NaH2PO4, 7M UREA, 1mM

797   DTT, (Reyes et al., 1997)). After the pellet had dissolved, the mixture was

798   diluted with CLIP lysis buffer to 1000 µl and an additional centrifugation was

799   performed. We found by Western Blotting that up to 50% of MATR3 protein is

800   insoluble by detergent without urea. Lysates were pooled (2ml total volume)

801   and incubated with 4 U/ml of RNase I and 2 µl antiRNase (1/1000, AM2690,

802   Thermo Fisher) at 37°C for 3 min, and centrifuged. We took care to prepare

803   the initial dilution of RNase in water, since we found that RNase I gradually

804   loses its activity when diluted in the lysis buffer. 1.5 ml of the supernatant was

805   then added to the beads and incubated at 4°C for 4 h. The rest of the iCLIP

806   protocol was identical to the published protocol (Huppertz et al., 2014).

807

## 808     **Mapping of iCLIP and eCLIP data**

809     MATR3 and PTBP1 iCLIP libraries were sequenced on Illumina HiSeq2

810     machines in a single-end manner with a read length of 50 nt. Before mapping

811     the reads, we removed adapter sequences using the FASTX toolkit version

812     0.7 and we discarded reads shorter than 24 nucleotides. Reads were then

813     mapped with the iCount suite to UCSC hg19/GRCh37 or mm9/NCBI37

814     genome assembly using Bowtie v2.0.5 allowing up to two mismatches and up

815     to 20 multiple hits. Unique and multiple mappers were separately analysed,

816     and to quantify binding to individual loci, only uniquely mapping reads were

817     used. Supplementary Table 1 lists the source and details including fil numbers

818     of all published iCLIP and HITS-CLIP data used within this study.

819     The eCLIP libraries were downloaded from ENCODE (Van Nostrand et al.,

820     2017, Sloan et al., 2016). Before mapping the reads, adapter sequences were

821     removed using Cutadapt v1.9.dev1 and reads shorter than 18 nucleotides

822     were dropped from the analysis. Reads were mapped with STAR v2.4.0i

823     (Dobin et al., 2013) to UCSC hg19/GRCh37 genome assembly. To quantify

824     binding to individual loci, only uniquely mapping reads were used. For

825     analysis of enrichment on repeat families, up to 20 multiple alignments were

826     allowed and fractional counts used.

827

## 828     **TEtranscript estimates of LINE family enrichments**

829     To consider both uniquely mapping and multimapping reads in estimating

830     binding to repeat (sub)families, we used the approach described in

831     TEtranscripts (Jin et al., 2015). In short, for eCLIP FASTQ files, adapters were

832     removed according to the ENCODE eCLIP standard operating procedure. For

833     iCLIP FASTQ files, barcodes were removed using the FASTX-Toolkit (v

834     0.0.14). For all files, reads aligning to rRNA or tRNA were removed by

835     aligning to custom rRNA and tRNA indices (human or mouse as appropriate)

836     using Bowtie2 (v. 2.2.9). The remaining reads were aligned to the appropriate

837    genome (GRCh38, Gencode V25 for human, and GRCm38, Gencode M13 for

838    mouse) using STAR (v. 2.5.2) with the addition of the parameters "--

839    `winAnchorMultimapNmax 100 --outFilterMultimapNmax 100`" as

840    recommended by TEtranscripts. For each CLIP dataset, TEtranscripts was

841    run using both stranded options (`--stranded reverse` and `--stranded`

842    `yes`) to obtain results for sense and antisense LINE binding.

843    RNAseq data from ENCODE was used as control, for eCLIP RNAseq of K562

844    and HEPG2 cells lines (ENCSR885DVH and ENCSR181ZG). For iCLIP

845    samples from mouse brain, we used P2 mouse brain from ENCODE. The

846    iCLIP data in mouse brain was produced from total mouse brain, so we

847    pooled the RNAseq of forebrain, midbrain and hindbrain, accession numbers

848    ENCSR723SZV, ENCSR255SDF and ENCSR749BAG (Sloan et al., 2016).

849

850    **Generation of RNAseq libraries and mapping with TopHat2 (human)**

851    Before library preparation, purified RNA was DNase I treated for a second

852    time and purified with the DNA-free kit (Ambion). To generate stranded

853    RNAseq libraries, we used the TruSeq stranded RNAseq library kit (Illumina)

854    according to manufacturer's recommendations; RNA was depleted of rRNA

855    using the RiboZero kit (Epicentre).

856    All libraries were sequenced on Illumina HiSeq2 machines in a single-end

857    manner with a read length of 100 nt. Before mapping the reads, adapter

858    sequences were removed using the FASTX toolkit version 0.7 and we

859    discarded reads shorter than 24 nucleotides. Reads were then mapped with

860    TopHat v2.0.5 (Kim et al., 2013) to UCSC hg19/GRCh37 genome assembly

861    using ENSEMBL version 72 gene annotation as reference, allowing up to two

862    mismatches and only using uniquely mapping hits. RNAseq data files of rRNA

863    depleted cytoplasmic and nuclear RNA from cells depleted of MATR3 and

864    PTBP1 are deposited on EBI ArrayExpress under the accession number

865    E-MTAB-6204.

866

**Generation of pAseq libraries and mapping**

867

868 To quantify polyA site usage, we used the QuantSeq mRNA 3' end
869 sequencing kit (Lexogen) according to manufacturer's recommendations. We
870 used both the forward and reverse library kit on two independent biological
871 replicates each (four replicates in total). Libraries were prepared from nuclear
872 RNA after individual or combined siRNA depletion of MATR3 and PTBP1/2.
873 All libraries were sequenced on Illumina HiSeq2 machines in a single-end
874 manner with a read length of 100 nt. polyA site usage was analysed with the
875 expressRNA platform. In short, reads were mapped with STAR v??? to UCSC
876 hg19/GRCh37 genome assembly, allowing up to ??? mismatches and ?only
877 using uniquely mapping hits?. pAseq raw data is deposited on ArrayExpress
878 at E-MTAB-6287.

879

**Mapping of published RNAseq with STAR (human)**

880

881 To test for any change in usage of LINE-derived exons upon depletion of the
882 NMD core factor UPF1, we made use of the data generated by Ge et al. and
883 in HEK293 cells, depleted of PTB, UPF1 or both proteins (Ge et al., 2016).
884 Raw sequencing data in FASTQ format was downloaded from SRA and
885 mapped with STAR v2.5.2a (Dobin et al., 2013) to UCSC hg19/GRCh37
886 genome assembly, allowing up to 10 mismatches and only using uniquely
887 mapping hits. Then we analysed the data using JunctionSeq (Hartley and
888 Mullikin, 2016) with ENSEMBL version 72 gene annotation as reference.

889

**Mapping of published RNAseq with STAR (mouse)**

890

891 To test for LINE-derived exon inclusion in mouse brain, we made use of the
892 data generated by Li et al. and Vuong et al. (Li et al., 2014, Vuong et al.,
893 2016). Raw sequencing data in FASTQ format was downloaded from SRA
894 and mapped with STAR v2.5.2a (Dobin et al., 2013) to UCSC
895 mm10/GRCm38 genome assembly, allowing up to 10 mismatches and only

31

896 using uniquely mapping hits. Then we analysed the data using JunctionSeq

897 (Hartley and Mullikin, 2016) with ENSEMBL version 72 gene annotation as

898 reference.

899 **Sequence motif analysis**

900 For PTBP1 motifs around iCLIP peaks, we used the strong binding motifs as

901 defined previously (Haberman et al., 2017), and counted their occurrence

902 around peak centres. To define enrichment, we divided the occurrence at

903 MATR3-dependent and independent peaks by the distribution across all other

904 PTBP1 peaks.

905 For motifs within antisense L1 elements, we used motifs described in the

906 literature; for PTBP1, TARDBP and hnRNPM their binding motifs were

907 validated in vitro and through functional studies (Gooding et al., 1998,

908 Oberstrass et al., 2005, Avendano-Vazquez et al., 2012, Ayala et al., 2005).

909 For all other proteins, we used RNAcompete motifs (Ray et al., 2013). The

910 number of motifs per 100nt gave a distribution for each motif, and we used

911 quartiles for each motif to describe gain or loss of motifs between evolutionary

912 groups. To obtain a false discovery rate of motif gain or loss, we generated an

913 empirical distribution of motif enrichments across groups. We compared the

914 change in Q1 and Q4 for each of the possible 256 tetramers, which resulted in

915 an approximately normal distribution. We then called motifs within the 2.5%

916 and 5% extremes as significant at FDR<0.05 and FDR<0.1.

917

918 **RNA maps**

919 All metaprofiles of iCLIP data and LINE sequence content around loci of

920 interest (also called RNAmaps) were drawn in R. Metaprofiles are normalised

921 to the number of input loci of each track, and data was smoothed using

922 binning as indicated in figure legends, using the *zoo* package. A generalised

923 script for generation of a metaprofile can be found at

924 https://github.com/JAttig/generalised-Rscripts.

925

**Annotation of known alternative exons**

For annotation of exons known to be alternatively spliced, we downloaded the 'knownAlt Events' and 'knownGene' from UCSC TableBrowser for hg19 on 28th March 2014. In addition, we downloaded the 'refGene' table on 23rd March 2017. The exons annotated by UCSC were collapsed within a gene to unique exonic ranges, and classified as alternative or constitutive exon as follows. All exons not annotated as alternative by UCSC and present in the RefSeq exon annotation with identical genomic coordinates were classified as constitutive, all other exons were considered alternative exons.

935

***De novo* identification of cryptic exons**

In order to predict exons from our RNAseq data, we ran Cufflinks (version 0.9.3, -min-isoform-fraction 0, Trapnell et al., 2012) on the collapsed reads from all cytoplasmic samples of our stranded RNAseq data and then extracted the exons of all predicted transcripts. After flattening the Cufflinks output to non-overlapping exonic bins, our Cufflinks prediction contained 671,956 exonic bins. However, we only considered exonic bins of at least 5 nucleotides in size. All exons with one or both splice sites residing within a LINE repeat (as annotated by RepeatMasker, (Smit et al., 1996-2010b)) were assigned as LINE-exons. In order to minimise noise, we only kept exons for analysis that were supported by at least one junction-spanning read (225,322 exonic bins). All exons that were not identical with exons annotated in UCSC gene annotation (hg19) were referred to as 'cryptic' (see also Supplementary Table 2) for complete breakdown of annotation of exonic bins). For readability, we refer to 'exonic bins' as 'exons' throughout the text.

951

**Analysis of differential gene expression and differential exon inclusion**

Analyses of differential gene expression were performed using DESeq2

954 (Anders and Huber, 2010, Love et al., 2014) with gene coordinates based on

955 ENSEMBL annotation (version 72). To combine the results from both siRNAs

956 targeting MATR3, we used a conditional thresholding approach, calling

957 expression changes as significant if they had an adjusted p-value < 0.01 in at

958 least one of the two depletion conditions and an adjusted p-value < 0.05 in the

959 other. Differential splicing was determined using DEXSeq (Anders et al.,

960 2012), and the two MATR3 depletion conditions were combined by conditional

961 thresholding accordingly.

962

963 **Analysis of exon inclusion in human tissues**

964 To analyse inclusion of exons across human tissues, we retrieved data on

965 mapped junctions from the V6p release of the GTEx consortium

966 (http://www.gtexportal.org/home/, (Consortium, 2015)). We used

967 UCSC/RefSeq annotation (see above) and isolated all LINE-derived exons as

968 well as Alu-exons. Then, we selected all exons from genes with at least one

969 Alu- or LINE-derived exon.  We identified junction-spanning reads to each of

970 these exons in a 2 nt grace window around the splice and used those to

971 identify the 5' and 3' splice site of the upstream and downstream exon. We

972 then calculated percent spliced in (PSI) index as the ratio of inclusion junction

973 reads (average of up+downstream junctions) to total junction reads  (average

974 of up+downstream junctions + skipping junctions), and inclusion within each

975 tissue as average of all samples. We only allowed a single exon inclusion

976 isoform across tissues (i.e. identical flanking exons) and choose the isoform

977 with more junction reads. To ensure sequencing depth and gene expression

978 were sufficient to calculate exon inclusion, we only used exons with at least

979 200 reads across the 8,555 samples (average of up+downstream junctions or

980 skipping junctions). If an exon was absent in any tissue, as judged by

981 absence of any junction spanning read and any read for the skipping junction,

982 it was treated as 'data not available' for this particular tissue. In total, we

983 covered 45,940 exons across 52 tissues and subtissues, which were adipose

984 tissue (subcutanoues and visceral omentum), adrenal glands, artery (aorta,

985    tibial and coronary artery), bladder, brain, breast, cervix (ecto- and endo-

986    cervix), colon (sigmoid and transverse), esophagus (mucosa, muscularis and

987    gastroesophageal junction), fallopian tube, heart (atrial appendage and left

988    ventricle), kidney (cortex), liver, lung, skeletal muscle, nerve tissue (amygdala,

989    anterior cingulate cortex, caudate basal ganglia, cerebellar.hemisphere,

990    cerebellum, cortex, frontal cortex, hippocampus, hypothalamus, nucleus

991    accumbens basal ganglia, putamen basal ganglia, cervical spinal cord,

992    substantia nigra, tibial), ovary, pancreas, pituitary, prostate, minor salivary

993    gland, small intestine terminal ileum), spleen, skin (suprapubic and lower leg),

994    stomach, thyroid, testis, uterus and vagina, as well as EBV transformed

995    lymphocytes and transformed fibroblasts. We did not use data from whole

996    blood, which had poor coverage on most genes. On top of the PSI index for

997    each tissue, we collated the data across tissues and computed the maximum

998    difference in PSI between the tissue(s) with highest inclusion and lowest

999    inclusion of each exon. Because testis is known to be a very promiscuously

1000    transcribed tissue (Soumillon et al., 2013) and accordingly showed many

1001    LINE-derived exons exclusively observed in the testis, we only included exons

1002    which showed at least 5% inclusion in any tissue, except testis.

1003

1004    **Classification of repeat element age by divergence or phylogenetic**

1005    **tracing**

1006    To compare the divergence of LINE insertions from their consensus

1007    sequence, we used the nucleotide difference / 1000nt, which is provided for

1008    each repeat element by the RepeatMasker table (hg19, Repeat Library

1009    20090604, (Smit et al., 1996-2010a)).

1010    For phylogenetic tracing, we tested for presence of orthologues positions with

1011    the UCSC Genome Browser LiftOver tool (Rosenbloom et al., 2015), using

1012    the respective all-chain BLASTZ files. Human and mouse LINE repeats from

1013    hg19 and mm9 RepeatMasker annotation were first lifted to hg38 and mm10.

1014    We then tested for the presence of each LINE repeat in the human and

1015    mouse lineage by retrieving orthologue genomic loci for the genomes of

1016 rhesus macaque (rheMac8), gorilla (gorGor5), mouse (mm10), rat (rn6), dog
1017 (canFam3) and cow (bosTau8). To curate the LiftOver results and safeguard
1018 against misannotation by errors in the genome lift, we cross-referenced for all
1019 liftover positions if the element overlaps with a LINE annotated by
1020 RepeatMasker for the respective genome, and only refer to the element as
1021 present in a species if at least 33% of the lifted genomic position are LINE-
1022 derived as annotated by RepeatMasker. All other elements are either
1023 'notLINE' if they were not identified by RepeatMasker, 'degenerate' if LiftOver
1024 reported them as 'partially-deleted', or 'absent' if LiftOver reported them as
1025 'deleted'. Elements from hg19 that were not 'present' in hg38 were discarded
1026 entirely. Then we converted the LiftOver annotation to phylogenetic groups
1027 after manual inspection of the liftover results in the following manner. We
1028 denoted elements as human- and primate-specific, which are 'absent' in all
1029 other species. We denoted additional elements as primate-specific, if they
1030 were either 'present', 'degenerate' or 'notLINE' in at least one of the two
1031 primate species, and 'absent' or 'notLINE' in all of the others. We denoted
1032 elements as specific for the euarchontoglires branch, if the element was
1033 'absent' or 'notLINE' in the two laurasiatherian species, and 'present' or
1034 'degenerate' in mouse or rat. The remaining elements were all lifted towards
1035 at least one of the two laurasiatherian species, and hence present in the last
1036 common ancestor of the species we surveyed. Elements present in one but
1037 absent in the other were denoted as found in 'one distant species', elements
1038 present in both as found in 'two distant species'. All remaining elements were
1039 either reported as degenerate in both species, or the liftover results were
1040 'unclear' (for example if the element was lifted to many species but did not
1041 overlap with the LINE annotation in any of those). In either case, we ignored
1042 the corresponding element for phylogenetic comparisons. Group sizes for the
1043 hg19 assembly were:

| | |
|---|---|
| Primate-specific LINE insertions | 516720 |
| Euarchontoglires-specific insertions | 64490 |
| One-distant species | 243610 |

| | |
|---|---|
| Two-distant species | 73965 |
| Sequence degenerated elements | 227587 |
| unclear liftover results | 274273 |

1044

**Statistics**

Whenever referred to in the text, *replicates* stands for biological replicates, defined as samples collected independently of one another in separated experiments. All experiments were done with biological replicates as indicated in Methods and Figure legends. In case of the iCLIP experiments from MATR3 or PTBP1 depleted cells, sequencing files were pooled across 2 biological replicates because coverage varied widely within them, and only the pooled data was used.

All statistical analyses were performed in the R software environment (version 3.1.3) or in GraphPad PRISM6. We made use of nonparametric tests in all statistical tests, since data distributions failed to conform with the assumption of normality and equal variance (homoscedasticity), assessed visually with qqnorm plots. Statistical tests are listed in figure legends. To compare multiple groups we used the Kruskal-Wallis Rank Sum test, and made pairwise comparisons with Dunn's test corrected according to Holm-Sidak, using functions implemented in the *stats* and the *dunn.test* (Dinno, 2017) R packages.

1084

## AUTHOR CONTRIBUTIONS

1085

1086    J.A., C.G., C.W.J.S and J.U. conceived the project and designed the
1087    experiments. F. A. supervised computational analysis. J.A., C.G. and A.S.
1088    performed experiments, and J.A., F.A., A.C., N.H. and W.E. performed
1089    computational analysis. J.A., F.A, C.S., N.L. and J.U. interpreted and
1090    conceptualised primary data.

1091

## DECLARATION OF INTERESTS

1092

1093    The author declare no competing interests.

1094

**REFERENCES**

ANDERS, S. & HUBER, W. 2010. Differential expression analysis for sequence count data. *Genome Biol,* 11**,** R106.

ANDERS, S., REYES, A. & HUBER, W. 2012. Detecting differential usage of exons from RNA-seq data. *Genome Res,* 22**,** 2008-17.

ATTIG, J., RUIZ DE LOS MOZOS, I., HABERMAN, N., WANG, Z., EMMETT, W., ZARNACK, K., KONIG, J. & ULE, J. 2016. Splicing repression allows the gradual emergence of new Alu-exons in primate evolution. *Elife,* 5**,** e19545.

AVENDANO-VAZQUEZ, S. E., DHIR, A., BEMBICH, S., BURATTI, E., PROUDFOOT, N. & BARALLE, F. E. 2012. Autoregulation of TDP-43 mRNA levels involves interplay between transcription, splicing, and alternative polyA site selection. *Genes Dev,* 26**,** 1679-84.

AYALA, Y. M., PANTANO, S., D'AMBROGIO, A., BURATTI, E., BRINDISI, A., MARCHETTI, C., ROMANO, M. & BARALLE, F. E. 2005. Human, Drosophila, and C.elegans TDP43: nucleic acid binding properties and splicing regulatory function. *J Mol Biol,* 348**,** 575-88.

BECK, C. R., COLLIER, P., MACFARLANE, C., MALIG, M., KIDD, J. M., EICHLER, E. E., BADGE, R. M. & MORAN, J. V. 2010. LINE-1 retrotransposition activity in human genomes. *Cell,* 141**,** 1159-70.

BELANCIO, V. P., HEDGES, D. J. & DEININGER, P. 2006. LINE-1 RNA splicing and influences on mammalian gene expression. *Nucleic Acids Res,* 34**,** 1512-21.

BROUHA, B., SCHUSTAK, J., BADGE, R. M., LUTZ-PRIGGE, S., FARLEY, A. H., MORAN, J. V. & KAZAZIAN, H. H., JR. 2003. Hot L1s account for the bulk of retrotransposition in the human population. *Proc Natl Acad Sci U S A,* 100**,** 5280-5.

COELHO, M. B., ATTIG, J., BELLORA, N., KONIG, J., HALLEGGER, M., KAYIKCI, M., EYRAS, E., ULE, J. & SMITH, C. W. 2015. Nuclear matrix protein Matrin3 regulates alternative splicing and forms overlapping regulatory networks with PTB. *EMBO J,* 34**,** 653-668.

CONSORTIUM, G. T. 2015. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. *Science,* 348**,** 648-60.

DAMIANOV, A., YING, Y., LIN, C. H., LEE, J. A., TRAN, D., VASHISHT, A. A., BAHRAMI-SAMANI, E., XING, Y., MARTIN, K. C., WOHLSCHLEGEL, J. A. & BLACK, D. L. 2016. Rbfox Proteins Regulate Splicing as Part of a Large Multiprotein Complex LASR. *Cell,* 165**,** 606-19.

DEININGER, P. L. & BATZER, M. A. 2002. Mammalian retroelements. *Genome Res,* 12**,** 1455-65.

DIGNAM, J. D., LEBOVITZ, R. M. & ROEDER, R. G. 1983. Accurate transcription initiation by RNA polymerase II in a soluble extract from isolated mammalian nuclei. *Nucleic Acids Res,* 11**,** 1475-89.

DINNO, A. 2017. dunn.test: Dunn's Test of Multiple Comparisons Using Rank Sums. R

  package version 1.3.4.

1143 DOBIN, A., DAVIS, C. A., SCHLESINGER, F., DRENKOW, J., ZALESKI, C.,
1144      JHA, S., BATUT, P., CHAISSON, M. & GINGERAS, T. R. 2013. STAR:
1145      ultrafast universal RNA-seq aligner. *Bioinformatics,* 29**,** 15-21.
1146 EOM, T., ZHANG, C., WANG, H., LAY, K., FAK, J., NOEBELS, J. L. &
1147      DARNELL, R. B. 2013. NOVA-dependent regulation of cryptic NMD
1148      exons controls synaptic protein levels after seizure. *Elife,* 2**,** e00178.
1149 GE, Z., QUEK, B. L., BEEMON, K. L. & HOGG, J. R. 2016. Polypyrimidine
1150      tract binding protein 1 protects mRNAs from recognition by the
1151      nonsense-mediated mRNA decay pathway. *Elife,* 5.
1152 GOODIER, J. L., CHEUNG, L. E. & KAZAZIAN, H. H., JR. 2012. MOV10 RNA
1153      helicase is a potent inhibitor of retrotransposition in cells. *PLoS Genet,*
1154      8**,** e1002941.
1155 GOODIER, J. L., CHEUNG, L. E. & KAZAZIAN, H. H., JR. 2013. Mapping the
1156      LINE1 ORF1 protein interactome reveals associated inhibitors of
1157      human retrotransposition. *Nucleic Acids Res,* 41**,** 7401-19.
1158 GOODING, C., ROBERTS, G. C. & SMITH, C. W. 1998. Role of an inhibitory
1159      pyrimidine element and polypyrimidine tract binding protein in
1160      repression of a regulated alpha-tropomyosin exon. *RNA,* 4**,** 85-100.
1161 HABERMAN, N., HUPPERTZ, I., ATTIG, J., KONIG, J., WANG, Z., HAUER,
1162      C., HENTZE, M. W., KULOZIK, A. E., LE HIR, H., CURK, T., SIBLEY,
1163      C. R., ZARNACK, K. & ULE, J. 2017. Insights into the design and
1164      interpretation of iCLIP experiments. *Genome Biol,* 18**,** 7.
1165 HAN, J. S., SZAK, S. T. & BOEKE, J. D. 2004. Transcriptional disruption by
1166      the L1 retrotransposon and implications for mammalian transcriptomes.
1167      *Nature,* 429**,** 268-74.
1168 HARTLEY, S. W. & MULLIKIN, J. C. 2016. Detection and visualization of
1169      differential splicing in RNA-Seq data with JunctionSeq. *Nucleic Acids*
1170      *Res,* 44**,** e127.
1171 HUANG, C. R., SCHNEIDER, A. M., LU, Y., NIRANJAN, T., SHEN, P.,
1172      ROBINSON, M. A., STERANKA, J. P., VALLE, D., CIVIN, C. I., WANG,
1173      T., WHEELAN, S. J., JI, H., BOEKE, J. D. & BURNS, K. H. 2010.
1174      Mobile interspersed repeats are major structural variants in the human
1175      genome. *Cell,* 141**,** 1171-82.
1176 HUPPERTZ, I., ATTIG, J., D'AMBROGIO, A., EASTON, L. E., SIBLEY, C. R.,
1177      SUGIMOTO, Y., TAJNIK, M., KONIG, J. & ULE, J. 2014. iCLIP:
1178      protein-RNA interactions at nucleotide resolution. *Methods,* 65**,** 274-87.
1179 IMBEAULT, M., HELLEBOID, P. Y. & TRONO, D. 2017. KRAB zinc-finger
1180      proteins contribute to the evolution of gene regulatory networks.
1181      *Nature,* 543**,** 550-554.
1182 JIN, Y., TAM, O. H., PANIAGUA, E. & HAMMELL, M. 2015. TEtranscripts: a
1183      package for including transposable elements in differential expression
1184      analysis of RNA-seq datasets. *Bioinformatics*.
1185 JURKA, J. 1998. Repeats in genomic DNA: mining and meaning. *Curr Opin*
1186      *Struct Biol,* 8**,** 333-7.
1187 KELLEY, D. R., HENDRICKSON, D. G., TENEN, D. & RINN, J. L. 2014.
1188      Transposable elements modulate human RNA abundance and splicing
1189      via specific RNA-protein interactions. *Genome Biol,* 15**,** 537.
1190 KIM, D., PERTEA, G., TRAPNELL, C., PIMENTEL, H., KELLEY, R. &
1191      SALZBERG, S. L. 2013. TopHat2: accurate alignment of
1192      transcriptomes in the presence of insertions, deletions and gene

1193        fusions. *Genome Biol,* 14**,** R36.
1194 KULA, A., GUERRA, J., KNEZEVICH, A., KLEVA, D., MYERS, M. P. &
1195        MARCELLO, A. 2011. Characterization of the HIV-1 RNA associated
1196        proteome identifies Matrin 3 as a nuclear cofactor of Rev function.
1197        *Retrovirology,* 8**,** 60.
1198 LEE, J. Y., JI, Z. & TIAN, B. 2008. Phylogenetic analysis of mRNA
1199        polyadenylation sites reveals a role of transposable elements in
1200        evolution of the 3'-end of genes. *Nucleic Acids Res,* 36**,** 5581-90.
1201 LI, Q., ZHENG, S., HAN, A., LIN, C. H., STOILOV, P., FU, X. D. & BLACK, D.
1202        L. 2014. The splicing regulator PTBP2 controls a program of embryonic
1203        splicing required for neuronal maturation. *Elife,* 3**,** e01201.
1204 LIN, Y., PROTTER, D. S., ROSEN, M. K. & PARKER, R. 2015. Formation and
1205        Maturation of Phase-Separated Liquid Droplets by RNA-Binding
1206        Proteins. *Mol Cell,* 60**,** 208-19.
1207 LING, J. P., CHHABRA, R., MERRAN, J. D., SCHAUGHENCY, P. M.,
1208        WHEELAN, S. J., CORDEN, J. L. & WONG, P. C. 2016. PTBP1 and
1209        PTBP2 Repress Nonconserved Cryptic Exons. *Cell Rep,* 17**,** 104-113.
1210 LOVE, M. I., HUBER, W. & ANDERS, S. 2014. Moderated estimation of fold
1211        change and dispersion for RNA-seq data with DESeq2. *Genome Biol,*
1212        15**,** 550.
1213 MEISCHL, C., BOER, M., AHLIN, A. & ROOS, D. 2000. A new exon created
1214        by intronic insertion of a rearranged LINE-1 element as the cause of
1215        chronic granulomatous disease. *Eur J Hum Genet,* 8**,** 697-703.
1216 MERKIN, J. J., CHEN, P., ALEXIS, M. S., HAUTANIEMI, S. K. & BURGE, C.
1217        B. 2015. Origins and impacts of new mammalian exons. *Cell Rep,* 10**,**
1218        1992-2005.
1219 MGI:1298379: International Mouse Phenotyping consortium database.
1220        Accessed last 03/10/2017.
1221 MODREK, B. & LEE, C. J. 2003. Alternative splicing in the human, mouse and
1222        rat genomes is associated with an increased frequency of exon
1223        creation and/or loss. *Nat Genet,* 34**,** 177-80.
1224 O'LEARY, M. A., BLOCH, J. I., FLYNN, J. J., GAUDIN, T. J.,
1225        GIALLOMBARDO, A., GIANNINI, N. P., GOLDBERG, S. L., KRAATZ,
1226        B. P., LUO, Z. X., MENG, J., NI, X., NOVACEK, M. J., PERINI, F. A.,
1227        RANDALL, Z. S., ROUGIER, G. W., SARGIS, E. J., SILCOX, M. T.,
1228        SIMMONS, N. B., SPAULDING, M., VELAZCO, P. M., WEKSLER, M.,
1229        WIBLE, J. R. & CIRRANELLO, A. L. 2013. The placental mammal
1230        ancestor and the post-K-Pg radiation of placentals. *Science,* 339**,** 662-
1231        7.
1232 OBERSTRASS, F. C., AUWETER, S. D., ERAT, M., HARGOUS, Y.,
1233        HENNING, A., WENTER, P., REYMOND, L., AMIR-AHMADY, B.,
1234        PITSCH, S., BLACK, D. L. & ALLAIN, F. H. 2005. Structure of PTB
1235        bound to RNA: specific binding and implications for splicing regulation.
1236        *Science,* 309**,** 2054-7.
1237 PAYNE, J. L. & WAGNER, A. 2015. Mechanisms of mutational robustness in
1238        transcriptional regulation. *Front Genet,* 6**,** 322.
1239 PETRYSZAK, R., KEAYS, M., TANG, Y. A., FONSECA, N. A., BARRERA, E.,
1240        BURDETT, T., FULLGRABE, A., FUENTES, A. M., JUPP, S.,
1241        KOSKINEN, S., MANNION, O., HUERTA, L., MEGY, K., SNOW, C.,
1242        WILLIAMS, E., BARZINE, M., HASTINGS, E., WEISSER, H.,

1243      WRIGHT, J., JAISWAL, P., HUBER, W., CHOUDHARY, J.,
1244      PARKINSON, H. E. & BRAZMA, A. 2016. Expression Atlas update--an
1245      integrated database of gene and protein expression in humans,
1246      animals and plants. *Nucleic Acids Res,* 44**,** D746-52.
1247 RAY, D., KAZAN, H., COOK, K. B., WEIRAUCH, M. T., NAJAFABADI, H. S.,
1248      LI, X., GUEROUSSOV, S., ALBU, M., ZHENG, H., YANG, A., NA, H.,
1249      IRIMIA, M., MATZAT, L. H., DALE, R. K., SMITH, S. A., YAROSH, C.
1250      A., KELLY, S. M., NABET, B., MECENAS, D., LI, W., LAISHRAM, R.
1251      S., QIAO, M., LIPSHITZ, H. D., PIANO, F., CORBETT, A. H.,
1252      CARSTENS, R. P., FREY, B. J., ANDERSON, R. A., LYNCH, K. W.,
1253      PENALVA, L. O., LEI, E. P., FRASER, A. G., BLENCOWE, B. J.,
1254      MORRIS, Q. D. & HUGHES, T. R. 2013. A compendium of RNA-
1255      binding motifs for decoding gene regulation. *Nature,* 499**,** 172-7.
1256 REED, R. 2000. Mechanisms of fidelity in pre-mRNA splicing. *Curr Opin Cell*
1257      *Biol,* 12**,** 340-5.
1258 REYES, J. C., MUCHARDT, C. & YANIV, M. 1997. Components of the human
1259      SWI/SNF complex are enriched in active chromatin and are associated
1260      with the nuclear matrix. *The Journal of cell biology,* 137**,** 263-74.
1261 ROSENBLOOM, K. R., ARMSTRONG, J., BARBER, G. P., CASPER, J.,
1262      CLAWSON, H., DIEKHANS, M., DRESZER, T. R., FUJITA, P. A.,
1263      GURUVADOO, L., HAEUSSLER, M., HARTE, R. A., HEITNER, S.,
1264      HICKEY, G., HINRICHS, A. S., HUBLEY, R., KAROLCHIK, D.,
1265      LEARNED, K., LEE, B. T., LI, C. H., MIGA, K. H., NGUYEN, N.,
1266      PATEN, B., RANEY, B. J., SMIT, A. F., SPEIR, M. L., ZWEIG, A. S.,
1267      HAUSSLER, D., KUHN, R. M. & KENT, W. J. 2015. The UCSC
1268      Genome Browser database: 2015 update. *Nucleic Acids Res,* 43**,**
1269      D670-81.
1270 ROT, G., WANG, Z., HUPPERTZ, I., MODIC, M., LENCE, T., HALLEGGER,
1271      M., HABERMAN, N., CURK, T., VON MERING, C. & ULE, J. 2017.
1272      High-Resolution RNA Maps Suggest Common Principles of Splicing
1273      and Polyadenylation Regulation by TDP-43. *Cell Rep,* 19**,** 1056-1067.
1274 SCHWAHN, U., LENZNER, S., DONG, J., FEIL, S., HINZMANN, B., VAN
1275      DUIJNHOVEN, G., KIRSCHNER, R., HEMBERGER, M., BERGEN, A.
1276      A., ROSENBERG, T., PINCKERS, A. J., FUNDELE, R., ROSENTHAL,
1277      A., CREMERS, F. P., ROPERS, H. H. & BERGER, W. 1998. Positional
1278      cloning of the gene for X-linked retinitis pigmentosa 2. *Nat Genet,* 19**,**
1279      327-32.
1280 SEMLOW, D. R. & STALEY, J. P. 2012. Staying on message: ensuring fidelity
1281      in pre-mRNA splicing. *Trends Biochem Sci,* 37**,** 263-73.
1282 SIBLEY, C. R., BLAZQUEZ, L. & ULE, J. 2016. Lessons from non-canonical
1283      splicing. *Nat Rev Genet,* 17**,** 407-21.
1284 SLOAN, C. A., CHAN, E. T., DAVIDSON, J. M., MALLADI, V. S., STRATTAN,
1285      J. S., HITZ, B. C., GABDANK, I., NARAYANAN, A. K., HO, M., LEE, B.
1286      T., ROWE, L. D., DRESZER, T. R., ROE, G., PODDUTURI, N. R.,
1287      TANAKA, F., HONG, E. L. & CHERRY, J. M. 2016. ENCODE data at
1288      the ENCODE portal. *Nucleic Acids Res,* 44**,** D726-32.
1289 SMIT, A., HUBLEY, R. & GREEN, P. 1996-2010a. RepeatMasker Open-3.0.
1290      http://www.repeatmasker.org.
1291 SMIT, A. F., HUBLEY, R. & GREEN, P. 1996-2010b. RepeatMasker Open-
1292      3.0. http://www.repeatmasker.org.

SOUMILLON, M., NECSULEA, A., WEIER, M., BRAWAND, D., ZHANG, X., GU, H., BARTHES, P., KOKKINAKI, M., NEF, S., GNIRKE, A., DYM, M., DE MASSY, B., MIKKELSEN, T. S. & KAESSMANN, H. 2013. Cellular source and mechanisms of high transcriptome complexity in the mammalian testis. *Cell Rep,* 3**,** 2179-90.

TAYLOR, M. S., LACAVA, J., MITA, P., MOLLOY, K. R., HUANG, C. R., LI, D., ADNEY, E. M., JIANG, H., BURNS, K. H., CHAIT, B. T., ROUT, M. P., BOEKE, J. D. & DAI, L. 2013. Affinity proteomics reveals human host factors implicated in discrete stages of LINE-1 retrotransposition. *Cell,* 155**,** 1034-48.

THOMAS, J. H. & SCHNEIDER, S. 2011. Coevolution of retroelements and tandem zinc finger genes. *Genome Res,* 21**,** 1800-12.

TIROSH, I., REIKHAV, S., SIGAL, N., ASSIA, Y. & BARKAI, N. 2010. Chromatin regulators as capacitors of interspecies variations in gene expression. *Mol Syst Biol,* 6**,** 435.

TRAPNELL, C., ROBERTS, A., GOFF, L., PERTEA, G., KIM, D., KELLEY, D. R., PIMENTEL, H., SALZBERG, S. L., RINN, J. L. & PACHTER, L. 2012. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc,* 7**,** 562-78.

VAN NOSTRAND, E. L., FREESE, P., PRATT, G. A., WANG, X., WEI, X., BLUE, S. M., DOMINGUEZ, D., CODY, N. A. L., OLSON, S., SUNDARARAMAN, B., XIAO, R., ZHAN, L., BAZILE, C., BENOIT BOUVRETTE, L. P., CHEN, J., DUFF, M. O., GARCIA, K., GELBOIN-BURKHART, C., HOCHMAN, A., LAMBERT, N. J., LI, H., NGUYEN, T. B., PALDEN, T., RABANO, I., SATHE, S., STANTON, R., LOUIE, A. L., AIGNER, S., BERGALET, J., ZHOU, B., SU, A., WANG, R., YEE, B. A., FU, X.-D., LECUYER, E., BURGE, C. B., GRAVELEY, B. & YEO, G. W. 2017. A Large-Scale Binding and Functional Map of Human RNA Binding Proteins. *bioRxiv*.

VORECHOVSKY, I. 2010. Transposable elements in disease-associated cryptic exons. *Hum Genet,* 127**,** 135-54.

VUONG, J. K., LIN, C. H., ZHANG, M., CHEN, L., BLACK, D. L. & ZHENG, S. 2016. PTBP1 and PTBP2 Serve Both Specific and Redundant Functions in Neuronal Pre-mRNA Splicing. *Cell Rep,* 17**,** 2766-2775.

WARD, M. C., WILSON, M. D., BARBOSA-MORAIS, N. L., SCHMIDT, D., STARK, R., PAN, Q., SCHWALIE, P. C., MENON, S., LUKK, M., WATT, S., THYBERT, D., KUTTER, C., KIRSCHNER, K., FLICEK, P., BLENCOWE, B. J. & ODOM, D. T. 2013. Latent regulatory potential of human-specific repetitive elements. *Mol Cell,* 49**,** 262-72.

WITTEN, J. T. & ULE, J. 2011. Understanding splicing regulation through RNA splicing maps. *Trends Genet,* 27**,** 89-97.

XING, Y. & LEE, C. 2006. Alternative splicing and RNA selection pressure--evolutionary consequences for eukaryotic genomes. *Nat Rev Genet,* 7**,** 499-509.

YEO, G. & BURGE, C. B. 2004. Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *J Comput Biol,* 11**,** 377-94.

YOSHIDA, K., NAKAMURA, A., YAZAKI, M., IKEDA, S. & TAKEDA, S. 1998. Insertional mutation by transposable element, L1, in the DMD gene results in X-linked dilated cardiomyopathy. *Hum Mol Genet,* 7**,** 1129-

1343        32.
1344    ZEITZ, M. J., MALYAVANTHAM, K. S., SEIFERT, B. & BEREZNEY, R. 2009.
1345        Matrin 3: chromosomal distribution and protein interactions. *J Cell*
1346        *Biochem,* 108**,** 125-33.
1347    ZHANG, Y., ROMANISH, M. T. & MAGER, D. L. 2011. Distributions of
1348        transposable elements reveal hazardous zones in mammalian introns.
1349        *PLoS Comput Biol,* 7**,** e1002046.
1350    ZHANG, Z. & CARMICHAEL, G. G. 2001. The fate of dsRNA in the nucleus: a
1351        p54(nrb)-containing complex mediates the nuclear retention of
1352        promiscuously A-to-I edited RNAs. *Cell,* 106**,** 465-75.
1353

1354

1355    **FIGURES AND FIGURE LEGENDS**

1356

**Figure 1**

**Figure 1: LINEs are binding platforms for a set of RBPs.**

(A) Estimate of abundance of L1-sequences in cytoplasmic and nuclear RNA fractions from HeLa, K562 and HepG2 cells. Strand-specific RNAseq was used to quantify abundance of L1 in sense and antisense (colored in orange and blue), relative to number of mapped reads. Data is split for libraries made from polyA-, polyA+ or rRNA-RNA. Data for K562 and HepG2 is from the ENCODE consortium. Data for HeLa is from replicates, and bargraph show mean ± s.d.m.

(B) Frequency of L1 repeat sequences among the bound RNA sequences of a panel of RBPs. For each RBP, all cDNAs recovered in an iCLIP or eCLIP experiment were counted if they mapped at least partially to a L1 element. Since e/iCLIP is strand-specific, binding to LINEs transcribed in sense or in antisense was quantified separately, coloured in orange and blue. The orange and blue lines indicate the average binding across all RBPs (median). The iCLIP data was derived either from HeLa cells or from HEK293 FlpIN cells, and the eCLIP data from K562 and HepG2 cells. This information and the full data set is available in Suppl. Table 1, together with the source of each data set. For visualisation, replicates were averaged and only data from one cell line is shown.

(C) Binding to introns of at least 7kb size was analysed in 100nt bins up to 5kb upstream and downstream of the exon, and quantified in percent relative to the total number of mapped reads. Data is shown for the first 100nt and as an average of bins 101-500nt, 501-2000nt and 2001-5000nt. A rank for deep intronic binding is given based on the average of the first 100nt of either splice site and average binding in the 20001-5000nt window.

46

# Figure 2

1391 **Figure 2: Binding of PTBP1 to antisense L1 elements is**

1392 **MATR3-dependent.**

1393 PTBP1 iCLIP was performed from HEK293T cells depleted of MATR3, PTBP1

1394 as well as controls.

1395 (A) TOP: $^{32}$P labelled RNA crosslinked to and co-precipitated with PTBP1

1396 under high RNase conditions. MIDDLE: To quantify the signal, grey

1397 pixel intensity if shown across the centre of each lane, analysed with

1398 ImageJ software. BOTTOM: The input lysate for the iCLIP experiment

1399 was probed for MATR3 and PTBP1 antibodies in a Western Blot to

1400 ensure reduced signal is not due to changes in protein abundance.

1401 Samples are the same as in the radiogram, but the gel image was cut

1402 to align them. Note replicates are shown in Fig. S2A.

1403 (B) PTBP1 binding peaks were identified from all iCLIP experiments, and

1404 classified according to susceptibility to MATR3 depletion as indicated

1405 based on moderated log2 fold changes. Binding peaks with a

1406 normalised count of less than 8 were ignored, indicated by the dotted

1407 line.

1408 (C) PTBP1 binding peaks susceptible to MATR3 depletion are shorter

1409 than those which are not.

1410 (D) MATR3 iCLIP is enriched around MATR3-dependent PTBP1 binding

1411 peaks.

1412 (E) Enrichment for high-affinity motifs around PTBP1 binding peaks.

1413 LEFT: all PTBP1 binding peaks show strong enrichment for PTBP

1414 binding motifs. RIGHT: MATR3-dependent PTBP1 binding peaks

1415 show enrichment in a 200nt region for high-affinity motifs above other

1416 PTBP1 binding peaks.

1417 (F) The overlap between the centre of PTBP1 binding peaks and different

1418 repeat classes was tested for antisense L1 elements, sense L2

1419 elements, and sense CT-/T-rich microsatellite repeats. Metaprofile

1420 shows percent of each class of clusters overlapping with each

1421 genomic element. MATR3-dependent binding peaks are more

1422 frequently derived from an antisense L1 element than MATR3-
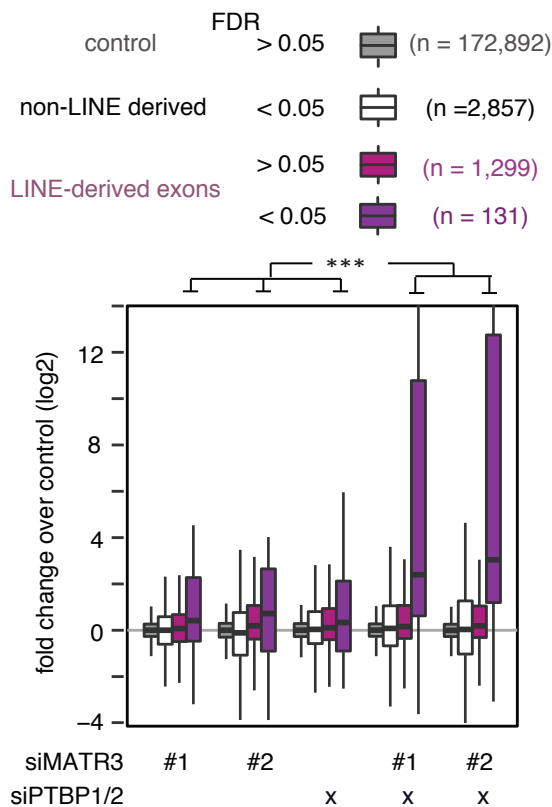
47

1423    independent once.

1424 (G) Protein-protein interaction between MATR3 and PTBP1 allows

1425    formation of a heteromeric complex on a substrate RNA with two

1426    ATGTT motifs *in vitro*. Recombinant PTBP1 (rPTBP1) and different

1427    MATR3 mutants (rMATR3) were crosslinked to the same RNA at

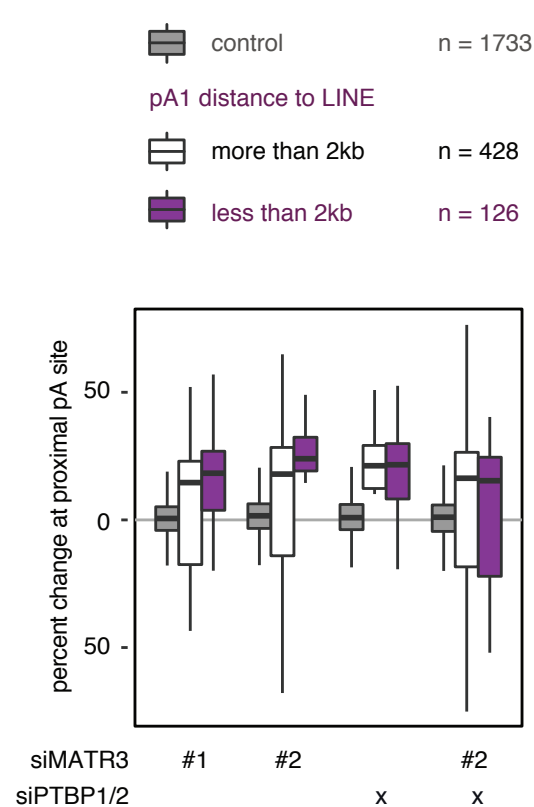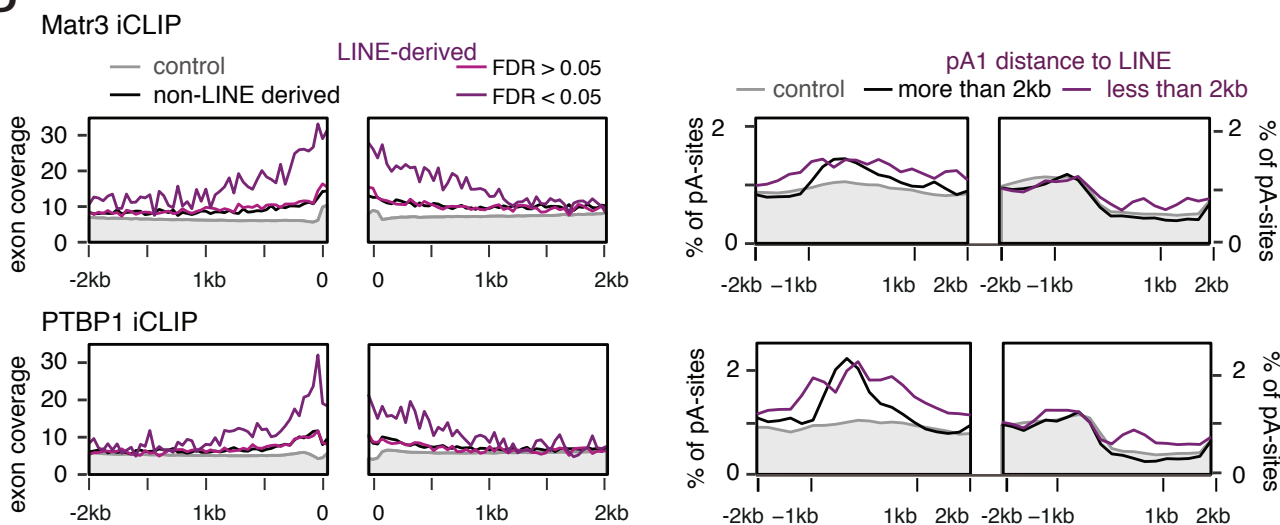1428    different MATR3 molarity (rPTBP1 at 0.5µM).

1429

**Figure 3**

1430  **Figure 3: MATR3 and PTBP1 repress usage of cryptic splice and**

1431  **polyA-sites in vicinity to LINEs.**

1432  (A)  The metaprofile shows the coverage of antisense L1 sequences in a

1433  ±2kb window flanking the splice sites and the proximal and distal

1434  polyA sites of MATR3/PTBP1/2 repressed events. Exon usage and

1435  polyA-site usage was analysed in cells depleted of MATR3 and

1436  PTBP1/2, individually or in combination, and events significantly

1437  increased in absence of either proteins were selected. Misregulated

1438  exons are alternative exons selected from a splice-array experiment

1439  (Coelho et al., 2015), polyA site pairs are from mRNA 3'end

1440  sequencing experiments. Controls are non-significant events site with

1441  no appreciable change (below 10%) and reflect the expected genomic

1442  frequency of L1 antisense sequence (shown in grey). Metaprofile was

1443  smoothed using 40 nucleotide bins.

1444  (B)  The transcriptome was *de novo* assembled from cells depleted of

1445  MATR3 and PTBP1/2, individually or in combination, in order to

1446  capture cryptic LINE-derived exons absent from microarrays. For

1447  each condition, the log2 fold changes of MATR3/PTBP1 regulated

1448  exons are plotted. Only events with at least one junction-spanning

1449  read were considered for analysis, and significant and non-significant

1450  LINE-derived exons are shown separately (at FDR < 0.05).

1451  Differences between the changes in exon abundance across groups

1452  were tested by Kruskal-Wallis Rank Sum test (p-value < $2.2e^{-16}$), and

1453  pairwise comparisons within each condition were tested with a two-

1454  sided Wilcoxon Rank Sum test, and corrected for multiple testing

1455  according to Bonferroni. Adjusted p-value indicated by *** was below

1456  0.0001. Whiskers are cut-off from the boxplot for visualisation, but

1457  data distribution extends to (ymax +24) in cells depleted of MATR3

1458  and PTBP1/2 simultaneously.

1459  (C)  Percent change in usage of the proximal polyA-sites (same as in A).

1460  Misregulated pA-sites are split into those within 2kb vicinity of a LINE
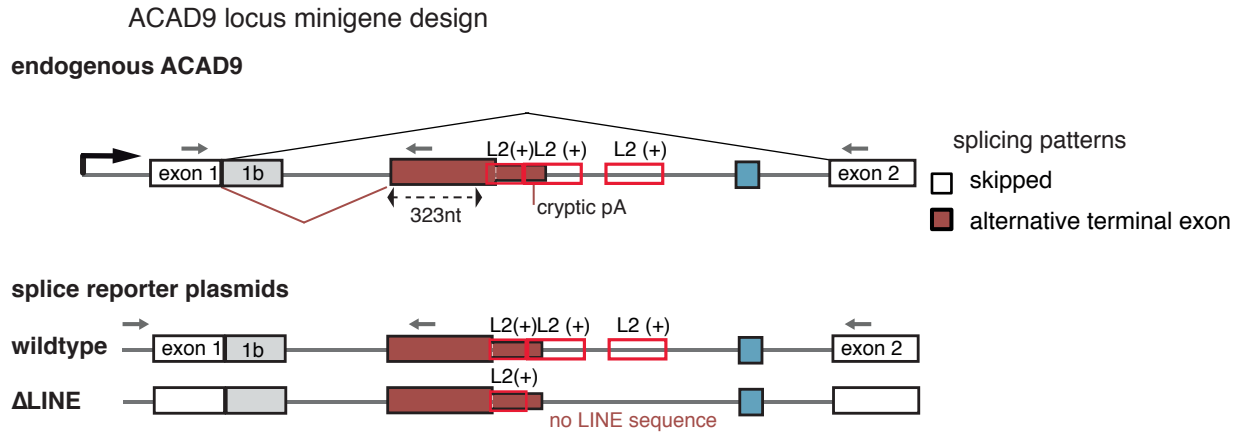
1461  and those which are not.

1462  (D)  Metaprofile of MATR3 and PTBP1 iCLIP binding across the splice and

1463  polyA sites ±2kb of the regulated event. Events were selected and

49

1464        grouped as in (B) and (C). iCLIP binding is presented as percentage

1465        of occupancy, and was smoothed using 40 nucleotide bins.

1466        Occupancy on non-regulated sites is shown as control (in grey).
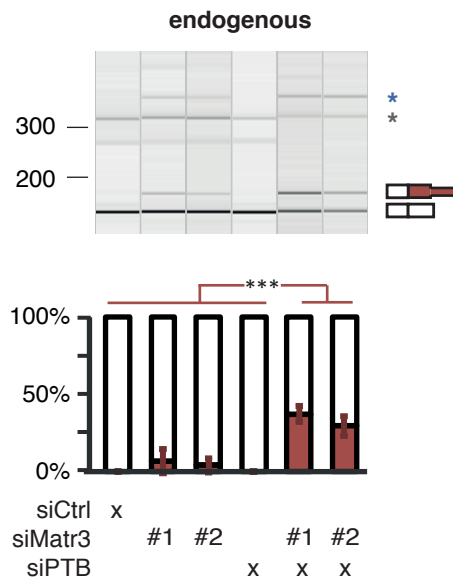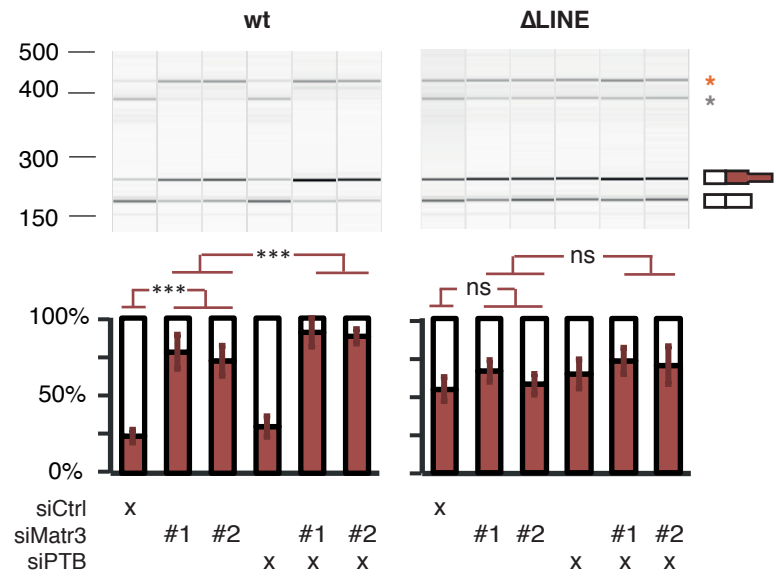
1467

## Figure 4

**A**

ACAD9 locus minigene design



**B**



**C**

1468 **Figure 4: Partial deletion of L2 sequences disrupts splicing repression**

1469 **of *ACAD9* by MATR3/PTBP1.**

1470 (A) Schematic illustrating the endogenous ACAD9 locus and the ACAD9
1471 splice reporter. The first two exons and the complete intron1 were
1472 cloned into a CMV driven reporter plasmid. In the ΔLINE splice
1473 reporter 499 base pairs of L2 sequence were replaced by non-
1474 repetitive sequence of intron2 of ACAD9.

1475 (B) The inclusion level of the LINE-proximal alternative terminal exon in
1476 endogenous ACAD9 was measured in total RNA of cells depleted of
1477 MATR3 and PTBP1/2 individually or in combination as well as
1478 controls. To test for significance, one-way ANOVA was used coupled
1479 with multiple comparison correction according to Tukey's HSD. ***
1480 indicates p-value below 0.001. Semi-quantitative RT-PCR analysis is
1481 averaged across three independent replicates, error bars represent
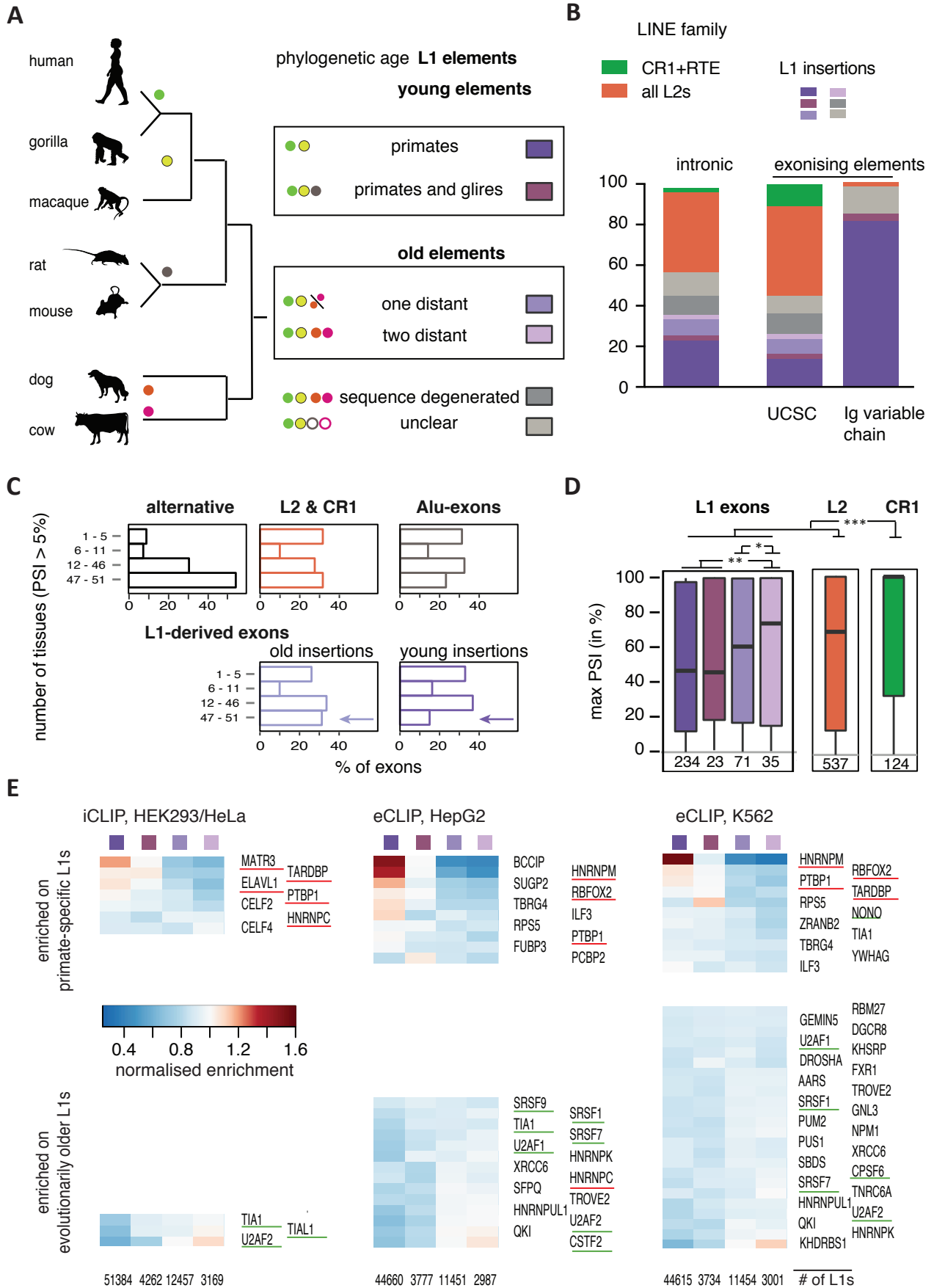1482 s.d.m.

1483 (C) The inclusion level of the LINE-derived exon was measured as in (B)
1484 in the wild-type and ΔLINE ACAD9 splice reporter.

1485 (B) and (C) Additional splice products are indicated by asterisks. These
1486 use the 5' splice site of exon 1b.

1487

1488

51

**Figure 5**

1489 **Figure 5: LINE-derived exons are a source of primate-specific alternative**

1490 **exons.**

1491 Percent splice index (PSI) was calculated in the GTEx panel of human tissue

1492 samples for LINE-derived exons annotated in UCSC (relative to the flanking

1493 exons). Inclusion levels range from 0 to 100%, showing no inclusion or full

1494 inclusion. If no support for expression of the flanking exons was found, the

1495 gene was assumed to be non-expressed.

1496 (A) The phylogenetic age of each LINE element in the human genome

1497 was mapped by comparison to the gorilla, rhesus macaque, mouse,

1498 rat, dog and cow genome assemblies using UCSC liftover genome

1499 alignments overlaid with RepeatMasker annotation (see Methods for

1500 details). Elements specific to the primate or euarchontoglires lineage

1501 are considered evolutionarily young elements, while elements present

1502 in cow and dog are considered old elements.

1503 (B) The phylogenetic age of a LINE element gave an estimate of the

1504 genomic age of each LINE-derived exon. UCSC annotated exons are

1505 generally of the youngest elements. Within UCSC, the Ig-encoding

1506 region (*abParts*) stands out with 1,152 out of 6,012 annotated LINE-

1507 derived exons, which are frequently primate-specific.

1508 (C) Exons derived from evolutionarily young L1 elements are rarely

1509 present across human tissue subtypes. We determined the number of

1510 tissues in which each exon was detectable (at PSI > 5%)  and

1511 compared repeat-derived exons to non-repeat derived known

1512 alternative exons.

1513 (D) Maximum inclusion in any tissue correlates with the genomic age of

1514 L1-derived exons. Significance was tested across groups by

1515 Kruskal-Wallis' Rank Sum test and pairwise comparisons by Dunn's

1516 test corrected according to Holm-Šidák. *, ** and *** indicate adjusted

1517 p-value was below 0.05, 0.01 and 0.001, respectively.

1518 (E) RBPs show preferences for binding to L1 elements of different

1519 evolutionary ages. The L1 elements with 10% highest coverage

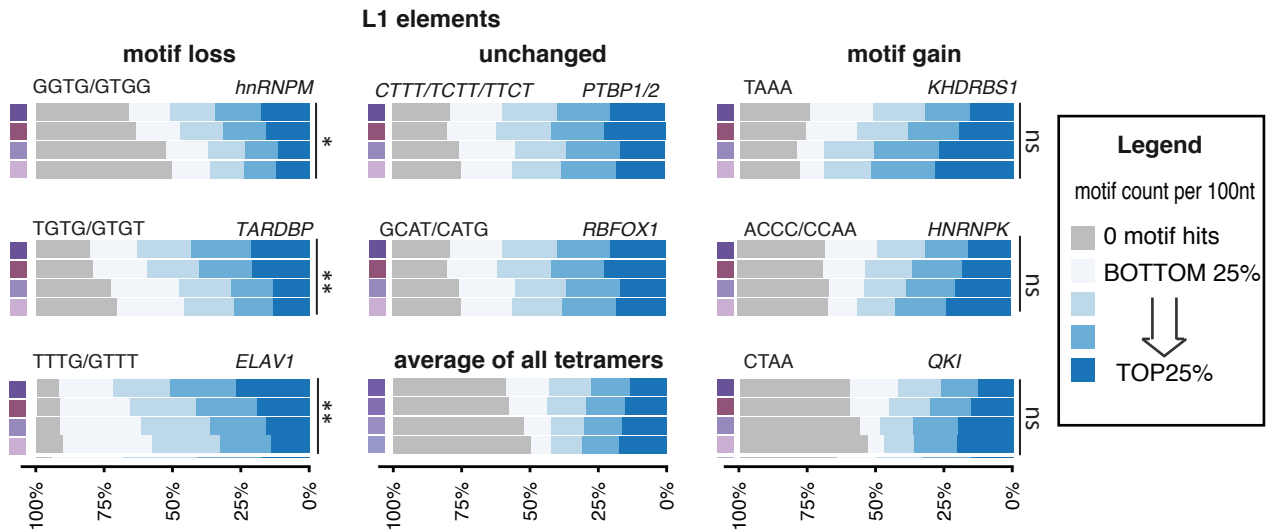1520 across any i/eCLIP data were used to calculate a normalised

52

1521    coverage for each RBP, and the number of L1 elements in each group

1522    is given at the bottom. Binding of each RBP was normalised by the

1523    sum of all RBPs within each cell line on an individual L1 element to

1524    obtain a relative binding estimate, and for visualisation of binding

1525    preference, normalised enrichment of each RBP was calculated by
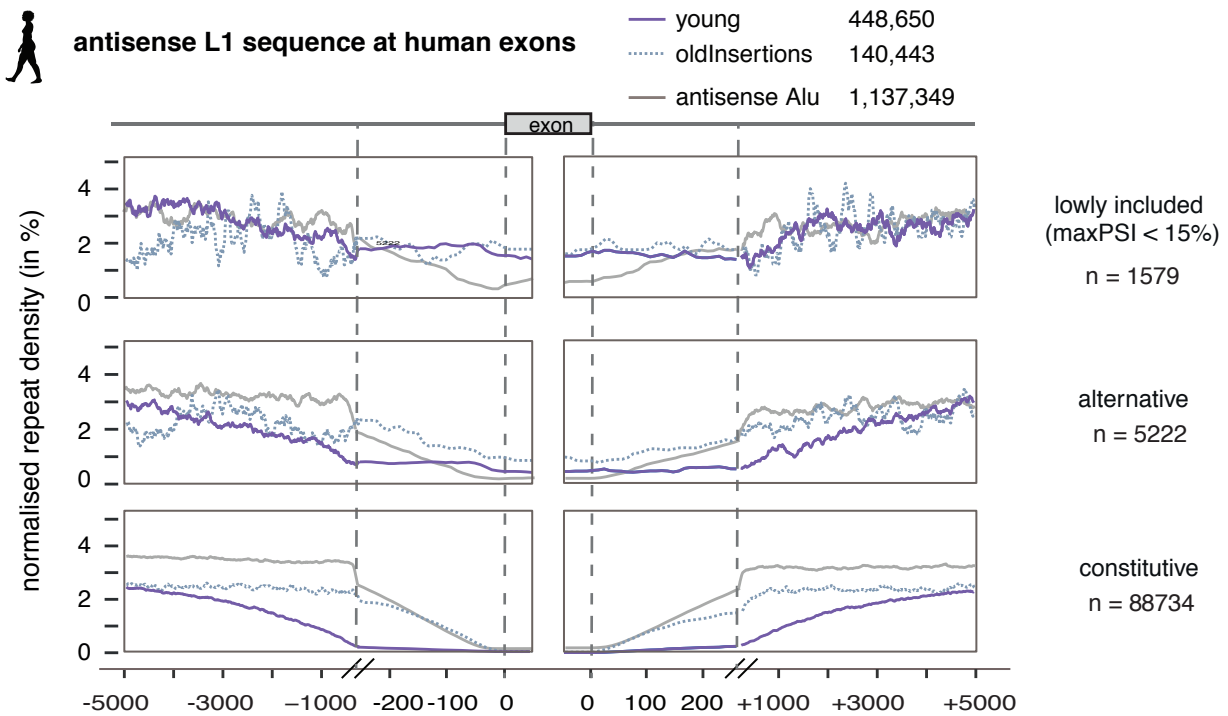
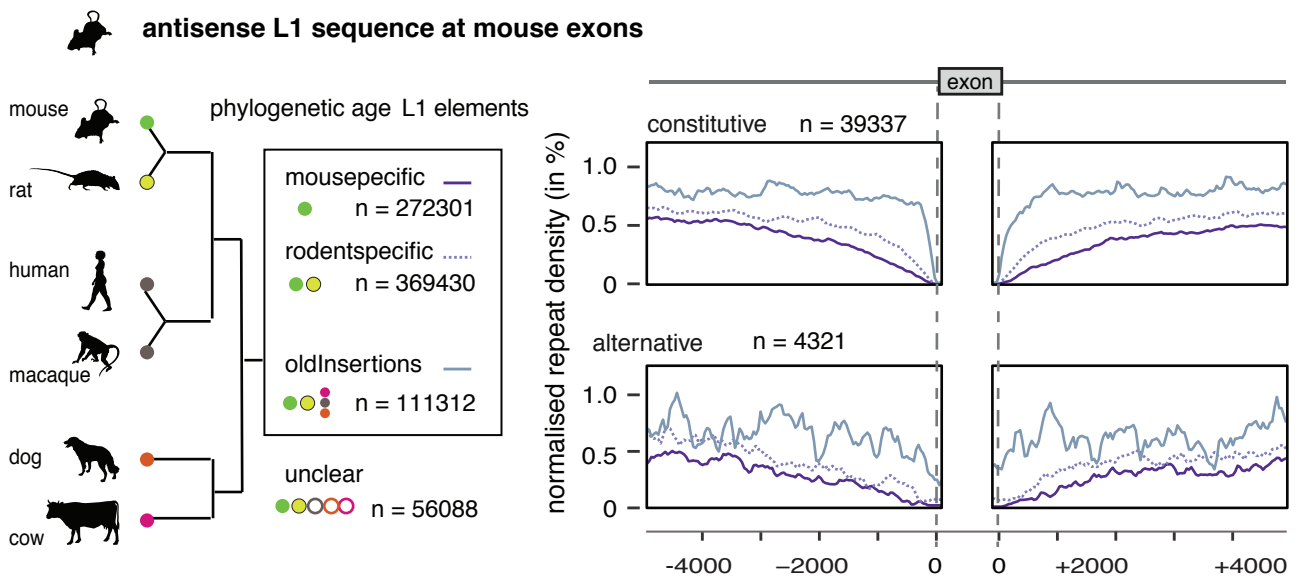1526    normalising to the mean.

1527

1528

## Figure 6

1529 **Figure 6: Young L1 elements are rich in splice repressor binding motifs**

1530 **and selected against at exons in a broad window.**

1531   (A)   The number of binding motifs associated with different RBPs is shown
1532        for antisense L1 sequences of different genomic age. Binding motifs
1533        of RBPs shown in Figure 5E were identified from literature where
1534        possible and searched for in antisense L1 elements. The genomic age
1535        of L1 elements is defined as in Figure 5A. Total motif count per 100nt
1536        was determined and categorised as quartiles (bottom to top 25% and
1537        0 motifs, see legend). For comparison, the average distribution of all
1538        possible tetramers is shown. Changes in motifs counts with
1539        evolutionary age of the elements were considered significant based
1540        on their empirical distribution (see Methods for details)., ** and *
1541        indicates FDR below 0.05 and  0.1, respectively; ns = not significant.
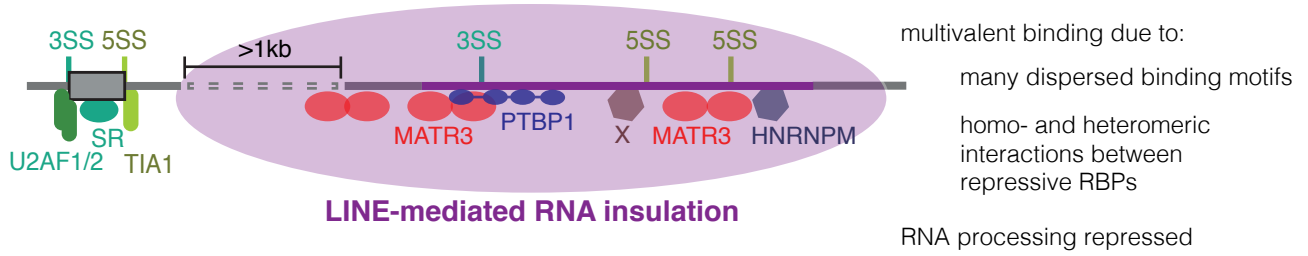
1542   (B)   Density profiles showing L1 antisense sequence 5kb upstream and
1543        downstream of human exons. L1s were split for evolutionary young
1544        and old insertions and repeat density is normalised to the total
1545        number of repeats in the two groups. For comparison, the primate-
1546        specific Alu insertions are shown. Exons were grouped by inclusion in
1547        human tissues from GTEx data into those which are more than 5%
1548        but less than 15% included in any tissue, those which are alternative
1549        and those which are constitutively included. To better present the
1550        repeat density around the splice sites, the x axis is cut at 250 nt to
1551        show a zoom-in of the 250nt flanking the exons.

1552   (C)   Density profiles showing L1 antisense sequence 5kb upstream and
1553        downstream of constitutive and alternative exons in the mouse. The
1554        genomic age of each L1 element in the mouse genome was mapped
1555        by comparison to the rat, rhesus macaque, human, dog and cow
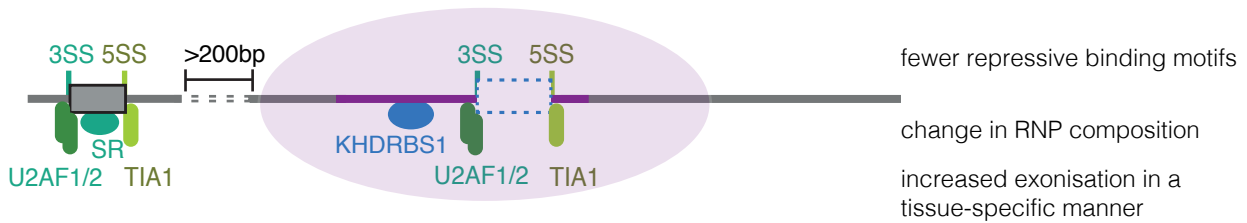1556        genome assemblies (see Methods for details).

1557

54

**Figure 7**

Young LINEs recruit repressive RBPs to insulate the LINE and surrounding RNA



LINE-mediated RNA insulation

multivalent binding due to:

many dispersed binding motifs

homo- and heteromeric
interactions between
repressive RBPs

RNA processing repressed

old LINEs are less repressed, and are a more common source of tissue-specific exons



fewer repressive binding motifs

change in RNP composition

increased exonisation in a
tissue-specific manner

1558  **Figure 7: LINE elements create a splice repressive zone that prevents**
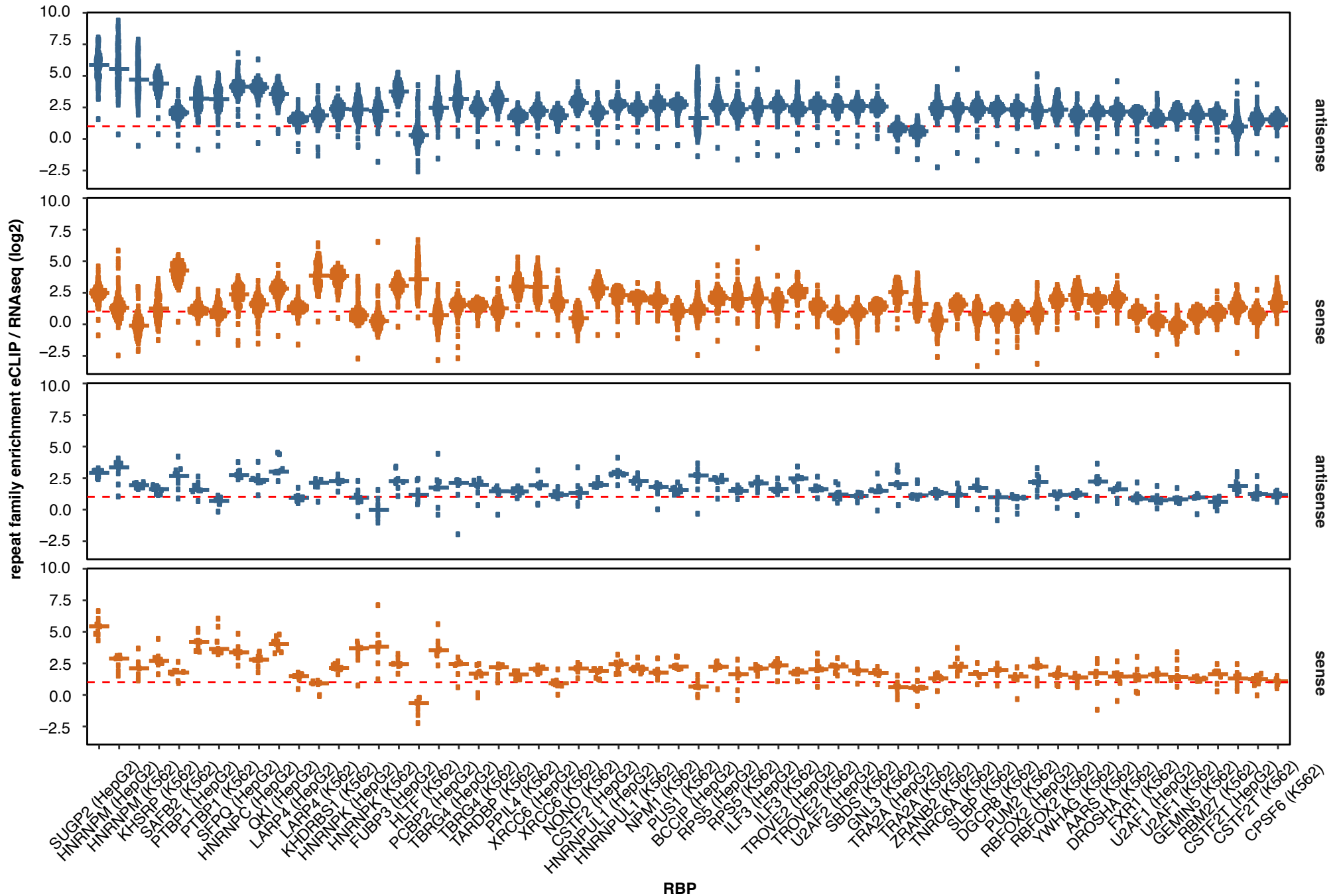
1559  **cryptic exonisation events**

1560  Consensus L1 elements are known to contain strong splice sites, but

1561  exonisation is rare and generally we observed exons from elements that are

1562  evolutionarily old. Evolutionarily young L1 insertions recruit a number of splice

1563  repressive proteins, including MATR3, PTBP1 and hnRNPM, as well as RBPs

1564  of yet unknown function (indicated as an X; but candidates are for instance

1565  BCCIP and SUGP2, see Figure 5C). These proteins recognise RNA motifs

1566  present within the L1 elements, which are diminished within evolutionary older

1567  L1s. The extent of splice-repressive proteins assembling on the L1s leads to

1568  selective pressure against young L1 insertions in a large proximity window of

1569  non-repeat derived exons. Hence evolutionary young LINEs insulate intronic

1570  regions from RNA processing. Evolutionarily older elements have a high

1571  probability of loosing binding sites of repressive RBPs. Hence, their

1572  exonisation is more common, but still largely tissue-specific.

1573

1574

55

**List of Supplementary files and Tables.**

**Figure S1. Related to Figure 1: Extended data for LINEs are binding platforms for a set of RBPs.**

**Figure S2. Related to Figure 2: Combinatorial binding of MATR3 and PTBP1 to the same LINEs.**

**Figure S3. Related to Figure 3: MATR3/PTBP1 repressed exons are frequently derived from LINEs or proximal to LINEs.**

**Figure S4. Related to Figure 4: Nonsense-mediated decay triggered by LINE-derived exons and depletion of ACAD9 expression following inclusion of a LINE-proximal exon.**

**Figure S5. Related to Figure 5: MATR3 and PTBP2 binds to mouse-specific L1 insertions and PTBP2 represses LINE-derived exon inclusion in the mouse brain.**

**Figure S6. Related to Figure 5: L1-derived exons are a source of primate-specific alternative exons with high tissue-specific variability.**

**Suppl. Table 1: Sources and references for iCLIP, eCLIP and RNAseq data used in this study and RBP binding motifs identified from literature.**

**Suppl. Table 2: Quantification of L1 and L2 sequences in iCLIP and eCLIP.**

**Suppl. Table 3:  Summary statistics of cryptic exon annotation from interleaving UCSC or ENSEMBL annotation and Cufflinks assembly.**

**Suppl. Table 4: Summary statistics of mRNA 3- end sequencing experiments.**

**Suppl. Table 5: Annotation derived from phylogenetic tracing of LINE elements in hg19.**

**Suppl. Table 6: Inclusion levels of 43583 UCSC annotated exons in 53 human tissue types.**

**Suppl. Table 7: Summary statistics of tetramer frequencies in antisense L1 sequences.**

Figure S1. Related to Figure 1

1606 **Figure S1. Related to Figure 1: Extended data for LINEs are binding**

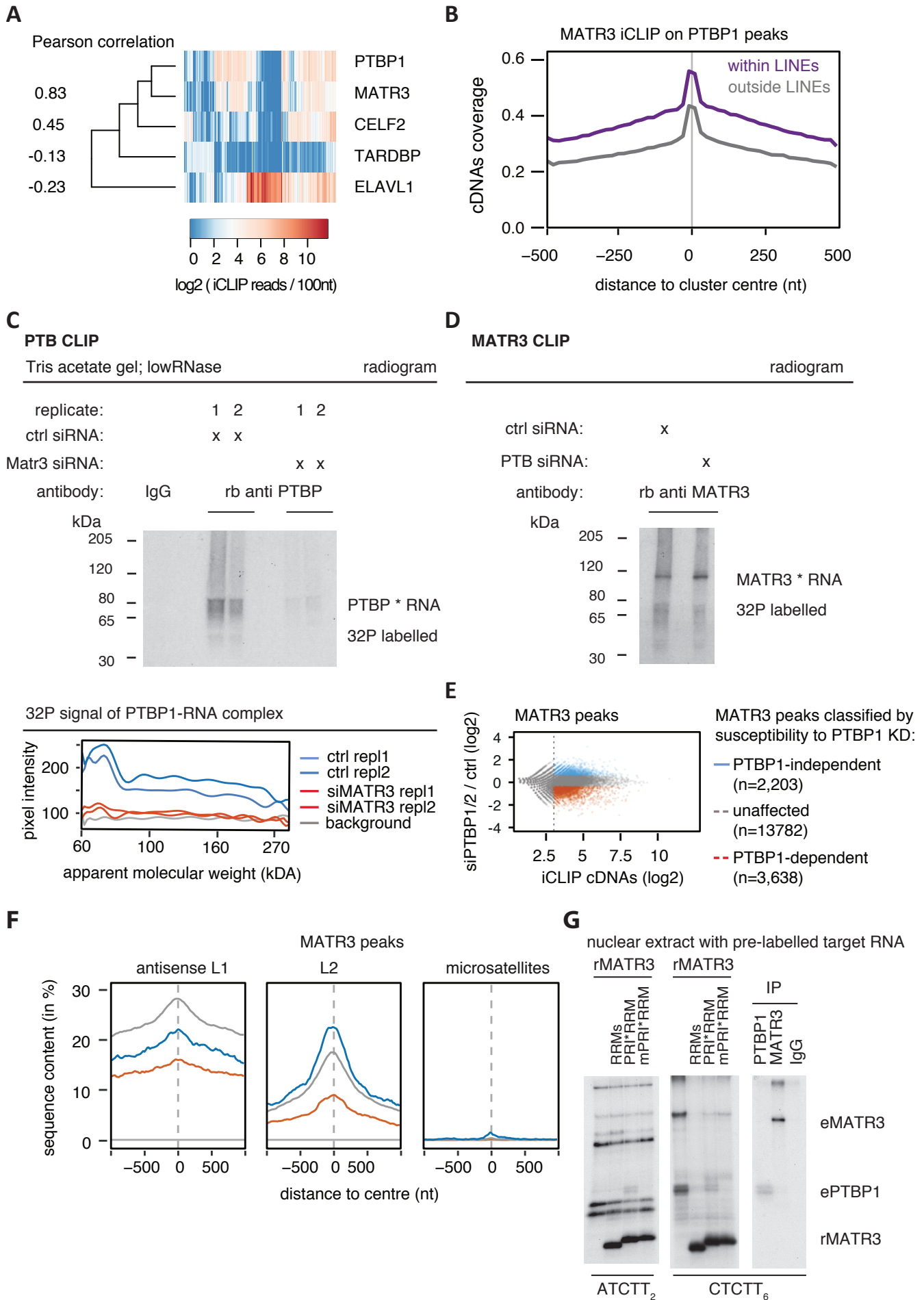1607 **platforms for a set of RBPs.**

1608 (A) TEtranscript (Jin et al., 2015) was used to estimate the enrichment of

1609 each subfamily of L1 and L2 repeats among the bound RNA

1610 sequences of a panel of RBPs, comparing the abundance in

1611 recovered eCLIP tags to the abundance in RNAseq reads. For each

1612 RBP, all 142 L1/L2 subfamilies (132 for L1, 10 for L2) were

1613 considered. Since eCLIP is strand-specific, binding to LINEs

1614 transcribed in sense or in antisense were quantified separately,

1615 coloured in red and blue. The cell lines used in each eCLIP

1616 experiment are indicated on the bottom.

1617

1618

1619

# Figure S2. Related to Figure 2



**A** Pearson correlation

**B** MATR3 iCLIP on PTBP1 peaks

**C** PTB CLIP
Tris acetate gel; lowRNase — radiogram

32P signal of PTBP1-RNA complex

**D** MATR3 CLIP — radiogram

**E** MATR3 peaks

MATR3 peaks classified by susceptibility to PTBP1 KD:
- PTBP1-independent (n=2,203)
- unaffected (n=13782)
- PTBP1-dependent (n=3,638)

**F** MATR3 peaks
antisense L1 | L2 | microsatellites

**G** nuclear extract with pre-labelled target RNA

1620 **Figure S2. Related to Figure 2: Combinatorial binding of MATR3 and**

1621 **PTBP1 to the same LINEs.**

1622 (A) For each RBP that showed considerable binding to LINE repeats in
1623    iCLIP (see B), we selected the 50 LINE repeats with strongest
1624    coverage (cDNAs per 100nt). For comparison we included TDP43,
1625    which showed little binding to LINE repeats. All iCLIP data selected
1626    was collected from HEK293 cells. The heatmap shows comparison of
1627    binding strength at this set of 214 LINE repeats, and the nearest
1628    neighbour analysis for each RBP. The values left to the dendrogram
1629    show the pearson correlation coefficient between all RBPs and
1630    PTBP1. Only LINEs with a minimal length of 50nt were considered to
1631    reduce the bias to short, highly expressed LINE repeats.

1632 (B) Metaprofile of iCLIP binding for MATR3 around iCLIP binding peaks of
1633    Celf2, Celf4, TDP43, HuR and PTBP1 within and outside of LINE
1634    repeats. The data was smoothed with 20nt bins.

1635 (C) HEK293T cells were transfected with siRNAs targeting MATR3,
1636    PTBP1 or scrambled controls, and 72 hours later labelled with 100μM
1637    4SU for 8 hours and cross-linked with 365nm UV light. The radiogram
1638    shows $^{32}$P labelled RNA crosslinked to and co-precipitated with
1639    PTBP1. Before immunoprecipitation, protein concentration was
1640    measured and equalised. The PTBP1 iCLIP was done under low
1641    RNase conditions (compare with Fig. 2A for high RNase condition).
1642    Replicate 1 and 2 are independent biological replicates processed in
1643    parallel.

1644 (D) $^{32}$P labelled RNA crosslinked to and co-precipitated with MATR3
1645    under equivalent conditions as in (C). The MATR3 iCLIP shown was
1646    done under high RNase conditions.

1647 (E) MATR3 binding peaks were identified from iCLIP experiments, and
1648    classified according to susceptibility to PTBP1 depletion as indicated
1649    based on moderated log2 fold change. Binding peaks with a
1650    normalised count of less than 8 were ignored, as indicated by the
1651    dotted line.

1652 (F) The overlap between the centre of MATR3 binding peaks and different

1653   repeat classes was tested for antisense L1 elements, sense L2

1654   elements, and sense CT-/T-rich microsatellite repeats. Metaprofile

1655   shows percent of each class of clusters overlapping with each
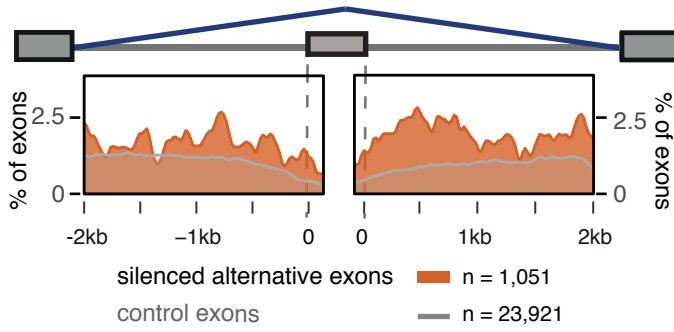
1656   genomic element.

1657  (G) Protein-protein interaction between MATR3 and PTBP1 allows

1658   recruitment of PTBP1 to a MATR3 bound RNA *in vitro*. Recombinant

1659   MATR3 mutants (rMATR3) and 32P labelled RNA probes were added

1660   to nuclear extracts from HeLa cells and UV-crosslinked. RNA

1661   substrates contained either two MATR3 or six PTBP1 RNA compete

1662   motifs motifs (ATCTT$_2$ and CTCTT$_6$). Crosslinking signals

1663   corresponding to endogenous PTBP1 (ePTBP1) and MATR3

1664   (eMATR3) were confirmed by immunoprecipitation.
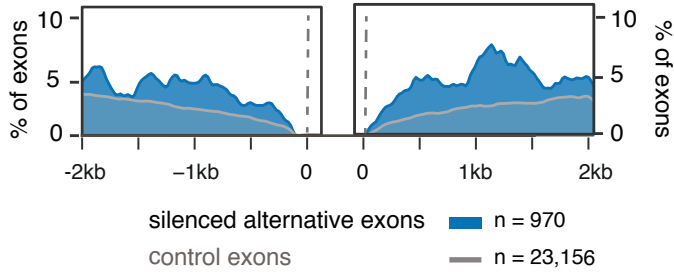
1665

62

**Figure S3 Related to Figure 3**



**A** sense L2 sequence content
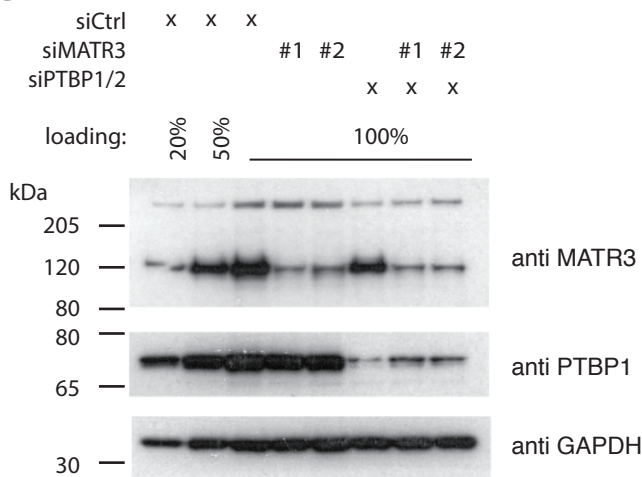
silenced alternative exons — n = 1,051
control exons — n = 23,921

**B** antisense L1 sequence - excluding LINE-derived exons

silenced alternative exons — n = 970
control exons — n = 23,156

**C**

**D** LINE-derived 3' and 5' splice sites

3SS 129 | 251 | 118 5SS

**E**
ΔsiPTB-ctrl  ΔsiMatr3-ctrl  ΔdKD-ctrl

change in exon inclusion (in %; measured by RT-PCR)

other | LINE-proximal | LINE-derived | LINE-derived terminal exon
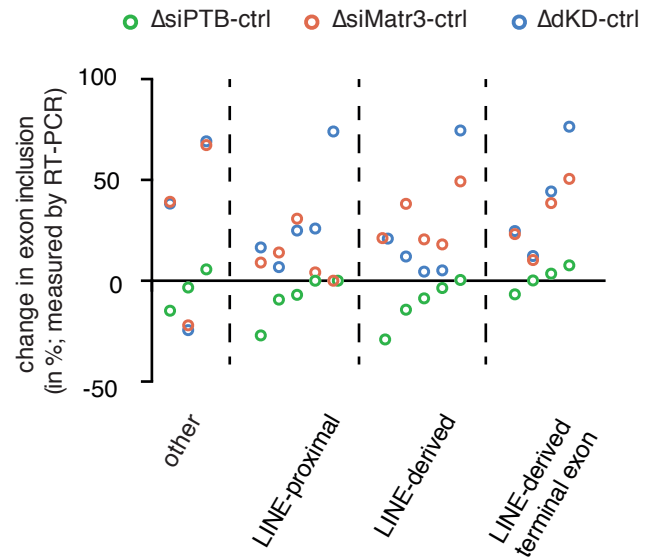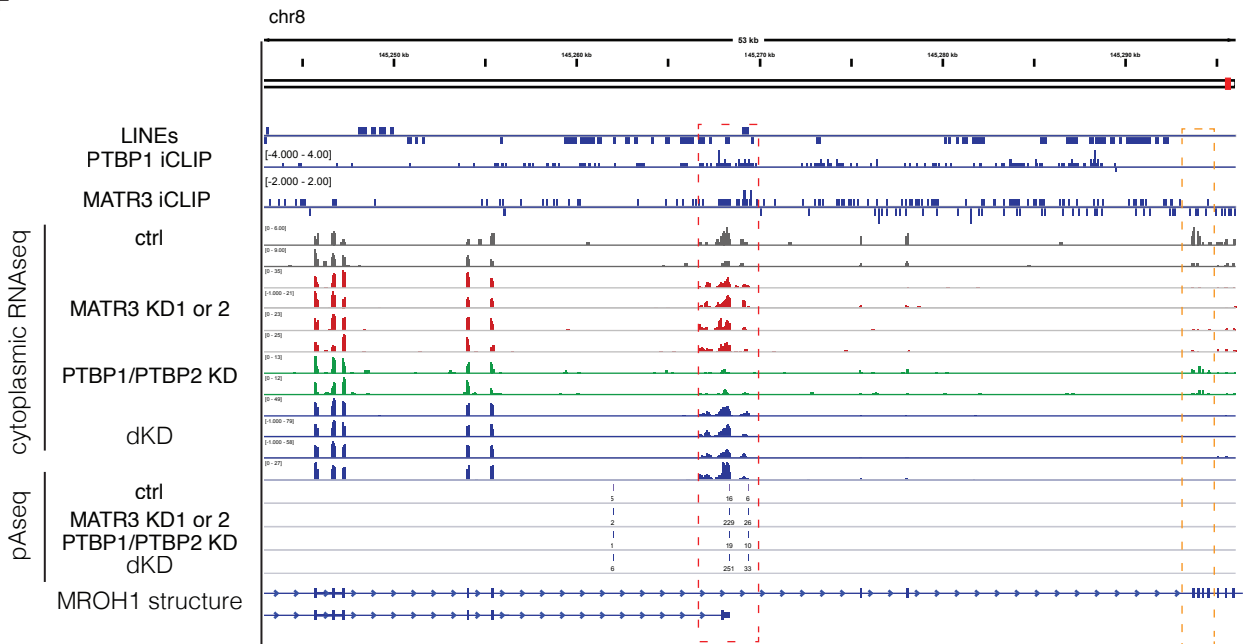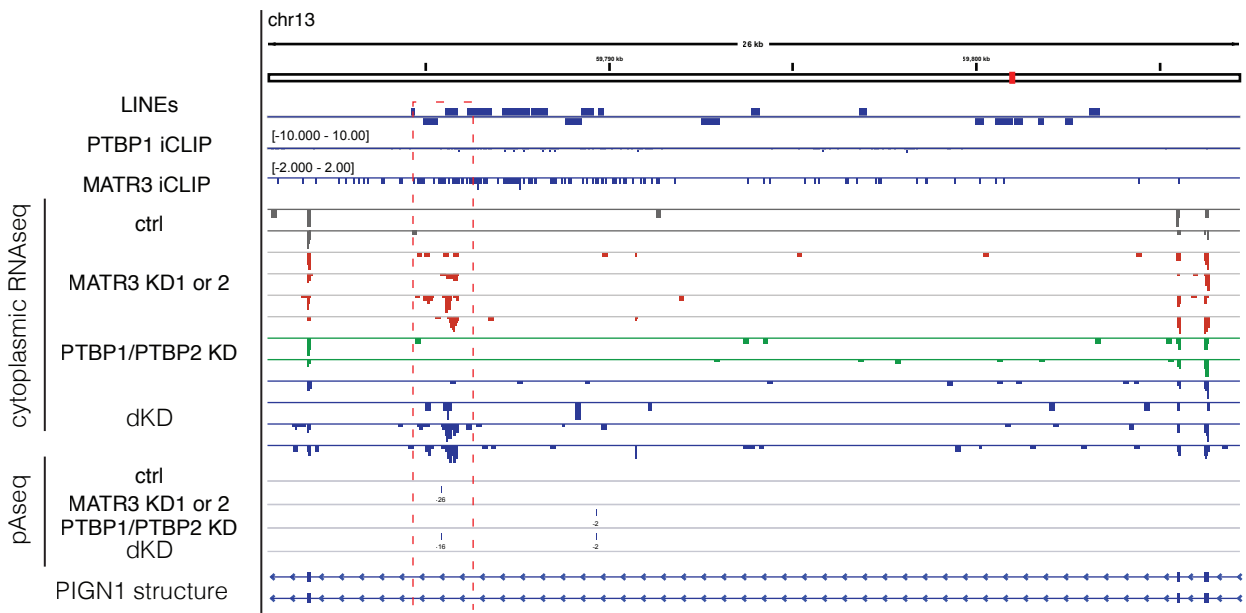
1666    **Figure S3. Related to Figure 3: MATR3/PTBP1 repressed exons are**

1667    **frequently derived from LINEs or proximal to LINEs.**

1668    (A)    The metaprofile shows the amount of sense L2 sequences flanking

1669            the splice sites of MATR3/PTBP1/2 repressed events. L2 sequences

1670            are particularly enriched towards the 3' splice site, and to a lesser

1671            extent than antisense L1 sequence.

1672    (B)    The metaprofile shows the amount of antisense L1 sequences

1673            flanking the splice sites of MATR3/PTBP1/2 repressed events, after

1674            removing all LINE-derived exons. The enrichment for L1 antisense

1675            sequence still persists (compare with Fig. 3A).

1676    (C)    Semi-quantitative Western blot showed efficient depletion of MATR3

1677            and PTBP1 in cells transfected with siRNAs against MATR3 or

1678            PTBP1 individually or in combination.

1679    (D)    The overlap of LINE-derived exons for which the 3' or 5' splice site is

1680            derived from a LINE element, only showing exons with

1681            junction-spanning reads on both sides (498 exons).

1682    (E)    Seventeen exons differentially included in MATR3 depleted cells were

1683            selected, and changes in exon inclusion were validated by RT-PCR.

1684            For each exon, the relative abundance of the isoform including the

1685            alternative exon was calculated compared to the exon exclusion

1686            isoform (conventional splicing pattern). The change between cells

1687            depleted of MATR3, PTBP1/2, or both simultaneously is shown, and

1688            exons are grouped by their positioning relative to the closest LINE

1689            element. Semi-quantitative RT-PCR analysis is averaged across three

1690            independent replicates.

1691    (F)    And (G) Examples of MATR3/PTBP1 repressed polyA sites. Genome

1692            browser tracks show position and orientation of LINE insertion

1693            (hg19/RepeatMasker annotation), PTBP1 and MATR3 iCLIP

1694            coverage, as well as tracks for RNAseq of cytoplasmic RNA and

1695            mRNA 3' end sequencing (pA-seq) from total RNA. All tracks are

1696            scaled appropriately to library size. (F) The MROH1 gene shows

1697            inclusion of additional exonic sequence and two different terminal

1698            exon isoforms in MATR3 depleted cells (highlighted by red dashed
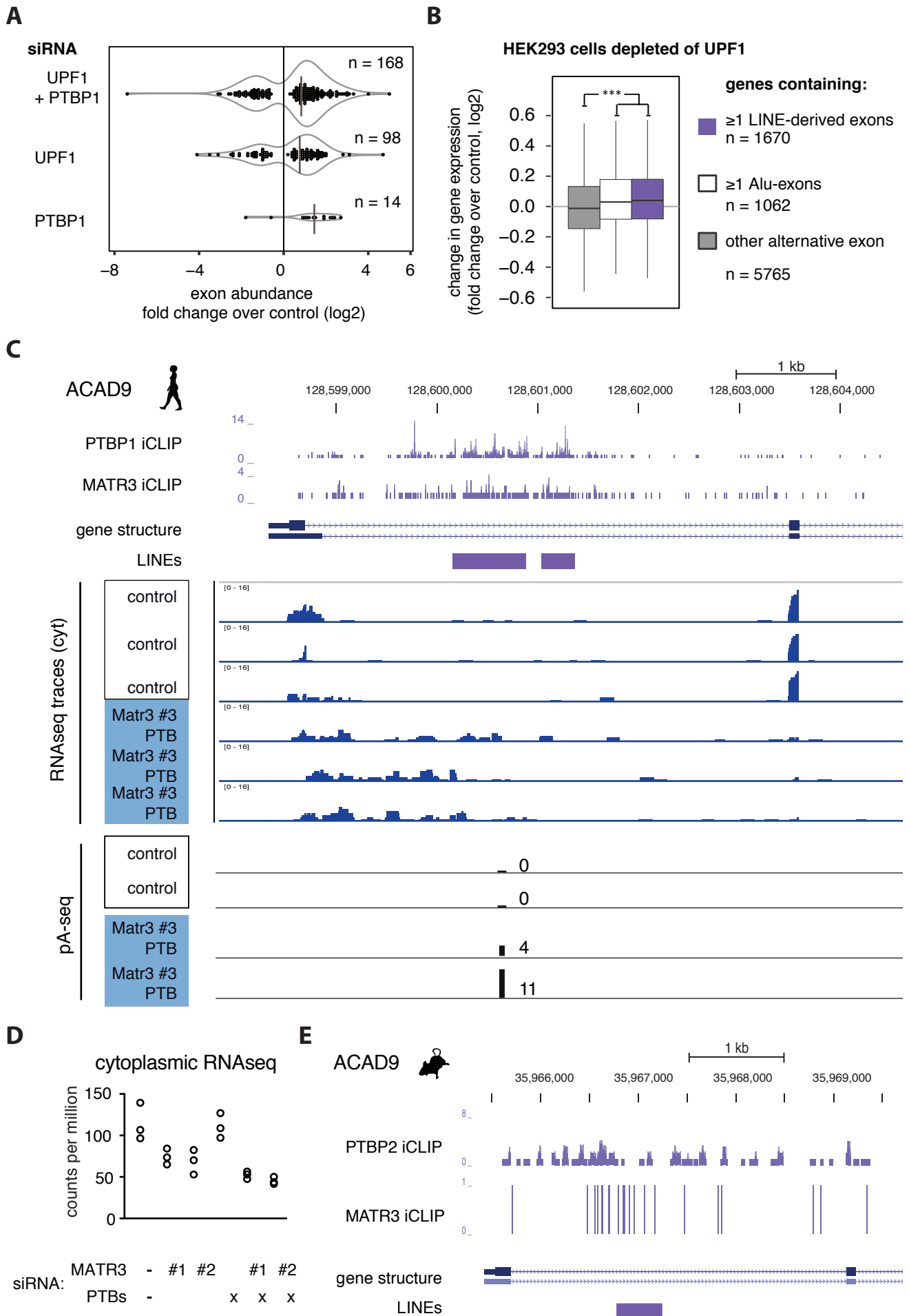
63

1699    lines). Usage of this alternative terminal exon results in gene

1700    truncation as seen by loss of expression downstream of it (highlighted

1701    by orange dashed lines). (G) The PIGN1 shows usage of a cryptic

1702    processing site resulting in a novel exon and a novel polyA site,

1703    derived from two antisense L1 insertions (highlighted by red dashed

1704    lines).

1705

1706

1707

# Figure S4. Related to Figure 4.

**A**

siRNA



exon abundance
fold change over control (log2)

**B**

HEK293 cells depleted of UPF1



genes containing:

■ ≥1 LINE-derived exons
n = 1670

□ ≥1 Alu-exons
n = 1062

▨ other alternative exon
n = 5765

**C**



**D**



**E**

1708 **Figure S4. Related to Figure 4: Nonsense-mediated decay triggered by**

1709 **LINE-derived exons and depletion of ACAD9 expression following**

1710 **inclusion of a LINE-derived exons.**

1711 (A) RNAseq data from Ge et al. on HEK293 cells depleted of PTBP1,
1712    UPF1 or both, was reanalysed with DEXSeq. The number of
1713    detectable LINE-derived exons and their change in abundance
1714    compared to control cells is shown. Consistent with the hypothesis
1715    that LINE-derived exons are repressed in wild-type cells by splicing
1716    repressors and through decay of the inclusion isoform, combined
1717    depletion of UPF1 and PTBP1 greatly increases the number of
1718    detectable LINE-derived exons.

1719 (B) The change in gene expression in UPF1-depleted cells over control is
1720    shown for genes that contained or did not contain one or more
1721    LINE-derived exons. As positive control, Alu-exon containing genes
1722    are shown since inclusion of Alu-exons frequently triggers NMD (Attig
1723    et al., 2016).

1724 (C) Genome browser tracks for PTBP1 and MATR3 iCLIP data from HeLa
1725    cells at the ACAD9 locus. Position of L2 insertions is annotated below
1726    the structure of annotated ACAD9 transcripts, and stranded RNAseq
1727    data from cytoplasmic RNA of HeLa cells depleted of
1728    MATR3/PTBP1/PTBP2 is shown. Below the position of a novel pA site
1729    within the second L2 repeat is shown, which is only detected in
1730    absence of MATR3/PTBP1/PTBP2.

1731 (D) Quantification of ACAD9 expression in single and combined depletion
1732    of MATR3 and PTBP1/2 from cytoplasmic RNAseq.

1733 (E) Genome browser tracks for PTBP2 and MATR3 on the mouse ACAD9
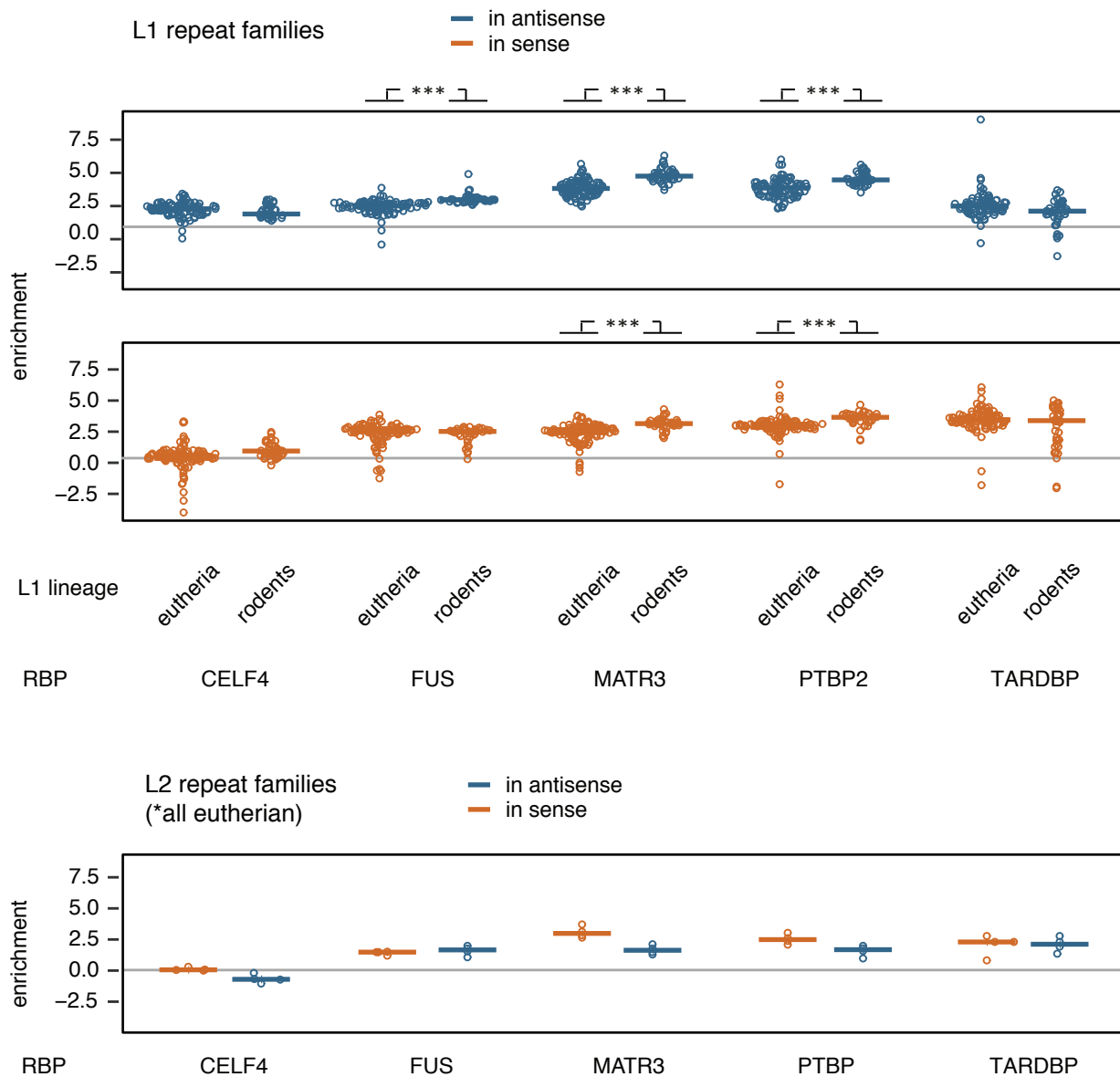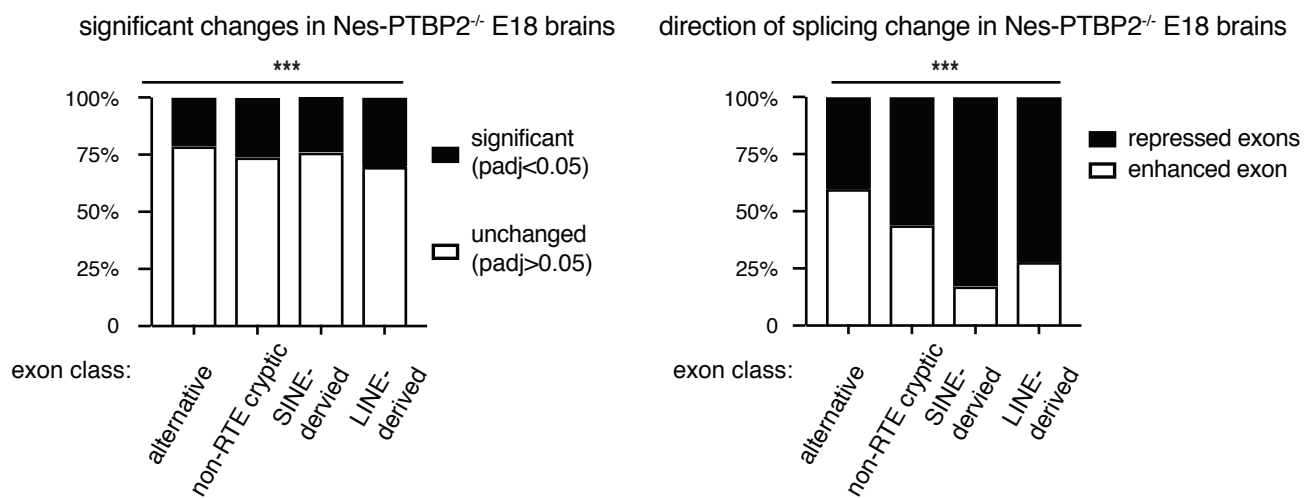1734    locus. In mouse, there is a single, 465bp long L2 insertion.

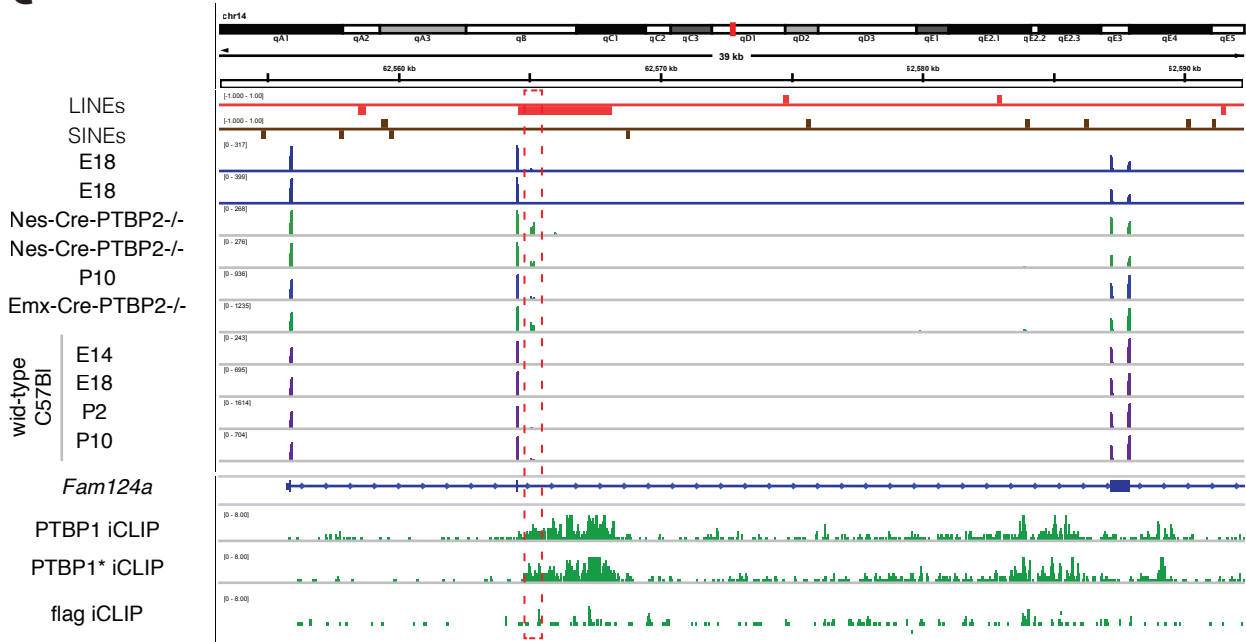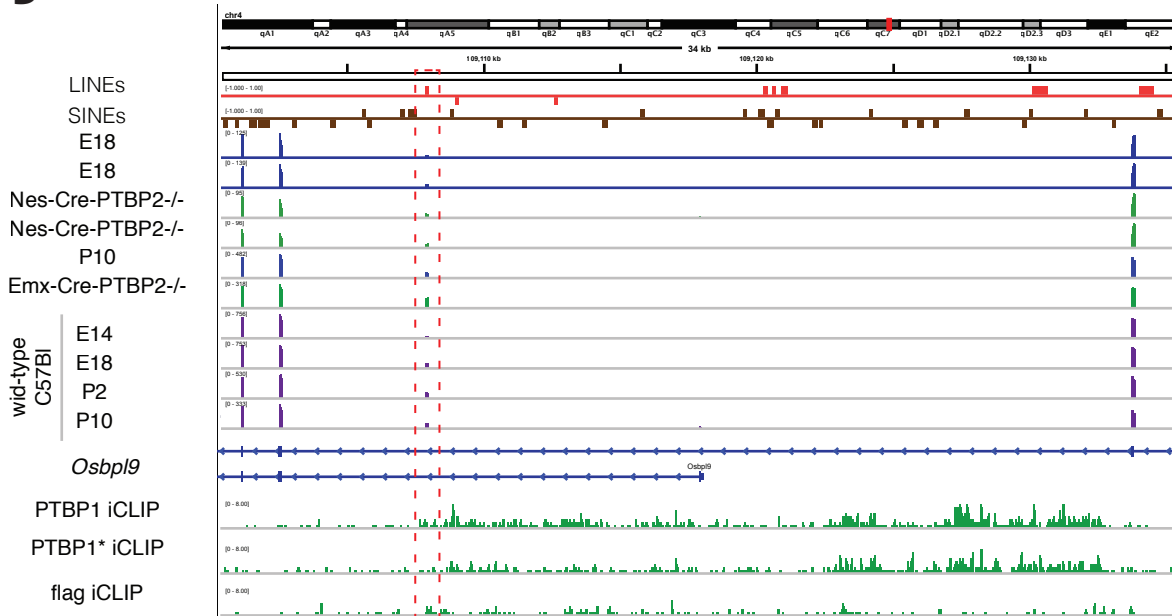1735

1736

1737

**Figure S5 Related to Figure 5**

**A**



**B**
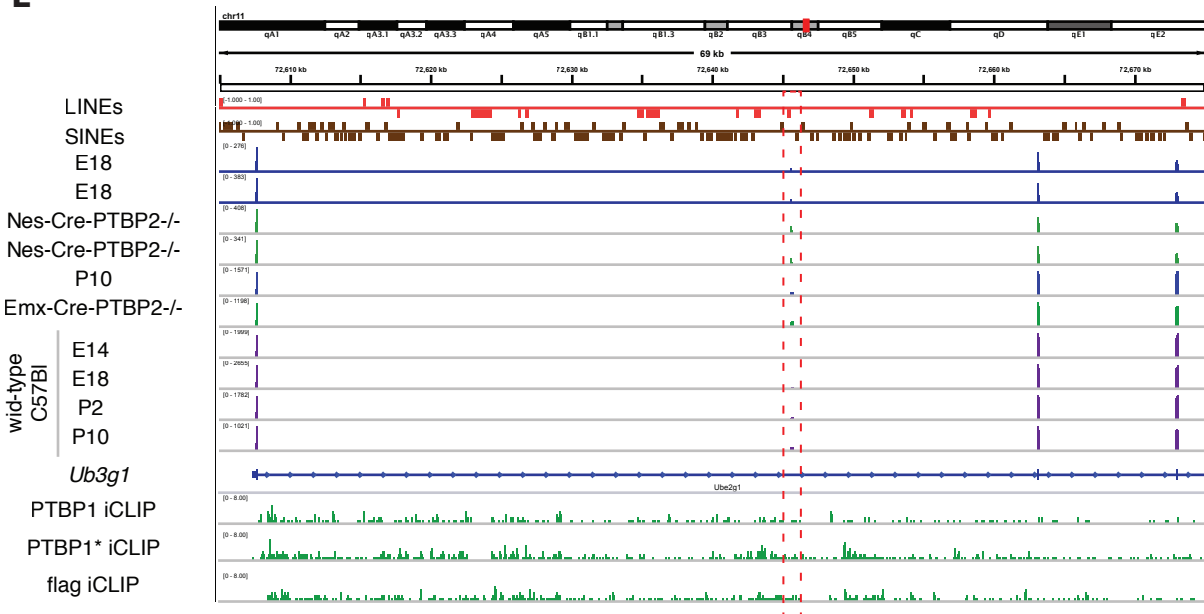
**Figure S5. Related to Figure 5: MATR3 and PTBP2 binds to mouse-specific L1 insertions and PTBP2 represses LINE-derived exon inclusion in the mouse brain.**
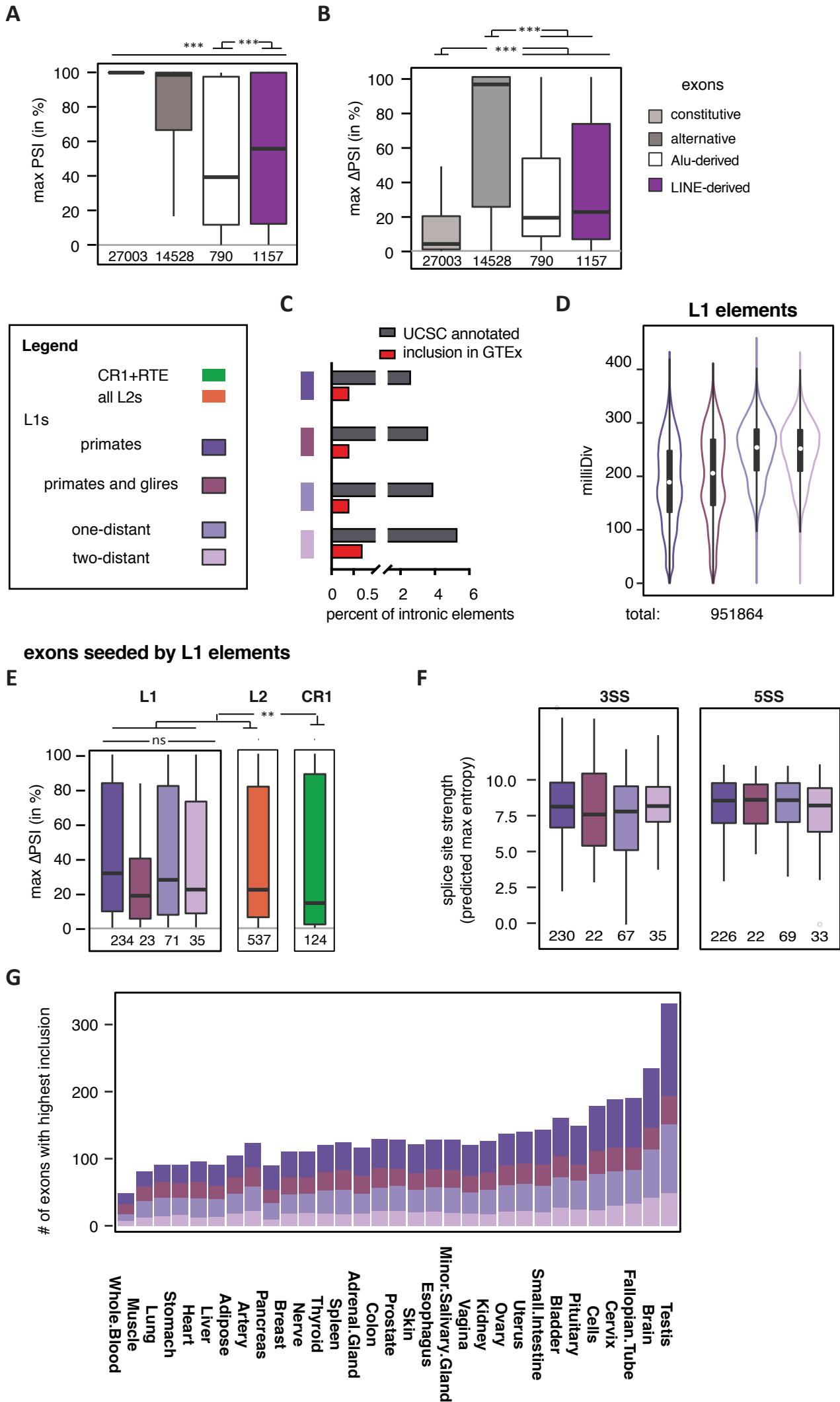
(A) TEtranscript (Jin et al., 2015) was used to estimate the enrichment of each subfamily of L1 and L2 repeats among the bound RNA sequences of a panel of RBPs, with CLIP data available for C57Bl mouse brain; comparing the abundance in recovered eCLIP tags to the abundance in RNAseq reads of P2. For each RBP, 133 repBase LINE subfamilies were considered (129 for L1, 4 for L2, (Jurka, 1998)). Since eCLIP is strand-specific, binding to LINEs transcribed in sense or in antisense were quantified separately, coloured in red and blue. Details and references of data sets are given in Supplementary Table 1.

(B) RNAseq data of PTBP2$^{-/-}$ knockout mouse brains from (Vuong et al., 2016) was re-analysed, and exons with significant differences in inclusion in Nes-Cre-PTBP$^{-/-}$ knockout mouse brains were stratified according to their relationship to retrotransposon repeats. LINE-derived exons were more likely to be mis-regulated than expected by chance ($\chi^2$ test), and PTBP2 acts primarily as repressor on LINE-derived exons. Number of exons in each group are: alternative exons n=8142, non-repeat derived cryptic exons n=33420, SINE derived exons n=459, LINE-derived exons n=308.

(C) to (E) RNAseq data of PTBP2$^{-/-}$ knockout mouse brains from (Vuong et al., 2016) compared to RNAseq data of C57/B6 wild-type mouse brain at different developmental stages (E10, E14, P2 and P10). LINE-derived exons were selected from the exon list of PTBP2 responsive exons provided by Vuong et al. in Suppl Table 3 and 5; that is they are selected to show a minimum 10% change in inclusion upon PTBP2 depletion. PTBP1 iCLIP was done with flag antibody in Rosa-PTBP1 transgenic mice, and PTBP1* iCLIP is PTBP1 iCLIP in PTBP2 knockouts (also from (Vuong et al., 2016)). Because of the high extent of intron-retention reads in mouse brain, only junction-spanning reads are shown. These exons are more included in

67

1773       postnatal brains than in foetal brain, suggesting PTBP2 suppresses

1774       exonisation in developing neurons but less in mature neurons. (C)

1775       Exon 3 of Fam124 is derived from a rodent specific L1 insertion. (D)

1776       Exon 5 of *Osbpl9* is derived from an old CR1 insertion conserved

1777       across mammalian lineages. (E) Exon 2 of *Ube2g1* is derived from an

1778       old HAL1 insertion conserved across mammalian lineages.

1779

1780

1781

# Figure S6. Related to Figure S1.



**A**

**B**

**C**

**D** L1 elements

**Legend**

CR1+RTE
all L2s

L1s

primates

primates and glires

one-distant

two-distant

exons seeded by L1 elements

**E**

**F** 3SS    5SS

**G**

1782 **Figure S6. Related to Figure 5: L1-derived exons are a source of primate-**

1783 **specific alternative exons with high tissue-specific variability.**

1784 Percent splice index (PSI) was calculated in the GTEx panel of human tissues

1785 for LINE-derived and Alu-derived exons, as well as all other exons of the

1786 same genes. All exons are annotated within UCSC and cross-referenced with

1787 RefSeq annotation. Inclusion levels range from 0 to 100%, showing no

1788 inclusion or full inclusion. If no support for expression of the flanking exons

1789 was found, the gene is assumed to be non-expressed. Genomic age of L1

1790 elements as defined in Figure 5A. Significance tests were done across groups

1791 by Kruskal-Wallis' test and pairwise comparisons were done with Dunn's test

1792 and corrected according to Holm-Šidák. ** and *** indicate adjusted p-value

1793 was below 0.01 and 0.001, respectively.

1794 (A) Maximum inclusion in any tissue was calculated for each exon, and

1795 the distribution is shown for LINE-derived exons, Alu-exons as well as

1796 non-repeat derived alternative and constitutive exons.

1797 (B) For all exons surveyed within the GTEx data, the difference in PSI

1798 between the tissues with highest and lowest inclusion was calculated

1799 as metric for tissue-specific inclusion.

1800 (C) Exons derived from old L1 insertions are most likely to form an exon

1801 based on UCSC annotation. Based on GTEx data, exons derived from

1802 old L1 insertions retained in primates, cow and dog, are most likely to

1803 be included in any of the tissue types covered.

1804 (D) The substitutions from L1 consensus families is shown for L1s

1805 grouped by genomic age. As expected, young elements show fewer

1806 substitutions from consensus then old elements.

1807 (E) Difference in PSI between tissues with highest and lowest inclusion

1808 for exons derived from L1 elements grouped by genomic age of the

1809 insertion, compared to exons derived from L2 and CR1 insertions.

1810 (F) Exons derived from L1 elements have strong splice sites irrespective

1811 of the genomic age of the insertion. The maximum entropy score of 5'

1812 and 3' splice sites of each exon was predicted based on nucleotide

1813 sequence (Yeo and Burge, 2004).

1814 (G) The number of L1-derived exons is shown for all primary tissues

1815        screened in the GTEx data, based on testing in which tissue an exon

1816        is most included. Exons are allowed to be counted multiple times if

1817        maximum inclusion was in multiple tissues, for instance because they

1818        are constitutive.

1819

1820

**Suppl Table 3**

## Cufflinks predicted exonic bins*

*Cufflinks prediction was flatted from gtf to gff for HTseqcounts.py

exonic bins of min 5nt length and sufficient coverage for DEXSeq analysis

referenced to UCSC

|  | total | truly novel exons | LINE-derived | novel LINE-derived |
|---|---|---|---|---|
| exonic bins | 264,169 | 26,939 | 1,430 | 634 |
| with min 1 splice site confirmed by junction-spanning read | 177,179 | 9,000 | 1,430 | 634 |
| # of genes | 12,929 | 4,455 | 1,065 | 499 |
| exonic bins in protein coding genes | 165,138 | 7,020 | 1,114 | 501 |
| # of genes | 11,582 | 3,689 | 899 | 86 |

referenced to ENSEMBL72

|  | total | truly novel exons | LINE-derived | novel LINE-derived |
|---|---|---|---|---|
| exonic bins | 264,169 | 6,627 | 1,430 | 257 |
| with min 1 splice site confirmed by junction-spanning read | 177,179 | 2,246 | 1,430 | 257 |
| # of genes | 12,929 | 1,305 | 1,065 | 188 |
| exonic bins in protein coding genes | 165,138 | 1,862 | 1,114 | 207 |
| # of genes | 11,582 | 1,233 | 899 | 178 |

LINE-derived exons

total: 1,430

|  | | constitutive | alternative | "cryptic exons" partial-overlap | truly novel |
|---|---|---|---|---|---|
| UCSC / ENSEMBL | | | | | |
| perfect overlap | | 266 | 318 | 375 | 40 |
| partial-overlap | | 0 | 3 | 972 | 130 |
| truly novel | | 0 | 0 | 16 | 34 |

## Suppl Table 4

changes in pA-site usage

| condition | proximal pA site usage | FDR | # | # within LINE repeats proximal | distal |
|---|---|---|---|---|---|
| Matr3 KD1 | up | < 0.05 | 257 | 9 | 6 |
| | | >= 0.05 | 417 | 24 | 24 |
| | down | < 0.05 | 329 | 26 | 5 |
| | | >= 0.05 | 428 | 54 | 22 |
| | neither | - | 2,605 | 150 | 63 |
| Matr3 KD2 | up | < 0.05 | 214 | 11 | 7 |
| | | >= 0.05 | 416 | 41 | 27 |
| | down | < 0.05 | 368 | 14 | 5 |
| | | >= 0.05 | 567 | 46 | 26 |
| | neither | - | 2,256 | 131 | 56 |
| PTB KD | up | < 0.05 | 219 | 5 | 4 |
| | | >= 0.05 | 420 | 28 | 21 |
| | down | < 0.05 | 401 | 18 | 4 |
| | | >= 0.05 | 516 | 52 | 28 |
| | neither | - | 2,332 | 139 | 65 |
| Matr3 KD1 + PTB KD | up | < 0.05 | 370 | 17 | 11 |
| | | >= 0.05 | 424 | 34 | 21 |
| | down | < 0.05 | 399 | 24 | 5 |
| | | >= 0.05 | 492 | 44 | 28 |
| | neither | - | 2,373 | 129 | 59 |