



## Ancient exapted transposable elements promote nuclear enrichment of human long noncoding RNAs

Joana Carlevaro-Fita, Taisia Polidori, Monalisa Das, et al.

*Genome Res.* published online December 26, 2018  
Access the most recent version at doi:[10.1101/gr.229922.117](https://doi.org/10.1101/gr.229922.117)

---

<b>P&lt;P</b>	Published online December 26, 2018 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by/4.0/">http://creativecommons.org/licenses/by/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .

---



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

1 *Ancient exapted transposable elements promote nuclear enrichment of human long*  
2 *noncoding RNAs*

3

4 Joana Carlevaro-Fita<sup>1,2,3,5</sup>

5 Taisia Polidori<sup>1,2,3,5</sup>

6 Monalisa Das<sup>1,2,5</sup>

7 Carmen Navarro<sup>4</sup>

8 Tatjana I. Zoller<sup>1,2</sup>

9 Rory Johnson<sup>1,2\*</sup>

10

11

12

13 1. Department for BioMedical Research (DBMR), University of Bern, 3008 Bern, Switzerland

14 2. Department of Medical Oncology, Inselspital, University Hospital and University of Bern, 3010 Bern,  
15 Switzerland

16 3. Graduate School of Cellular and Biomedical Sciences, University of Bern, Bern, Switzerland

17 4. Department of Computer Science and Artificial Intelligence, University of Granada, Spain

18 5. Equal contribution

19

20 \*Correspondence: [rory.johnson@dbmr.unibe.ch](mailto:rory.johnson@dbmr.unibe.ch)

21 Keywords: Transposable element; subcellular localization; long noncoding RNA; lncRNA; evolution;  
22 exaptation.

23

24

25 **Abstract**

26

27 **The sequence domains underlying long noncoding RNA (lncRNA) activities, including their**  
28 **characteristic nuclear enrichment, remain largely unknown. It has been proposed that these**  
29 **domains can originate from neofunctionalised fragments of transposable elements (TEs), otherwise**  
30 **known as RIDLs (Repeat Insertion Domains of Long Noncoding RNA), although just a handful**  
31 **have been identified. It is challenging to distinguish functional RIDL instances against a numerous**  
32 **genomic background of neutrally-evolving TEs. We here show evidence that a subset of TE types**  
33 **experience evolutionary selection in the context of lncRNA exons. Together these comprise an**  
34 **enrichment group of 5374 TE fragments in 3566 loci. Their host lncRNAs tend to be functionally**  
35 **validated and associated with disease. This RIDL group was used to explore the relationship**  
36 **between TEs and lncRNA subcellular localisation. Using global localisation data from ten human**  
37 **cell lines, we uncover a dose-dependent relationship between nuclear/cytoplasmic distribution, and**  
38 **evolutionarily-conserved L2b, MIRb and MIRc elements. This is observed in multiple cell types,**  
39 **and is unaffected by confounders of transcript length or expression. Experimental validation with**  
40 **engineered transgenes shows that these TEs drive nuclear enrichment in a natural sequence**  
41 **context. Together these data reveal a role for TEs in regulating the subcellular localisation of**  
42 **lncRNAs.**

## 43 Introduction

44           The human genome contains many thousands of long noncoding RNAs (lncRNAs), of which at  
45 least a fraction are likely to have evolutionarily-selected biological functions (Ulitsky and Bartel 2013).  
46 Our current working hypothesis is that, similar to proteins, lncRNA functions are encoded in primary  
47 sequence through “domains”, or discrete elements that mediate specific aspects of lncRNA activity. Such  
48 activities range from molecular interactions to subcellular localisation (Guttman and Rinn 2012; Mercer  
49 and Mattick 2013; Johnson and Guigó 2014). Experimental support for this domain model is beginning to  
50 emerge (Marín-Béjar et al. 2017). Mapping domains in a comprehensive manner is thus a key step  
51 towards the understanding and prediction of lncRNA functions.

52           One possible source of lncRNA domains are transposable elements (TEs) (Johnson and Guigó  
53 2014). TEs are known to have been major contributors to genomic evolution through the insertion and  
54 neofunctionalisation of sequence fragments – a process known as *exaptation* (Bourque 2009)(Feschotte  
55 2008). This process has contributed to the evolution of diverse features in genomic DNA, including  
56 transcriptional regulatory motifs (Bourque et al. 2008; Johnson et al. 2006), microRNAs (Roberts et al.  
57 2014), gene promoters (Faulkner et al.; Huda et al. 2011), and splice sites (Lev-Maor et al. 2003; Sela et  
58 al. 2007).

59           We recently proposed that exaptation also takes place in the context of lncRNAs, with TEs  
60 contributing pre-formed functional domains. We termed these “RIDLs” – *Repeat Insertion Domains of*  
61 *Long noncoding RNAs* (Johnson and Guigó 2014). As RNA, TEs are known to interact with a rich variety  
62 of proteins, meaning that in the context of lncRNA they could plausibly act as protein-docking sites  
63 (Blackwell et al. 2012). Diverse evidence also points to repetitive sequences forming intermolecular  
64 Watson-Crick RNA:RNA and RNA:DNA hybrids (Johnson and Guigó 2014; Gong and Maquat 2011;  
65 Holdt et al. 2013). However, it is likely that *bona fide* RIDLs represent a small minority of the many  
66 exonic TEs, with the remainder being phenotypically-neutral “passengers”.

67           A small but growing number of RIDLs have been described, reviewed in (Johnson and Guigó  
68 2014). These are found in lncRNAs with clearly-demonstrated functions, including the X Chromosome  
69 silencing transcript *XIST* (Elisaphenko et al. 2008), the oncogene *ANRIL* (Holdt et al. 2013) and the  
70 regulatory antisense *Uchl1AS* (Carrieri et al. 2012). In each case, domains of repetitive origin are necessary  
71 for a defined function: the structured A-repeat of *XIST*, of retroviral origin, recruits the PRC2 silencing  
72 complex (Elisaphenko et al. 2008); Watson-Crick hybridisation between RNA and DNA *Alu* elements  
73 recruits *ANRIL* to target genes (Holdt et al. 2013); a SINEB2 repeat in *Uchl1AS* increases translational  
74 rate of its sense mRNA (Carrieri et al. 2012). In parallel, transcriptome-wide maps of lncRNA-linked TEs

75 have shown how TEs have contributed extensively to lncRNA gene evolution (Kelley and Rinn 2012;  
76 Kapusta et al. 2013)(Schmitt et al. 2016)(Hezroni et al. 2015). However, there has been no attempt to  
77 enrich these maps for RIDLs with evidence of selected functions in the context of mature lncRNA  
78 molecules.

79 Subcellular localisation, and the domains controlling it, are crucial determinants of lncRNA  
80 functions (reviewed in (Chen 2016)). For example, transcriptional regulatory lncRNAs must be located in  
81 the nucleus and chromatin, whereas those regulating microRNAs or translation should be present in the  
82 cytoplasm (Zhang et al. 2014b). Although higher nuclear/cytoplasmic ratios are a hallmark of lncRNAs, a  
83 large population of cytoplasmic transcripts also exists (Mukherjee et al. 2017) (Carlevaro-Fita et al. 2016;  
84 Derrien et al. 2012)(Cabili et al. 2015; Mas-Ponte et al. 2017)(Benoit Bouvrette et al. 2018). If lessons  
85 learned from mRNA are also valid for lncRNAs, then short sequence motifs recognised by RNA binding  
86 proteins (RBPs) will be an important localisation-regulatory mechanism (Martin and Ephrussi 2009). This  
87 was recently demonstrated for the *BORG* lncRNA, where a pentameric motif was shown to mediate  
88 nuclear retention (Zhang et al. 2014a). Similarly, multiple copies of the 156 bp RRD repeat motif mediate  
89 nuclear enrichment of the *FIRRE* lncRNA, through binding to hnRNPU (Hacisuleyman et al. 2016a)  
90 (Hacisuleyman et al. 2014). Another study implicated an inverted pair of *Alu* elements in nuclear  
91 retention of *lincRNA-P21* (Chillón and Pyle 2016). This raises the possibility that, by “copying and  
92 pasting” generic RNA motifs, RIDLs could fine-tune lncRNA localisation at a global scale.

93 The aim of the present study is to create a human transcriptome-wide catalogue of putative  
94 RIDLs. Supporting its relevance, lncRNAs carrying these RIDLs are enriched for functional genes.  
95 Finally, we provide *in silico* and experimental evidence that certain RIDL types, derived from ancient  
96 transposable elements, promote the nuclear enrichment of their host transcripts.

## 97 **Results**

98           The objective of this study is to create a map of repeat insertion domains of long noncoding  
99 RNAs (RIDLs) and link them to lncRNA functions. We hypothesise that RIDLs could confer such  
100 functions through interactions with DNA, RNA or protein molecules (Johnson and Guigó 2014) (Figure  
101 1A).

102           Any attempt to map RIDLs must deal with two challenges. First, that they will likely represent a  
103 small minority amongst many phenotypically-neutral “passenger” transposable elements (TEs) in lncRNA  
104 exons (Figure 1B). Second, many TE instances may be under evolutionary selection, but for functions  
105 executed at the *DNA level* (eg transcription factor binding sites, enhancer elements), rather than the RNA  
106 level (Bassett et al. 2014)

107           Therefore, it is necessary to identify RIDLs by some signature of selection that is specific for a  
108 mature RNA product using an appropriate background model. In this study we use three types of such  
109 signatures: exonic enrichment, strand bias (with respect to host gene), and exon-specific evolutionary  
110 conservation (Figure 1B). To estimate background, we utilise intronic TEs, since they should mirror any  
111 biases of TE distribution across the genome but are not incorporated into mature lncRNA transcripts.

112           Resulting RIDL predictions should be considered as “enrichment groups”, due to high rates of  
113 false positive predictions, and all downstream analyses should be interpreted accordingly.

114

### 115 **A map of exonic transposable elements in GENCODE v21 lncRNAs**

116           Our first aim was to create a comprehensive map of transposable elements (TEs) within the exons  
117 of GENCODE v21 human lncRNAs (Figure 2A). Altogether 5,520,018 distinct TE insertions were  
118 intersected with 48684 exons from 26414 transcripts of 15877 GENCODE v21 lncRNA genes, resulting  
119 in 46474 exonic TE insertions in lncRNA (Figure 1B). 13121 lncRNA genes (82.6%) carry at least one  
120 exonic TE fragment in one or more of their mature transcripts.

121           We also created a reference dataset with 31,004 GENCODE lncRNA introns, resulting in 562,640  
122 intron-overlapping TE fragments (Figure 2A). Comparing intronic and exonic TE data, we see that  
123 lncRNA exons are depleted for TE insertions: 29.2% of exonic nucleotides are of TE origin, compared to  
124 43.4% of intronic nucleotides (Figure 2B), similar to previous studies (Kapusta et al. 2013). This may  
125 reflect generalised selection against disruption of functional lncRNA transcripts by TEs. The exonic  
126 depletion of TEs in lncRNAs is less pronounced than for protein-coding loci, whereas the intronic TE  
127 density of both is similar to the whole-genome average.

128

129 Contribution of transposable elements to lncRNA gene structures

130 TEs have contributed widely to both coding and noncoding gene structures by the insertion of  
131 elements such as promoters, splice sites and termination sites (Sela et al. 2007). We next classified  
132 inserted TEs by their contribution to lncRNA gene structure (Figure 2C,D). It should be borne in mind  
133 that this analysis is dependent on the accuracy of underlying GENCODE annotations, which are often  
134 incomplete at 5' and 3' ends (Lagarde et al. 2017). Altogether 4993 (18.9%) transcripts' promoters lie  
135 within a TE, most often those of *Alu*, L1 and ERVL-MaLR classes (Figure 2E). 7497 (28.4%) lncRNA  
136 transcripts are terminated by a TE, most commonly by L1, *Alu*, ERVL-MaLR classes. 8494 lncRNA  
137 splice sites (32.2%) are of TE origin, and 2681 entire exons are fully contributed by TEs (10.1%) (Figure  
138 2E). These observations support known contributions of TEs to gene structural features (Sela et al. 2007).  
139 Nevertheless, the most frequent case is represented by 22,031 TEs that lie completely within an exon and  
140 do not overlap any splice junction (“inside”).

141

142 Evidence for selection on certain exonic transposable element types

143 This exonic TE map represents the starting point for the identification of RIDLs, defined as the  
144 subset of TEs with evidence for functionality in the context of mature lncRNAs. In this and subsequent  
145 analyses, TEs are grouped by type as defined by *RepeatMasker* (Smit, AFA, Hubley, R & Green). We  
146 utilise three distinct sources of evidence for selection on TEs: exonic enrichment, strand bias and  
147 evolutionary conservation (Figure 1B).

148 We first asked whether particular TE types are enriched in lncRNA exons, compared to intronic  
149 sequence (Kelley and Rinn 2012). Thus, we calculated the ratio of exonic / intronic sequence coverage by  
150 TEs (Figure 3A). We found enrichment >2-fold for numerous repeat types, including endogenous  
151 retrovirus classes (HERVE-int, HERVK9-int, HERV3-int, LTR12D) in addition to others such as  
152 ALR/Alpha, BSR/Beta and REP522. A number of simple repeats are also enriched in lncRNA, including  
153 GC-rich repeats. A weaker but more generalized trend of enrichment is also observed for various MLT  
154 repeat classes. These findings are consistent with previous analyses by Kelley and Rinn using whole  
155 genome, rather than introns, as background (Kelley and Rinn 2012). Similarly, both studies agree in  
156 finding no difference in *Alu* density between lncRNA exons and intergenic / intronic DNA.

157 Despite their overall abundance throughout the genome, presently-active LINE1 elements are  
158 relatively depleted in lncRNA exons (Figure 3A). It is possible that this reflects selection against

159 disruption to normal gene expression, where numerous weak polyadenylation signals lead to premature  
160 transcription termination when the LINE1 element lies on the same strand as the overlapping gene  
161 (Perepelitsa-Belancio and Deininger 2003). Other explanations may be low transcriptional processivity  
162 exhibited by the LINE1 ORF2 in the sense strand (Perepelitsa-Belancio and Deininger 2003), or else  
163 epigenetic silencing effects (Hollister and Gaut 2009).

164 As a second source of evidence for selection, we searched for TE types displaying a strand  
165 preference relative to host lncRNA (Johnson and Guigó 2014). We were conscious of a major source of  
166 bias: as shown above, many TSS and splice sites of lncRNA are contributed by TEs, and such cases  
167 would lead to artefactual strand bias. To avoid this, we ignored any TEs that overlap an exon-intron  
168 boundary. We calculated the relative strand overlap of all remaining TEs in lncRNA exons. Statistical  
169 significance was assessed by randomisation, with significance defined at  $P < 0.001$ , corresponding to a  
170 false discovery rate (FDR) below 5% (similar cutoffs apply to subsequent analyses, more details may be  
171 found in Materials and Methods) (Figure 3B). In lncRNA exons, a number of TE types are enriched in  
172 either sense or antisense, dominated by LINE1 family members, possibly for the reasons mentioned  
173 above. Other significantly enriched TE types include LTR78, MLT1B, and MIRc (Figure 3B).

174 To test the specificity of this exonic strand bias, we performed equivalent analysis using introns.  
175 Although intronic strand bias is weaker, we did detect a modest yet statistically-significant depletion of  
176 same-strand TE insertions (Supplemental Figure S1). This is especially true for LINE1 elements, possibly  
177 for aforementioned reasons. In contrast to exons, almost no TE types were significantly enriched on the  
178 same-strand in introns.

179 To test for TE type-specific conservation, we turned to two sets of predictions of evolutionarily-  
180 conserved elements. First, the widely-used phastCons conserved elements, based on phylogenetic hidden  
181 Markov model (Siepel et al. 2005) calculated separately on primate, placental mammal and vertebrate  
182 alignments; second, the more recent “Evolutionarily Conserved Structures” (ECS) set (Smith et al. 2013).  
183 Importantly, the phastCons regions are defined based on sequence conservation alone, while the ECS are  
184 defined by phylogenetic analysis of RNA structure evolution.

185 To look for evidence of evolutionary conservation on exonic TEs, we calculated the fraction of  
186 nucleotides overlapped by evolutionarily-conserved genomic elements, and compared to the equivalent  
187 fraction for intronic TEs of the same type. To assess statistical significance, we again used positional  
188 randomisation (see inset in Figure 3C). This pipeline was applied independently to the phastCons  
189 (placental mammal shown in Figure 3C, primate and vertebrate in Supplemental Figure S1B,C) and ECS  
190 (Supplemental Figure S1D) data. The majority of TE types do not exhibit signatures of conservation (grey

191 points). However, for each conservation type, the method detects significant conservation for a minority  
192 of TE types (Figure 3C). This enrichment disappeared when phastCons elements were positionally  
193 randomised (Supplemental Figure S2A). It is unlikely that overlap with protein-coding loci biases the  
194 results, since equivalent analyses using intergenic lncRNAs yielded similar candidate RIDLs  
195 (Supplemental Figure S2B). A similar analysis was performed using protein-coding exons, and although a  
196 number of significantly-conserved TEs were identified, they display limited overlap with those from  
197 lncRNAs (Supplemental Figure S2C). We also found a small number of TEs depleted for signatures of  
198 conservation in lncRNA exons, namely the young *AluSz*, *AluSx* and *AluJb* (phastCons) and LIM4c and  
199 *AluSx1* (ECS) (coloured orange in Figure 3C and Supplemental Figure S1). The cause of this depletion is  
200 unclear, although one explanation is enrichment of conservation in intronic TEs due to RNA-independent  
201 regulatory roles as observed previously (Su et al. 2014).

202 All the selection evidence is summarised in Figure 3D. As might be expected, one observes a  
203 high degree of concordance in candidate TEs identified by the three phastCons methods, in addition to a  
204 smaller number with both phastCons and ECS evidence, including L2b and MIRb. This is not surprising  
205 given the distinct methodologies used to infer conservation. Less concordance is observed between  
206 conservation, enrichment, and strand bias candidates, although some TEs are identified by multiple  
207 methods, such as MIRc (strand bias and ECS).

208

### 209 An annotation of RIDLs

210 We next combined all TE classes with evidence of functionality into a draft annotation of RIDLs.  
211 This annotation combined altogether 99 TE types with at least one type of selection evidence. For each  
212 TE / evidence pair, only those TE instances satisfying that evidence were included. In other words, if  
213 MIRb elements were found to be associated with vertebrate phastCons elements, then *only* those instances  
214 of exonic MIRb elements overlapping such an element would be included in the RIDL annotation, and all  
215 other exonic MIRs would be excluded. This operation was performed for all three phastCons element  
216 types, ECS elements and strand-bias. An example is *CCATI* lncRNA oncogene: it carries three exonic  
217 MIR elements, of which one is defined as a RIDL based on its overlapping a phastCons element (Figure  
218 4A).

219 After removing redundancy, the final RIDL annotation consists of 5374 elements, located within  
220 3566 distinct lncRNA genes (Figure 3D). These represent 12% (5374/46474) of all exonic TE fragments.  
221 The most predominant TE families are MIR and L2 repeats, representing 2329 and 1143 RIDLs (Figure  
222 4B). The majority of both are defined based on evolutionary evidence (Figure 4B, Supplemental Figure

223 S3). In contrast, RIDLs composed by ERV1, low complexity, satellite and simple repeats families are  
224 more frequently identified due to exonic enrichment (Figure 4B). The entire RIDL annotation is available  
225 in Supplemental File S1.

226 It is important to consider this RIDL annotation as an “enrichment group”, with a greater  
227 proportion of functional TEs than when using the entire exonic TE set. Using introns as a reference, we  
228 conservatively estimate the fraction of true positive predictions to range from 12% (strand bias) to 40%  
229 (phastCons primate) and 78% (exonic enrichment) (Supplemental Figure S4).

230 We also examined the evolutionary history of RIDLs. Using 6-mammal alignments, their depth of  
231 evolutionary conservation could be inferred (Supplemental Figure S5). 12% of instances appear to be  
232 Great Ape-specific, with no orthologous sequence beyond chimpanzee. 47% are primate-specific, while  
233 the remaining 40% are identified in at least one non-primate mammal. The wide timeframe for  
234 appearance of RIDLs is consistent with the wide diversity of TE types, from ancient MIR elements to  
235 presently-active LINE1 (Jurka et al. 1995; Smith et al. 2013; Konkel et al. 2010).

236 Instances of genomic TE insertions typically represent a fragment of the full consensus sequence.  
237 We hypothesised that particular regions of the TE consensus will be important for RIDL activity,  
238 introducing selection for these regions that would distinguish them from unselected, intronic copies. To  
239 test this, we compared insertion profiles of RIDLs to intronic instances, for each TE type, and used the  
240 correlation coefficient (CC) as a quantitative measure of similarity (Figure 4C and Supplemental File S2).  
241 For 17 cases, a  $CC < 0.9$  points to possible selective forces acting on RIDL insertions. An example is the  
242 macrosatellite SST1 repeat where RIDL copies in 41 lncRNAs show a strong preference inclusion of the  
243 3' end, in contrast to the general 5' preference observed in introns (Figure 4C). This suggests a possible  
244 functional relevance for the 1000-1500 nt region of the SST1 consensus.

245 To assess whether RIDLs experience purifying evolutionary selection in modern humans, we  
246 analysed the derived allele frequency (DAF) spectrum of their overlapping SNPs (Supplemental Figure  
247 S6) (Haerty and Ponting 2013)(Tan et al. 2017). This showed that RIDLs (orange bars) have a greater  
248 proportion of rare ( $DAF < 0.1$ ) alleles compared to other TEs in exons (green bars) or introns (turquoise  
249 bars) of the same lncRNAs, and indeed compared to non-RIDL exonic nucleotides (black bars). These  
250 differences fail to reach statistical significance, possibly due to small sample sizes. Overall these data are  
251 consistent with RIDLs experiencing an elevated rate of purifying evolutionary selection in modern  
252 humans compared to nearby neutral sequence, although larger datasets will be required before this can be  
253 stated conclusively.

254

## 255 RIDL-carrying lncRNAs are enriched for functions and disease roles

256 We next looked for evidence to support the RIDL annotation by investigating the properties of  
257 their host lncRNAs. We first asked whether RIDLs are randomly distributed amongst lncRNAs, or else  
258 non-randomly clustered in a smaller number of genes. Figure 4D shows that the latter is the case, with a  
259 significant deviation of RIDLs from a random distribution. These lncRNAs carry a mean of 1.15 RIDLs /  
260 kb of exonic sequence (median: 0.84 RIDLs/kb) (Supplemental Figure S7).

261 Are RIDL-lncRNAs more likely to be functional? To address this, we compared lncRNA genes  
262 carrying one or more RIDLs, to a length-matched set of control lncRNAs (Figure 4E, Supplemental  
263 Figure S8). We observed that RIDL-lncRNAs are (1) over-represented in the reference database for  
264 functional lncRNAs, lncRNAdb (Quek et al. 2015), (2) are enriched in associations with cancer and other  
265 diseases, and (3) enriched in their exons for trait/disease-associated SNPs. In order to estimate the impact  
266 of carrying RIDLs on the functional-associated outcomes mentioned above, while controlling for  
267 potential biases from conservation and length, we performed multiple logistic regression analysis. In each  
268 case, the overlap with RIDL-lncRNAs was positive and statistically significant (Figure 4F). However we  
269 did not observed any difference in mean or maximum expression of RIDL-lncRNAs to length matched  
270 controls across ten tissues of the Human Body Map RNA-seq dataset (Supplemental Figure S9).

271 In addition to *CCAT1* (Figure 4A) (Nissan et al. 2012) there are a number of deeply-studied  
272 RIDL-containing genes. *XIST*, the X Chromosome silencing RNA contains seven internal RIDL elements.  
273 As we pointed out previously (Johnson and Guigó 2014) these include an array of four similar pairs of  
274 MIRc / L2b repeats. The prostate cancer-associated *UCA1* gene has a transcript isoform promoted from  
275 an LTR7c, as well as an additional internal RIDL, thereby making a potential link between cancer gene  
276 regulation and RIDLs. The *TUG1* gene, involved in neuronal differentiation, contains highly  
277 evolutionarily-conserved RIDLs including Charlie15k and MLT1K elements (Johnson and Guigó 2014).  
278 Other RIDL-containing lncRNAs include *MEG3*, *MEG9*, *SNHG5*, *ANRIL*, *NEAT1*, *CARMEN1* and  
279 *SOX2OT*. *LINC01206*, located adjacent to *SOX2OT*, also contains numerous RIDLs. A full list can be  
280 found in Supplemental File S3.

281

## 282 Correlation between RIDLs and subcellular localisation of host transcript

283 The location of a lncRNA within the cell is of key importance to its molecular function (Cabili et  
284 al. 2015; Derrien et al. 2012)(Mas-Ponte et al. 2017), therefore we next investigated whether RIDLs  
285 might regulate lncRNA localisation (Zhang et al. 2014a; Hacısuleyman et al. 2016b)(Chillón and Pyle

286 2016) (Figure 5A). Using subcellular RNA-seq data based on 10 ENCODE cell lines (Djebali et al.  
287 2012), we calculated the relative nuclear/cytoplasmic localisation in  $\log_2$  units, or “Relative Concentration  
288 Index” (RCI) (Mas-Ponte et al. 2017). Using this dataset, we tested each of the 99 RIDL types for  
289 association with localisation of their host transcript.

290 After correcting for multiple hypothesis testing using the Benjamini-Hochberg method  
291 (Benjamini and Hochberg 1995), this approach identified four distinct RIDL types: L1PA16, L2b, MIRb  
292 and MIRc (Figure 5B). For example, 44 lncRNAs carrying L2b RIDLs have 6.9-fold higher relative  
293 nuclear/cytoplasmic ratio in IMR-90 cells, and this tendency is observed in six different cell types (Figure  
294 5B,C).

295 The degree of nuclear localisation increases in lncRNAs as a function of the number of RIDLs  
296 (L1PA16, L2b, MIRb and MIRc) they carry (Figure 5D). We also found a significant relationship  
297 between GC-rich elements and cytoplasmic enrichment across three independent cell samples. The GC-  
298 rich-containing lncRNAs have between 2 and 2.3-fold higher relative expression in the cytoplasm of these  
299 cells (Supplemental Figure S10).

300 We were curious whether this relationship with localisation is only a property of RIDLs, or  
301 conversely, holds true when considering any instances of L1PA16, L2b, MIRb and MIRc. Indeed when  
302 the preceding analysis was repeated with unfiltered TE instances, the latter was observed (Supplemental  
303 Figure S11). However, the strength of the effect was consistently lower than for RIDLs (Supplemental  
304 Figure S12). This difference between RIDLs and unfiltered TEs supports both the usefulness of the RIDL  
305 identification method, and the idea that RIDLs are under selection as a result of their effect on  
306 localisation.

307 We were concerned that two un-modelled confounding factors that positively correlated with TE  
308 number could explain the observed data: transcript length and whole-cell gene expression. To address  
309 this, we performed multiple linear regression for localisation with explanatory variables of RIDL number,  
310 transcript length and whole-cell expression (Figure 5E). Such a model accounts independently for each  
311 variable, enabling one to eliminate confounding effects. Training such models for each cell type / RIDL  
312 pair, we observed positive and statistically-significant contributions for RIDL number in most cases. We  
313 also observed weaker but significant contributions from transcript length and whole-cell expression terms,  
314 indicating that our intuition was correct that these factors influence localisation independent of RIDLs  
315 (Supplemental Figure S13A,B). We drew similar conclusions from equivalent analyses using partial  
316 correlation analysis (Supplemental Figure S13C). In summary, observed RIDLs correlate with lncRNA  
317 localisation even when controlling for other factors.

318           Given that L2b and MIR elements predate human-mouse divergence, we attempted to perform  
319 similar analyses in mouse cells. However given that just two equivalent datasets are available at present  
320 (Bahar Halpern et al. 2015)(Tan et al. 2015), as well as the relatively low number of annotated lncRNAs  
321 in mouse, we were unable to draw statistically-robust conclusions regarding the evolutionary conservation  
322 of this phenomenon.

323

#### 324 Intra-gene correlation between RIDLs and subcellular localisation

325           LncRNA gene loci are often composed of multiple, differentially-spliced transcript isoforms that  
326 partially differ in their mature sequence. We reasoned that differential inclusion of RIDL-containing  
327 exons should give rise to differences in localisation amongst transcripts from the same gene locus. In  
328 other words, for RIDL-lncRNA gene loci having multiple transcript isoforms, those isoforms *with* a RIDL  
329 should display greater nuclear enrichment than those isoforms *without* a RIDL (Figure 6, left panel).

330           We tested this individually for each cell type. For every appropriate RIDL-lncRNA locus  
331 (numbers shown inside boxplot), we calculated the difference in the mean of the localisation between  
332 their RIDL and non-RIDL isoforms (Figure 6, right panel). For every cell line, the median difference was  
333 positive, indicating that RIDL-carrying transcript isoforms are more nuclear enriched than their non-  
334 RIDL cousins from the same gene locus. Given our *a priori* hypothesis that RIDLs promote nuclear  
335 enrichment, statistical significance was tested by comparison to zero using a 1-sided *t*-test. Altogether  
336 these data point to a consistent correlation between the presence of certain exonic TE elements, L1PA16,  
337 L2b, MIRb and MIRc, and the nuclear enrichment of their host lncRNA.

338

#### 339 RIDLs play a causative role in lncRNA nuclear localisation

340           To more directly test whether RIDLs play a causative role in nuclear localisation, we designed an  
341 experimental approach to quantify the effect of exonic TEs on localisation of a transfected lncRNA. We  
342 selected three lncRNAs, based on: (i) presence of L2b, MIRb and MIRc RIDLs; (ii) moderate expression;  
343 (iii) nuclear localisation, as inferred from RNA-seq (Figure 7A,B and Supplemental Figure S14). Nuclear  
344 localisation of these candidates could be validated in HeLa cells using qRT-PCR (Figure 7C).

345           We formulated an assay to compare the localisation of transfected lncRNAs carrying wild-type  
346 RIDLs, and mutated versions where the RIDL sequence was randomised without altering sequence  
347 composition (“Mutant”) (Figure 7D, full sequences available in Supplemental File S4). Wild-type and  
348 Mutant lncRNAs were transfected into cultured cells and their localisation evaluated by fractionation.

349 qRT-PCR primers were designed to distinguish transfected Wild-Type and Mutant transcripts from  
350 endogenously-expressed copies. Transgenes were typically expressed in a range of 0.2- to 10-fold  
351 compared to their endogenous transcripts (Supplemental Figure S15). Fractionation purity was verified by  
352 Western blotting (Figure 7E) and qRT-PCR (Figure 7F), and stringent DNase-treatment ensured that  
353 plasmid DNA made negligible contributions to our results (Supplemental Figure S16).

354         With this setup, we compared the nuclear/cytoplasmic localisation of lncRNAs with and without  
355 exonic RIDL sequences (Figure 7F). We observed a potent and consistent impact of RIDLs on  
356 nuclear/cytoplasmic localisation in HeLa cells: for all three candidates, loss of RIDL sequence resulted in  
357 relocalisation of the host transcript from nucleus to cytoplasm (Figure 7F, upper panel). We repeated  
358 these experiments in another cell line, A549, and observed similar, albeit less pronounced, effects (Figure  
359 7F, lower panel). This difference may be due to the less nuclear localisation of the endogenous transcripts  
360 in A549 (Supplemental Figure S17). To summarise, exonic L2b, MIRb and MIRc elements promote the  
361 nuclear enrichment of host lncRNAs.

## 362 Discussion

363           Recent years have seen a rapid increase in the number of annotated lncRNAs. However, our  
364 understanding of their molecular functions, and how such functions are encoded in primary RNA  
365 sequences, lag far behind. Two recent conceptual developments offer hope for resolving the sequence-  
366 function code of lncRNAs: First, the idea that the subcellular localisation of lncRNAs is a readily  
367 quantifiable characteristic that holds important clues to function; Second, that the abundant transposable  
368 element (TE) content of lncRNAs may contribute to functionality.

369           In this study, we have linked these two ideas, by showing evidence that certain TEs can drive the  
370 nuclear enrichment of lncRNAs. A global correlation analysis of TEs and RNA localisation data revealed  
371 a handful of TEs, most notably LINE2b, MIRb and MIRc, which positively and significantly correlate  
372 with the degree of nuclear/cytoplasmic localisation of their host transcripts. This correlation is observed  
373 in multiple cell types, and scales with the number of TEs present. A causative link was established  
374 experimentally, confirming that the indicated TEs are sufficient for a two- to four-fold increase in  
375 nuclear/cytoplasmic localisation. There are two principal explanations for this phenomenon. First, an  
376 “active” process whereby TEs are recognised by a cellular transport pathway, as demonstrated for *Alus* by  
377 Lubelsky and Ulitsky (Lubelsky and Ulitsky 2018). Second, a “passive” process where TEs destabilise  
378 transcripts leading to a concentration gradient from nucleus to cytoplasm. Although future studies will  
379 examine this question in detail, the fact that we do not observe a constant difference in steady-state levels  
380 in TE/mutated transgenes, would be more consistent with the active model.

381           These data support the hypothesis that exonic TE elements can act as functional lncRNA  
382 domains. In this “RIDL hypothesis”, transposable elements are co-opted by natural selection to form  
383 “Repeat Insertion Domains of lncRNA”, that is, fragments of sequence that confer adaptive advantage  
384 through some change in the activity of their host lncRNA. We proposed that RIDLs may serve as binding  
385 sites for proteins or other nucleic acids, and indeed a growing body of evidence supports this (reviewed in  
386 (Johnson and Guigó 2014)). In the context of localisation, RIDLs could mediate nuclear retention through  
387 hybridisation to complementary repeats in genomic DNA or through their described interactions with  
388 nuclear proteins (Kelley et al. 2014). In the course of this study we bioinformatically identified five  
389 candidate proteins (HNRNPU, HNRNPH2, HuR, KHDRBS1, TARDBP), however we could not find  
390 evidence that they contribute to RIDL-lncRNA localisation. Identification of any proteins that mediate  
391 RIDLs’ localisation activity may be achieved in future through pulldown approaches (Marín-Béjar and  
392 Huarte 2015).

393 The localisation RIDLs discovered – MIR and LINE2 - are both ancient and contemporaneous,  
394 being active prior to the mammalian radiation (Cordaux and Batzer 2009). Both have previously been  
395 associated with acquired roles in the context of genomic DNA, but not to our knowledge in RNA (Jjingo  
396 et al. 2014)(Johnson et al. 2006). Although the evolutionary history of lncRNAs remains an active area of  
397 research and accurate dating of lncRNA gene birth is challenging, it appears that the majority of human  
398 lncRNAs were born after the mammalian radiation (Hezroni et al. 2015)(Hezroni et al. 2017)(Necsulea et  
399 al. 2014)(Washietl et al. 2014). This would mean that MIR and LINE2 RIDLs were pre-existing  
400 sequences that were exapted by newly-born lncRNAs, corresponding to the “latent” exaptation model  
401 proposed by Feschotte and colleagues (Chuong et al. 2017). However it is also possible that for other  
402 cases the reverse could be true – a pre-existing lncRNA exapts a newly-inserted TE. Given that nuclear  
403 retention is at odds with the primary needs of natural TE transcripts to be exported to the cytoplasm, we  
404 propose that the observed nuclear localisation activity is a more modern feature of L2b/MIR RIDLs,  
405 which is unrelated to their original roles.

406 Our approach for identifying localisation-regulating RIDLs has advantages over previous studies  
407 (Lubelsky and Ulitsky 2018; Haciosuleyman et al. 2016c) in terms of its genome-wide scale. However an  
408 unavoidable consequence of our use of evolutionary conservation as a filter, is that it likely biases our  
409 analysis against recently-evolved TEs such as *Alus*. It remains entirely possible that modern TEs also  
410 influence lncRNA localisation, but cannot be detected using the signals of selection that we have  
411 employed. On the other hand, MIRb and MIRc were only identified in one cell type each. We expect this  
412 reflects low sensitivity of the statistical screen, rather than cell-type specificity alone, because (i) in a  
413 focussed re-analysis (Supplemental Figure S11) the effect was observed in multiple cells, and (ii)  
414 experimental validation confirmed it in two independent cell types (Figure 7F).

415 This is further supported by the recent study of Lubelsky and Ulitsky, who performed an  
416 experimental screen for localisation motifs in 37 nuclear-enriched lncRNAs, and identified *AluSx* as a  
417 nuclear-localisation element (Lubelsky and Ulitsky 2018). These 37 lncRNAs are enriched for RIDLs  
418 (62% of Lubelsky lncRNAs contain at least one RIDL, compared to 22% for other Gencode v21  
419 lncRNAs,  $P=4 \times 10^{-6}$ , Fisher’s exact test), as well as for the three localisation RIDLs identified here (L2b,  
420 MIRb, MIRc: 32% vs 9%,  $P=3 \times 10^{-4}$ ) (Supplemental Figure S18A). Although our bioinformatic screen  
421 did not identify *AluSx*, a naive unfiltered re-analysis of our data supports Lubelsky’s experimental finding  
422 that *AluSx*-carrying lncRNAs tend to be more nuclear across multiple cell types (Supplemental Figure  
423 S18B). Together, these considerations open the possibility that other localisation-controlling TE types  
424 may await discovery.

425 More generally, the RIDL predictions showed rather low concordance between the various  
426 selection evidence used (Figure 3D). This likely reflects a number of factors: young evolutionary age of  
427 some of the most common TEs, generally low statistical power due to large background of neutral TEs  
428 and multiple hypothesis testing, and false positives due to TEs that promote transcription or splicing of  
429 lncRNAs. However it is worthy of note that validated candidates L2b, MIRb and MIRc are all implicated  
430 by multiple, independent evidence sources (Figure 3D).

431 This work marks a step in the ongoing efforts to map the domains of lncRNAs. Previous studies  
432 have utilised a variety of approaches, from integrating experimental protein-binding data (Van Nostrand  
433 et al. 2016)(Hu et al. 2017)(Li et al. 2014), to evolutionarily-conserved segments (Smith et al.  
434 2013)(Seemann et al. 2017). Previous maps of TEs have highlighted their profound roles in lncRNA gene  
435 evolution (Kapusta et al. 2013)(Kelley and Rinn 2012)(Hezroni et al. 2015). However, the present RIDL  
436 annotation stands apart in attempting to identify the subset of TEs with evidence for selection. We hope  
437 that this RIDL map will prove a resource for future studies to better understand functional domains of  
438 lncRNAs. Although various evidence suggests that the RIDL annotation is a useful enrichment group of  
439 functional TE elements, it contains a substantial false positive (and likely also false negative) rates that  
440 will have to be improved in future.

441 This study may help to explain a longstanding and unexplained property of lncRNAs: their  
442 nuclear enrichment (Derrien et al. 2012)(Ulitsky and Bartel 2013). Although they are readily detected in  
443 the cytoplasm, lncRNAs general tendency is to have higher nuclear/cytoplasmic ratios compared to  
444 mRNAs (Clark et al. 2012)(Ulitsky and Bartel 2013)(Derrien et al. 2012)(Mas-Ponte et al. 2017). This is  
445 true across various human and mouse cell types. Although this may partially be explained by decreased  
446 stability (Mukherjee et al. 2017), it is likely that RNA sequence motifs also contribute to nuclear  
447 localisation (Chillón and Pyle 2016)(Zhang et al. 2014a). Here we show that this is the case, and that the  
448 enrichment of certain RIDL types in lncRNA mature sequences is likely to be a major contributor to  
449 lncRNA nuclear retention. In contrast, the far lower exonic content of TEs in protein-coding mRNAs may  
450 help explain their greater cytoplasmic abundance (Kapusta and Feschotte 2014). Indeed, even within the  
451 cytoplasm, there is evidence that TE content may also influence the efficiency with which lncRNAs are  
452 trafficked to the translation machinery (Carlevaro-Fita et al. 2016). Together, this evidence may reflect  
453 unknown cellular quality control mechanisms that vet RNAs based on their TE content, tending to retain  
454 TE-rich sequences (including lncRNAs or incorrectly processed mRNAs) in the nucleus, and promote the  
455 cytoplasmic export and ribosomal loading of canonical TE-poor mRNAs.

456 In summary therefore, we have made available a first annotation of selected RIDLs in lncRNAs,  
457 and described a new paradigm for TE-derived fragments as drivers of nuclear localisation in lncRNAs.

## 458 Materials and Methods

459 All operations were carried out on human genome version GRCh38/hg38, unless stated otherwise.

460

### 461 Exonic TE curation

462 *RepeatMasker* annotations were downloaded from the UCSC Genome Browser (version hg38) on  
463 December 31st 2014 (Smit, AFA, Hubley, R & Green), and GENCODE version 21 lncRNA annotations  
464 in GTF format were downloaded from [www.encodegenes.org](http://www.encodegenes.org) (Harrow et al. 2012). Annotations were  
465 not filtered further. The ‘*transposon.profiler*’ script, largely based on BEDTools’ *intersect* and *merge*  
466 functionalities (Quinlan and Hall 2010), was used to annotate exonic and intronic TEs of the given gene  
467 annotation (Supplemental Code). Exons of all transcripts belonging to the given gene annotation were  
468 merged, henceforth referred to as “exons”. The set of introns was curated by subtracting the merged  
469 exonic sequences from the full gene spans, and only retaining those introns that belonged to a single gene.  
470 Intronic regions were assigned the strand of the host gene.

471 The *RepeatMasker* annotation file was intersected with exons and classified into one of 6  
472 categories: TSS (transcription start site), overlapping the first exonic nucleotide of the first exon; splice  
473 acceptor, overlapping exon 5’ end; splice donor, overlapping exon 3’ end; internal, residing within an  
474 exon and not overlapping any intronic sequence; encompassing, where an entire exon lies within the TE;  
475 TTS (transcription termination site), overlapping the last nucleotide of the last exon. In every case, the  
476 TEs are separated by strand relative to the host gene: + where both gene and TE are annotated on the  
477 same strand, otherwise -. The result is the “Exonic TE Annotation” (Supplemental File S5).

478

### 479 RIDL identification

480 Using this Exonic TE Annotation, we identified the subset of individual TEs with evidence for  
481 functionality. For certain analysis, an Intronic TE Annotation was also employed, being the output for the  
482 equivalent intron annotation described above. Three different types of evidence were used: enrichment,  
483 strand bias and evolutionary conservation.

484 In enrichment analysis, the exon/intron ratio of the fraction of nucleotide coverage by each repeat  
485 type was calculated. Any repeat type with >2-fold exon/intron ratio was considered as a candidate. All  
486 exonic TE instances belonging to such TE types are defined as RIDLs.

487 In strand bias analysis, a subset of Exonic TE Annotation was used, being the set of non-splice  
 488 junction crossing TE instances (“noSJ”). This additional filter was employed to guard against false  
 489 positive enrichments for TEs known to provide splice sites (Sela et al. 2007; Lev-Maor et al. 2003). For  
 490 all TE instances, the “relative strand” was calculated: positive, if the annotated TE strand matches that of  
 491 the host transcript; negative, if not. Then for every TE type, the ratio of relative strand sense/antisense  
 492 was calculated. Statistical significance was calculated empirically: entire gene structures were randomly  
 493 re-positioned in the genome using *BEDTools shuffle*, and the intersection with the entire *RepeatMasker*  
 494 annotation was re-calculated. For each iteration, sense/antisense ratios were calculated for all TE types. A  
 495 TE type was considered to have significant strand bias, if its true ratio exceeded (positively) all of 1000  
 496 simulations. All exonic instances of these TE types that also have the same strand orientation to the host  
 497 transcript are defined as RIDLs. On the other hand, after inspection of the data, we decided to exclude  
 498 TEs with significant antisense enrichment. This is because most instances were from the LINE1 class,  
 499 which are known to interfere with gene expression when falling on the same strand (Perepelitsa-Belancio  
 500 and Deininger 2003). Therefore, we considered it likely that observed antisense enrichment is simply an  
 501 artefact of selection against insertion on the same strand, and in the interests of controlling the false  
 502 positive prediction rate, decided to exclude these cases.

503 In evolutionary analysis, four different annotations of evolutionarily-conserved regions were  
 504 treated similarly, using unfiltered Exonic TE Annotations. Primate, Placental Mammal and Vertebrate  
 505 phastCons elements based on 46-way alignments were downloaded as BED files from UCSC Genome  
 506 Browser (Siepel et al. 2005), while the ECS conserved regions from obtained from Supplemental Data of  
 507 Smith et al (Smith et al. 2013) (see Supplemental File S6 for summary). Because at the time of analysis,  
 508 phastCons elements were only available for hg19 genome build, we mapped them to hg38 using *LiftOver*  
 509 utility (Hinrichs et al. 2006). For each TE type we calculated the exonic/intronic conservation ratio. To do  
 510 this we used *IntersectBED* (Quinlan and Hall 2010) to overlap exonic locations with TEs, and calculate  
 511 the total number of nucleotides overlapping. We performed a similar operation for intronic regions. Then  
 512 for each TE type, we calculated the ratio of conserved TE nucleotides for exons compared to introns:

$$513 \text{ Relative exonic-intronic conservation (REIC)} = (C_e / (C_e + N_e)) / (C_i / (C_i + N_i))$$

514 Where  $C$  is conserved TE nucleotides,  $N$  is non-conserved TE nucleotides, and subscripts  $e$  and  $i$  denote  
 515 exonic and intronic, respectively. Note that, because it calculates fractional overlap of TEs by conserved  
 516 elements, REIC normalises for different lengths of exons and introns (Supplemental Figure S19).

517 To estimate the background, the conserved element BED files were positionally randomized 1000 times  
 518 using *BEDTools shuffle*, each time recalculating REIC. We considered to be significantly conserved those

519 TE types where the true REIC was greater or less than every one of 1000 randomised REIC values. All  
520 exonic instances of these TE types that also intersect the appropriate evolutionarily conserved element are  
521 defined as RIDLs. This approach of shuffling conserved elements displayed no apparent bias in the length  
522 of TEs it identifies (Supplemental Figure S2D). We also tested an alternative approach for estimating  
523 significance whereby conserved elements were held constant, and TEs were positionally randomised.  
524 While there was a significant overlap in identified candidate RIDLs, this method displayed a bias towards  
525 longer TEs (Supplemental Figure S2D), and therefore was not employed further.

526 We chose to randomise conserved elements, rather than TEs because the former are enriched in  
527 lncRNA exons (Pegueroles and Gabaldón 2016). Thus, using randomised TEs to estimate background  
528 REIC would lead to overestimation of exonic TE conservation, and hence underestimate the rate of  
529 conservation of TEs in real data.

530 All RIDL predictions were then merged using *mergeBED* and any instances with length <10 nt  
531 were discarded. The outcome, a BED format file with coordinates for hg38, is found in Supplemental File  
532 S1.

533 False discovery rates (FDR) were estimated for RIDL predictions. TE type FDR estimates were  
534 based on shuffling simulations described above. Empirical *p*-values for true data were estimated  
535 according to  $P=(\text{rank in distribution})/(1 + \text{number of simulations})$ . For significant cases, where the true  
536 value exceeded all  $n=1000$  simulations, this value was conservatively defined to be  $P=0.001$ . These  
537 empirical *p*-values were then converted to FDR using the R command “p.adjust” with “fdr” setting  
538 (Rackham et al. 2011; R Core Team). Accordingly, empirical significance cutoff ( $P<0.001$ ) mentioned in  
539 the main text corresponds to the following FDR values: Strand bias: 0.027; Vertebrate phastCons: 0.013;  
540 Placental phastCons: 0.014; Primate phastCons: 0.009; ECS: 0.034. This analysis is conservative, since  
541 empirical *p*-values of candidates are rounded up in every case to 0.001.

542 FDR rates were also estimated at the element level. Here, the set of significant TEs were grouped  
543 for each evidence type. Then, the frequency of overlap of these TEs with the evidence type was compared  
544 for lncRNA exons and introns. This data is shown in Supplemental Figure S4.

545

#### 546 RIDL orthology analysis

547 In order to assess evolutionary history of RIDLs, we used chained alignments of human to chimp  
548 (hg19ToPanTro4), macaque (hg19ToRheMac3), mouse (hg19ToMm10), rat (hg19ToRn5), and cow

549 (hg19ToBosTau7). Due to availability of chain files, RIDL coordinates were first converted from hg38 to  
 550 hg19. Orthology was defined by *LiftOver* utility used at default settings (Hinrichs et al. 2006).

551

#### 552 Derived allele frequency (DAF) analysis

553 We used allele frequencies from African population provided by the 1000 Genomes Project (The 1000  
 554 Genomes Project Consortium, 2015), as performed previously by (Haerty and Ponting 2013). DAF was  
 555 determined for human common SNPs from dbSNP (build 150) (Sherry et al. 2001) for every group  
 556 analysed. Ancestral repeats (AR) were defined as human repeats (excluding simple repeats) intersecting at  
 557 least 1 nucleotide of mouse repeats defined by LiftOver, and falling within 5kb of but not overlapping  
 558 RIDL-containing genes.

559

#### 560 Comparing RIDL-carrying lncRNAs versus other lncRNAs

561 In order to test for functional enrichment amongst lncRNAs hosting RIDLs, we tested for statistical  
 562 enrichment of the following traits in RIDL- carrying lncRNAs compared to other lncRNAs (see below)  
 563 by Fisher's exact test:

564 A) Functionally-characterised lncRNAs: lncRNAs from GENCODE v21 that are present in lncRNADB  
 565 (Quek et al. 2015).

566 B) Disease-associated genes: lncRNAs from GENCODE v21 that are present in at least in one of the  
 567 following databases or public sets: lncRNADisease (Chen et al. 2013), lnc2Cancer (Ning et al.  
 568 2016), Cancer lncRNA Census (CLC) (Carlevaro-Fita et al. bioRxiv doi: 10.1101/152769)

569 C) GWAS SNPs: We collected SNPs from the NHGRI-EBI Catalog of published genome-wide  
 570 association studies (Welter et al. 2014; Hindorff et al. 2009) (<https://www.ebi.ac.uk/gwas/home>).

571 We intersected its coordinates with lncRNA exons coordinates.

572 For defining a comparable set of “other lncRNAs” we sampled from the rest of GENCODE v21 a set of  
 573 lncRNAs matching RIDL-lncRNAs’ exonic length distribution (Supplemental Figure S8). We performed  
 574 sampling using *matchDistribution* script: <https://github.com/julienlag/matchDistribution>. In order to  
 575 simultaneously control for both conservation and length, we performed multiple logistic regression  
 576 analysis using glm R function (R Core Team), with the following structure:

577 Functional-association outcome ~ RIDLs + transcript length + exonic conservation

578 Where functional-association outcome indicates A, B and C traits defined above; RIDLs indicates the  
 579 number of RIDL instances in the host gene; transcript length indicates the projected exonic length;

580 conservation indicates the percent of exonic lncRNA nucleotides overlapping the union of primate,  
581 placental mammal and vertebrate phastCons elements. We did not find evidence for multicollinearity in  
582 any case (variance inflation factors (VIF) <1.1). We used the ‘VIF’ command from the R package ‘fmsb’  
583 (Nakazawa 2018).

584

#### 585 Subcellular localisation analysis

586 Processed RNA-seq data from human cell fractions were obtained from ENCODE in the form of  
587 RPKM (reads per kilobase per million mapped reads) quantified against the GENCODE v19 annotation  
588 (Djebali et al. 2012; Mas-Ponte et al. 2017). Only transcripts common to both the v21 and v19  
589 annotations were considered. For the following analysis only one transcript per gene was considered,  
590 defined as the one with largest number of exons. Nuclear/cytoplasmic ratio expression for each transcript  
591 was defined as (nuclear poly(A)+ RPKM)/(cytoplasmic poly(A)+ RPKM), and only transcripts having  
592 non-zero values (at Irreproducible Discovery Rate (IDR) between samples < 1) in both were considered.  
593 These ratios were  $\log_2$ -transformed, to yield the Relative Concentration Index (RCI) (Mas-Ponte et al.  
594 2017). For each RIDL type and cell type in turn, the nuclear/cytoplasmic ratio distribution of RIDL-  
595 containing to non-RIDL-containing lncRNAs was compared using Wilcoxon test. Only RIDLs having at  
596 least three expressed transcripts in at least one cell type were tested. Resulting *p*-values were globally  
597 adjusted to False Discovery Rate using the Benjamini-Hochberg method (Benjamini and Hochberg 1995).

598

#### 599 Multiple linear regression and partial correlation analysis

600 Linear models were created in R using the “lm” function (R Core Team) , at the level of lncRNA  
601 transcripts with the form:

602 localisation ~ RIDL + transcript length + expression

603 Localisation refers to nuclear/cytoplasmic RCI; RIDL denotes the number of instances of a given RIDL in  
604 a transcript; expression denotes the whole cell expression level, as inferred from RNA-seq in units of  
605 RPKM. Equivalent partial correlation analyses were performed, using the R “pcor.test” function from the  
606 ‘ppcor’ package (Spearman correlation) (Kim 2015), correlating RCI with RIDL number, while  
607 controlling for transcript length and expression. We checked all regression models for multicollinearity by  
608 searching for variance inflation factors (VIF) using the ‘VIF’ command from the R package ‘fmsb’  
609 (Nakazawa 2018). In no case did VIF exceed 1.1, below values raising concern of multicollinearity (>4).

610

611 Cell lines and reagents

612 Human cervical cancer cell line HeLa and human lung cancer cell line A549 were cultured in  
613 Dulbecco's Modified Eagle's medium (Sigma-Aldrich, # D5671) supplemented with 10% FBS and 1%  
614 penicillin streptomycin at 37°C and 5 % CO<sub>2</sub>. Anti-GAPDH antibody (Sigma-Aldrich, # G9545) and anti-  
615 histone H3 antibody (Abcam, # ab24834) were used for Western blot analysis.

616

617 Gene synthesis and cloning of lncRNAs

618 The three lncRNA sequences (RP11-5407, LINC00173, RP4-806M20.4) containing wild-type  
619 RIDLs, and corresponding mutated versions where RIDL sequence has been randomised ("Mutant"),  
620 were synthesized commercially (BioCat GmbH). For each gene locus, only one transcript contained the  
621 RIDL(s), and was chosen for experimental study. The sequences were cloned into pcDNA 3.1 (+) vector  
622 within the *NheI* and *XhoI* restriction enzyme sites. The clones were checked by restriction digestion and  
623 Sanger sequencing. The sequence of the wild-type and mutant clones are provided in Supplemental File  
624 S4.

625

626 lncRNA transfection and sub-cellular fractionation

627 Wild-type and mutant lncRNA clones for each tested gene were transfected independently in  
628 separate wells of a 6-well plate. Transfections and subsequent analysis were repeated as biological  
629 replicates (four for HeLa, four for A549), defined as transfections performed on different days with  
630 different cell passages. Transfections were carried out with 2 µg of total plasmid DNA in each well using  
631 Lipofectamine 2000. 48 h post-transfection, cells from each well were harvested, pooled and re-seeded  
632 into a 10 cm dish and allowed to grow till 100% confluence. Expression of transgenes was checked by qRT-  
633 PCR using specific primers, and found to typically be several-fold greater than endogenous copies (HeLa)  
634 or from 0.2- to 1-fold (A549) (Supplemental Figure S15).

635 The nuclear and cytoplasmic fractionation was carried out as described previously (Suzuki et al.  
636 2010) with minor modifications. In brief, cells from 10 cm dishes were harvested by scraping and washed  
637 with 1x ice-cold PBS. For fractionation, cell pellet was re-suspended in 900 µl of ice-cold 0.1% NP-40 in  
638 PBS and triturated 7 times using a p1000 micropipette. 300 µl of the cell lysate was saved as the whole  
639 cell lysate. The remaining 600 µl of the cell lysate was centrifuged for 30 sec on a table top centrifuge and  
640 the supernatant was collected as "cytoplasmic fraction". 300 µl from the cytoplasmic supernatant was  
641 kept for RNA isolation and the remaining 300 µl was saved for protein analysis by Western blot. The

642 pellet containing the intact nuclei was washed with 1 ml of 0.1% NP-40 in PBS. The nuclear pellet was  
643 re-suspended in 200  $\mu$ l 1X PBS and subjected to a quick sonication of 3 pulses with 2 sec ON-2 sec OFF  
644 to lyse the nuclei and prepare the “nuclear fraction”. 100  $\mu$ l of nuclear fraction was saved for RNA  
645 isolation and the remaining 100  $\mu$ l was kept for Western blot.

646

#### 647 RNA isolation and real time PCR

648 The RNA from each nuclear and cytoplasmic fraction was isolated using Quick-RNA MiniPrep kit  
649 (ZYMO Research, # R1055). The RNAs were subjected to on-column DNase I treatment and clean up  
650 using the manufacturer’s protocol. For A549 samples, additional units of DNase were employed, due to  
651 residual signal in –RT samples. The RNA from each fraction was converted to cDNA using GoScript  
652 reverse transcriptase (Promega, # A5001) and random hexamer primers. The expression of each of the  
653 individual transcripts was quantified by qRT-PCR (Applied Biosystems® 7500 Real-Time) using  
654 indicated primers (Supplemental File S7) and GoTaq qPCR master mix (Promega, # A6001). In order to  
655 distinguish expression of transfected wild-type genes from endogenous copies, we designed forward  
656 primers against a transcribed region of the expression vector backbone. Human *GAPDH* mRNA and  
657 *MALAT1* lncRNA were used as cytoplasmic and nuclear markers, respectively. The absence of  
658 contaminating plasmid DNA in cDNA was checked for all samples using qPCR (see Supplemental Figure  
659 S16 for a representative example).

660

#### 661 Western Blotting

662 The protein concentration of each of the fractions was determined, and equal amounts of protein (50  $\mu$ g)  
663 from whole cell lysate, cytoplasmic fraction, and nuclear fraction were resolved on 12 % Tris-glycine  
664 SDS-polyacrylamide gels and transferred onto polyvinylidene fluoride (PVDF) membranes (VWR, #  
665 1060029). Membranes were blocked with 5% skimmed milk and incubated overnight at 4°C with anti-  
666 GAPDH antibody as a cytoplasmic marker and anti p-histone H3 antibody as nuclear marker. Membranes  
667 were washed with PBS-T (1X PBS with 0.1 % Tween 20) followed by incubation with HRP-conjugated  
668 anti-rabbit or anti-mouse secondary antibodies respectively. The bands were detected using  
669 SuperSignal™ West Pico chemiluminescent substrate (Thermo Fisher Scientific, # 34077).

670

#### 671 Software availability

672 “*transposon.profiler*”, is available on Github at [https://github.com/gold-lab/shared\\_scripts](https://github.com/gold-lab/shared_scripts) and in  
673 Supplemental Code.

674 **Acknowledgements**

675           We wish to thank Roderic Guigó (CRG), Marc Friedlaender (SciLife Lab) and Marta Melé  
676 (Harvard) for many helpful discussions. Roberta Esposito (DBMR) and Samir Ounzain (CHUV)  
677 contributed valuable suggestions regarding experimental design and analysis. Julien Lagarde (CRG)  
678 kindly provided help in gene sampling analysis. Carlos Pulido (DBMR) assisted with RNA-seq analysis,  
679 and Reza Sodaie (CRG) helped with combinatorial analysis of TEs. We acknowledge Deborah Re  
680 (DBMR), Silvia Roesselet (DBMR) and Marianne Zahn (Inselspital) for administrative support. CN is  
681 supported by grants TIN-2013-41990-R and DPI-2017-84439-R from the Spanish Ministry of Economy,  
682 Industry and Competitiveness (MINECO). This research was funded by the NCCR “RNA & Disease”  
683 funded by the Swiss National Science Foundation, and by the Medical Faculty of the University and  
684 University Hospital of Bern.

685 **References**

686

687 Bahar Halpern K, Caspi I, Lemze D, Levy M, Landen S, Elinav E, Ulitsky I, Itzkovitz S. 2015. Nuclear  
688 Retention of mRNA in Mammalian Tissues. *Cell Rep* **13**: 2653–2662.

689 Bassett AR, Akhtar A, Barlow DP, Bird AP, Brockdorff N, Duboule D, Ephrussi A, Ferguson-Smith AC,  
690 Gingeras TR, Haerty W, et al. 2014. Considerations when investigating lncRNA function in vivo.  
691 *Elife* **3**: e03058.

692 Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful  
693 Approach to Multiple Testing. *J R Stat Soc Ser B* **57**: 289–300.

694 Benoit Bouvrette LP, Cody NAL, Bergalet J, Lefebvre FA, Diot C, Wang X, Blanchette M, Lécuyer E.  
695 2018. CeFra-seq reveals broad asymmetric mRNA and noncoding RNA distribution profiles in  
696 *Drosophila* and human cells. *RNA* **24**: 98–113.

697 Blackwell BJ, Lopez MF, Wang J, Krastins B, Sarracino D, Tollervey JR, Dobke M, Jordan IK, Lunyak  
698 V V. 2012. Protein interactions with piALU RNA indicates putative participation of retroRNA in  
699 the cell cycle, DNA repair and chromatin assembly. *Mob Genet Elements* **2**: 26–35.

700 Bourque G. 2009. Transposable elements in gene regulation and in the evolution of vertebrate genomes.  
701 *Curr Opin Genet Dev* **19**: 607–612.

702 Bourque G, Leong B, Vega VB, Chen X, Lee YL, Srinivasan KG, Chew J-L, Ruan Y, Wei C-L, Ng HH,  
703 et al. 2008. Evolution of the mammalian transcription factor binding repertoire via transposable  
704 elements. *Genome Res* **18**: 1752–62.

705 Cabili MN, Dunagin MC, McClanahan PD, Biaesch A, Padovan-Merhar O, Regev A, Rinn JL, Raj A.  
706 2015. Localization and abundance analysis of human lncRNAs at single-cell and single-molecule  
707 resolution. *Genome Biol* **16**: 20.

708 Carlevaro-Fita J, Rahim A, Guigó R, Vardy LA, Johnson R. 2016. Cytoplasmic long noncoding RNAs  
709 are frequently bound to and degraded at ribosomes in human cells. *RNA* **22**: 867–82.

710 Carrieri C, Cimatti L, Biagioli M, Beugnet A, Zucchelli S, Fedele S, Pesce E, Ferrer I, Collavin L,  
711 Santoro C, et al. 2012. Long non-coding antisense RNA controls Uchl1 translation through an  
712 embedded SINEB2 repeat. *Nature* **491**: 454–7.

713 Chen G, Wang Z, Wang D, Qiu C, Liu M, Chen X, Zhang Q, Yan G, Cui Q. 2013. LncRNADisease: a  
714 database for long-non-coding RNA-associated diseases. *Nucleic Acids Res* **41**: D983–D986.

715 Chen L-L. 2016. Linking Long Noncoding RNA Localization and Function. *Trends Biochem Sci*.

716 Chillón I, Pyle AM. 2016. Inverted repeat elements in the human lincRNA-p21 adopt a conserved  
717 secondary structure that regulates RNA function. *Nucleic Acids Res* gkw599.

718 Chuong EB, Elde NC, Feschotte C. 2017. Regulatory activities of transposable elements: from conflicts

719 to benefits.

720 Clark MB, Johnston RL, Inostroza-Ponta M, Fox AH, Fortini E, Moscato P, Dinger ME, Mattick JS.  
721 2012. Genome-wide analysis of long noncoding RNA stability. *Genome Res* **22**: 885–98.

722 Cordaux R, Batzer MA. 2009. The impact of retrotransposons on human genome evolution. *Nat Rev*  
723 *Genet* **10**: 691–703.

724 Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D, Merkel A,  
725 Knowles DG, et al. 2012. The GENCODE v7 catalog of human long noncoding RNAs: analysis of  
726 their gene structure, evolution, and expression. *Genome Res* **22**: 1775–89.

727 Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,  
728 Schlesinger F, et al. 2012. Landscape of transcription in human cells. *Nature* **489**: 101–108.

729 Elisaphenko EA, Kolesnikov NN, Shevchenko AI, Rogozin IB, Nesterova TB, Brockdorff N, Zakian SM.  
730 2008. A Dual Origin of the Xist Gene from a Protein-Coding Gene and a Set of Transposable  
731 Elements.

732 Faulkner GJ, Kimura Y, Daub CO, Wani S, Plessy C, Irvine KM, Schroder K, Cloonan N, Steptoe AL,  
733 Lassmann T, et al. The regulated retrotransposon transcriptome of mammalian cells.

734 Feschotte C. 2008. Transposable elements and the evolution of regulatory networks. *Nat Rev Genet* **9**:  
735 397–405.

736 Gibbs RA, Boerwinkle E, Doddapaneni H, Han Y, Korchina V, Kovar C, Lee S, Muzny D, Reid JG, Zhu  
737 Y, et al. 2015. A global reference for human genetic variation. *Nature* **526**: 68–74.

738 Gong C, Maquat LE. 2011. lncRNAs transactivate STAU1-mediated mRNA decay by duplexing with 3’  
739 UTRs via Alu elements. *Nature* **470**: 284–8.

740 Guttman M, Rinn JL. 2012. Modular regulatory principles of large non-coding RNAs. *Nature* **482**: 339–  
741 46.

742 Hacisuleyman E, Goff LA, Trapnell C, Williams A, Henao-Mejia J, Sun L, McClanahan P, Hendrickson  
743 DG, Sauvageau M, Kelley DR, et al. 2014. Topological organization of multichromosomal regions  
744 by the long intergenic noncoding RNA Firre. *Nat Struct Mol Biol* **21**: 198–206.

745 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL. 2016a. Function and evolution of local repeats in the  
746 Firre locus. *Nat Commun* **7**: 11021.

747 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL, Koning AP de, Gu W, Castoe TA, Batzer MA, Pollock  
748 DD, Wicker T, et al. 2016b. Function and evolution of local repeats in the Firre locus. *Nat Commun*  
749 **7**: 11021.

750 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL, Koning AP de, Gu W, Castoe TA, Batzer MA, Pollock  
751 DD, Wicker T, et al. 2016c. Function and evolution of local repeats in the Firre locus. *Nat Commun*  
752 **7**: 11021.

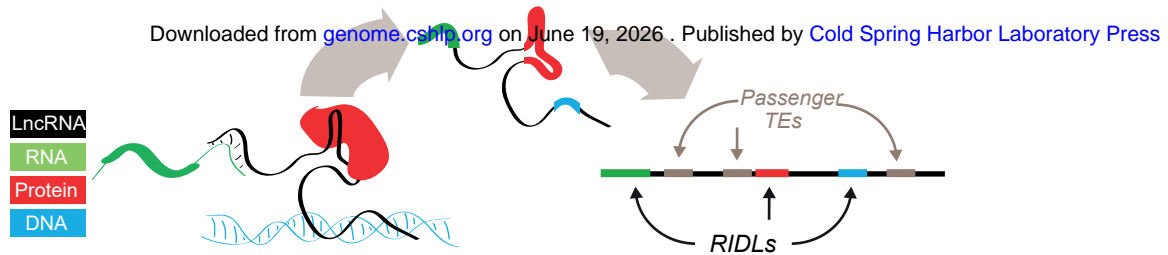
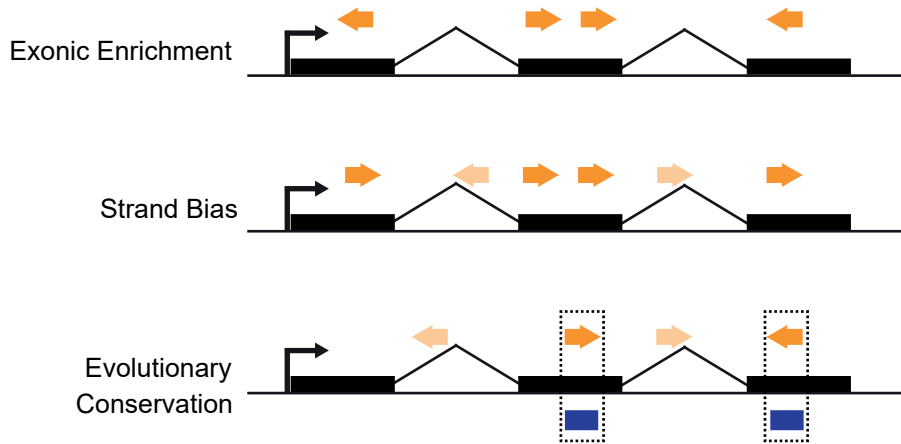
- 753 Haerty W, Ponting CP. 2013. Mutations within lncRNAs are effectively selected against in fruitfly but  
754 not in human. *Genome Biol* **14**: R49.
- 755 Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D,  
756 Zadissa A, Searle S, et al. 2012. GENCODE: The reference human genome annotation for The  
757 ENCODE Project. *Genome Res* **22**: 1760–1774.
- 758 Hezroni H, Ben-Tov Perry R, Meir Z, Housman G, Lubelsky Y, Ulitsky I. 2017. A subset of conserved  
759 mammalian long non-coding RNAs are fossils of ancestral protein-coding genes. *Genome Biol* **18**:  
760 162.
- 761 Hezroni H, Koppstein D, Schwartz MG, Avrutin A, Bartel DP, Ulitsky I. 2015. Principles of Long  
762 Noncoding RNA Evolution Derived from Direct Comparison of Transcriptomes in 17 Species. *Cell*  
763 *Rep* **11**: 1110–1122.
- 764 Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, Manolio TA. 2009. Potential  
765 etiologic and functional implications of genome-wide association loci for human diseases and traits.  
766 *Proc Natl Acad Sci* **106**: 9362–9367.
- 767 Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, Diekhans M, Furey TS,  
768 Harte RA, Hsu F, et al. 2006. The UCSC Genome Browser Database: update 2006. *Nucleic Acids*  
769 *Res* **34**: D590–D598.
- 770 Holdt LM, Hoffmann S, Sass K, Langenberger D, Scholz M, Krohn K, Finstermeier K, Stahringer A,  
771 Wilfert W, Beutner F, et al. 2013. Alu Elements in ANRIL Non-Coding RNA at Chromosome 9p21  
772 Modulate Atherogenic Cell Functions through Trans-Regulation of Gene Networks. *PLoS Genet* **9**.
- 773 Hollister JD, Gaut BS. 2009. Epigenetic silencing of transposable elements: a trade-off between reduced  
774 transposition and deleterious effects on neighboring gene expression. *Genome Res* **19**: 1419–28.
- 775 Hu B, Yang Y-CT, Huang Y, Zhu Y, Lu ZJ. 2017. POSTAR: a platform for exploring post-transcriptional  
776 regulation coordinated by RNA-binding proteins. *Nucleic Acids Res* **45**: D104–D114.
- 777 Huda A, Bowen NJ, Conley AB, Jordan IK. 2011. Epigenetic regulation of transposable element derived  
778 human gene promoters. *Gene* **475**: 39–48.
- 779 Jjingo D, Conley AB, Wang J, Mariño-Ramírez L, Lunyak V V, Jordan IK. 2014. Mammalian-wide  
780 interspersed repeat (MIR)-derived enhancers and the regulation of human gene expression. *Mob*  
781 *DNA* **5**: 14.
- 782 Johnson R, Guigó R. 2014. The RIDL hypothesis: transposable elements as functional domains of long  
783 noncoding RNAs. *RNA* **20**: 959–76.
- 784 Johnson R, W.B. D, B.T. L, J.R. A, S.L. G, M. M, G.J. W, C.M. M, B.T. L, R. R. 2006. Identification of  
785 the REST regulon reveals extensive transposable element-mediated binding site duplication. *Nucleic*  
786 *Acids Res* **34**: 3862–3877.

- 787 Jurka J, Zietkiewicz E, Labuda D. 1995. Ubiquitous mammalian-wide interspersed repeats (MIRs) are  
788 molecular fossils from the mesozoic era. *Nucleic Acids Res* **23**: 170–5.
- 789 Kapusta A, Feschotte C. 2014. Volatile evolution of long noncoding RNA repertoires: mechanisms and  
790 biological implications. *Trends Genet* **30**: 439–452.
- 791 Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, Yandell M, Feschotte C. 2013.  
792 Transposable elements are major contributors to the origin, diversification, and regulation of  
793 vertebrate long noncoding RNAs. ed. H.E. Hoekstra. *PLoS Genet* **9**: e1003470.
- 794 Kelley D, Rinn J. 2012. Transposable elements reveal a stem cell-specific class of long noncoding RNAs.  
795 *Genome Biol* **13**: R107.
- 796 Kelley DR, Hendrickson DG, Tenen D, Rinn JL. 2014. Transposable elements modulate human RNA  
797 abundance and splicing via specific RNA-protein interactions. *Genome Biol* **15**: 537.
- 798 Kim S. 2015. ppcor: Partial and Semi-Partial (Part) Correlation.
- 799 Konkel MK, Walker JA, Batzer MA. 2010. LINEs and SINEs of primate evolution. *Evol Anthropol* **19**:  
800 236–249.
- 801 Lagarde J, Uszczynska-Ratajczak B, Carbonell S, Pérez-Lluch S, Abad A, Davis C, Gingeras TR,  
802 Frankish A, Harrow J, Guigo R, et al. 2017. High-throughput annotation of full-length long  
803 noncoding RNAs with capture long-read sequencing. *Nat Genet* **49**: 1731–1740.
- 804 Lev-Maor G, Sorek R, Shomron N, Ast G. 2003. The Birth of an Alternatively Spliced Exon: 3' Splice-  
805 Site Selection in Alu Exons. *Science (80- )* **300**.
- 806 Li J-H, Liu S, Zhou H, Qu L-H, Yang J-H. 2014. starBase v2.0: decoding miRNA-ceRNA, miRNA-  
807 ncRNA and protein–RNA interaction networks from large-scale CLIP-Seq data. *Nucleic Acids Res*  
808 **42**: D92–D97.
- 809 Lubelsky Y, Ulitsky I. 2018. Sequences enriched in Alu repeats drive nuclear localization of long RNAs  
810 in human cells. *Nature* **555**: 107–111.
- 811 Marín-Béjar O, Huarte M. 2015. RNA Pulldown Protocol for In Vitro Detection and Identification of  
812 RNA-Associated Proteins. In *Methods in molecular biology (Clifton, N.J.)*, Vol. 1206 of, pp. 87–95.
- 813 Marín-Béjar O, Mas AM, González J, Martínez D, Athie A, Morales X, Galduroz M, Raimondi I, Grossi  
814 E, Guo S, et al. 2017. The human lncRNA LINC-PINT inhibits tumor cell invasion through a highly  
815 conserved sequence element. *Genome Biol* **18**: 202.
- 816 Martin KC, Ephrussi A. 2009. mRNA localization: gene expression in the spatial dimension. *Cell* **136**:  
817 719–30.
- 818 Mas-Ponte D, Carlevaro-Fita J, Palumbo E, Hermoso Pulido T, Guigo R, Johnson R. 2017. LncAtlas  
819 database for subcellular localisation of long noncoding RNAs. *RNA* rna.060814.117.
- 820 Mercer TR, Mattick JS. 2013. Structure and function of long noncoding RNAs in epigenetic regulation.

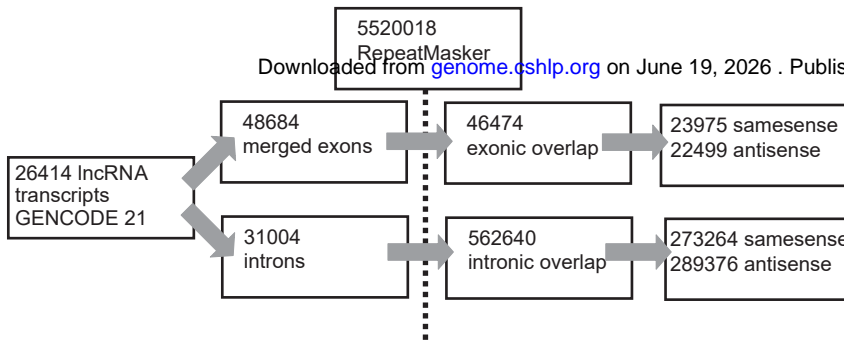
- 821 *Nat Publ Gr* **20**.
- 822 Mukherjee N, Calviello L, Hirsekorn A, de Pretis S, Pelizzola M, Ohler U. 2017. Integrative classification  
823 of human coding and noncoding genes through RNA metabolism profiles. *Nat Struct Mol Biol* **24**:  
824 86–96.
- 825 Nakazawa M. 2018. fmsb: Functions for Medical Statistics Book with some Demographic Data.
- 826 Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Zeller U, Baker JC, Grützner F, Kaessmann  
827 H. 2014. The evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* **505**:  
828 635–40.
- 829 Ning S, Zhang J, Wang P, Zhi H, Wang J, Liu Y, Gao Y, Guo M, Yue M, Wang L, et al. 2016.  
830 Lnc2Cancer: a manually curated database of experimentally supported lncRNAs associated with  
831 various human cancers. *Nucleic Acids Res* **44**: D980–D985.
- 832 Nissan A, Stojadinovic A, Mitrani-Rosenbaum S, Halle D, Grinbaum R, Roistacher M, Bochem A,  
833 Dayanc BE, Ritter G, Gomceli I, et al. 2012. Colon cancer associated transcript-1: A novel RNA  
834 expressed in malignant and pre-malignant human tissues. *Int J Cancer* **130**: 1598–1606.
- 835 Pegueroles C, Gabaldón T. 2016. Secondary structure impacts patterns of selection in human lncRNAs.  
836 *BMC Biol* **14**: 60.
- 837 Perepelitsa-Belancio V, Deininger P. 2003. RNA truncation by premature polyadenylation attenuates  
838 human mobile element activity. *Nat Genet* **35**: 363–366.
- 839 Quek XC, Thomson DW, Maag JL V, Bartonicek N, Signal B, Clark MB, Gloss BS, Dinger ME. 2015.  
840 lncRNADB v2.0: expanding the reference database for functional long noncoding RNAs. *Nucleic*  
841 *Acids Res* **43**: D168-73.
- 842 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.  
843 *Bioinformatics* **26**: 841–2.
- 844 Rackham O, Shearwood A-MJ, Mercer TR, Davies SMK, Mattick JS, Filipovska A. 2011. Long  
845 noncoding RNAs are generated from the mitochondrial genome and regulated by nuclear-encoded  
846 proteins. *RNA* **17**: 2085–93.
- 847 R Core Team. R: A Language and Environment for Statistical Computing.
- 848 Roberts JT, Cardin SE, Borchert GM. 2014. Burgeoning evidence indicates that microRNAs were  
849 initially formed from transposable element sequences. *Mob Genet Elements* **4**: e29255.
- 850 Schmitt AM, Chang HY, Abdelmohsen K, Panda A, Kang MJ, Xu J, Selimyan R, Yoon JH, Martindale  
851 JL, De S, et al. 2016. Long Noncoding RNAs in Cancer Pathways. *Cancer Cell* **29**: 452–463.
- 852 Seemann SE, Mirza AH, Hansen C, Bang-Berthelsen CH, Garde C, Christensen-Dalsgaard M,  
853 Torarinsson E, Yao Z, Workman CT, Pociot F, et al. 2017. The identification and functional  
854 annotation of RNA structures conserved in vertebrates. *Genome Res* **27**: 1371–1383.

- 855 Sela N, Mersch B, Gal-Mark N, Lev-Maor G, Hotz-Wagenblatt A, Ast G. 2007. Comparative analysis of  
856 transposed element insertion within human and mouse genomes reveals Alu's unique role in shaping  
857 the human transcriptome. **8**.
- 858 Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP: the  
859 NCBI database of genetic variation. *Nucleic Acids Res* **29**: 308–311.
- 860 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier  
861 LW, Richards S, et al. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and  
862 yeast genomes. *Genome Res* **15**: 1034–50.
- 863 Smit, AFA, Hubley, R & Green P. RepeatMasker Open-4.0. <<http://www.repeatmasker.org>>.
- 864 Smith MA, Gesell T, Stadler PF, Mattick JS. 2013. Widespread purifying selection on RNA structure in  
865 mammals. *Nucleic Acids Res* **41**: 8220–8236.
- 866 Su M, Han D, Boyd-Kirkup J, Yu X, Han J-DJ. 2014. Evolution of Alu elements toward enhancers. *Cell*  
867 *Rep* **7**: 376–85.
- 868 Suzuki K, Bose P, Leong-Quong RYY, Fujita DJ, Riabowol K. 2010. REAP: A two minute cell  
869 fractionation method. *BMC Res Notes* **3**: 294.
- 870 Tan JY, Sirey T, Honti F, Graham B, Piovesan A, Merkschlager M, Webber C, Ponting CP, Marques  
871 AC. 2015. Extensive microRNA-mediated crosstalk between lincRNAs and mRNAs in mouse  
872 embryonic stem cells. *Genome Res* **25**: 655–666.
- 873 Tan JY, Smith AAT, Ferreira da Silva M, Matthey-Doret C, Rueedi R, Sönmez R, Ding D, Kutalik Z,  
874 Bergmann S, Marques AC. 2017. cis-Acting Complex-Trait-Associated lincRNA Expression  
875 Correlates with Modulation of Chromosomal Architecture. *Cell Rep* **18**: 2280–2288.
- 876 Ulitsky I, Bartel DP. 2013. lincRNAs: Genomics, Evolution, and Mechanisms. *Cell* **154**: 26–46.
- 877 Van Nostrand EL, Pratt GA, Shishkin AA, Gelboin-Burkhart C, Fang MY, Sundararaman B, Blue SM,  
878 Nguyen TB, Surka C, Elkins K, et al. 2016. Robust transcriptome-wide discovery of RNA-binding  
879 protein binding sites with enhanced CLIP (eCLIP). *Nat Methods* **13**: 508–514.
- 880 Washietl S, Kellis M, Garber M. 2014. Evolutionary dynamics and tissue specificity of human long  
881 noncoding RNAs in six mammals. *Genome Res* **24**: 616–28.
- 882 Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T,  
883 Hindorff L, et al. 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations.  
884 *Nucleic Acids Res* **42**: D1001–D1006.
- 885 Zhang B, Gunawardane L, Niazi F, Jahanbani F, Chen X, Valadkhan S. 2014a. A Novel RNA Motif  
886 Mediates the Strict Nuclear Localization of a Long Noncoding RNA. *Mol Cell Biol* **34**: 2318–2329.
- 887 Zhang K, Shi Z-M, Chang Y-N, Hu Z-M, Qi H-X, Hong W. 2014b. The ways of action of long non-  
888 coding RNAs in cytoplasm and nucleus. *Gene* **547**: 1–9.

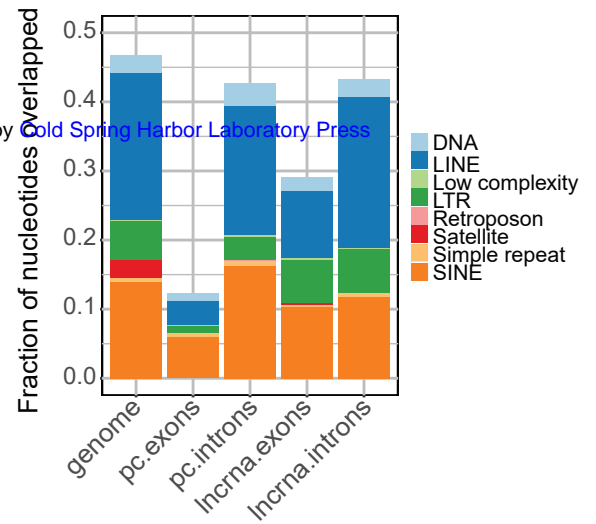


**A****B**

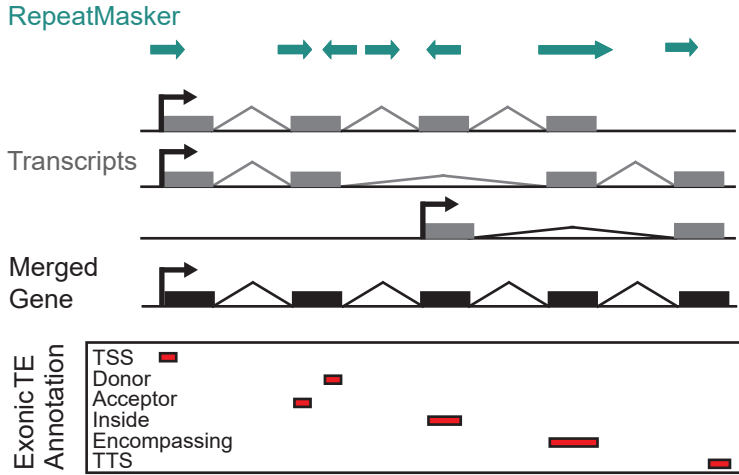
**A**



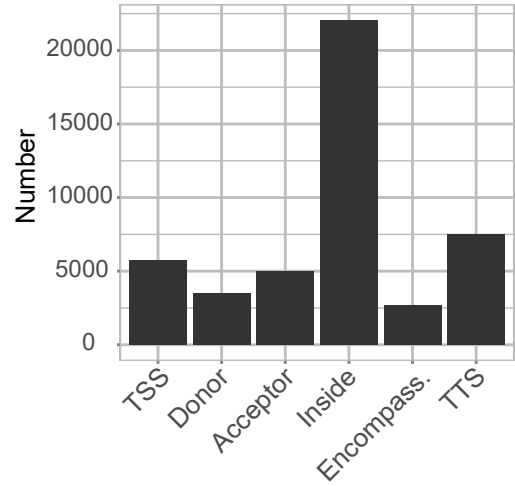
**B**



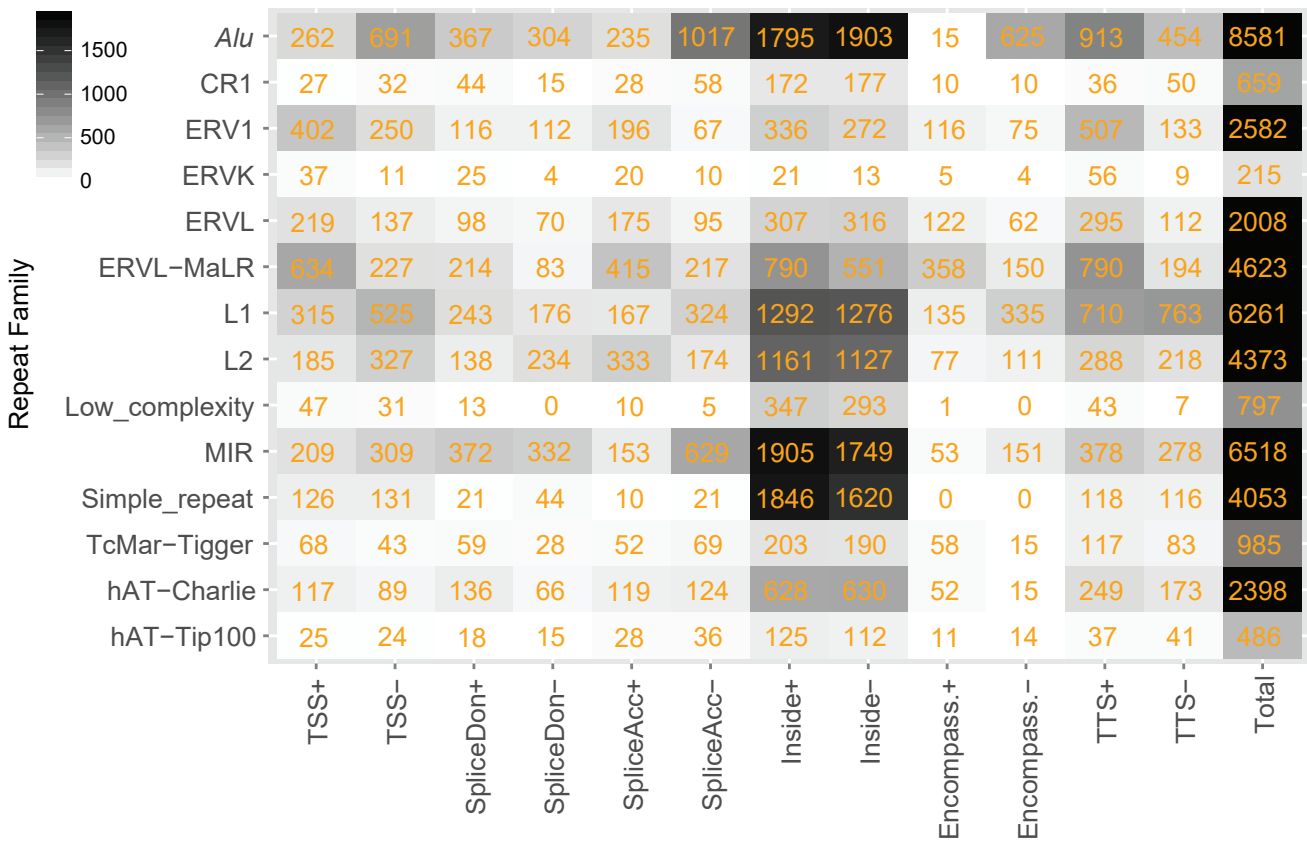
**C**



**D**



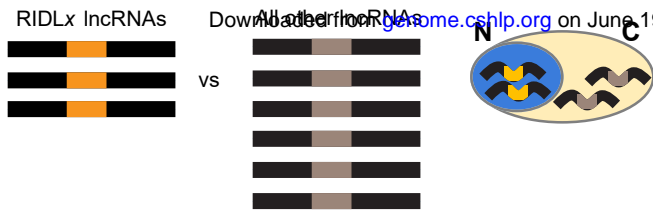
**E**



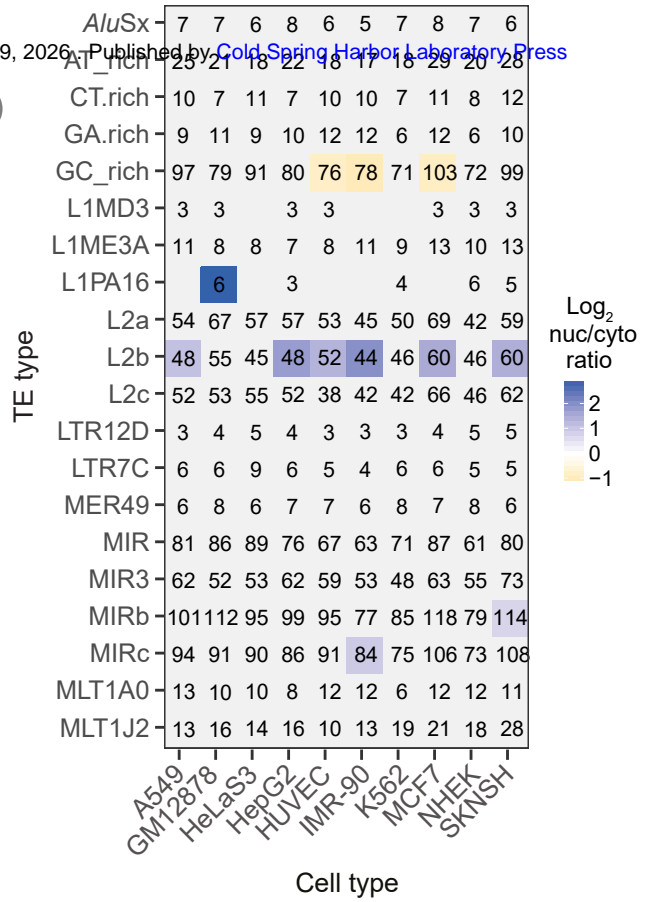




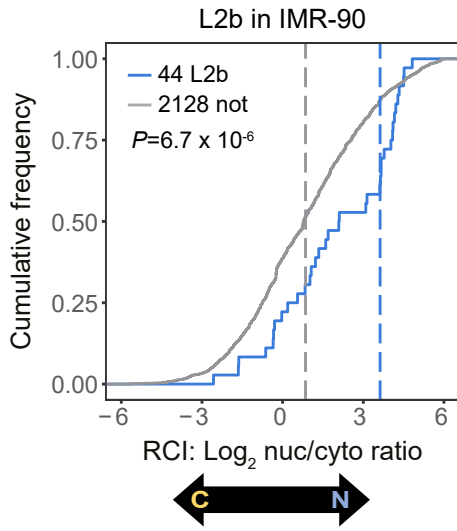
**A**



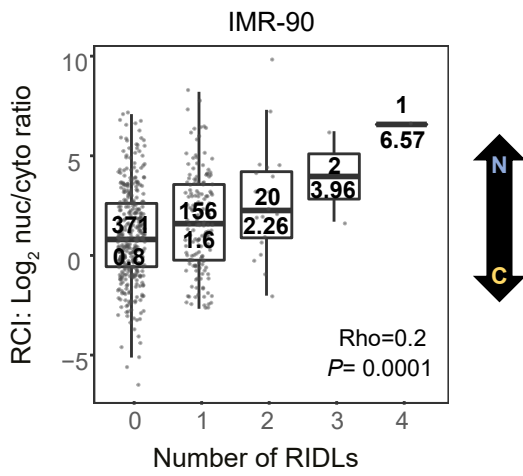
**B**



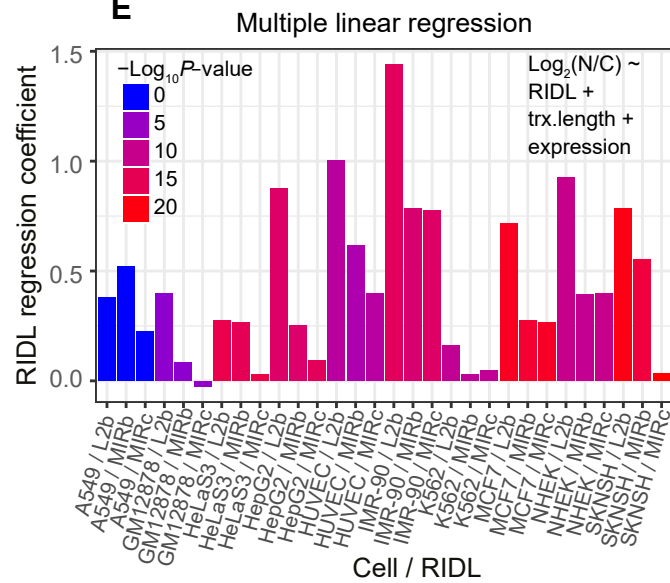
**C**

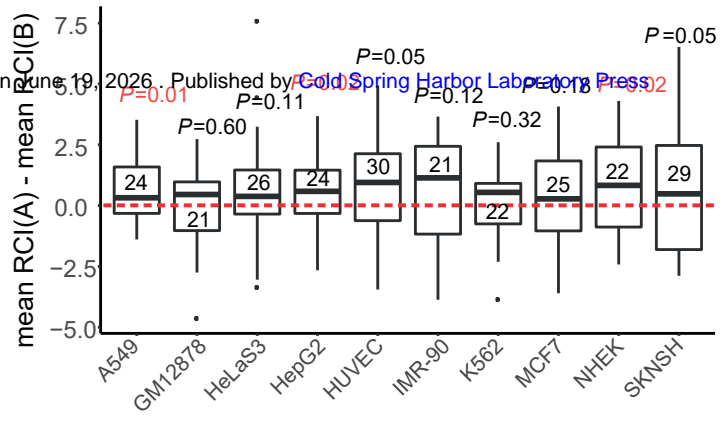
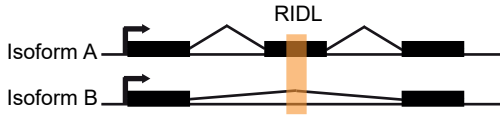


**D**



**E**





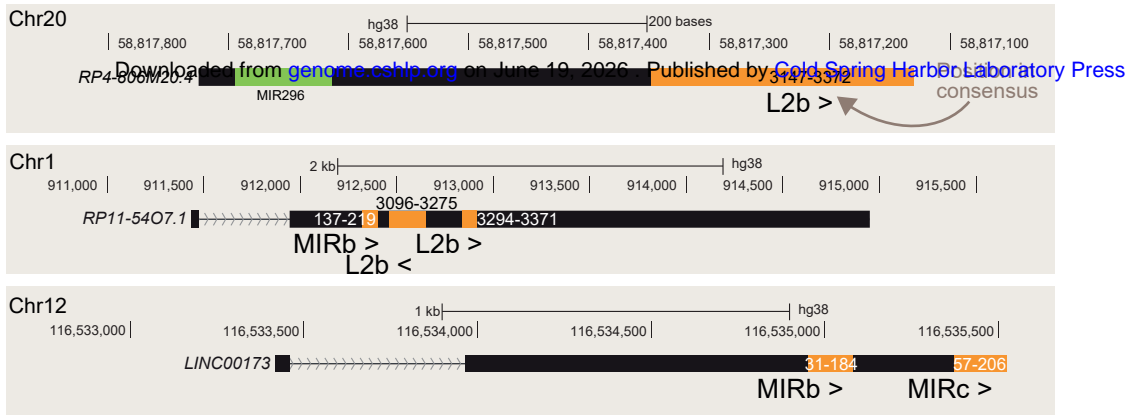
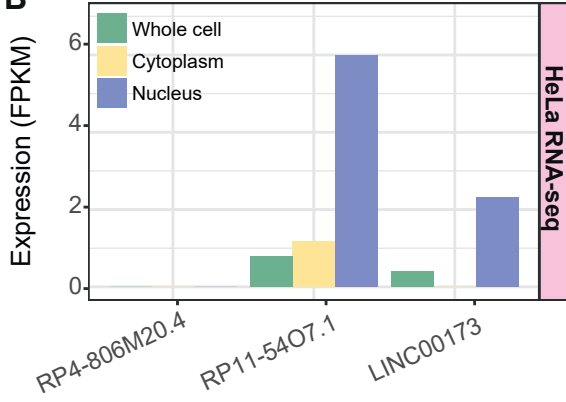
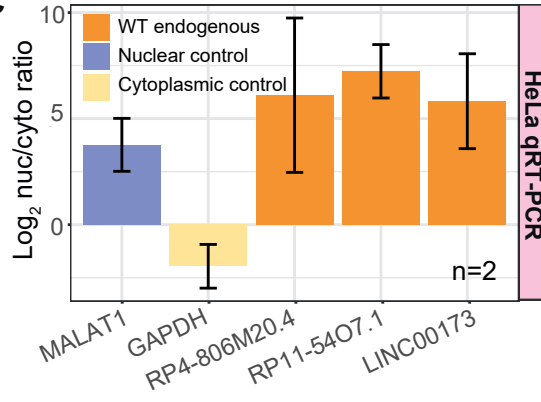
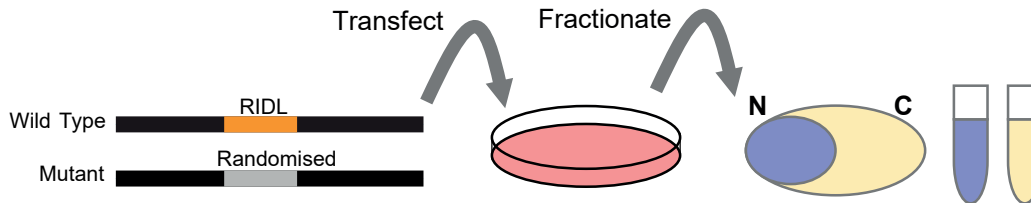
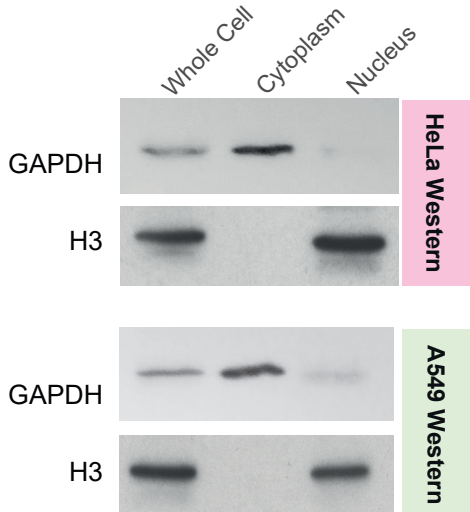
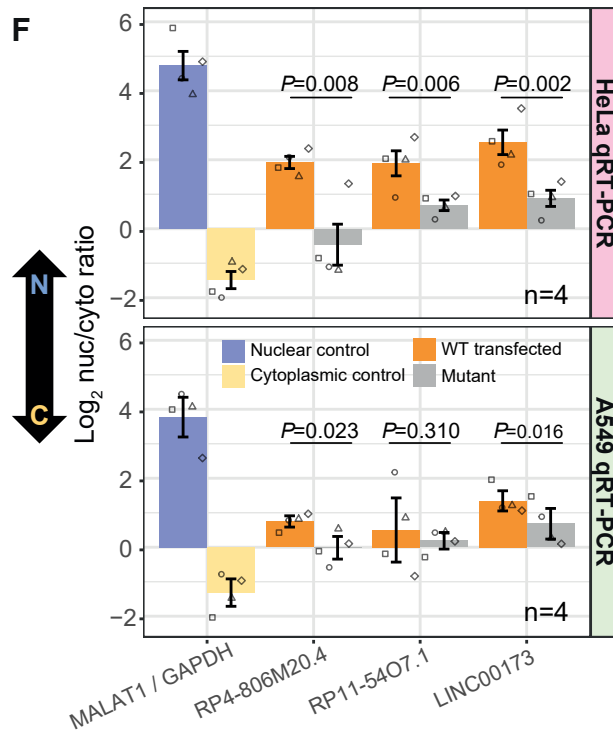
**A****B****C****D****E****F**

Figure Legends:**Figure 1: Repeat insertion domains of lncRNAs (RIDLs).**

(A) In the “Repeat Insertion Domain of lncRNAs” (RIDL) model, exonic-inserted fragments of transposable elements contain pre-formed protein-binding (red), RNA-binding (green) or DNA-binding (blue) activities, that contribute to functionality of the host lncRNA (black). RIDLs are likely to be a small minority of exonic TEs, coexisting with large numbers of non-functional “passengers” (grey). (B) RIDLs (dark orange arrows) will be distinguished from passenger TEs by signals of selection, including: (1) simple enrichment in exons; (2) a preference for residing on a particular strand relative to the host transcript; (3) elevated evolutionary conservation in exons compared to introns. Selection might be identified by comparing exonic TEs to a neutral population, for example those residing in lncRNA introns (light coloured arrows).

**Figure 2: An exonic transposable element annotation with the GENCODE v21 lncRNA catalogue**

(A) Statistics for the exonic TE annotation process using GENCODE v21 lncRNAs. (B) The fraction of nucleotides overlapped by TEs for lncRNA exons and introns, protein-coding introns and exons (“pc”), and the whole genome. (C) Overview of the annotation process. The exons of all transcripts within a lncRNA gene annotation are merged. Merged exons are intersected with the RepeatMasker TE annotation. Intersecting TEs are classified into one of six categories (bottom panel) according to the gene structure with which they intersect, and the relative strand of the TE with respect to the gene: “TSS”, overlapping the transcription start site; “Donor”, splice donor site; “Acceptor”, splice acceptor site; “Inside”, the TE boundaries both lie within the exon; “Encompassing”, the exon boundaries both lie within the TE; “TTS”, the transcription termination site. (D) Summary of classification breakdown for exonic TE annotation. (E) Classification of TE classes in exonic TE annotation. Numbers indicate instances of each type. +/- indicate the relative strand of the TE with respect to lncRNA transcript.

**Figure 3: Evidence for selection on transposable elements in lncRNA exons.**

(A) Figure shows, for every TE type, the enrichment of per nucleotide coverage in exons compared to introns (*y* axis) and overall exonic nucleotide coverage (*x* axis). Enriched TE types (at a 2-fold cutoff) are shown in blue. (B) As for (A), but this time the *y* axis records the ratio of nucleotide coverage in sense vs antisense configuration. “Sense” here is defined as sense of TE annotation relative to the overlapping exon. Similar results for lncRNA introns may be found in Supplemental Figure S1. Significantly-enriched TE types are shown in blue. Statistical significance was estimated by a randomisation procedure, and significance is defined at an uncorrected empirical *p*-value < 0.001 (See Material and Methods). (C) As for (A), but here the *y* axis records the ratio of per-nucleotide overlap by phastCons mammalian-conserved elements for exons vs introns. Similar results for three other measures of evolutionary conservation may be found in Supplemental Figure S1. Significantly-enriched TE types are shown in blue. Statistical significance was estimated by a randomisation procedure, and significance is defined at an uncorrected empirical *p*-value < 0.001 (See Material and Methods). An example of significance estimation is shown in the inset: the distribution shows the exonic/intronic conservation ratio for 1000 simulations. The green arrow shows the true value, in this case for MLT1A0 type. (D) Summary of TE types with evidence of exonic selection. Six distinct evidence types are shown in rows, and TE types in columns. On the right are summary statistics for (i) the number of unique TE types identified by each method, and (ii) the number of instances of exonic

TEs from each type with appropriate selection evidence. The latter are henceforth defined as “RIDLs”.

**Figure 4: Annotated RIDLs and RIDL-lncRNAs.**

(A) Example of a RIDL-lncRNA gene: *CCAT1*. Of note is that although several exonic TE instances are identified (grey), including three separate MIR elements, only one is defined a RIDL (orange) due to overlap of a conserved element. (B) Breakdown of RIDL instances by TE family and evidence sources. (C) Insertion profile of SST1 RIDLs (blue) and intronic insertions (red). *x* axis shows the entire consensus sequence of SST1. *y* axis indicates the frequency with which each nucleotide position is present in the aggregate of all insertions. “*CC*”: Spearman correlation coefficient of the two profiles. “RIDLs” / “Intronic TEs” indicate the numbers of individual insertions considered for RIDLs / intronic insertions, respectively. (D) Number of lncRNAs (*y* axis) carrying the indicated number of RIDL (*x* axis) given the true distribution (black) and randomized distribution (red). The 95% confidence interval was computed empirically, by randomly shuffling RIDLs across the entire lncRNA annotation. (E) Percentage of RIDL-lncRNAs, and a length-matched set of non-RIDL lncRNAs, which are present in disease- and cancer-associated lncRNA databases (see Materials and Methods), in the lncRNAdb database of functional lncRNAs (functional characterisation), or contain at least one trait/disease-associated SNP in an exonic region (GWAS SNP overlap). Numbers denote gene counts. (F) Plot shows regression coefficients for the “RIDL” term in the indicated multiple logistic regression model using the same measures of functionality than in (E). Colours indicate the associated *p*-value. These values assess the correlation between RIDL number and measures of functionality of their host transcript, while accounting for transcript length (trx length) and conservation.

**Figure 5: Correlation between RIDLs and host lncRNA nuclear/cytoplasmic localisation.**

(A) Outline of in silico screen for localisation-regulating RIDLs. For each RIDL-type / cell-type combination, the nuclear/cytoplasmic localisation of RIDL-lncRNAs is compared to all other detected lncRNAs. (B) Results of an in silico screen. Rows: RIDL types; Columns: Cell types. Significant RIDL-cell type combinations are coloured (Benjamini-Hochberg corrected *p*-value < 0.01; Wilcoxon test). Colour scale indicates the nuclear/cytoplasmic ratio mean of RIDL-lncRNAs. Numbers in cells indicate the number of considered RIDL-lncRNAs. Analyses were performed using a single representative transcript isoform from each gene locus, being that with the greatest number of exons. (C) The nuclear/cytoplasmic localization of lncRNAs carrying L2b RIDLs in IMR-90 cells. Blue indicates lncRNAs carrying ≥1 RIDLs, grey indicates all other detected lncRNAs (“not”). Dashed lines represent medians. Significance was calculated using Wilcoxon test (*P*). (D) The nuclear/cytoplasmic ratio of lncRNAs as a function of the number of RIDLs that they carry (L1PA16, L2b, MIRb, MIRc). Correlation coefficient (Rho) and corresponding *p*-value (*P*) were calculated using Spearman correlation, two-sided test. In each box, upper value indicates the number of lncRNAs, and lower value the median. (E) Plot shows regression coefficients for the “RIDL” term in the indicated linear model using L2b, MIRb and MIRc RIDLs (see Methods). Colours indicate the associated *p*-value. These values assess the correlation between RIDL number and nuclear/cytoplasmic localisation ( $\text{Log}_2(N/C)$ ) of their host transcript, while accounting for possible confounding factors of transcript length (trx.length) or whole-cell expression levels (expression).

**Figure 6: RIDLs correlate with differential localisation of lncRNA transcripts from the same locus.**

Distribution of differences between RCI mean of transcripts with nuclear RIDL (mean RCI(A)) and RCI mean of transcripts without nuclear RIDL (mean RCI(B)). A positive value indicates that RIDL-carrying transcripts are more nuclear-enriched than non-RIDL transcripts. Data were calculated individually for every gene that has ≥1 RIDL-transcript and ≥1 non-RIDL transcript expressed in a given cell line. Numbers inside the boxplots indicate the number of gene loci analysed for each cell line. Horizontal

bar indicates the median. Here “nuclear RIDL” refers to L1AP16, MIRb, MIRc and L2b. *P*-values obtained from one-sided *t*-test are shown (in red when  $P < 0.05$ ).

**Figure 7: Disruption of RIDLs results in lncRNA relocalisation from nucleus to cytoplasm.**

(A) Structures of candidate RIDL-lncRNAs. Orange indicates RIDL positions. For each RIDL, numbers indicate the position within the TE consensus, and its orientation with respect to the lncRNA is indicated by arrows (“>” for same strand, “<” for opposite strand). (B) Expression of the three lncRNA candidates as inferred from HeLa RNA-seq (40). (C) Nuclear/cytoplasmic localisation of endogenous candidate lncRNA copies in wild-type HeLa cells, as measured by qRT-PCR. (D) Experimental design. (E) The purity of HeLa and A549 subcellular fractions was assessed by Western blotting against specific markers. GAPDH / Histone H3 proteins are used as cytoplasmic / nuclear markers, respectively. (F) Nuclear/cytoplasmic localisation of transfected candidate lncRNAs in HeLa (upper panel) and A549 (lower panel). GAPDH/MALAT1 are used as cytoplasmic/nuclear controls, respectively. N indicates the number of biological replicates (values from all replicates are plotted, each replicate is represented by a different dot shape), and error bars represent standard error of the mean. *P*-values for paired *t*-test (1 tail) are shown.