

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

3

4 Joana Carlevaro-Fita^{1,2,3,5}

5 Taisia Polidori^{1,2,3,5}

6 Monalisa Das^{1,2,5}

7 Carmen Navarro⁴

8 Tatjana I. Zoller^{1,2}

9 Rory Johnson^{1,2*}

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

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

42 **Introduction**

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

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

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

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

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

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

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

95 **Results**

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

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

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

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

111

112 **A map of exonic transposable elements in GENCODE v21 lncRNAs**

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

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

125

126 Contribution of transposable elements to lncRNA gene structures

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

138

139 Evidence for selection on certain exonic transposable element types

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

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

153 Despite their overall abundance throughout the genome, presently-active LINE1 elements are
154 relatively depleted in lncRNA exons (Figure 3A). It is possible that this reflects selection against disruption
155 to normal gene expression, where numerous weak polyadenylation signals lead to premature transcription

156 termination when the LINE1 element lies on the same strand as the overlapping gene (Perepelitsa-Belancio
157 and Deininger 2003). Other explanations may be low transcriptional processivity exhibited by the LINE1
158 ORF2 in the sense strand (Perepelitsa-Belancio and Deininger 2003), or else epigenetic silencing effects
159 (Hollister and Gaut 2009).

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

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

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

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

188 TE types (Figure 3C). This enrichment disappeared when phastCons elements were positionally randomised
189 (Supplemental Figure S2A). It is unlikely that overlap with protein-coding loci biases the results, since
190 equivalent analyses using intergenic lncRNAs yielded similar candidate RIDLs (Supplemental Figure S2B).
191 A similar analysis was performed using protein-coding exons, and although a number of significantly-
192 conserved TEs were identified, they display limited overlap with those from lncRNAs (Supplemental
193 Figure S2C). We also found a small number of TEs depleted for signatures of conservation in lncRNA
194 exons, namely the young *AluSz*, *AluSx* and *AluJb* (phastCons) and L1M4c and *AluSx1* (ECS) (coloured
195 orange in Figure 3C and Supplemental Figure S1). The cause of this depletion is unclear, although one
196 explanation is enrichment of conservation in intronic TEs due to RNA-independent regulatory roles as
197 observed previously (Su et al. 2014).

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

204

205 An annotation of RIDLs

206 We next combined all TE classes with evidence of functionality into a draft annotation of RIDLs.
207 This annotation combined altogether 99 TE types with at least one type of selection evidence. For each TE
208 / evidence pair, only those TE instances satisfying that evidence were included. In other words, if MIRb
209 elements were found to be associated with vertebrate phastCons elements, then *only* those instances of
210 exonic MIRb elements overlapping such an element would be included in the RIDL annotation, and all
211 other exonic MIRs would be excluded. This operation was performed for all three phastCons element types,
212 ECS elements and strand-bias. An example is *CCAT1* lncRNA oncogene: it carries three exonic MIR
213 elements, of which one is defined as a RIDL based on its overlapping a phastCons element (Figure 4A).

214 After removing redundancy, the final RIDL annotation consists of 5374 elements, located within
215 3566 distinct lncRNA genes (Figure 3D). These represent 12% (5374/46474) of all exonic TE fragments.
216 The most predominant TE families are MIR and L2 repeats, representing 2329 and 1143 RIDLs (Figure
217 4B). The majority of both are defined based on evolutionary evidence (Figure 4B, Supplemental Figure
218 S3). In contrast, RIDLs composed by ERV1, low complexity, satellite and simple repeats families are more

219 frequently identified due to exonic enrichment (Figure 4B). The entire RIDL annotation is available in
220 Supplemental File S1.

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

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

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

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

248

249 RIDL-carrying lncRNAs are enriched for functions and disease roles

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

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

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

274

275 Correlation between RIDLs and subcellular localisation of host transcript

276 The location of a lncRNA within the cell is of key importance to its molecular function (Derrien et
277 al. 2012; Cabili et al. 2015)(Mas-Ponte et al. 2017), therefore we next investigated whether RIDLs might
278 regulate lncRNA localisation (Zhang et al. 2014a; Hacısuleyman et al. 2016b)(Chillón and Pyle 2016)
279 (Figure 5A). Using subcellular RNA-seq data based on 10 ENCODE cell lines (Djebali et al. 2012), we
280 calculated the relative nuclear/cytoplasmic localisation in log₂ units, or “Relative Concentration Index”

281 (RCI) (Mas-Ponte et al. 2017). Using this dataset, we tested each of the 99 RIDL types for association with
282 localisation of their host transcript.

283 After correcting for multiple hypothesis testing using the Benjamini-Hochberg method (Benjamini
284 and Hochberg 1995), this approach identified four distinct RIDL types: L1PA16, L2b, MIRb and MIRc
285 (Figure 5B). For example, 44 lncRNAs carrying L2b RIDLs have 6.9-fold higher relative
286 nuclear/cytoplasmic ratio in IMR90 cells, and this tendency is observed in six different cell types (Figure
287 5B,C).

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

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

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

310 Given that L2b and MIR elements predate human-mouse divergence, we attempted to perform
311 similar analyses in mouse cells. However given that just two equivalent datasets are available at present
312 (Bahar Halpern et al. 2015)(Tan et al. 2015), as well as the relatively low number of annotated lncRNAs in

313 mouse, we were unable to draw statistically-robust conclusions regarding the evolutionary conservation of
314 this phenomenon.

315

316 Intra-gene correlation between RIDLs and subcellular localisation

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

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

330

331 RIDLs play a causative role in lncRNA nuclear localisation

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

337 We formulated an assay to compare the localisation of transfected lncRNAs carrying wild-type
338 RIDLs, and mutated versions where the RIDL sequence was randomised without altering sequence
339 composition (“Mutant”) (Figure 7D, full sequences available in Supplemental File S4). Wild-type and
340 Mutant lncRNAs were transfected into cultured cells and their localisation evaluated by fractionation. qRT-
341 PCR primers were designed to distinguish transfected Wild-Type and Mutant transcripts from
342 endogenously-expressed copies. Transgenes were typically expressed in a range of 0.2- to 10-fold compared
343 to their endogenous transcripts (Supplemental Figure S15). Fractionation purity was verified by Western

344 blotting (Figure 7E) and qRT-PCR (Figure 7F), and stringent DNase-treatment ensured that plasmid DNA
345 made negligible contributions to our results (Supplemental Figure S16).

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

354 **Discussion**

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

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

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

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

398 Our approach for identifying localisation-regulating RIDLs has advantages over previous studies
399 (Lubelsky and Ulitsky 2018; Haciosuleyman et al. 2016c) in terms of its genome-wide scale. However an
400 unavoidable consequence of our use of evolutionary conservation as a filter, is that it likely biases our
401 analysis against recently-evolved TEs such as *Alus*. It remains entirely possible that modern TEs also
402 influence lncRNA localisation, but cannot be detected using the signals of selection that we have employed.
403 On the other hand, MIRb and MIRc were only identified in one cell type each. We expect this reflects low
404 sensitivity of the statistical screen, rather than cell-type specificity alone, because (i) in a focussed re-
405 analysis (Supplemental Figure S12) the effect was observed in multiple cells, and (ii) experimental
406 validation confirmed it in two independent cell types (Figure 7F).

407 This is further supported by the recent study of Lubelsky and Ulitsky, who performed an
408 experimental screen for localisation motifs in 37 nuclear-enriched lncRNAs, and identified *AluSx* as a
409 nuclear-localisation element (Lubelsky and Ulitsky 2018). These 37 lncRNAs are enriched for RIDLs (62%
410 of Lubelsky lncRNAs contain at least one RIDL, compared to 22% for other Gencode v21 lncRNAs, $P=4e-$
411 6 , Fisher exact test), as well as for the three localisation RIDLs identified here (L2b, MIRb, MIRc: 32% vs
412 9%, $P=3e-4$) (Supplemental Figure S18A). Although our bioinformatic screen did not identify *AluSx*, a
413 naive unfiltered re-analysis of our data supports Lubelsky’s experimental finding that *AluSx*-carrying
414 lncRNAs tend to be more nuclear across multiple cell types (Supplemental Figure S18B). Together, these
415 considerations open the possibility that other localisation-controlling TE types may await discovery.

416 More generally, the RIDL predictions showed rather low concordance between the various
417 selection evidence used (Figure 3D). This likely reflects a number of factors: young evolutionary age of

418 some of the most common TEs, generally low statistical power due to large background of neutral TEs and
419 multiple hypothesis testing, and false positives due to TEs that promote transcription or splicing of
420 lncRNAs. However it is worthy of note that validated candidates L2b, MIRb and MIRc are all implicated
421 by multiple, independent evidence sources (Figure 3D).

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

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

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

449 **Materials and Methods**

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

451

452 **Exonic TE curation**

453 *RepeatMasker* annotations were downloaded from the UCSC Genome Browser (version hg38) on
454 December 31st 2014, and GENCODE version 21 lncRNA annotations in GTF format were downloaded
455 from www.genecodegenes.org (Harrow et al. 2012). Annotations were not filtered further. The
456 ‘*transposon.profiler*’ script, largely based on BEDTools’ *intersect* and *merge* functionalities (Quinlan and
457 Hall 2010), was used to annotate exonic and intronic TEs of the given gene annotation (Supplemental File
458 S5). Exons of all transcripts belonging to the given gene annotation were merged, henceforth referred to as
459 “exons”. The set of introns was curated by subtracting the merged exonic sequences from the full gene
460 spans, and only retaining those introns that belonged to a single gene. Intronic regions were assigned the
461 strand of the host gene.

462 The *RepeatMasker* annotation file was intersected with exons and classified into one of 6
463 categories: TSS (transcription start site), overlapping the first exonic nucleotide of the first exon; splice
464 acceptor, overlapping exon 5’ end; splice donor, overlapping exon 3’ end; internal, residing within an exon
465 and not overlapping any intronic sequence; encompassing, where an entire exon lies within the TE; TTS
466 (transcription termination site), overlapping the last nucleotide of the last exon. In every case, the TEs are
467 separated by strand relative to the host gene: + where both gene and TE are annotated on the same strand,
468 otherwise -. The result is the “Exonic TE Annotation” (Supplemental File S6).

469

470 **RIDL identification**

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

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

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

494 In evolutionary analysis, four different annotations of evolutionarily-conserved regions were
495 treated similarly, using unfiltered Exonic TE Annotations. Primate, Placental Mammal and Vertebrate
496 phastCons elements based on 46-way alignments were downloaded as BED files from UCSC Genome
497 Browser (Siepel et al. 2005), while the ECS conserved regions from obtained from Supplemental Data of
498 Smith et al (Smith et al. 2013) (see Supplemental File S7 for summary). Because at the time of analysis,
499 phastCons elements were only available for hg19 genome build, we mapped them to hg38 using *LiftOver*
500 utility. For each TE type we calculated the exonic/intronic conservation ratio. To do this we used
501 *IntersectBED* (Quinlan and Hall 2010) to overlap exonic locations with TEs, and calculate the total number
502 of nucleotides overlapping. We performed a similar operation for intronic regions. Then for each TE type,
503 we calculated the ratio of conserved TE nucleotides for exons compared to introns:

$$504 \text{ Relative exonic-intronic conservation (REIC)} = (C_e / (C_e + N_e)) / (C_i / (C_i + N_i))$$

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

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

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

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

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

524 False discovery rates (FDR) were estimated for RIDL predictions. TE type FDR estimates were
525 based on shuffling simulations described above. Empirical p -values for true data were estimated according
526 to $P=(\text{rank in distribution})/(1 + \text{number of simulations})$. For significant cases, where the true value exceeded
527 all $n=1000$ simulations, this value was conservatively defined to be $P=0.001$. These empirical p -values were
528 then converted to FDR using the `r` command “`p.adjust`” with “`fdr`” setting. Accordingly, empirical
529 significance cutoff ($P<0.001$) mentioned in the main text corresponds to the following FDR values: Strand
530 bias: 0.027; Vertebrate phastCons: 0.013; Placental phastCons: 0.014; Primate phastCons: 0.009; ECS:
531 0.034. This analysis is conservative, since empirical p -values of candidates are rounded up in every case to
532 0.001.

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

536

537 RIDL orthology analysis

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

540 (hg19ToBosTau7). Due to availability of chain files, RIDL coordinates were first converted from hg38 to
541 hg19. Orthology was defined by *LiftOver* utility used at default settings.

542

543 Derived allele frequency (DAF) analysis

544 We used allele frequencies from African population provided by the 1000 Genomes Project (Auton et al.
545 2015), as performed previously by (Haerty and Ponting 2013). DAF was determined for human common
546 SNPs from dbSNP (build 150) (Sherry et al. 2001) for every group analysed. Ancestral repeats (AR) were
547 defined as human repeats (excluding simple repeats) intersecting at least 1 nucleotide of mouse repeats
548 defined by LiftOver, and falling within 5kb of but not overlapping RIDL-containing genes.

549

550 Comparing RIDL-carrying lncRNAs versus other lncRNAs

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

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

556 B) Disease-associated genes: lncRNAs from GENCODE v21 that are present in at least in one of the
557 following databases or public sets: lncRNADisease (Chen et al. 2013), lnc2Cancer (Ning et al.
558 2016), Cancer lncRNA Census (CLC) (Carlevaro-Fita et al. bioRxiv doi: 10.1101/152769)

559 C) GWAS SNPs: We collected SNPs from the NHGRI-EBI Catalog of published genome-wide
560 association studies (Welter et al. 2014; Hindorff et al. 2009) (<https://www.ebi.ac.uk/gwas/home>). We
561 intersected its coordinates with lncRNA exons coordinates.

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

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

568 Where functional-association outcome indicates A, B and C traits defined above; RIDLs indicates the
569 number of RIDL instances in the host gene; transcript length indicates the projected exonic length;
570 conservation indicates the percent of exonic lncRNA nucleotides overlapping the union of primate,

571 placental mammal and vertebrate phastCons elements. We did not find evidence for multicollinearity in
572 any case (variance inflation factors (VIF) <1.1).

573

574 Subcellular localisation analysis

575 Processed RNA-seq data from human cell fractions were obtained from ENCODE in the form of
576 RPKM (reads per kilobase per million mapped reads) quantified against the GENCODE v19 annotation
577 (Djebali et al. 2012; Mas-Ponte et al. 2017). Only transcripts common to both the v21 and v19 annotations
578 were considered. For the following analysis only one transcript per gene was considered, defined as the one
579 with largest number of exons. Nuclear/cytoplasmic ratio expression for each transcript was defined as
580 (nuclear polyA+ RPKM)/(cytoplasmic polyA+ RPKM), and only transcripts having non-zero values (at
581 Irreproducible Discovery Rate (IDR) between samples < 1) in both were considered. These ratios were
582 log₂-transformed, to yield the Relative Concentration Index (RCI) (Mas-Ponte et al. 2017). For each RIDL
583 type and cell type in turn, the nuclear/cytoplasmic ratio distribution of RIDL-containing to non-RIDL-
584 containing lncRNAs was compared using Wilcoxon test. Only RIDLs having at least three expressed
585 transcripts in at least one cell type were tested. Resulting *p*-values were globally adjusted to False Discovery
586 Rate using the Benjamini-Hochberg method (Benjamini and Hochberg 1995).

587

588 Multiple linear regression and partial correlation analysis

589 Linear models were created in R using the “lm” package, at the level of lncRNA transcripts with
590 the form:

591 localisation ~ RIDL + transcript length + expression

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

599

600 Cell lines and reagents

601 Human cervical cancer cell line HeLa and human lung cancer cell line A549 were cultured in
602 Dulbecco's Modified Eagle's medium (Sigma-Aldrich, # D5671) supplemented with 10% FBS and 1%
603 penicillin streptomycin at 37°C and 5 % CO₂. Anti-GAPDH antibody (Sigma-Aldrich, # G9545) and anti-
604 histone H3 antibody (Abcam, # ab24834) were used for Western blot analysis.

605

606 Gene synthesis and cloning of lncRNAs

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

613

614 lncRNA transfection and sub-cellular fractionation

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

623 The nuclear and cytoplasmic fractionation was carried out as described previously (Suzuki et al. 2010)
624 with minor modifications. In brief, cells from 10 cm dishes were harvested by scraping and washed with
625 1x ice-cold PBS. For fractionation, cell pellet was re-suspended in 900 µl of ice-cold 0.1% NP40 in PBS
626 and triturated 7 times using a p1000 micropipette. 300 µl of the cell lysate was saved as the whole cell
627 lysate. The remaining 600 µl of the cell lysate was centrifuged for 30 sec on a table top centrifuge and the
628 supernatant was collected as "cytoplasmic fraction". 300 µl from the cytoplasmic supernatant was kept for
629 RNA isolation and the remaining 300 µl was saved for protein analysis by western blot. The pellet
630 containing the intact nuclei was washed with 1 ml of 0.1% NP40 in PBS. The nuclear pellet was re-
631 suspended in 200 µl 1X PBS and subjected to a quick sonication of 3 pulses with 2 sec ON-2 sec OFF to

632 lyse the nuclei and prepare the “nuclear fraction”. 100 µl of nuclear fraction was saved for RNA isolation
633 and the remaining 100 µl was kept for western blot.

634

635 RNA isolation and real time PCR

636 The RNA from each nuclear and cytoplasmic fraction was isolated using Quick-RNA MiniPrep kit
637 (ZYMO Research, # R1055). The RNAs were subjected to on-column DNase I treatment and clean up
638 using the manufacturer’s protocol. For A549 samples, additional units of DNase were employed, due to
639 residual signal in –RT samples. The RNA from each fraction was converted to cDNA using GoScript
640 reverse transcriptase (Promega, # A5001) and random hexamer primers. The expression of each of the
641 individual transcripts was quantified by qRT-PCR (Applied Biosystems® 7500 Real-Time) using indicated
642 primers (Supplemental File S8) and GoTaq qPCR master mix (Promega, # A6001). In order to distinguish
643 expression of transfected wild-type genes from endogenous copies, we designed forward primers against a
644 transcribed region of the expression vector backbone. Human *GAPDH* mRNA and *MALAT1* lncRNA were
645 used as cytoplasmic and nuclear markers, respectively. The absence of contaminating plasmid DNA in
646 cDNA was checked for all samples using qPCR (see Supplemental Figure S16 for a representative
647 example).

648

649 Western Blotting

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

658

659 Software availability

660 “*transposon.profiler*”, is available on Github at https://github.com/gold-lab/shared_scripts and in
661 Supplemental File S5.

662 **Acknowledgements**

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

673 **References**

674

675 Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S,
676 McVean GA, Abecasis GR. 2015. A global reference for human genetic variation. *Nature* **526**: 68–
677 74.

678 Bahar Halpern K, Caspi I, Lemze D, Levy M, Landen S, Elinav E, Ulitsky I, Itzkovitz S. 2015. Nuclear
679 Retention of mRNA in Mammalian Tissues. *Cell Rep* **13**: 2653–2662.
680 <http://www.ncbi.nlm.nih.gov/pubmed/26711333> (Accessed October 19, 2017).

681 Bassett AR, Akhtar A, Barlow DP, Bird AP, Brockdorff N, Duboule D, Ephrussi A, Ferguson-Smith AC,
682 Gingeras TR, Haerty W, et al. 2014. Considerations when investigating lncRNA function in vivo.
683 *Elife* **3**: e03058. <http://www.ncbi.nlm.nih.gov/pubmed/25124674> (Accessed January 16, 2018).

684 Benjamini Y, Hochberg Y. 1995. Controlling the False Discovery Rate: A Practical and Powerful
685 Approach to Multiple Testing. *J R Stat Soc Ser B* **57**: 289–300.
686 <https://www.jstor.org/stable/2346101> (Accessed August 25, 2017).

687 Benoit Bouvrette LP, Cody NAL, Bergalet J, Lefebvre FA, Diot C, Wang X, Blanchette M, Lécuyer E.
688 2018. CeFra-seq reveals broad asymmetric mRNA and noncoding RNA distribution profiles in
689 *Drosophila* and human cells. *RNA* **24**: 98–113. <http://www.ncbi.nlm.nih.gov/pubmed/29079635>
690 (Accessed January 8, 2018).

691 Blackwell BJ, Lopez MF, Wang J, Krastins B, Sarracino D, Tollervey JR, Dobke M, Jordan IK, Lunyak
692 V V. 2012. Protein interactions with piALU RNA indicates putative participation of retroRNA in
693 the cell cycle, DNA repair and chromatin assembly. *Mob Genet Elements* **2**: 26–35.
694 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3383447/pdf/mge-2-26.pdf> (Accessed August 25,
695 2017).

696 Bourque G. 2009. Transposable elements in gene regulation and in the evolution of vertebrate genomes.
697 *Curr Opin Genet Dev* **19**: 607–612.
698 <http://www.sciencedirect.com/science/article/pii/S0959437X09001725> (Accessed August 25, 2017).

699 Bourque G, Leong B, Vega VB, Chen X, Lee YL, Srinivasan KG, Chew J-L, Ruan Y, Wei C-L, Ng HH,
700 et al. 2008. Evolution of the mammalian transcription factor binding repertoire via transposable
701 elements. *Genome Res* **18**: 1752–62. <http://www.ncbi.nlm.nih.gov/pubmed/18682548> (Accessed
702 August 25, 2017).

703 Cabili MN, Dunagin MC, McClanahan PD, Biaesch A, Padovan-Merhar O, Regev A, Rinn JL, Raj A.
704 2015. Localization and abundance analysis of human lncRNAs at single-cell and single-molecule
705 resolution. *Genome Biol* **16**: 20. <http://genomebiology.com/2015/16/1/20> (Accessed January 30,
706 2015).

- 707 Carlevaro-Fita J, Rahim A, Guigó R, Vardy LA, Johnson R. 2016. Cytoplasmic long noncoding RNAs
708 are frequently bound to and degraded at ribosomes in human cells. *RNA* **22**: 867–82.
709 <http://www.ncbi.nlm.nih.gov/pubmed/27090285> (Accessed April 26, 2016).
- 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. <http://www.ncbi.nlm.nih.gov/pubmed/23064229>
713 (Accessed February 10, 2015).
- 714 Chen G, Wang Z, Wang D, Qiu C, Liu M, Chen X, Zhang Q, Yan G, Cui Q. 2013. LncRNADisease: a
715 database for long-non-coding RNA-associated diseases. *Nucleic Acids Res* **41**: D983–D986.
716 <http://www.ncbi.nlm.nih.gov/pubmed/23175614> (Accessed December 24, 2016).
- 717 Chen L-L. 2016. Linking Long Noncoding RNA Localization and Function. *Trends Biochem Sci*.
718 <http://www.ncbi.nlm.nih.gov/pubmed/27499234> (Accessed August 25, 2016).
- 719 Chillón I, Pyle AM. 2016. Inverted repeat *Alu* elements in the human lincRNA-p21 adopt a conserved
720 secondary structure that regulates RNA function. *Nucleic Acids Res* gkw599.
721 <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkw599> (Accessed September 5, 2016).
- 722 Chuong EB, Elde NC, Feschotte C. 2017. Regulatory activities of transposable elements: from conflicts
723 to benefits. <https://www.nature.com/nrg/journal/v18/n2/pdf/nrg.2016.139.pdf> (Accessed August 25,
724 2017).
- 725 Clark MB, Johnston RL, Inostroza-Ponta M, Fox AH, Fortini E, Moscato P, Dinger ME, Mattick JS.
726 2012. Genome-wide analysis of long noncoding RNA stability. *Genome Res* **22**: 885–98.
727 <http://genome.cshlp.org/content/22/5/885.long> (Accessed June 3, 2014).
- 728 Cordaux R, Batzer MA. 2009. The impact of retrotransposons on human genome evolution. *Nat Rev*
729 *Genet* **10**: 691–703. [https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2884099/pdf/nihms-](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2884099/pdf/nihms-201920.pdf)
730 [201920.pdf](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2884099/pdf/nihms-201920.pdf) (Accessed August 29, 2017).
- 731 Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, Guernec G, Martin D, Merkel A,
732 Knowles DG, et al. 2012. The GENCODE v7 catalog of human long noncoding RNAs: analysis of
733 their gene structure, evolution, and expression. *Genome Res* **22**: 1775–89.
734 <http://genome.cshlp.org/content/22/9/1775.long> (Accessed May 23, 2014).
- 735 Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,
736 Schlesinger F, et al. 2012. Landscape of transcription in human cells. *Nature* **489**: 101–108.
737 <http://www.nature.com/doi/10.1038/nature11233> (Accessed April 20, 2017).
- 738 Elisaphenko EA, Kolesnikov NN, Shevchenko AI, Rogozin IB, Nesterova TB, Brockdorff N, Zakian SM.
739 2008. A Dual Origin of the Xist Gene from a Protein-Coding Gene and a Set of Transposable
740 Elements.

741 <http://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0002521&type=printable>
742 (Accessed August 25, 2017).

743 Faulkner GJ, Kimura Y, Daub CO, Wani S, Plessy C, Irvine KM, Schroder K, Cloonan N, Steptoe AL,
744 Lassmann T, et al. The regulated retrotransposon transcriptome of mammalian cells.
745 <https://www.nature.com/ng/journal/v41/n5/pdf/ng.368.pdf> (Accessed August 25, 2017).

746 Feschotte C. 2008. Transposable elements and the evolution of regulatory networks. *Nat Rev Genet* **9**:
747 397–405. <http://www.ncbi.nlm.nih.gov/pubmed/18368054> (Accessed September 1, 2017).

748 Gong C, Maquat LE. 2011. lncRNAs transactivate STAU1-mediated mRNA decay by duplexing with 3'
749 UTRs via Alu elements. *Nature* **470**: 284–8.
750 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3073508&tool=pmcentrez&rendertype=](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3073508&tool=pmcentrez&rendertype=abstract)
751 [abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3073508&tool=pmcentrez&rendertype=abstract) (Accessed March 6, 2015).

752 Guttman M, Rinn JL. 2012. Modular regulatory principles of large non-coding RNAs. *Nature* **482**: 339–
753 46. <http://www.ncbi.nlm.nih.gov/pubmed/22337053> (Accessed May 11, 2017).

754 Hacisuleyman E, Goff LA, Trapnell C, Williams A, Henao-Mejia J, Sun L, McClanahan P, Hendrickson
755 DG, Sauvageau M, Kelley DR, et al. 2014. Topological organization of multichromosomal regions
756 by the long intergenic noncoding RNA Firre. *Nat Struct Mol Biol* **21**: 198–206.
757 <http://www.nature.com/articles/nsmb.2764> (Accessed January 16, 2018).

758 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL. 2016a. Function and evolution of local repeats in the
759 Firre locus. *Nat Commun* **7**: 11021. <http://www.nature.com/doi/10.1038/ncomms11021>
760 (Accessed January 16, 2018).

761 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL, Koning AP de, Gu W, Castoe TA, Batzer MA, Pollock
762 DD, Wicker T, et al. 2016b. Function and evolution of local repeats in the Firre locus. *Nat Commun*
763 **7**: 11021. <http://www.nature.com/doi/10.1038/ncomms11021> (Accessed November 20, 2016).

764 Hacisuleyman E, Shukla CJ, Weiner CL, Rinn JL, Koning AP de, Gu W, Castoe TA, Batzer MA, Pollock
765 DD, Wicker T, et al. 2016c. Function and evolution of local repeats in the Firre locus. *Nat Commun*
766 **7**: 11021. <http://www.nature.com/doi/10.1038/ncomms11021> (Accessed November 20, 2016).

767 Haerty W, Ponting CP. 2013. Mutations within lncRNAs are effectively selected against in fruitfly but
768 not in human. *Genome Biol* **14**: R49. [http://genomebiology.biomedcentral.com/articles/10.1186/gb-](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-5-r49)
769 [2013-14-5-r49](http://genomebiology.biomedcentral.com/articles/10.1186/gb-2013-14-5-r49) (Accessed October 23, 2018).

770 Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, Aken BL, Barrell D,
771 Zadissa A, Searle S, et al. 2012. GENCODE: The reference human genome annotation for The
772 ENCODE Project. *Genome Res* **22**: 1760–1774.
773 <http://genome.cshlp.org/cgi/doi/10.1101/gr.135350.111> (Accessed November 20, 2016).

774 Hezroni H, Ben-Tov Perry R, Meir Z, Housman G, Lubelsky Y, Ulitsky I. 2017. A subset of conserved

775 mammalian long non-coding RNAs are fossils of ancestral protein-coding genes. *Genome Biol* **18**:
776 162. <http://www.ncbi.nlm.nih.gov/pubmed/28854954> (Accessed September 2, 2017).

777 Hezroni H, Koppstein D, Schwartz MG, Avrutin A, Bartel DP, Ulitsky I. 2015. Principles of Long
778 Noncoding RNA Evolution Derived from Direct Comparison of Transcriptomes in 17 Species. *Cell*
779 *Rep* **11**: 1110–1122. <http://linkinghub.elsevier.com/retrieve/pii/S2211124715004106> (Accessed
780 September 1, 2017).

781 Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, Manolio TA. 2009. Potential
782 etiologic and functional implications of genome-wide association loci for human diseases and traits.
783 *Proc Natl Acad Sci* **106**: 9362–9367. <http://www.ncbi.nlm.nih.gov/pubmed/19474294> (Accessed
784 August 30, 2017).

785 Holdt LM, Hoffmann S, Sass K, Langenberger D, Scholz M, Krohn K, Finstermeier K, Stahringer A,
786 Wilfert W, Beutner F, et al. 2013. Alu Elements in ANRIL Non-Coding RNA at Chromosome 9p21
787 Modulate Atherogenic Cell Functions through Trans-Regulation of Gene Networks. *PLoS Genet* **9**.
788 <http://journals.plos.org/plosgenetics/article/file?id=10.1371/journal.pgen.1003588&type=printable>
789 (Accessed August 25, 2017).

790 Hollister JD, Gaut BS. 2009. Epigenetic silencing of transposable elements: a trade-off between reduced
791 transposition and deleterious effects on neighboring gene expression. *Genome Res* **19**: 1419–28.
792 <http://www.ncbi.nlm.nih.gov/pubmed/19478138> (Accessed January 19, 2018).

793 Hu B, Yang Y-CT, Huang Y, Zhu Y, Lu ZJ. 2017. POSTAR: a platform for exploring post-transcriptional
794 regulation coordinated by RNA-binding proteins. *Nucleic Acids Res* **45**: D104–D114.
795 <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkw888> (Accessed September 2,
796 2017).

797 Huda A, Bowen NJ, Conley AB, Jordan IK. 2011. Epigenetic regulation of transposable element derived
798 human gene promoters. *Gene* **475**: 39–48.
799 <http://www.sciencedirect.com/science/article/pii/S0378111910004762> (Accessed August 27, 2017).

800 Jjingo D, Conley AB, Wang J, Mariño-Ramírez L, Lunyak V V, Jordan IK. 2014. Mammalian-wide
801 interspersed repeat (MIR)-derived enhancers and the regulation of human gene expression. *Mob*
802 *DNA* **5**: 14. <http://mobilejournal.biomedcentral.com/articles/10.1186/1759-8753-5-14> (Accessed
803 September 2, 2017).

804 Johnson R, Guigó R. 2014. The RIDL hypothesis: transposable elements as functional domains of long
805 noncoding RNAs. *RNA* **20**: 959–76.
806 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4114693&tool=pmcentrez&rendertype=
807 abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4114693&tool=pmcentrez&rendertype=abstract) (Accessed April 29, 2015).

808 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

809 the REST regulon reveals extensive transposable element-mediated binding site duplication. *Nucleic*
810 *Acids Res* **34**: 3862–3877. <http://nar.oxfordjournals.org/lookup/doi/10.1093/nar/gkl525> (Accessed
811 August 25, 2017).

812 Jurka J, Zietkiewicz E, Labuda D. 1995. Ubiquitous mammalian-wide interspersed repeats (MIRs) are
813 molecular fossils from the mesozoic era. *Nucleic Acids Res* **23**: 170–5.
814 <http://www.ncbi.nlm.nih.gov/pubmed/7870583> (Accessed August 29, 2017).

815 Kapusta A, Feschotte C. 2014. Volatile evolution of long noncoding RNA repertoires: mechanisms and
816 biological implications. *Trends Genet* **30**: 439–452.
817 <http://linkinghub.elsevier.com/retrieve/pii/S0168952514001346> (Accessed August 25, 2017).

818 Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, Yandell M, Feschotte C. 2013.
819 Transposable elements are major contributors to the origin, diversification, and regulation of
820 vertebrate long noncoding RNAs. ed. H.E. Hoekstra. *PLoS Genet* **9**: e1003470.
821 <http://dx.plos.org/10.1371/journal.pgen.1003470> (Accessed August 25, 2017).

822 Kelley D, Rinn J. 2012. Transposable elements reveal a stem cell-specific class of long noncoding RNAs.
823 *Genome Biol* **13**: R107. <http://genomebiology.com/2012/13/11/R107> (Accessed May 25, 2014).

824 Kelley DR, Hendrickson DG, Tenen D, Rinn JL. 2014. Transposable elements modulate human RNA
825 abundance and splicing via specific RNA-protein interactions. *Genome Biol* **15**: 537.
826 <http://www.ncbi.nlm.nih.gov/pubmed/25572935> (Accessed September 2, 2017).

827 Konkel MK, Walker JA, Batzer MA. 2010. LINEs and SINEs of primate evolution. *Evol Anthropol* **19**:
828 236–249. <http://www.ncbi.nlm.nih.gov/pubmed/25147443> (Accessed August 29, 2017).

829 Lagarde J, Uszczyńska-Ratajczak B, Carbonell S, Pérez-Lluch S, Abad A, Davis C, Gingeras TR,
830 Frankish A, Harrow J, Guigo R, et al. 2017. High-throughput annotation of full-length long
831 noncoding RNAs with capture long-read sequencing. *Nat Genet* **49**: 1731–1740.
832 <http://www.ncbi.nlm.nih.gov/pubmed/29106417> (Accessed January 3, 2018).

833 Lev-Maor G, Sorek R, Shomron N, Ast G. 2003. The Birth of an Alternatively Spliced Exon: 3' Splice-
834 Site Selection in Alu Exons. *Science (80-)* **300**.
835 <http://science.sciencemag.org/content/300/5623/1288/tab-pdf> (Accessed August 25, 2017).

836 Li J-H, Liu S, Zhou H, Qu L-H, Yang J-H. 2014. starBase v2.0: decoding miRNA-ceRNA, miRNA-
837 ncRNA and protein–RNA interaction networks from large-scale CLIP-Seq data. *Nucleic Acids Res*
838 **42**: D92–D97. <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkt1248> (Accessed
839 September 2, 2017).

840 Lubelsky Y, Ulitsky I. 2018. Sequences enriched in Alu repeats drive nuclear localization of long RNAs
841 in human cells. *Nature* **555**: 107–111. <http://www.ncbi.nlm.nih.gov/pubmed/29466324> (Accessed
842 March 3, 2018).

- 843 Marín-Béjar O, Huarte M. 2015. RNA Pulldown Protocol for In Vitro Detection and Identification of
844 RNA-Associated Proteins. In *Methods in molecular biology (Clifton, N.J.)*, Vol. 1206 of, pp. 87–95
845 <http://www.ncbi.nlm.nih.gov/pubmed/25240889> (Accessed March 11, 2018).
- 846 Marín-Béjar O, Mas AM, González J, Martínez D, Athie A, Morales X, Galduroz M, Raimondi I, Grossi
847 E, Guo S, et al. 2017. The human lncRNA LINC-PINT inhibits tumor cell invasion through a highly
848 conserved sequence element. *Genome Biol* **18**: 202. <http://www.ncbi.nlm.nih.gov/pubmed/29078818>
849 (Accessed January 8, 2018).
- 850 Martin KC, Ephrussi A. 2009. mRNA localization: gene expression in the spatial dimension. *Cell* **136**:
851 719–30. <http://www.ncbi.nlm.nih.gov/pubmed/19239891> (Accessed August 29, 2017).
- 852 Mas-Ponte D, Carlevaro-Fita J, Palumbo E, Hermoso Pulido T, Guigo R, Johnson R. 2017. LncATLAS
853 database for subcellular localisation of long noncoding RNAs. *RNA* rna.060814.117.
854 <http://rnajournal.cshlp.org/lookup/doi/10.1261/rna.060814.117> (Accessed April 10, 2017).
- 855 Mercer TR, Mattick JS. 2013. Structure and function of long noncoding RNAs in epigenetic regulation.
856 *Nat Publ Gr* **20**. <https://www.nature.com/nsmb/journal/v20/n3/pdf/nsmb.2480.pdf> (Accessed August
857 25, 2017).
- 858 Mukherjee N, Calviello L, Hirsekorn A, de Pretis S, Pelizzola M, Ohler U. 2017. Integrative classification
859 of human coding and noncoding genes through RNA metabolism profiles. *Nat Struct Mol Biol* **24**:
860 86–96. <http://www.nature.com/doi/10.1038/nsmb.3325> (Accessed September 2, 2017).
- 861 Necsulea A, Soumillon M, Warnefors M, Liechti A, Daish T, Zeller U, Baker JC, Grützner F, Kaessmann
862 H. 2014. The evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* **505**:
863 635–40. <http://www.ncbi.nlm.nih.gov/pubmed/24463510> (Accessed May 11, 2017).
- 864 Ning S, Zhang J, Wang P, Zhi H, Wang J, Liu Y, Gao Y, Guo M, Yue M, Wang L, et al. 2016.
865 Lnc2Cancer: a manually curated database of experimentally supported lncRNAs associated with
866 various human cancers. *Nucleic Acids Res* **44**: D980–D985. [https://academic.oup.com/nar/article-
867 lookup/doi/10.1093/nar/gkv1094](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv1094) (Accessed August 30, 2017).
- 868 Nissan A, Stojadinovic A, Mitrani-Rosenbaum S, Halle D, Grinbaum R, Roistacher M, Bochem A,
869 Dayanc BE, Ritter G, Gomceli I, et al. 2012. Colon cancer associated transcript-1: A novel RNA
870 expressed in malignant and pre-malignant human tissues. *Int J Cancer* **130**: 1598–1606.
871 <http://www.ncbi.nlm.nih.gov/pubmed/21547902> (Accessed August 29, 2017).
- 872 Pegueroles C, Gabaldón T. 2016. Secondary structure impacts patterns of selection in human lncRNAs.
873 *BMC Biol* **14**: 60. <http://www.ncbi.nlm.nih.gov/pubmed/27457204> (Accessed January 8, 2018).
- 874 Perepelitsa-Belancio V, Deininger P. 2003. RNA truncation by premature polyadenylation attenuates
875 human mobile element activity. *Nat Genet* **35**: 363–366.
876 <http://www.nature.com/doi/10.1038/ng1269> (Accessed August 25, 2017).

- 877 Quek XC, Thomson DW, Maag JL V, Bartonicek N, Signal B, Clark MB, Gloss BS, Dinger ME. 2015.
878 lncRNADB v2.0: expanding the reference database for functional long noncoding RNAs. *Nucleic*
879 *Acids Res* **43**: D168-73. <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gku988>
880 (Accessed December 24, 2016).
- 881 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
882 *Bioinformatics* **26**: 841–2. <http://www.ncbi.nlm.nih.gov/pubmed/20110278> (Accessed August 25,
883 2017).
- 884 Roberts JT, Cardin SE, Borchert GM. 2014. Burgeoning evidence indicates that microRNAs were
885 initially formed from transposable element sequences. *Mob Genet Elements* **4**: e29255.
886 <http://www.tandfonline.com/doi/abs/10.4161/mge.29255> (Accessed August 25, 2017).
- 887 Schmitt AM, Chang HY, Abdelmohsen K, Panda A, Kang MJ, Xu J, Selimyan R, Yoon JH, Martindale
888 JL, De S, et al. 2016. Long Noncoding RNAs in Cancer Pathways. *Cancer Cell* **29**: 452–463.
889 <http://linkinghub.elsevier.com/retrieve/pii/S1535610816300927> (Accessed June 28, 2016).
- 890 Seemann SE, Mirza AH, Hansen C, Bang-Berthelsen CH, Garde C, Christensen-Dalsgaard M,
891 Torarinsson E, Yao Z, Workman CT, Pociot F, et al. 2017. The identification and functional
892 annotation of RNA structures conserved in vertebrates. *Genome Res* **27**: 1371–1383.
893 <http://www.ncbi.nlm.nih.gov/pubmed/28487280> (Accessed September 2, 2017).
- 894 Sela N, Mersch B, Gal-Mark N, Lev-Maor G, Hotz-Wagenblatt A, Ast G. 2007. Comparative analysis of
895 transposed element insertion within human and mouse genomes reveals Alu’s unique role in shaping
896 the human transcriptome. **8**. [https://genomebiology.biomedcentral.com/track/pdf/10.1186/gb-2007-](https://genomebiology.biomedcentral.com/track/pdf/10.1186/gb-2007-8-6-r127?site=genomebiology.biomedcentral.com)
897 [8-6-r127?site=genomebiology.biomedcentral.com](https://genomebiology.biomedcentral.com/track/pdf/10.1186/gb-2007-8-6-r127?site=genomebiology.biomedcentral.com) (Accessed August 25, 2017).
- 898 Sherry ST, Ward M-H, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP: the
899 NCBI database of genetic variation. *Nucleic Acids Res* **29**: 308–311.
900 <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/29.1.308> (Accessed October 23, 2018).
- 901 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier
902 LW, Richards S, et al. 2005. Evolutionarily conserved elements in vertebrate, insect, worm, and
903 yeast genomes. *Genome Res* **15**: 1034–50. <http://www.ncbi.nlm.nih.gov/pubmed/16024819>
904 (Accessed August 25, 2017).
- 905 Smith MA, Gesell T, Stadler PF, Mattick JS. 2013. Widespread purifying selection on RNA structure in
906 mammals. *Nucleic Acids Res* **41**: 8220–8236.
907 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3783177&tool=pmcentrez&rendertype=](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3783177&tool=pmcentrez&rendertype=abstract)
908 [abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3783177&tool=pmcentrez&rendertype=abstract) (Accessed April 20, 2015).
- 909 Su M, Han D, Boyd-Kirkup J, Yu X, Han J-DJ. 2014. Evolution of Alu elements toward enhancers. *Cell*
910 *Rep* **7**: 376–85. <http://linkinghub.elsevier.com/retrieve/pii/S2211124714001892> (Accessed October

911 18, 2017).

912 Suzuki K, Bose P, Leong-Quong RYY, Fujita DJ, Riabowol K. 2010. REAP: A two minute cell
913 fractionation method. *BMC Res Notes* **3**: 294. <http://www.ncbi.nlm.nih.gov/pubmed/21067583>
914 (Accessed August 28, 2017).

915 Tan JY, Sirey T, Honti F, Graham B, Piovesan A, Merkenschlager M, Webber C, Ponting CP, Marques
916 AC. 2015. Extensive microRNA-mediated crosstalk between lncRNAs and mRNAs in mouse
917 embryonic stem cells. *Genome Res* **25**: 655–666. <http://www.ncbi.nlm.nih.gov/pubmed/25792609>
918 (Accessed October 19, 2017).

919 Tan JY, Smith AAT, Ferreira da Silva M, Matthey-Doret C, Rueedi R, Sönmez R, Ding D, Kutalik Z,
920 Bergmann S, Marques AC. 2017. cis-Acting Complex-Trait-Associated lincRNA Expression
921 Correlates with Modulation of Chromosomal Architecture. *Cell Rep* **18**: 2280–2288.
922 <http://www.ncbi.nlm.nih.gov/pubmed/28249171> (Accessed January 4, 2018).

923 Ulitsky I, Bartel DP. 2013. lincRNAs: Genomics, Evolution, and Mechanisms. *Cell* **154**: 26–46.
924 <http://linkinghub.elsevier.com/retrieve/pii/S0092867413007599> (Accessed November 20, 2016).

925 Van Nostrand EL, Pratt GA, Shishkin AA, Gelboin-Burkhart C, Fang MY, Sundararaman B, Blue SM,
926 Nguyen TB, Surka C, Elkins K, et al. 2016. Robust transcriptome-wide discovery of RNA-binding
927 protein binding sites with enhanced CLIP (eCLIP). *Nat Methods* **13**: 508–514.
928 <http://www.ncbi.nlm.nih.gov/pubmed/27018577> (Accessed September 2, 2017).

929 Washietl S, Kellis M, Garber M. 2014. Evolutionary dynamics and tissue specificity of human long
930 noncoding RNAs in six mammals. *Genome Res* **24**: 616–28.
931 <http://www.ncbi.nlm.nih.gov/pubmed/24429298> (Accessed May 25, 2014).

932 Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, Klemm A, Flicek P, Manolio T,
933 Hindorff L, et al. 2014. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations.
934 *Nucleic Acids Res* **42**: D1001–D1006. [https://academic.oup.com/nar/article-](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkt1229)
935 [lookup/doi/10.1093/nar/gkt1229](https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkt1229) (Accessed August 30, 2017).

936 Zhang B, Gunawardane L, Niazi F, Jahanbani F, Chen X, Valadkhan S. 2014a. A Novel RNA Motif
937 Mediates the Strict Nuclear Localization of a Long Noncoding RNA. *Mol Cell Biol* **34**: 2318–2329.
938 <http://mcb.asm.org/cgi/doi/10.1128/MCB.01673-13> (Accessed November 20, 2016).

939 Zhang K, Shi Z-M, Chang Y-N, Hu Z-M, Qi H-X, Hong W. 2014b. The ways of action of long non-
940 coding RNAs in cytoplasm and nucleus. *Gene* **547**: 1–9.
941 <http://www.ncbi.nlm.nih.gov/pubmed/24967943> (Accessed March 17, 2015).

942

A

bioRxiv preprint doi: <https://doi.org/10.1101/189753>; this version posted November 27, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



B

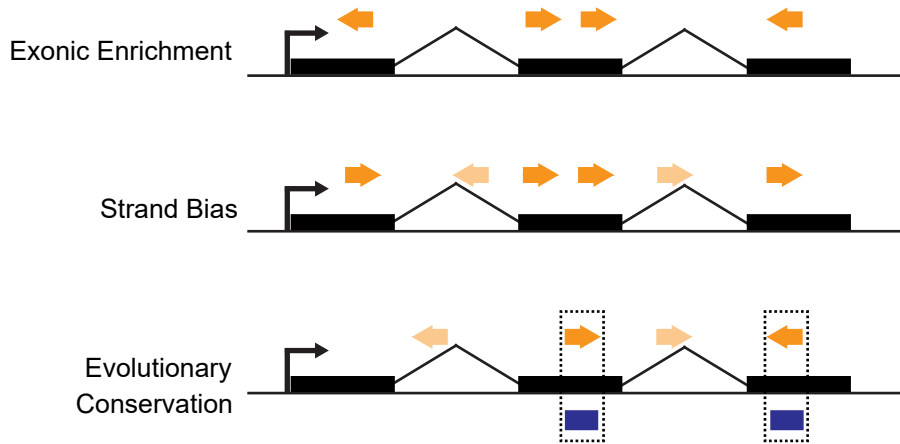
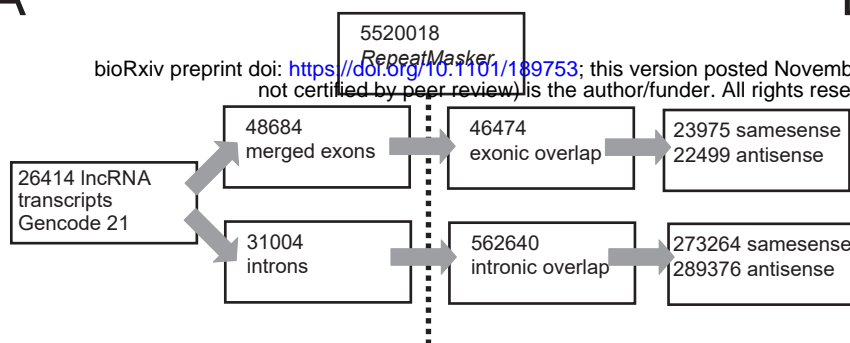


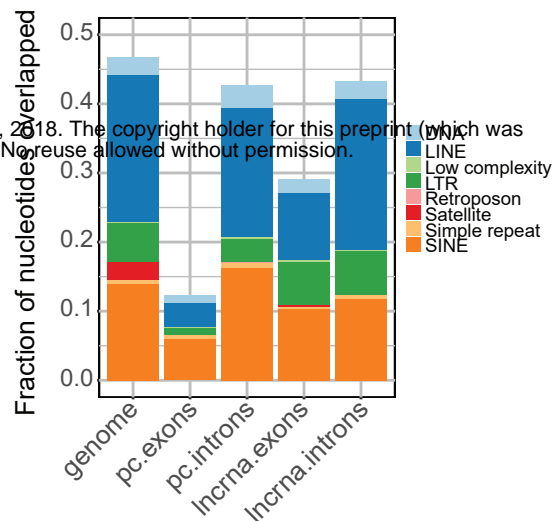
Figure 1: Repeat insertion domains of lncRNAs (RIDLs).

(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 introns. Selection might be identified by comparing exonic TEs to a neutral population, for example those residing in lncRNA introns (light coloured arrows).

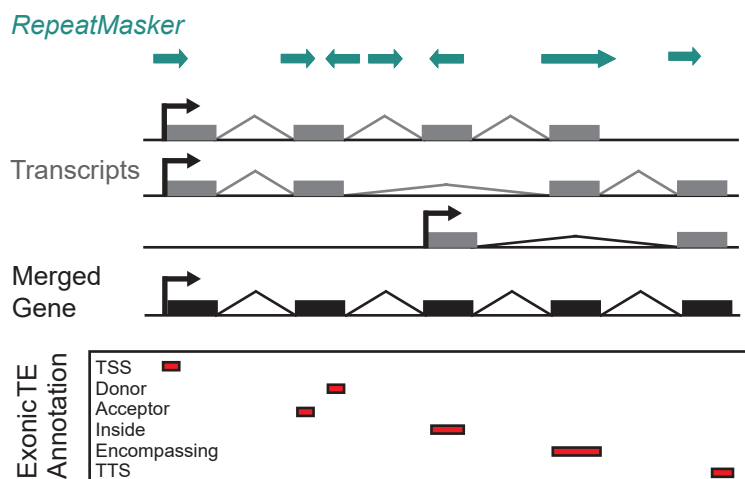
A



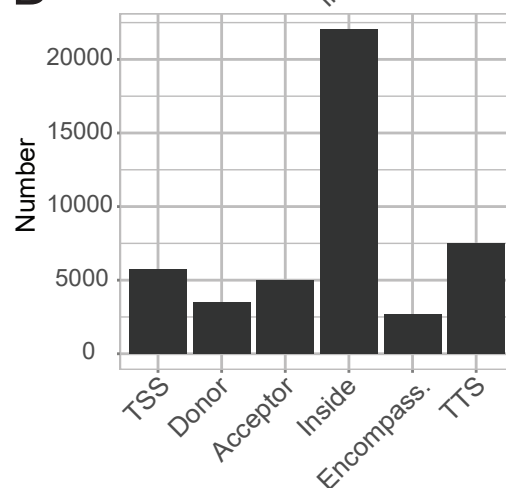
B



C



D



E

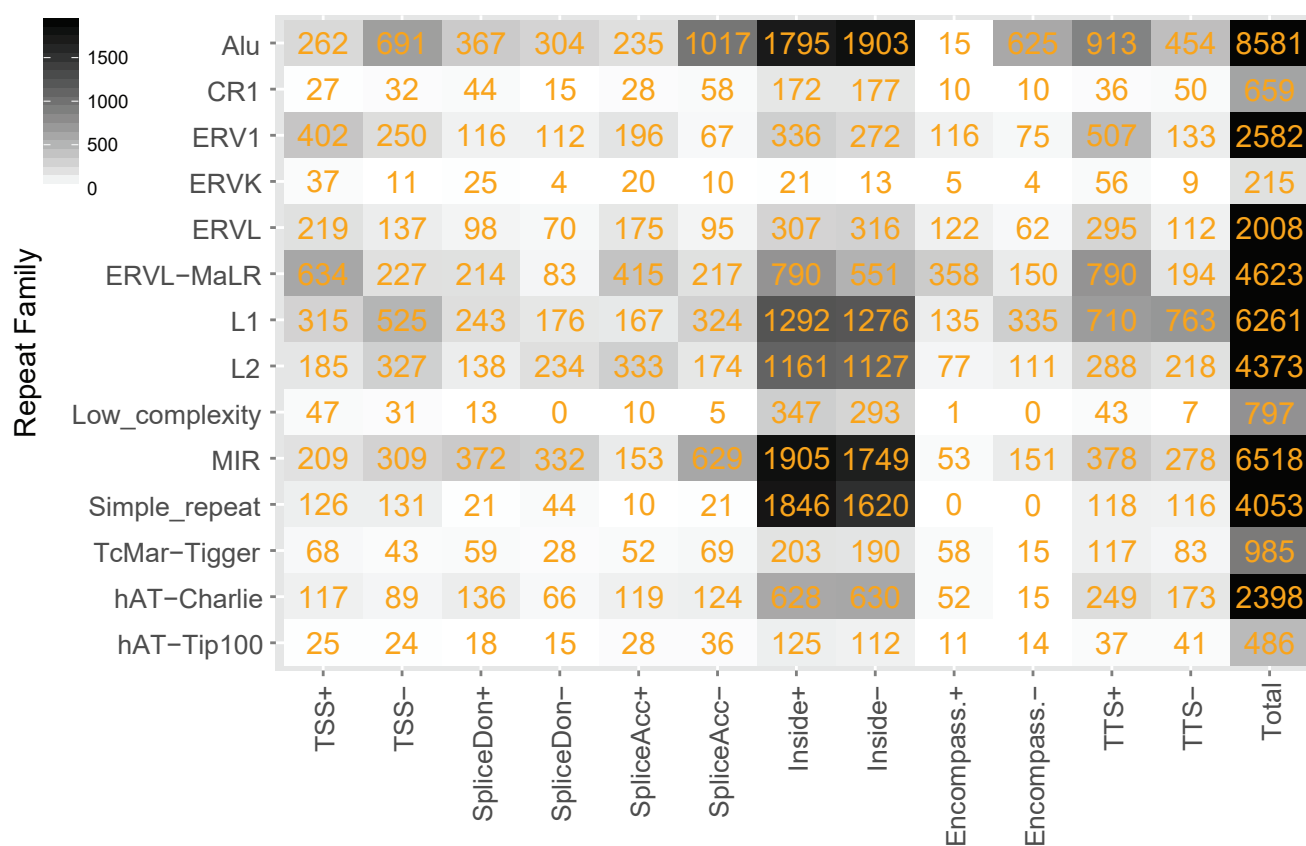
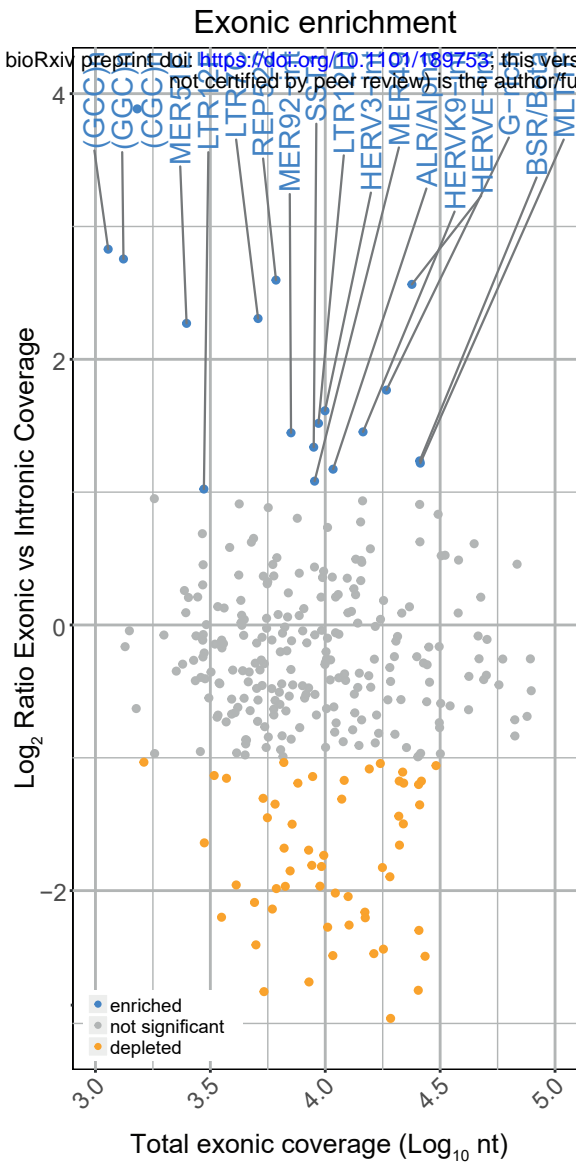


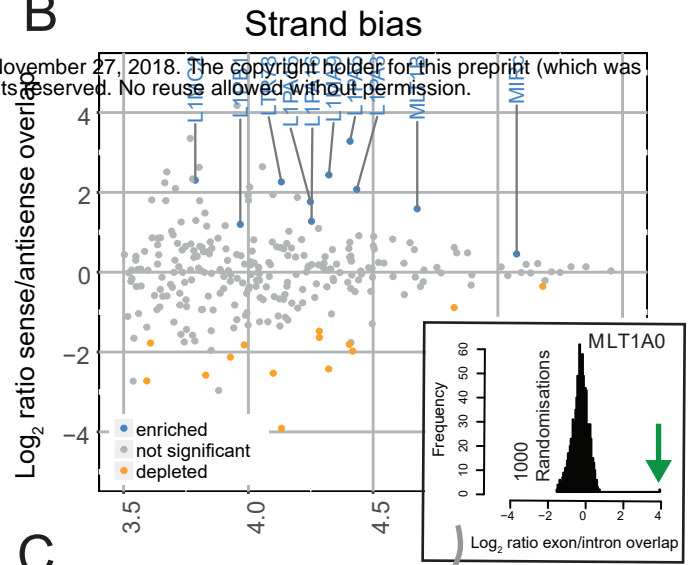
Figure 2: An exonic transposable element annotation with the GENCODE v21 lncRNA catalogue.

(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.

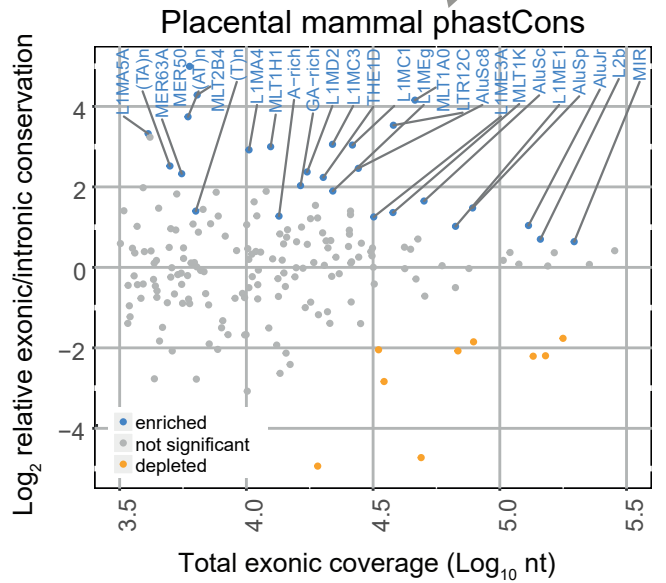
A



B



C



D

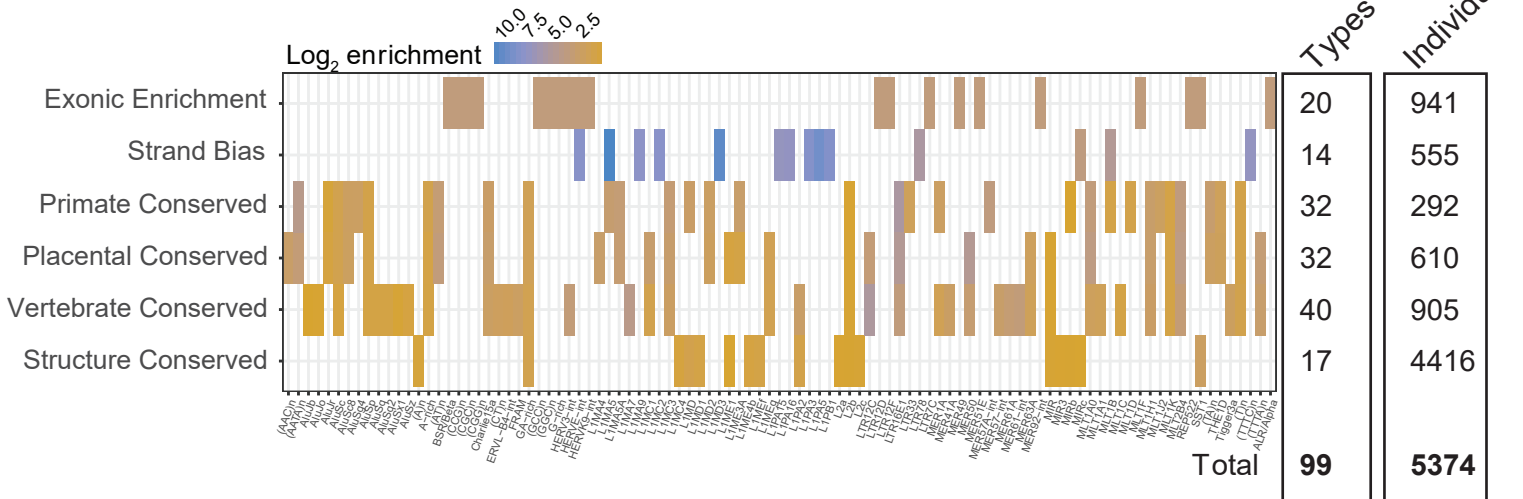


Figure 3: Evidence for selection on transposable elements in lncRNA exons.

(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 Supplementary 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 Supplementary 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”.

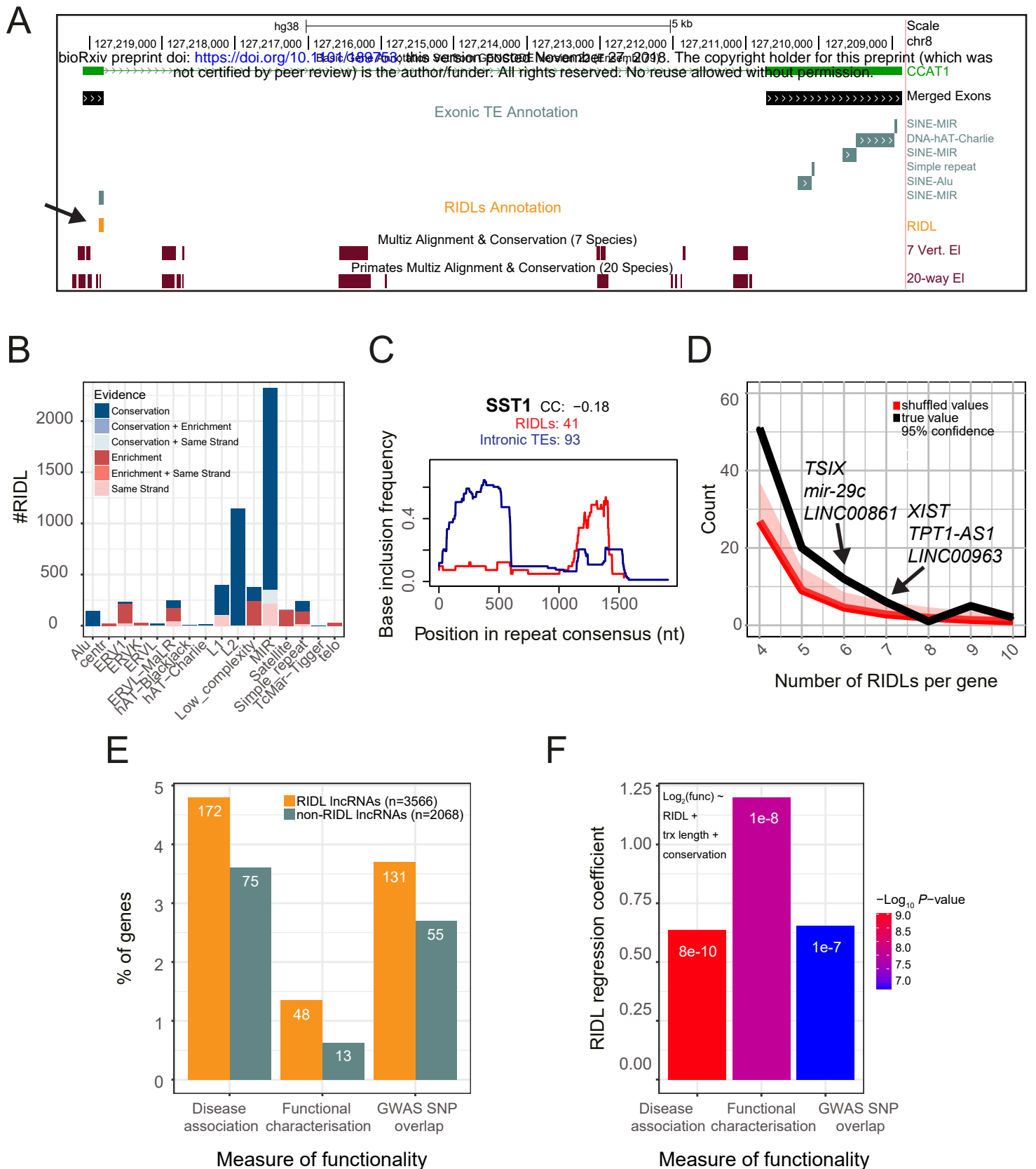


Figure 4: Annotated RIDLs and RIDL-lncRNAs.

(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 (36), or contain at least one trait/disease-associated SNP in an exonic region. 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.

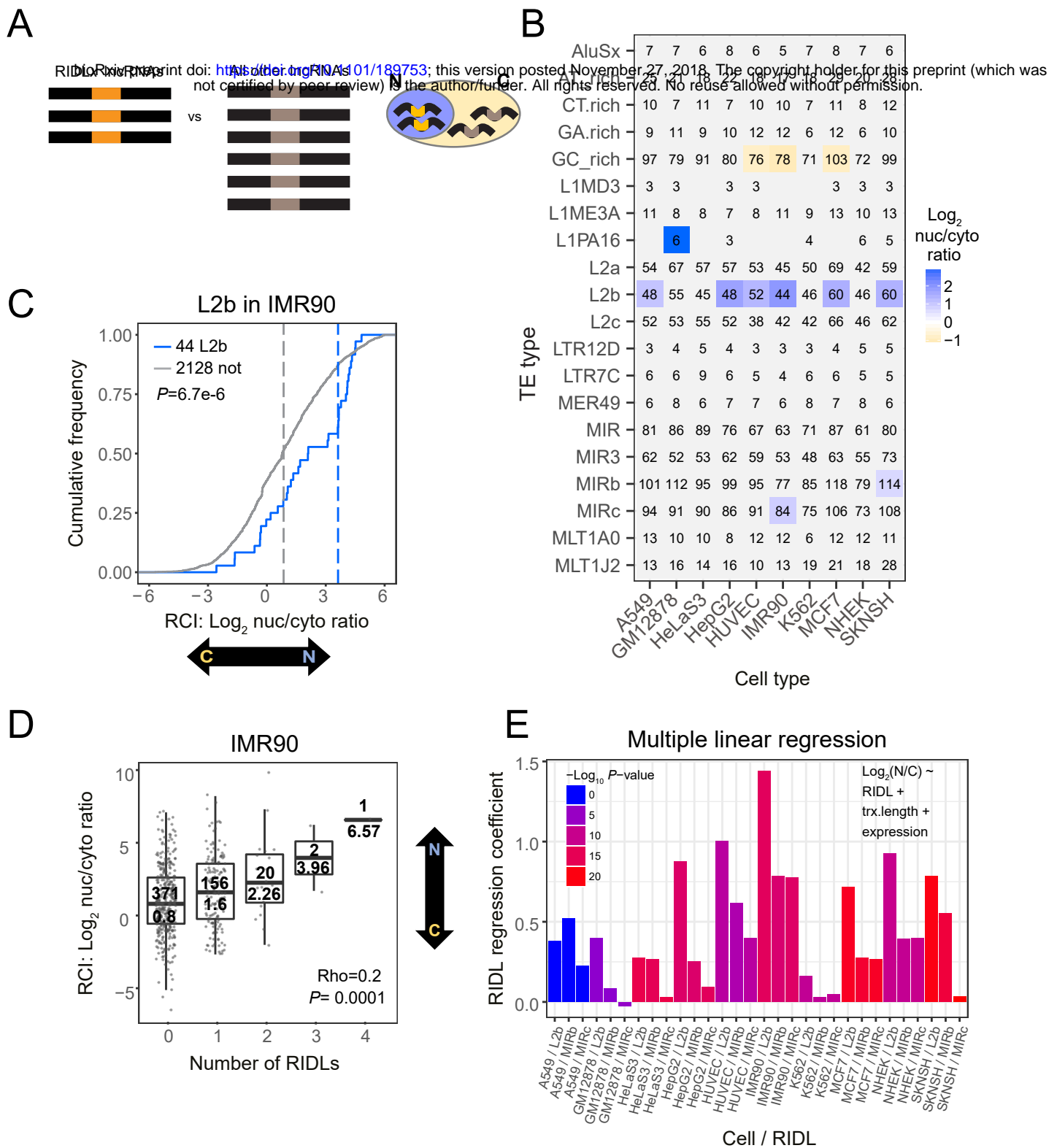


Figure 5: Correlation between RIDLs and host lncRNA nuclear/cytoplasmic localisation.

(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 IMR90 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 ($\text{Log}_2(\text{N/C})$) of their host transcript, while accounting for possible confounding factors of transcript length (trx.length) or whole-cell expression levels (expression).

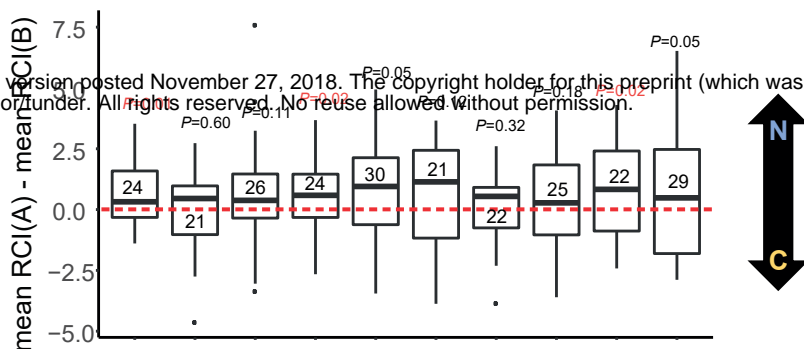
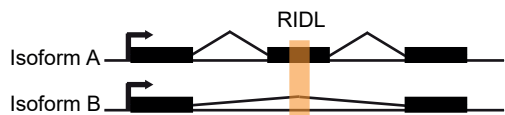
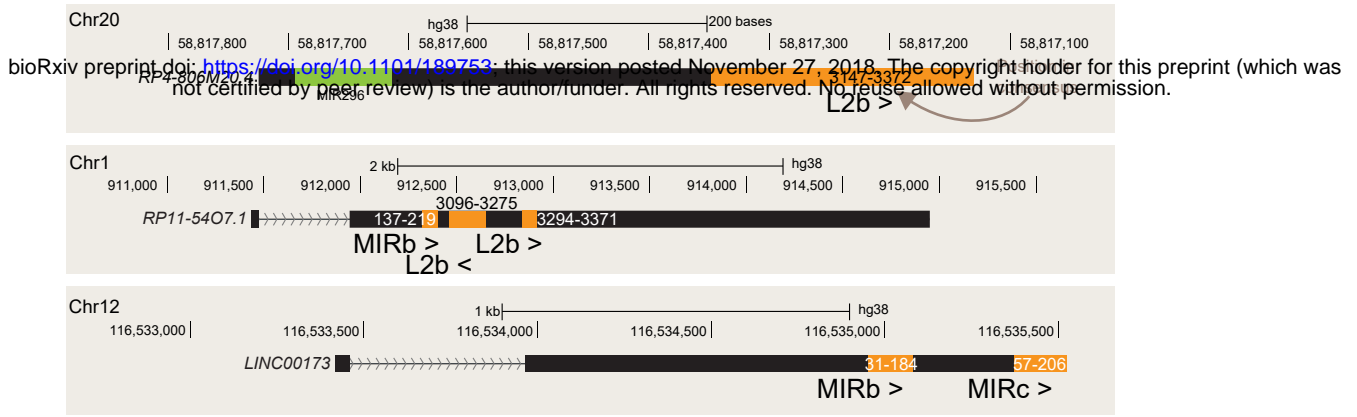


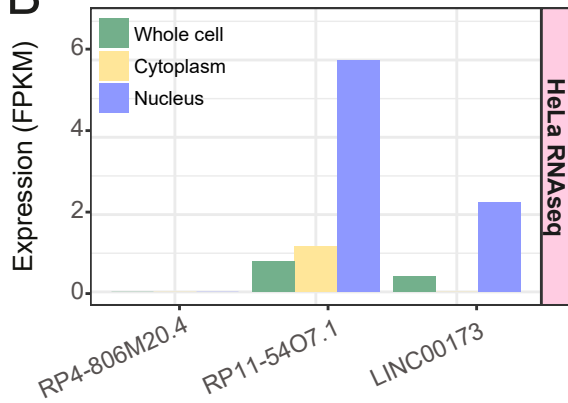
Figure 6: RIDLs correlate with differential localisation of lncRNA transcripts from the same locus.

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 $P < 0.05$).

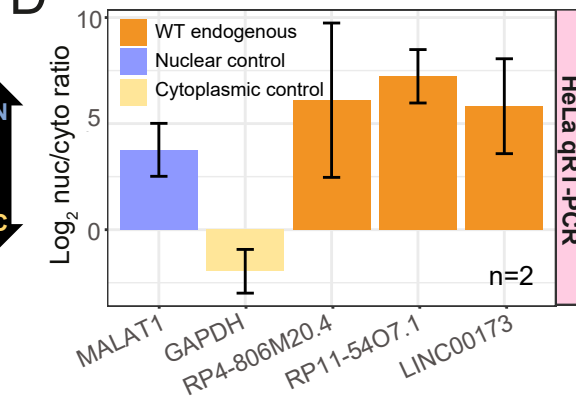
A



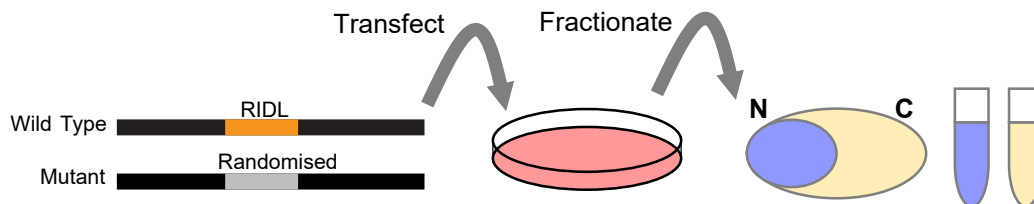
B



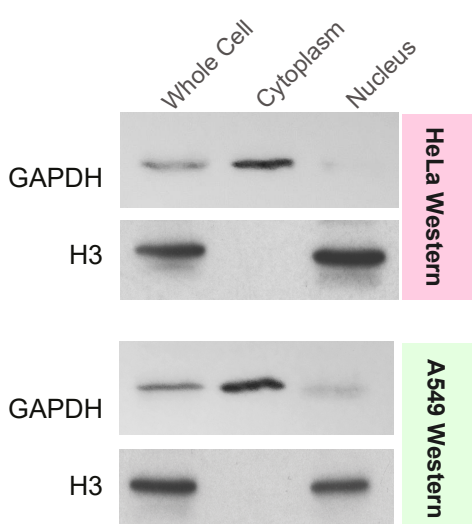
D



D



E



F

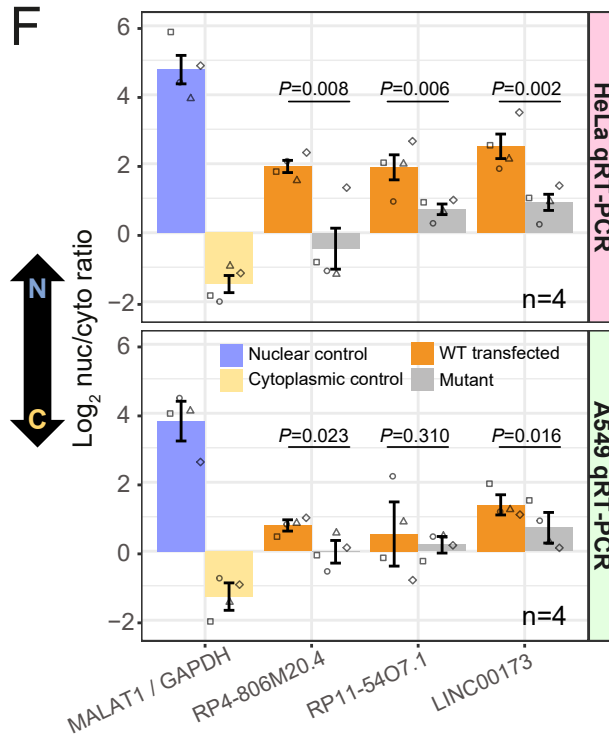


Figure 7: Disruption of RIDLs results in lncRNA relocalisation from nucleus to cytoplasm.

(A) Structures of candidate RIDL-ncRNAs. 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 RNAseq (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.