1

# Mosaic cis-regulatory evolution drives transcriptional partitioning of HERVH endogenous retrovirus in the human embryo

4

5    Thomas A. Carter[1], Manvendra Singh[1], Gabrijela Dumbović[2,3], Jason D. Chobirko[1],

6    John L. Rinn[2], Cédric Feschotte[1*]

7    *: corresponding author

8

9    1) Department of Molecular Biology and Genetics, Cornell University, Ithaca, NY 14850, USA

10   2) Department of Biochemistry, University of Colorado Boulder, Boulder, CO 80309, USA

11   3) Present address: Max Planck Institute of Immunobiology and Epigenetics, Freiburg, Germany

## **Abstract**

The human endogenous retrovirus type-H (HERVH) family is expressed in the preimplantation embryo. A subset of these elements are specifically transcribed in pluripotent stem cells where they appear to exert regulatory activities promoting self-renewal and pluripotency. How HERVH elements achieve such transcriptional specificity remains poorly understood. To uncover the sequence features underlying HERVH transcriptional activity, we performed a phyloregulatory analysis of the long terminal repeats (LTR7) of the HERVH family, which harbor its promoter, using a wealth of regulatory genomics data. We found that the family includes at least 8 previously unrecognized subfamilies that have been active at different timepoints in primate evolution and display distinct expression patterns during human embryonic development. Notably, nearly all HERVH elements transcribed in ESCs belong to one of the youngest subfamilies we dubbed LTR7up. LTR7 sequence evolution was driven by complex mutational processes, including multiple recombination events between subfamilies, that led to transcription factor binding motif modules characteristic of each subfamily. Using a reporter assay, we show that one such motif, a predicted SOX2/3 binding site unique to LTR7up, is essential for robust promoter activity in induced pluripotent stem cells. Together these findings illuminate the mechanisms by which HERVH diversified its expression pattern during evolution to colonize distinct cellular niches within the human embryo.

## Introduction

Transposable elements (TEs) are genomic parasites that use the host cell machinery for their own propagation. To propagate in the host genome, they must generate new insertions in germ cells or their embryonic precursors, as to be passed on to the next generation (Charlesworth and Langley, 1986; Cosby et al., 2019; Haig, 2016). To this end, many TEs have evolved stage-specific expression in germ cells or early embryonic development (Faulkner et al., 2009; Fort et al., 2014; Göke et al., 2015; Miao et al., 2020; Urusov et al., 2011). But how does this precise control of TE expression evolve?

Many endogenous retroviruses (ERVs) are known to exhibit highly stage-specific expression during early embryonic development (Chang et al., 2021; Göke et al., 2015; Hermant and Torres-Padilla, 2021; Peaston et al., 2004; Svoboda et al., 2004). ERVs are derived from exogenous retroviruses with which they share the same prototypical structure with two long terminal repeats (LTRs) flanking an internal region encoding products promoting their replication (Eickbush and Malik, 2002). There are hundreds of ERV families and subfamilies in the human genome, each associated to unique LTR sequences (Kojima, 2018; Vargiu et al., 2016). Each family has infiltrated the germline at different evolutionary timepoints and have achieved various levels of genomic amplification (Bannert and Kurth, 2004; Vargiu et al., 2016). One of the most abundant families is HERVH, a family derived from a gamma retrovirus that first entered the genome of the common ancestor of apes, Old World monkeys, and New World monkeys more than 40 million years ago (mya) (Goodchild et al., 1993; Izsvák et al., 2016; Mager and Freeman, 1995).

There are four subfamilies of HERVH elements currently recognized in the Dfam (Storer et al., 2021) and Repbase (Bao et al., 2015; Kojima, 2018) databases and annotated in the reference human genome based on distinct LTR consensus sequences: LTR7 (formerly known as Type I), 7b (Type II), 7c, and 7y (Type Ia) (Bao et al., 2015; Goodchild et al., 1993; Jern et al., 2005, 2004). Additional subdivisions of HERVH elements were also proposed based on phylogenetic analysis and structural variation of their internal gene sequences (Gemmell et al., 2019; Jern et al., 2005, 2004). However, all HERVH elements are currently annotated in the human genome using a single

62  consensus sequence for the internal region (HERVH_int) and the aforementioned four

63  LTR subfamilies.

64  HERVH has been the focus of extensive genomic investigation for its high level of RNA

65  expression in human embryonic stem cells (ESCs) and induced pluripotent stem cells

66  (iPSCs) (Fort et al., 2014; Gemmell et al., 2015; Izsvák et al., 2016; Kelley and Rinn,

67  2012; Loewer et al., 2010; Römer et al., 2017; Santoni et al., 2012; Zhang et al., 2019).

68  Several studies showed that family-wide HERVH knockdown results in the loss of

69  pluripotency of human ESC and reduced reprogramming efficiency of somatic cells to

70  iPSC (Lu et al., 2014; Ohnuki et al., 2014; Wang et al., 2014). Others reported similar

71  phenotypes with the knockdown of individual HERVH-derived RNAs such as those

72  produced from the *lincRNA-RoR* and *ESRG* loci (Loewer et al., 2010; Wang et al.,

73  2014) or the deletion of individual HERVH loci acting as boundaries for topological

74  associated domains (Zhang et al., 2019). These results converge on the notion that

75  HERVH products (RNA or proteins) exert some modulatory effect on the cellular

76  homeostasis of pluripotent stem cells. However, it is important to emphasize that

77  different HERVH knockdown constructs produced variable results and inconsistent

78  phenotypes (Lu et al., 2014; Wang et al., 2014; Zhang et al., 2019), and a recent

79  knockout experiment of the most highly transcribed locus (*ESRG*) failed to recapitulate

80  its previous knockdown phenotype (Takahashi et al., 2021). Despite intense study,

81  which expressed HERVH loci, if any, are necessary for the maintenance of pluripotency

82  remain unclear.

83  The mechanisms regulating the transcription of HERVH also remain poorly understood.

84  RNA-seq analyses have established that HERVH expression in human ESCs, iPSCs,

85  and the pluripotent epiblast can be attributed to a relatively small subset of loci

86  (estimated between 83 and 209) driven by LTR7 (sensu stricto) sequences (Göke et al.,

87  2015; Wang et al., 2014; Zhang et al., 2019). The related 7y sequences are known to

88  be expressed in the pluripotent epiblast of human embryos (Göke et al., 2015) and a

89  distinct subset of elements associated with 7b and 7y sequences are expressed even

90  earlier in development at the onset of embryonic genome activation (Göke et al., 2015).

91  These observations suggest that the HERVH family is composed of subsets of elements

92   expressed at different timepoints during embryonic development and that these

93   expression patterns reflect, at least in part, the unique cis-regulatory activities of their

94   LTRs. While it has been reported that several transcription factors (TFs) bind and

95   activate HERVH LTRs, including the pluripotency factors OCT4, NANOG, SP1, and

96   SOX2 (Göke et al., 2015; Ito et al., 2017; Kelley and Rinn, 2012; Kunarso et al., 2010;

97   Ohnuki et al., 2014; Pontis et al., 2019; Santoni et al., 2012), it remains unclear how TF

98   binding contributes to the differential expression of HERVH subfamilies and why only a

99   minority of HERVH are robustly transcribed in pluripotent stem cells and embryonic

100  development.

101  To shed light on these questions, we focused this study on the cis-regulatory evolution

102  of LTR7 elements. We use a "phyloregulatory" approach combining phylogenetic

103  analyses and regulatory genomics to investigate the sequence determinants underlying

104  the partitioning of expression of HERVH/LTR7 subfamilies during early embryonic

105  development.

# Results

## LTR7 consists of 8 previously undefined subfamilies

We began our investigation by examining the sequence relationships of the four LTR7 subfamilies currently recognized in the human genome: LTR7 sensu stricto (748 proviral copies; 711 solo LTRs), 7b (113; 524), 7c (24; 223), and 7y (77; 77). We built a maximum likelihood phylogenetic tree from a multiple sequence alignment of a total of 781 5' LTR and 1073 solo LTR sequences of near complete length (>350 bp) representing all intact LTR subfamilies extracted from the RepeatMasker output of the hg38 human reference assembly. While 7b and 7y sequences cluster, as expected, into clear monophyletic clades with relatively short internode distances and little subclade structure, sequences from the 7c and LTR7 subfamilies were much more heterogeneous and formed many subclades (Fig. 1A). Notably, sequences annotated as LTR7 were split into distinct monophyletic clades indicative of previously unrecognized subfamilies within that group. The branch length separating some of these LTR7 subclades were longer from one another than they were from those falling within the 7b, 7c, and 7y clades, indicating that they represent subfamilies as different from each other as those previously recognized (Fig. 1A).

We next sought to classify LTR7 elements more finely by performing a phylogenetic analysis using a multiple sequence alignment of all intact LTR7 sequences (>350 bp) along with the consensus sequences for the other LTR7 subfamilies for reference. We defined high-confidence subfamilies as those forming a clade supported by >95% ultrafast bootstrap (UFbootstrap) and internal branches >0.015 (1.5 nucleotide substitutions per 100 bp) separating subgroup nodes. Based on these criteria, LTR7 elements could be divided into 8 subfamilies (Fig. 1B).

While long internal branches with high UFbootstrap support separate LTR7 subfamilies, intra-subfamily internal branches with >95% UFbootstrap support were shorter (<0.015), suggesting that each subfamily was the product of a rapid burst of amplification of a
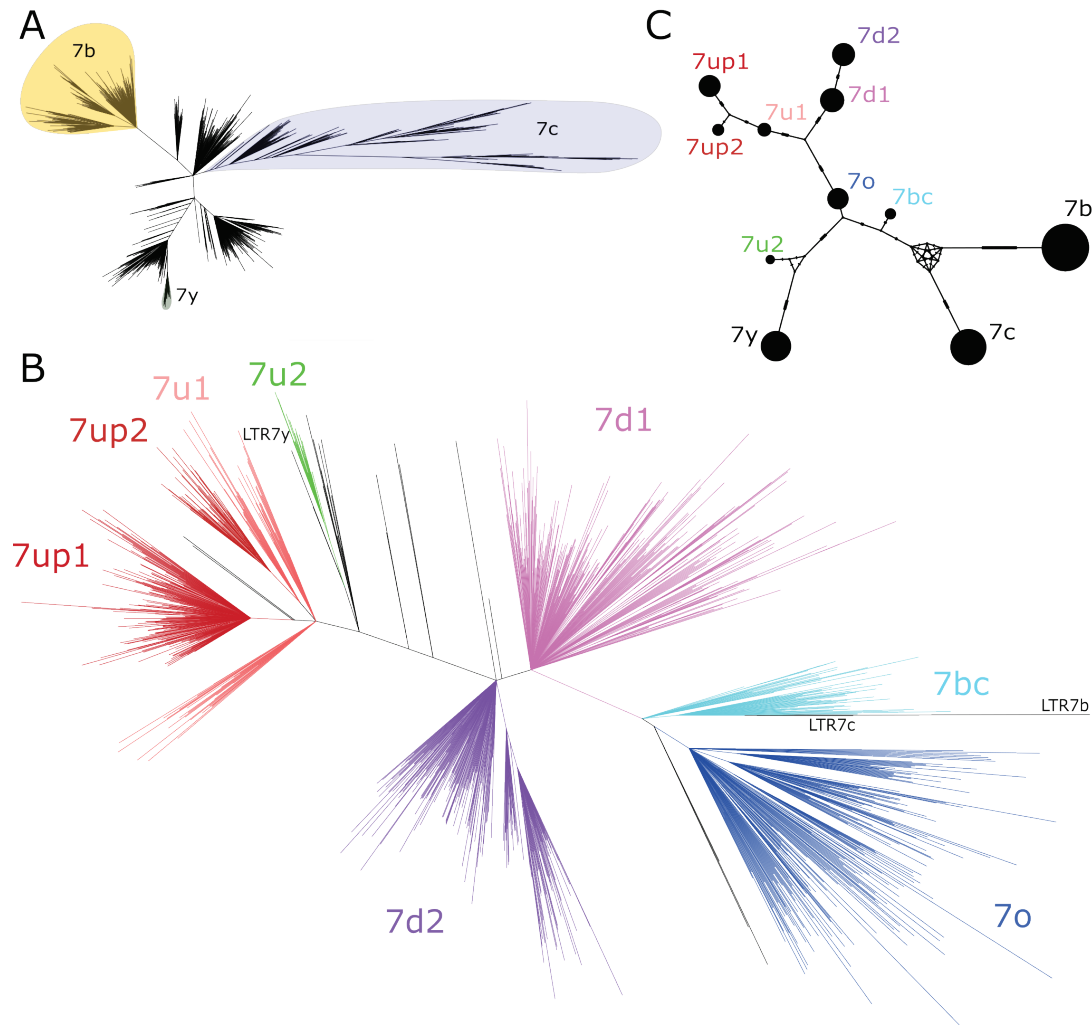
6

Fig. 1: Phylogenetic analysis of LTR7 sequences. A) Unrooted phylogeny of all solo and 5' LTR7 sequences. All nodes with UFbootstraps >0.95, >10 member insertions, and >1.5 substitutions / 100 bp (~6 base pairs) are grouped and colored (see methods). Previously listed consensus sequences from 7b/c/y were included in the alignment and are shown in black. B) Unrooted phylogeny of all solo and 5' LTR7 subfamilies from 1a, 7b, 7c, and 7y. Colors denote clades consisting of previously annotated 7b, 7c, and 7y with >95% concordance. C) Median joining network analysis of all LTR7 and related majority rule consensus sequences. Ticks indicate the number of SNPs at non-gaps between consensus sequences. The size of circles is proportional to the number of members in each subfamily. Only LTR7 insertions that met filtering requirements (see methods) are included while 7b/c/y counts are from dfam.

133    common ancestor. To approximate the sequence of these ancestral elements we

134    generated majority-rule consensus sequence for each of the 8 newly defined LTR7

135    subfamilies (7o, 7bc, 7up, etc.). The consensus sequences were deposited at

136    www.dfam.org.

137    To investigate the evolutionary relationships among the newly defined and previously

138    known LTR7 subfamilies, we conducted a median-joining network analysis (Leigh and

139 Bryant, 2015) of their consensus sequences (Fig. 1C). The network analysis provides

140 additional information on the relationships between subfamilies and approximates the

141 shortest and most parsimonious paths between them (Bandelt et al., 1999; Cordaux et

142 al., 2004; Posada and Crandall, 2001). The results place 7o in a central position from

143 which two major lineages are derived. One lineage led to two sub-lineages, formed by

144 7up1, 7up2, and 7u1 (with 7up1 and 7up2 being most closely related) and by 7d1 and

145 7d2. The other lineage emanating from 7o rapidly split into two sub-lineages; one gave

146 rise to 7u2 and then to 7y and the other gave rise to 7bc which is connected to the two

147 more diverged subfamilies 7b and 7c (Fig. 1C). Together these results indicate that the

148 LTRs of HERVH elements can be divided into additional subfamilies than those

149 previously recognized.

150

**151 The age of LTR7 subfamilies suggests three major waves of HERVH propagation**

152 The genetic differences between LTR7 subfamilies suggest that they may have been

153 active at different evolutionary timepoints. To examine this, we used reciprocal *liftover*

154 analysis to infer the presence/absence of each human LTR7 locus across five other

155 primate genomes. Insertions shared at orthologous genomic position across a set of

156 species are deemed to be ancestral to these species and thus can be inferred to be at

157 least as old as the divergence time of these species (Johnson, 2019).

158 The results of this cross-species analysis indicate that LTR7 subfamilies have been

159 transpositionally active at different timepoints in the primate lineage (Fig. 3A). The
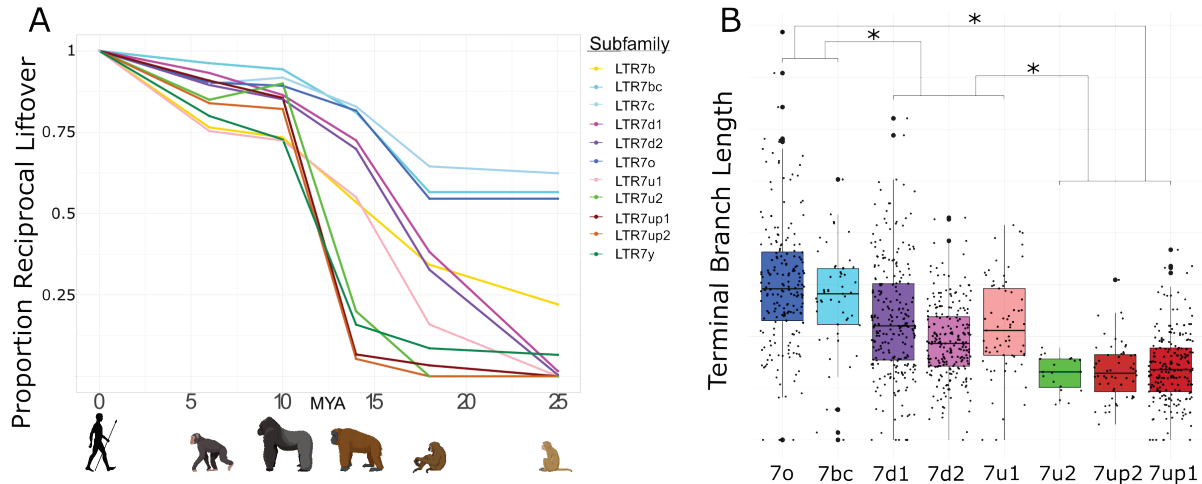
Fig. 2: Age analysis of LTR7 subfamilies. A) Proportion of a given subfamily that have 1:1 orthologous insertions between human and other primate species. LTR7 subfamilies are from trees in Figs. 1a and 2a; 7b/c/y subfamilies are from RepeatMasker annotations. Non-human primates are spaced out on the X axis in accordance with their approximate divergence times to the human lineage. B) Terminal branch lengths of all LTR7 insertions from Fig. 1a. Groups with similar liftover profiles were merged for statistical testing (see methods). Differences with padj<1e-15 are denoted with * (Wilcox rank-sum test with Bonferroni correction).

160  subfamilies 7o, 7bc, and 7c are the oldest since the majority of their insertions are found

161  at orthologous position in rhesus macaque, an Old World Monkey (OWM).  These three

162  subfamilies share similar evolutionary trajectories, with most of their proliferation

163  occurring prior to the split of OWM and hominoids, ~25 mya (Fig. 2a). Members of the

164  7b subfamily (the most numerous, 637 solo and full-length insertions) appear to be

165  overall younger, since only 22% of the human 7b elements could be lifted over to

166  rhesus macaque and the vast majority appeared to have inserted between 10 and 20

167  mya (Fig. 2A, Figure supplement 1). Only 5 of the 550 elements in the 7d1 and 7d2

168  subfamilies could be retrieved in rhesus macaque, but ~30% were shared with gibbon

169  and ~75% were shared with orangutan. Thus, these two subfamilies are largely

170  hominoid-specific and achieved most of their proliferation prior to the split of African and

171  Asian great apes ~14 mya (Fig. 3a). Members of the 7u1 subfamily also emerged in the

172  hominoid ancestor, but the majority (55%) of 7u1 elements present in the human

173  genome inserted after the split of gibbons in the great ape ancestor, between 14 and 20

174  mya. Thus, the 7b, 7d1/2, and 7u1 subfamilies primarily amplified during the same

175  evolutionary window, 14 to 20 mya.

176     The 7up1/2, 7y, and 7u2 subfamilies represent the youngest in the human genome, with

177     most of their proliferation occurring between ~10 and ~14 mya, in the ancestor of

178     African great apes (Fig. 3A). Based on these results, these subfamilies seem to have

179     experienced a burst of transposition after the divergence of African and Asian great

180     apes but before the split of the pan/homo and gorilla lineages. For example, only 14 of

181     the 208 (6.7%) human 7up1 elements can be retrieved in orangutan, but 178 (85.6%)

182     can be found in gorilla. These data indicate that the three youngest LTR7 subfamilies

183     mostly expanded in the ancestor of African great apes (Fig. 2C).

184     As an independent dating method, we used the terminal branch length separating each

185     insertion from its nearest node in Fig. 1B (Fig. 2B). Here, the terminal branch lengths

186     are proportional to nucleotide divergence accumulated after insertion and can thus

187     approximate each insertion's relative age. This method largely corroborated the results

188     of the *liftover* analysis and revealed three age groups among LTR7 subfamilies

189     characterized by statistically different mean branch lengths (p(adj)< 1e-15; Wilcox rank-

190     sum test). By contrast, we found no statistical difference between the mean branch

191     length of the subfamilies within these three age groups, suggesting that they were

192     concomitantly active. Taken together, our dating analyses distinguish 3 major waves of

193     HERV propagation: an older wave 25-40 mya involving 7c, 7o, and 7bc elements, an

194     intermediate wave 9-20 mya involving 7b, 7d1/2 and 7u1, and a most recent wave 4-10

195     mya implicating primarily 7up1/2, 7u2 and 7y elements.

196

197     **Only LTR7up shows robust transcription in human ESC and iPSC**

198     Our data thus far indicate that LTR7 is composed of genetically and evolutionarily

199     distinct subfamilies. Because a subset of HERVH elements linked to LTR7 were

200     previously reported to be transcribed in pluripotent stem cells (human ESCs and

201     iPSCs), we wondered whether this activity was restricted to one or several of the LTR7

202     subfamilies newly defined herein. To investigate this, we performed a "phyloregulatory"

203     analysis, where we layered locus-specific regulatory data obtained from publicly

204     available genome-wide assays in ESCs (mostly from the H1 cell line, see methods) for

205     each LTR insertion on top of a phylogenetic tree depicting their evolutionary

206  relationship. We called an individual LTR7 insertion as positive for a given feature if

207  there is overlap between the coordinates of the LTR and that of a peak called for this

208  mark (see methods). We predicted that if transcriptional activity was an ancestral

209  property of a given subfamily, evidence of transcription and "activation" marks should be

210  clustered within the cognate clade. Alternatively, if transcription and activation marks

211  were to be distributed throughout the tree, it would indicate that LTR7 transcriptional

212  activity in pluripotent cells was primarily driven by post-insertional changes or context-

213  specific effects. Differences in the proportion of positive insertions for a given mark

214  between LTR7 subfamilies were tested using a chi-square test with Bonferroni

215  correction. Unless otherwise noted, all proportions compared thereafter were

216  significantly different (padj< 0.05).

217  The results (Fig. 3A) show that HERVH elements inferred to be "highly expressed"

218  (fpkm > 2) based on RNA-seq analysis (Wang et al., 2014) were largely confined to two

219  closely related subfamilies, 7up1 and 7up2, together referred to as 7up hereafter.

220  Indeed, we estimated that 33% of 7up elements (88 loci) are highly expressed

221  according to RNA-seq compared with only 2% of highly expressed elements from all

222  other subfamilies combined (17 loci). Nascent RNA mapping using GRO-seq data

223  (Estarás et al., 2015) recapitulated this trend with 22% of 7up loci with visible signal

224  (Figure supplement 2), compared with only 4% of other LTR7 loci (Fig. 3D, Figure

225  supplement 2). Half of the loci displaying GRO-seq signal (53/96) also showed evidence

226  of mature RNA product (supp. file 1). Thus, HERVH transcriptional activity in H1 ESCs

227  is largely limited to loci driven by 7up sequences.

228  As previously noted from ChIP-seq data (Ohnuki et al., 2014), we found that KLF4

229  binding is a strong predictor of transcriptional activity: KLF4 ChIP-seq peaks overlap

230  91% of 7up loci and KLF4 binding is strongly enriched for the 7up subfamilies relative to

231  other subfamilies (Fig. 3A,B,D). NANOG binding is also enriched for 7up (97.7% of loci

232  overlap ChIP-seq peaks) but is observed to varying degrees at other LTR7 loci that do

233  not show evidence of active transcription based on GRO-seq and/or RNA-seq (85% of

234  7u1 loci, 32% 7d1, 45% 7d2, 13% 7o, 8.7% 7bc, and 0% of 7u2). Other TFs with known

235  roles in pluripotency are also enriched at 7up loci, such as SOX2 (32% LTR7up, 1-3%

236  all other LTR7), FOXP1(49%, 0-4.3%), and FOXA1(28%, 0-1.4%). In fact, FOXA1 binds

237  only a single non-7up insertion in our dataset, making it the most exclusive feature of

238  7up loci among the TFs examined in this analysis. In contrast, OCT4 binds merely 12%

239  of 7up loci (see supp. file 8 for full statistical analysis of all marks).

240  Congruent with having generally more TF binding and transcriptional activity, 7up loci

241  also have a propensity to be decorated by H3K4me3, a mark of active promoters (76%

242  LTR7up vs 19% all others) and the broader activity mark H3K27ac (89% vs 48%) (Fig.

243  3A,B). By contrast, H3K4me1, a mark typically associated with low POLII loading as

244  seen at enhancers as opposed to promoters, is spread rather evenly throughout the tree

245  of LTR7 sequences (26% vs. 18%) (Fig. 2A,B). Thus, promoter marks are primarily

246  restricted to 7up loci, but a broader range of 7up loci display putative enhancer marks.

247  Taken together, our phyloregulatory analysis suggests that strong promoter activity in

248  ESCs is restricted to 7up elements.

249

**250  Differential activation, rather than repression, explain the differential**

**251  transcriptional activity of LTR7 subfamilies in ESCs**

252  The pattern described above could be explained by two non-mutually exclusive

253  hypotheses: (i) 7up elements (most likely their progenitor) have acquired unique

254  sequences (TF binding sites, TFBS) that promote Pol II recruitment and active

255  transcription, and/or (ii) they somehow escape repressive mechanisms that actively

256  target the other subfamilies, preventing their transcription. For instance, 7up elements

257  may lack sequences targeted by transcriptional repressors such as KRAB-Zinc Finger

258  proteins (KZFP) that silence the other subfamilies in ESCs. KZFP are well-known for

259  binding TEs in a subfamily-specific manner where they nucleate inheritable epigenetic

260  silencing (Ecco et al., 2017; Jacobs et al., 2014; Wolf et al., 2020; Yang et al., 2017)

261  and several KZFPs are known to be capable of binding LTR7 loci (Imbeault et al.,

262  2017). To examine whether KZFPs may differentially bind to LTR7 subfamilies, we

263  analyzed the loading of the corepressor KAP1 and the repressive histone mark
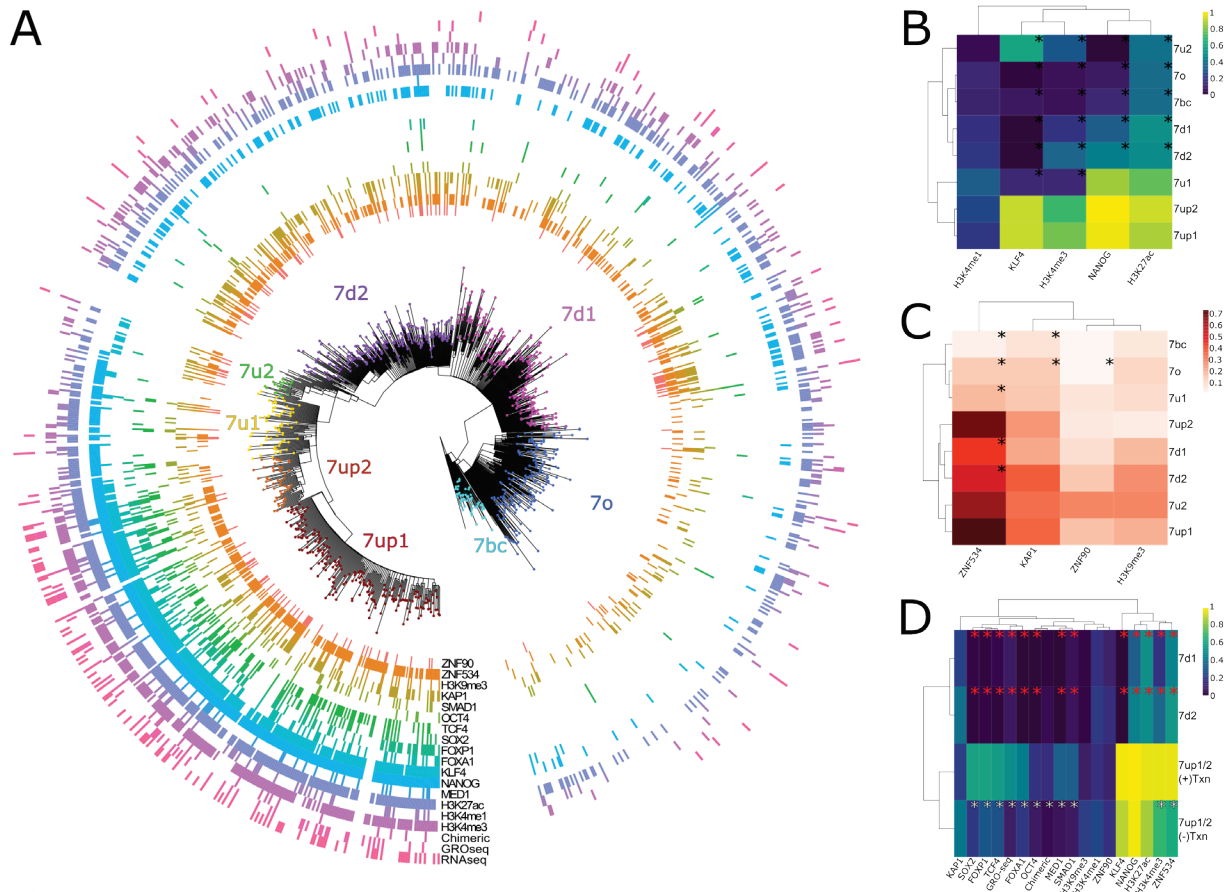
Fig. 3: Phyloregulatory analysis of LTR7. A) "Phyloregulatory" map of LTR7. The phylogenetic analysis to derive the circular tree is the same as for the tree in Fig. 1A but rooted on the 7b consensus. Subfamilies defined in Fig. 1 are denoted with dotted colored tips. Positive regulatory calls for each insertion are shown as tick marks of different colors and no tick mark indicates a negative call. All marks are derived from ESC except for ZNF90 and ZNF534, which are derived from ChIP-exo data after overexpression of these factors in HEK293 cells (see methods) B) Heatmap of major activation and repression profiles. Proportions indicate the proportion of each group positive for a given characteristic. Trees group LTR7 subfamilies on regulatory signature, not sequence similarity. Asterisks denote statistical differences between given group and 7up1 (padj> 0.05 Wilcox rank-sum with Bonferroni correction). C) Heatmap done in similar fashion to Fig. 3B but for repression marks. D) Heatmap of transcribed (>2 fpkm) and untranscribed 7up1/2 (<2 fpkm) and all 7d1/2. Red asterisks denote statistical differences between 7d1/2 and 7up1 (padj< 0.05 chi-square Bonferroni correction). White asterisks denote differences between transcribed and untranscribed LTR7up.

H3K9me3 typically deposited through the KZFP/KAP1 complex, across the LTR7

phylogeny using ChIP-seq data previously generated for ESCs (Imbeault et al., 2017;

Theunissen et al., 2016). We found that KAP1 and H3K9me3 loading were neither

enriched nor depleted for 7up elements relative to other subfamilies (Fig. 3A,C). Overall,

there were no significant differences in the level of H3K9me3 marking across

subfamilies and the only difference in KAP1 binding was a slight but significant

13

270  depletion for 7bc and 7o compared to all other subfamilies including 7up (14% vs. 35% -

271  padj< 0.05 chi-square Bonferroni correction). Furthermore, KAP1 and H3K9me3 loading

272  were found in similar proportions in expressed and unexpressed 7up elements (padj>

273  0.05) (Fig. 2C). This was also the case for CpG methylation, whose presence was not

274  differential between subfamilies (padj> 0.05 Wilcox rank-sum with Bonferroni correction)

275  (Figure supplement 2). Thus, KAP1 binding and repressive marks at LTR7 in ESCs

276  poorly correlate with their transcriptional activity and differential repression is unlikely to

277  explain the differential promoter activity of LTR7 subfamilies in ESCs.

278  We also examined the binding profile of ZNF534 and ZNF90, two KZFPs previously

279  reported to be enriched for binding LTR7 elements using ChIP-exo data in human

280  embryonic kidney 293 cells (Imbeault et al., 2017), in order to examine whether they

281  bind a particular subset of elements in our LTR7 phylogeny. We found that while ZNF90

282  bound all LTR7 subfamilies to a similar extent, ZNF534 preferentially bound members of

283  the 7up subfamily (72% of LTR7up vs. 34-53% of non-LTR7up). However, ZNF534

284  binding in 293 cells did not correlate with transcriptional activity of 7up elements in

285  ESCs nor with KAP1 binding or H3K9me3 deposition in these cells (Fig. 3A,D). In other

286  words, there was no significant enrichment for ZNF534 binding within untranscribed 7up

287  elements nor depletion within the 7up elements we inferred to be highly transcribed in

288  ESCs. These observations could simply reflect the fact that ZNF534 itself is not highly

289  expressed in ESCs (Figure supplement 3) and do not preclude that ZNF534 represses

290  7up in other cellular contexts or cell types. Collectively these data suggest that

291  differential LTR binding of KZFP/KAP1 across subfamilies cannot readily explain their

292  differential regulatory activities in ESCs. Thus, differential activation is the most likely

293  driver for the promoter activity of 7up elements in ESCs.

294  To determine which factors are associated and potentially determinant for 7up promoter

295  activity, we compared the set of "highly expressed" 7up loci to 7up loci which are

296  apparently poorly expressed, using 7d1/d2 as outgroups (Fig. 3D). While known

297  regulators of LTR7 transcription, KLF4 and NANOG, are enriched for binding to 7up

298  elements, their occupancy alone cannot distinguish transcribed from untranscribed 7up

299  loci (Fig. 3D). Thus, other factors must contribute to the transcriptional activation of 7up

300 elements. Our analysis of pluripotent transcriptional activators SOX2, FOXA1, FOXP1,

301 OCT4, TCF4, and SMAD1 (Boyer et al., 2005; Chambers and Smith, 2004; Niwa, 2007)

302 binding profiles show that all of these TFs are enriched in robustly transcribed 7up loci

303 compared to non-transcribed loci (Fig. 3D). Intriguingly, when overexpressed in HEK293

304 cells, the potential KZFP repressor ZNF534 preferably binds ESC-transcribed 7up over

305 untranscribed 7up, suggesting that ZNF534 may suppress transcription-competent 7up

306 in cellular contexts where this factor is expressed.

307 Together these data suggest that differential repression cannot explain the differential

308 promoter activity of LTR7 subfamilies in ESCs but rather that highly expressed LTR7up

309 loci are preferentially bound by a cocktail of transcriptional activators that are less

310 prevalent on poorly-expressed loci.

311

312 **Inter-element recombination and intra-element duplication drove LTR7 sequence**

313 **evolution**

314 The data presented above suggest that the transcriptional activity of 7up in ESCs

315 emerged from the gain of a unique combination of TFBS. To identify sequences unique

316 to 7up relative to its closely related subfamilies, we aligned the consensus sequences of

317 the newly defined LTR7 subfamilies and those of 7b/c/y consensus sequences. This

318 multiple sequence alignment revealed blocks of sequences that tend to be highly

319 conserved across subfamilies, only diverging by a few SNPs, while other regions

320 showed insertion/deletion (indel) segments specific to one or a few subfamilies (Fig.

321 4A). These indels resulted in substantial gain and loss of DNA between closely related

322 subfamilies, with the longest consensus (7y) having a length of 472-bp and the shortest

323 (7o) a mere 365-bp. These observations suggest that segmental rearrangements have

324 played an important role in the evolution of LTR7 sequences.

325 Upon closer scrutiny, we noticed that the indels characterizing some of the subfamilies

326 were at odds with the evolutionary relationship of the subfamilies defined by overall

327 phylogenetic and network analyses. This was particularly obvious in segments we
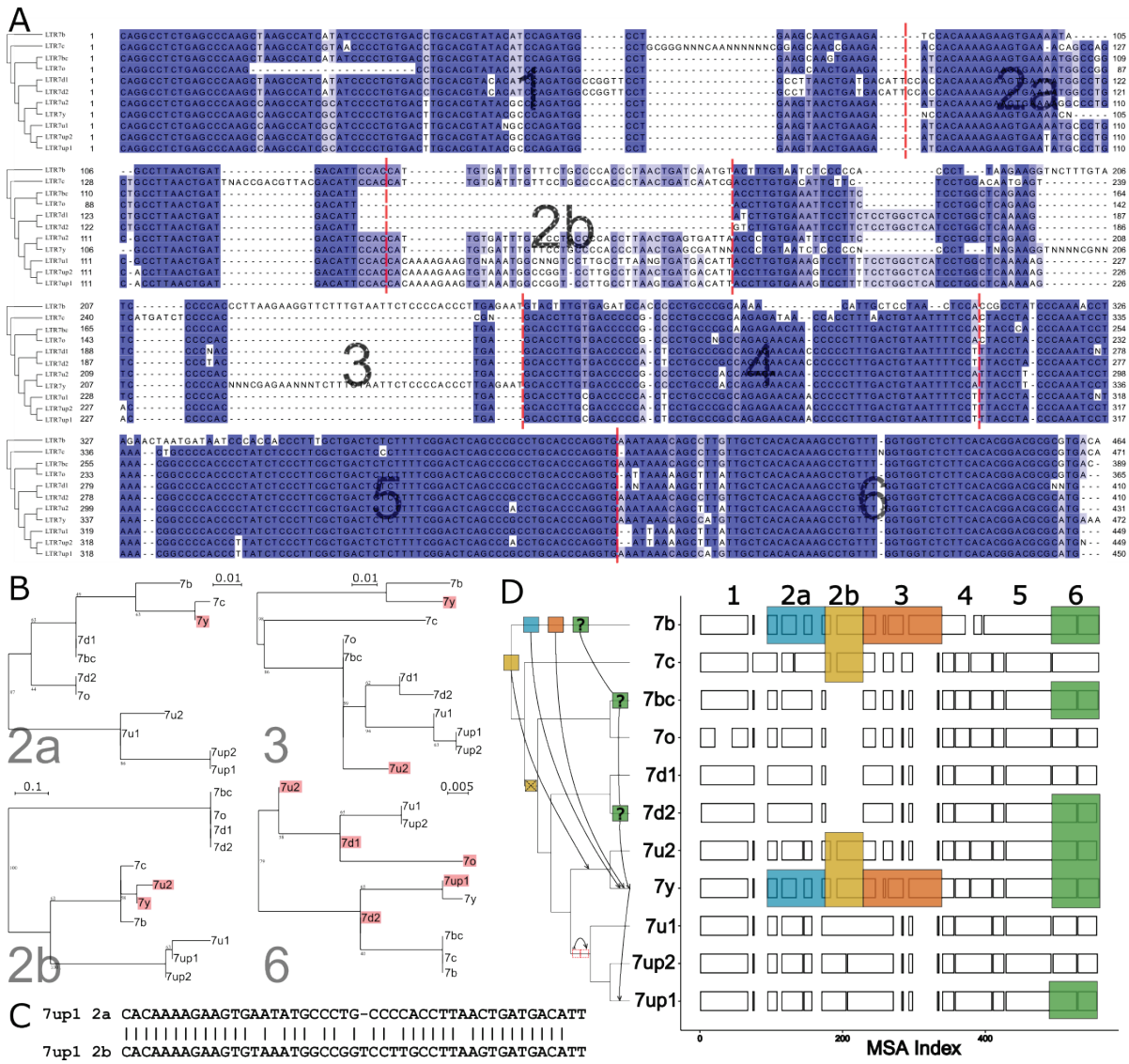
Fig. 4: Modular block evolution of LTR7 subfamilies. A) A multiple sequence alignment of LTR7 subfamily consensus sequences. The phylogenetic topology from Figure 1 is shown on the left. The MSA is broken down into sequence blocks (red lines) with differential patterns of relationships. B) Parsimony trees from Fig. 4a sequence blocks. Subfamilies whose blocks do not match the overall phylogeny are highlighted in red. Bootstrap values >0 are shown. C) Blastn alignment of LTR7up1 block 2a and 2b. D) A multiple sequence alignment of majority-rule consensus sequences from each LTR7 subfamily detailing shared structure. Blocks show aligned sequence; gaps represent absent sequence. Colored sections identify putative phylogeny-breaking events. Recombination events whose directionality can be inferred (via aging) are shown with blocks and arrows on the cladogram. Recombination events with multiple possible routes are denoted with "?". The deletion of 2b is denoted on the cladogram with a red "X"; the duplication of 2a is denoted with 2 red rectangles.

328    termed block 2b (where 7y and 7u2 share a large insertion with 7b and 7c) and block 3

329    (where 7y and 7b share a large insertion). This led us to carefully examine the multiple

330    sequence alignment of the LTR7 consensus sequences to identify indels with different

331    patterns of inter-subfamily relationships. Based on this analysis, we defined seven

332    sequence blocks shared by a different subset of subfamilies, pointing at relationships

333    that were at odds with the overall phylogeny of the LTR7 subfamilies (Fig. 4A-B). These

334    observations suggested that some of the blocks have been exchanged between LTR7

335    subfamilies through recombination events.

336    To systematically test if recombination events between elements drove the evolution of

337    LTR7 subfamilies, we generated parsimony trees for each block of consensus

338    sequences and looked for incongruences with the overall consensus phylogeny. We

339    found a minimum of 6 instances of clades supported in the block parsimony trees that

340    were incongruent with those supported by the overall phylogeny (Fig. 4B,D).

341    We also found some blocks evolved via tandem duplication. Notably, block 2b was

342    absent from 7d1/2 and 7bc/o but present in all other subfamilies. However, block 2b

343    from 7b, 7c, 7u2, and 7y aligned poorly with block 2b from 7up and 7u1. Instead, block

344    2b from 7up/u1 2b was closely related (~81% nucleotide similarity) to block 2a from the

345    same subfamilies (Fig. 4D), suggestive that it arose via tandem duplication in the

346    common ancestor of these subfamilies. To further clarify the evolutionary history of the

347    2a-2b duplication, we aligned all 2a and 2b blocks from all subfamilies and generated a

348    parsimony tree (Figure supplement 4). This analysis indicated that the 2b block from

349    7up/u1 most closely resembles the 2a block from 7d.

350    The results above suggest that the evolution of HERVH was characterized by extensive

351    diversification of LTR sequences through a mixture of point mutations, indels, and

352    recombination events.

353

**354    HERVH subfamilies show distinct expression profiles in the preimplantation**

**355    embryo**

356    We hypothesized that the mosaic pattern of LTR sequence evolution described above

357    gave rise to TFBS combinations unique to each family that drove shifts in HERVH

358    expression during early embryogenesis. To test this, we aimed to reanalyze the

359    expression profiles of newly defined LTR7 subfamilies during early human

360    embryogenesis and correlate these patterns with the acquisition of embryonic TF

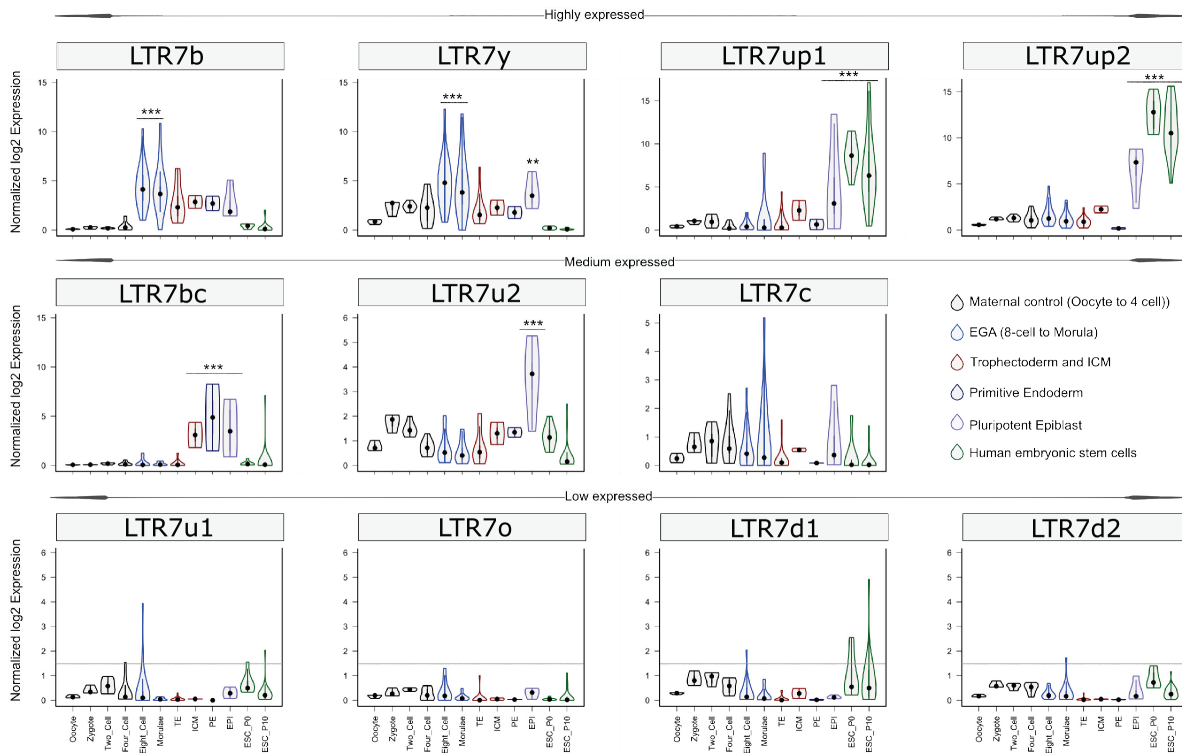361    binding motifs within each of the subfamilies.



Fig. 5: Expression profile of LTR7 subfamilies in human preimplantation embryonic lineages and ESCs. The solid dots and lines encompassing the violins represent the median and quartiles of single cellular RNA expression. The color scheme is based on embryonic stages, defined as maternal control of early embryos (Oocytes, Zygote, 2-cell and 4-cell stage), EGA (8-cell and Morula), inner cell mass (ICM), trophectoderm (TE), epiblast (EPI) and primitive endoderm (PE) from the blastocyst, and ESCs at passages 0 and 10.

362    To perform this analysis, we first reannotated the hg38 reference genome assembly

363    using Repeatmasker with a custom library consisting of the consensus sequences for

364    the 8 newly defined LTR7 subfamilies plus newly generated consensus sequences for

365    7b, 7c, and 7y subfamilies redefined from the phylogenetic analysis presented in Fig. 1B

366    (Figure supplement 5) (see methods). Our newly generated Repeatmasker annotations

367    (supp. file 2) did not drastically differ from previous annotations of LTR7 and 7c, where

368    90% and 86% of insertions, respectively, were concordant with the old Repeatmasker

369    annotations (though LTR7 insertions were now assigned to one of the 8 newly defined

370    subfamilies). 7y and 7b annotations, however, shifted significantly. Only 33% of

371    previously annotated 7y reannotated concordantly with 53% now being annotated as

372    7u2 and only 52% of 7b reannotated concordantly, with 22% now annotated as 7y.

18

373    These shifts can be largely explained by the fact that 7u2 and 7y are closely related

374    (Fig. 1A-C) and 7y and 7b share a great deal of sequence through recombination

375    events (Fig. 4B-C).

376    Next we used the newly generated Repeatmasker annotations to examine the RNA

377    expression profiles of the different LTR7 subfamilies using scRNA-seq data from human

378    pre-implantation embryos and RNA-seq data from human ESCs (Blakeley et al., 2015;

379    Tang et al., 2010) (see methods).

380    As expected, we found that the 7up subfamilies were highly expressed in the pluripotent

381    epiblast and in ESCs (Fig. 5). 7up expression was highly specific to these pluripotent

382    cell types, with little to no transcription at earlier developmental time points. As

383    previously observed (Göke et al., 2015), the 7b subfamily exhibited expression at the 8-

384    cell and morula stages, coinciding with EGA (Fig. 5). Another remarkable expression

385    pattern was that of 7u2 which was restricted to the pluripotent epiblast (Fig. 5).

386    Interestingly, the 7y subfamily combined the expression of 7b and 7u2 (8-cell and

387    morula plus epiblast), perhaps reflecting the acquisition of sequence blocks from both

388    subfamilies (Fig. 4B-C). Despite very similar sequence and age (Fig. 1, Fig. 2, Fig. 4A),

389    7bc and 7o elements show stark contrast in their expression profiles. 7o elements show

390    no significant transcription at any time point in early development, while 7bc elements

391    display RNA expression throughout the blastocyst, including trophectoderm and inner

392    cell mass, primitive endoderm, and pluripotent epiblast (Fig. 5). Previous expression

393    analysis of the oldest LTR7 subfamily, 7c, did not find robust stage-specific expression

394    (Göke et al., 2015). Our analysis revealed that some 7c elements display moderate

395    RNA expression at various developmental stages (Fig. 5). This pattern may reflect the

396    relatively high level of sequence heterogeneity within this subfamily (Fig. 1).

397    In summary, our analysis indicates that LTR7 subfamilies have distinct but partially

398    overlapping expression profiles during human early embryonic development that appear

399    to mirror their complex history of sequence diversification.

400

**A predicted SOX2/3 motif unique to 7up is required for transcriptional activity in pluripotent stem cells**

We hypothesized that differences in embryonic transcription among LTR7 subfamilies were driven by the gain and loss of TF binding motifs, and that one or more of these mutations led to 7up's pluripotent-specific transcription. To find TF motifs enriched within each LTR7 subfamily relative to the others, we performed an unbiased motif enrichment analysis using the program HOMER to calculate enrichment scores of known TF motifs within each segmental block defined in Fig. 4A in a pairwise comparison of each subfamily against each of the other subfamilies (see methods). The results yielded a slew of TF motifs enriched for each subfamily relative to the others (see Fig. 6A for 7up1 and enrichment for all HERVH subfamilies in supp. files 3,4). These results suggested that each LTR7 subfamily possesses a unique repertoire of TF binding motifs, which could explain their differential expression during embryonic development.

Next, we sought to pinpoint mutational events responsible for the gain of TF motifs responsible for the unique expression of 7up in ESC. The single most striking motif distinguishing the 7up clade from the others was a SOX2/3 motif which coincided with an 8-bp insertion in block 2b (Fig. 6A,B). Note this motif (and insertion) was also present in 7u1, the closest relative to 7up (Fig. 4C), but absent in all other subfamilies (Fig. 6B).

We hypothesized that the 8-bp insertion provided a binding motif for SOX2 and/or SOX3 contributing to 7up promoter activity in ESCs. Indeed, SOX2 and SOX3 bind a highly similar motif (Bergsland et al., 2011; Heinz et al., 2010), activate an overlapping set of genes and play a redundant function in pluripotency (Corsinotti et al., 2017; Niwa et al., 2016; Wang et al., 2012). In addition, we observed that both SOX2 and SOX3 are expressed in human ESCs but SOX3 was more highly and more specifically expressed in ESCs (Figure supplement 6A,C). While SOX3 binding has not been profiled in human ESCs, ChIP-seq data available for SOX2 indicated that it binds preferentially 7up in a region coinciding with the 8-bp motif (Fig. 6B). Together these observations suggest that 7up promoter activity in ESCs might be conferred in part by the gain of a SOX2/3 motif in block 2b.
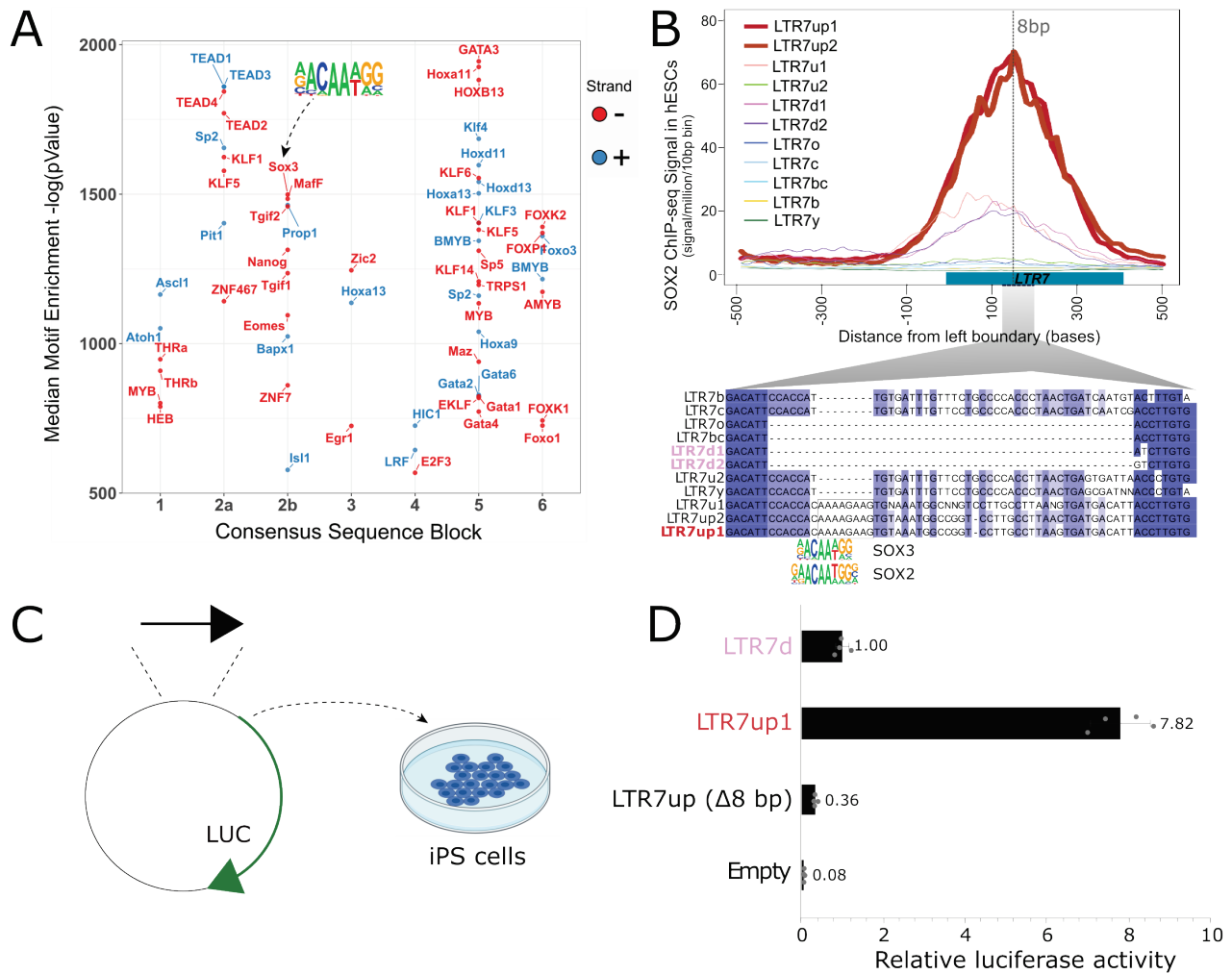
Fig. 6: An 8-bp insertion, SOX2/3 binding site necessary for LTR7up transcription. A) (log) p-values >500 for HOMER motifs enriched in 7up1 insertion's sequence blocks vs the same blocks from other insertions from other HERVH subfamilies are shown. B) Line plots show SOX2 ChIP-seq signal at LTR7 subfamily loci in human ESCs. Signal from genomic loci was compiled relative to position 0. The 7up/u1 8bp insertion position is shown with a dotted line. Region 2b harboring SOX2/3 TFBS is detailed below. C) Scheme of DNA fragments cloned into pGL3-basic vector driving luciferase gene expression (LUC) with identified SOX2/3 motifs. 3 constructs were analyzed: Entire LTR7up (7up1), 7d1/2 consensus sequence (approximate ancestral sequence for all LTR7d) and LTR7up with 8 nucleotides deleted (LTR7up (Δ8bp - AAAAGAAG)) (see panel B). D) Normalized relative luciferase activity of tested fragments compared to LTR7 down; n = 4 measurements; bars, means across replicates; error bars, standard deviation of the mean, dots, individual replicates.

431    To experimentally test this prediction, we used a luciferase reporter to assay promoter

432    activity of three different LTR7 sequences in iPSCs (see methods). The first consisted

433    of the full-length 7d consensus sequence (predicted to be inactive in iPSCs), the second

434    contained the full-length 7up1 consensus (predicted to be active) and the third used the

435    same 7up1 consensus sequence but lacking the 8-bp motif unique to 7up1/2 and 7u1

436    elements overlapping the SOX2/3 motif (Fig. 6B,C).  The results of the assays revealed

21

437    that the 7d construct exhibited, as predicted, only weak promoter activity in iPSC

438    compared to the empty vector (Fig. 6D), while the 7up1 construct had much stronger

439    promoter activity, driving on average 7.8-fold more luciferase expression than 7d and

440    100-fold more than the empty vector (Fig. 6D). Strikingly, the promoter activity was

441    essentially abolished in the 7up1 construct lacking the 8-bp motif, which drove minimal

442    luciferase expression (on average, 3-fold less than LTR7d and 20-fold less than the

443    intact LTR7up sequence). These results demonstrate that the 8-bp motif in 7up1 is

444    necessary for robust promoter activity in iPSCs, likely by providing a SOX2/3 binding

445    site essential for this activity.

## Discussion

The HERVH family has been the subject of intense investigation for its transcriptional and regulatory activities in human pluripotent stem cells. These studies often have treated the entire family as one homogenous, monophyletic entity and it has remained generally unclear which loci are transcribed and potentially important for pluripotency. This is in part because HERVH/LTR7 is an abundant and young family which poses technical challenges to interrogate the activity of individual loci and design experiments targeting specific members of the family (Chuong et al., 2017; Lanciano and Cristofari, 2020). Here, we applied a 'phyloregulatory' approach that integrates regulatory genomics data to a phylogenetic analysis of LTR7 sequences to reveal several new insights into the origin, evolution, and transcription of HERVH elements. In brief, our results show that: (i) LTR7 is a polyphyletic group composed of at least eight monophyletic subfamilies; (ii) these subfamilies have distinct evolutionary histories and transcriptional profiles in human embryos and a single and relatively small subgroup (~264 loci), LTR7up, exhibits robust promoter activity in ESC; (iii) LTR7 evolution is characterized by the gain, loss, and exchange of cis-regulatory modules likely underlying their transcriptional partitioning during early embryonic development.


### Phyloregulatory analysis of LTR7 disentangles the cis-regulatory evolution of HERVH

Previous studies have treated LTR7 *sensu stricto* insertions as equivalent representatives of their subfamilies (Bao et al., 2015; Gemmell et al., 2019; Göke et al., 2015; Izsvák et al., 2016; Storer et al., 2021; Wang et al., 2014; Zhang et al., 2019). While some of these studies were able to detect differential transcriptional partitioning between LTR7, LTR7y, and LTR7b (Göke et al., 2015), the amalgamating of LTR7 loci limited the ability to detect transcriptional variations among LTR7 and to identify key sequence differences responsible for divergent transcription patterns. Our granular parsing of LTR7 elements and their phyloregulatory profiling has revealed striking genetic, regulatory, and evolutionary differences amongst these sequences.

23

475     Importantly, a phylogeny based on the coding sequence (RVT domain) of HERVH

476     provided less granularity to separate the subfamilies than the LTR sequences (Figure

477     supplement 7). The classification of new subfamilies within LTR7 enabled us to discover

478     that they have distinct expression profiles during early embryonic development (Fig. 5)

479     that were previously obscured by their aggregation into a single group of elements. For

480     example, the 7u2 subfamily is, to our knowledge, the first subfamily of human TEs

481     reported to have preimplantation expression exclusively in the epiblast.

482     It has been observed for some time that only a small subset of HERVH elements are

483     expressed in ESCs (Gemmell et al., 2019; Göke et al., 2015; Ohnuki et al., 2014;

484     Santoni et al., 2012; Schön et al., 2001; Wang et al., 2014; Zhang et al., 2019). Some

485     have attributed this property to variation in the internal region of HERVH, context-

486     dependent effects (local chromatin or cis-regulatory environment) and/or age (Gemmell

487     et al., 2019; Zhang et al., 2019). Our results provide an additional, perhaps simpler

488     explanation: we found that HERVH elements expressed in ESCs are almost exclusively

489     driven by two closely related subfamilies of LTR7 (7up) that emerged most recently in

490     hominoid evolution. We identified one 8-bp sequence motif overlaps a predicted

491     SOX2/3 binding site unique to 7up that is required for promoter activity in pluripotent

492     stem cells. These results highlight that the primary sequence of the LTR plays an

493     important role in differentiating and diversifying HERVH expression during human

494     embryonic development.

495     The phyloregulatory approach outlined in this study could be applied to illuminate the

496     regulatory activities of LTR elements in other cellular contexts. In addition to

497     embryogenesis, subsets of LTR7 and LTR7y elements are known to be upregulated in

498     oncogenic states (Babaian and Mager, 2016; Glinsky, 2015; Kong et al., 2019; Yu et al.,

499     2013). It would be interesting to explore whether these activities can be linked to the

500     gain of specific TFBS using the new LTR7 annotations and regulatory information

501     presented herein. Other human LTR families, such as MER41, LTR12C, or LTR13 have

502     been previously identified as enriched for particular TF binding and cis-regulatory

503     activities in specific cellular contexts (Chuong et al., 2016; Deniz et al., 2020; Ito et al.,

504     2017; Krönung et al., 2016; Sundaram et al., 2014). In each case, TF binding

505  enrichment was driven by a relatively small subset of loci within each family. We

506  suspect that some of the intrafamilial differences in TF binding and cis-regulatory

507  activity may be caused by unrecognized subfamily structure and subfamily-specific

508  combinations of TFBS, much like we observe for LTR7.

509

**Recombination as a driver of LTR cis-regulatory evolution**

510

511  Recombination is a common and important force in the evolution of exogenous RNA

512  viruses (Jetzt et al., 2000; Pérez-Losada et al., 2015; Simon-Loriere and Holmes, 2011)

513  and endogenous retroviruses (Vargiu et al., 2016). Traditional models of recombination

514  describe recombination occurring due to template switching during reverse transcription,

515  a process that requires the co-packaging of RNA genomes, a feature of retroviruses

516  and some retrotransposons (Lai, 1992; Matsuda and Garfinkel, 2009). Previous studies

517  proposed that the HERVH family had undergone inter-element recombination events of

518  both its coding genes (Mager and Freeman, 1987; Vargiu et al., 2016) and LTR

519  (Goodchild et al., 1993). Specifically, it was inferred that recombination event between

520  Type I LTR (i.e., LTR7) and Type II LTR (LTR7b) led to the emergence of Type Ia

521  (LTR7y).

522  Our findings of extensive sequence block exchange between 7y and 7b (Fig. 4D) are

523  consistent with these inferences. Furthermore, our division of HERVH into at least 11

524  subfamilies, rather than the original trio (Type I, II, Ia), and systematic analysis of

525  recombination events (Fig. 4) suggest that recombination has occurred between

526  multiple lineages of elements and has been a pervasive force underlying LTR

527  diversification. We identified a minimum of six recombination events spanning 20 million

528  years of primate evolution (see Fig. 4D and summary model in Fig. 7). The coincidence

529  of recombination events with changes in expression profiles (Fig. 7) suggests that these

530  events were instrumental to the diversification of HERVH embryonic expression. The

531  hybrid origin and subsequent burst of amplification of LTR7 subfamilies (Fig. 1,2)

532  suggest they expanded rapidly after shifting their transcriptional profiles. The

533  coincidence of niche colonization with a burst in transposition leads us to speculate that

534  these shifts in expression were foundational to the formation and successful expansion

535    of new HERVH subfamilies. It would be interesting to explore whether inter-element

536    recombination has also contributed to the evolution of other LTR subfamilies and the

537    diversification of their expression patterns.

538    Previous work has highlighted the role of TEs, and LTRs in particular, in donating built-

539    in cis-regulatory sequences promoting the evolutionary rewiring of mammalian

540    transcriptional networks (Chuong et al., 2017; Feschotte, 2008; Hermant and Torres-

541    Padilla, 2021; Jacques et al., 2013; Rebollo et al., 2012; Sundaram and Wysocka, 2020;

542    Thompson et al., 2016). We show that recombination provides another layer to this

543    idea, where combinations of TFBS can be mixed-and-matched, then mobilized and

544    propagated, further accelerating the diversification of these regulatory DNA elements.

545    As HERVH expanded and diversified, its newly evolved cis-regulatory modules became

546    confined to specific host lineages (Fig. 2). Thus, it is possible that the formation of new



Fig. 7: Model of LTR7 subfamily evolution. Estimated LTR7 subfamily transpositional activity in mya are listed with corresponding approximate primate divergence times (bottom). The positioning and duration of transpositional activity are based on analysis from Fig. 3b. The grey connections between subfamilies indicate average tree topology which is driven by overall pairwise sequence similarity. Dashed lines indicate likely recombination events which led to the founding of new subfamilies. Stage-specific expression profiles from Fig. 5a are detailed to the right of each corresponding branch.

547     LTR via recombination and their subsequent amplification catalyzed cis-regulatory

548     divergence across primate species.

549

**LTR evolution enabled HERVH's colonization of different niches in the human**

551     **embryo**

552     Our evolutionary analysis reveals that multiple HERVH subfamilies were

553     transpositionally active in parallel during the past ~25 my of primate evolution (Fig. 2,7).

554     This is in stark contrast to the pattern of LINE1 evolution in primates, which is

555     characterized by a single subfamily being predominantly active at any given time (Khan

556     et al., 2006). We hypothesize that the ability of HERVH to colonize multiple cellular

557     niches underlie this difference. Indeed, we observe that concurrently active HERVH

558     subfamilies are transcribed at different developmental stages, such as 7up and 7u2

559     being transcribed in the pluripotent epiblast at the same time that 7y and the youngest

560     7b were transcribed at the 8 cell and morula stages (Fig. 7). We posit that this

561     partitioning allowed multiple HERVH subfamilies to amplify in parallel without causing

562     overt genome instability and cell death during embryonic development.

563     Niche diversification may have also enabled HERVH to evade cell-type-specific

564     repression by host-encoded factors such as KZFPs. KZFPs are thought to emerge and

565     adapt during evolution to silence specific TE subfamilies in a cell-type specific manner

566     (Bruno et al., 2019; Cosby et al., 2019; Ecco et al., 2017; Imbeault et al., 2017). For

567     example, there is evidence that the progenitors of the currently active L1HS subfamily

568     became silenced in human ESCs via KZFP targeting, but evaded that repression and

569     persisted in that niche through the deletion of the KZFP binding site (Jacobs et al.,

570     2014). HERVH may have persisted through another evasive strategy: changing their

571     TFBS repertoire to colonize niches lacking their repressors. To silence all LTR7, any

572     potential HERVH-targeting KZFP would need to gain expression in multiple cellular

573     contexts. For example, one potential repressor, ZNF534, binds a wide range of LTR7

574     sequences, but is particularly enriched at 7up in HEK293 cells (Fig. 3A,D). Our analysis

575     shows that ZNF534 is most highly expressed in the morula, but dips in human ESC

576     (Figure supplement 3). Thus, ZNF534 may repress 7up at earlier stages of development

577 but is apparently unable to suppress 7up transcription in pluripotent stem cells. If true,

578 this scenario would illustrate how LTR diversification facilitated HERVH persistence in

579 the face of KZFP coevolution. Further investigation is needed to explore the interplay

580 between KZFPs and HERVH subfamilies during primate evolution.

581

582 **Implications for stem cell and regenerative biology**

583 Lastly, our findings may provide new opportunities for stem cell research and

584 regenerative medicine. Our data on 7up reinforces previous findings (Corsinotti et al.,

585 2017; Wang et al., 2012) that place SOX2/3 as central players in pluripotency.

586 Furthermore, our analysis identified a set of TFs whose motifs are uniquely enriched in

587 different LTR7 subfamilies with distinct expression patterns in early embryonic cells,

588 which may enable a functional discriminatory analysis of the role of these TFs in each

589 cell type. HERVH/LTR7 has been used as a marker for human pluripotency (Ohnuki et

590 al., 2014; Santoni et al., 2012; Wang et al., 2014), and recent work has revealed that

591 HERVH/LTR7-positive cells may be more amenable to differentiation, and are therefore

592 referred to as "primed" cells (Göke et al., 2015; Theunissen et al., 2016). However,

593 primed cells are not as promising for regenerative medicine as so-called "naïve" cells

594 (Nichols and Smith, 2009), which are less differentiated and resemble cells from late

595 morula to epiblast, or so-called "formative" cells, which most closely resemble cells from

596 the early post-implantation epiblast (Kalkan and Smith, 2014; Kinoshita et al., 2021;

597 Rossant and Tam, 2017). Of relevance to this issue is our finding that elements of the

598 7u2 subfamily are highly and exclusively expressed in the pluripotent epiblast in vivo

599 (Fig. 5), but weakly so in H1 ESC, which consists of a majority of primed cells and a

600 minority of naïve or formative cells (Gafni et al., 2013). Thus, it might be possible to

601 develop a LTR7u2-driven reporter system to mark and purify naïve or formative cells

602 from an heterogenous ESC population. Similarly, a MERVL LTR-GFP transgene has

603 been used in mouse to purify rare 2-cell-like totipotent cells where this LTR is

604 specifically expressed amidst mouse ESCs in culture (Hermant and Torres-Padilla,

605 2021; Macfarlan et al., 2012).

606     In conclusion, our study highlights the modular cis-regulatory evolution of an

607     endogenous retrovirus which has facilitated its transcriptional partitioning in early

608     embryogenesis. We believe that phyloregulatory dissection of endogenous retroviral

609     LTRs has the potential to further our understanding of the evolution, impact, and

610     applications of these elements in a broad range of biomedical areas.

## **Acknowledgements:**

## **Methods**

622

### HERVH LTR sequence identification:

623

624 All HERVH-int and accompanying LTRs (LTR7, 7b, 7c, and 7y) were extracted from

625 masked (RepeatMasker version 4.0.5 repeat Library 20140131 - (Smit et al., 2013))

626 GRCh38/hg38 (alt chromosomes removed). All annotated HERVH-int and HERVH LTR

627 were run through OneCodeToFindThemAll.pl (Bailly-Bechet et al., 2014) followed by

628 rename_mergedLTRelements.pl (Thomas et al., 2018) to identify solo and full-length

629 HERVH insertions. 5' LTRs from full-length insertions >4kb were combined with and

630 solo LTRs. LTRs >350bp were considered for future analysis.

### Multiple sequence alignment, phylogenetic tree generation, and LTR7 subdivision:

631

632 All HERVH LTRs (Fig. 1A – supp. file 5) or only LTR7s (Fig. 1B – supp. file 6) were

633 aligned with mafft –auto (Nakamura et al., 2018) strategy: FFT-NS-2/Progressive

634 method followed by PRANK (Löytynoja and Goldman, 2010) with options -showanc -

635 support -njtree -uselogs -prunetree -prunedata -F -showevents. Uninformative structural

636 variations were removed with Trimal (Capella-Gutierrez et al., 2009) with option -gt

637 0.01.

638 To visualize inter-insertion relationships, the MSA was input into IQtree with options -nt

639 AUTO -m MFP -bb 6000 -asr -minsup .95 (Chernomor et al., 2016). This only displays

640 nodes with ultrafast (UF) bootstrap support >0.95.

641 Clusters of >10 insertions sharing a node with UFbootstrap support that were separated

642 from other insertions by internal branch lengths >0.015 (1.5subs / 100 bp) were defined

643 as belonging to a new bona fide LTR7 subfamily (Fig. 1B).

### LTR7 consensus generation and network analysis:

644

645 Majority rule (51%) was used to generate each LTR7 subfamily at nodes described in

646 Fig. 1. Positions without majority consensus are listed as "N". Majority rule consensus

647 sequences were aligned with MUSCLE in SEAVIEW (supp. file 7) (Edgar, 2004; Gouy

648 et al., 2010). Alignment was visualized with Jalview2 (Waterhouse et al., 2009)(Fig. 4A)

649 and ggplot2 (Fig. 4).

31

650    Non-gap SNPs from the muscle alignment were used to construct a median-joining

651    network (Bandelt et al., 1999) with POPART (Leigh and Bryant, 2015).

652    <u>Reverse Transcriptase Domain extraction, alignment, and tree generation:</u>

653    The reverse transcriptase (RT) domain was extracted from HERVH-int consensus via

654    repbrowser (Fernandes et al., 2020):

655    CACCCTTACCCCGCTCAATGCCAATATCCCATCCCACAGCATGCTTTAAAAGGATT

656    AAAGCCTGTTATCACTCGCCTGCTACAGCATGGCCTTTTAAAGCCTATAAACTCTCC

657    TTACAATTCCCCCATTTTACCTGTCCTAAAACCAGACAAGCCTTACAAGTTAGTTCA

658    GGATCTGTGCCTTATCAACCAAATTGTTTTGCCTATCCACCCCATGGTGCCAAACC

659    CATATACTCTCCTATCCTCAATACCTCCCTCCACAACCCATTATTCTGTTCTGGATC

660    TCAAACATGCTTTCTTTACTATTCCTTTGCACCCTTCATCCCAGCCTCTCTTCGCTTT

661    CACTTGGA

662    This sequence was blated (best hit) against all annotated HERVH-int in the human

663    genome and matches were extracted. Corresponding LTR7 subdivision annotations

664    from figure 1 were matched with these HERVH-int RT domains. Mafft alignment and

665    IQTree generation were done identically to the Mafft and IQTree run for the LTRs (see

666    corresponding methods section).

667    <u>Peak calling:</u>

668    ChIP-seq datasets representing transcription factors (TFs), histone modifications, and

669    regulatory complexes in human embryonic stem cells and differentiated cells were

670    retrieved from GSE61475 (38 distinct TFs and histone modifications), GSE69647

671    (H3K27Ac, POU5F1, MED1 and CTCF), GSE117395 (H3K27Ac, H3K9Me3, KLF4, and

672    KLF17), and GSE78099 (An array of KRAB-ZNFs and TRIM28) (Imbeault et al., 2017).

673    ZNFs enriched in LTR7 binding (ZNF90, ZNF534, ZNF75, ZNF69B, ZNF257, ZNF57,

674    and ZNF101) from HEK293 peaks were all evaluated, but only ZNF90 and ZNF534

675    bound >100 LTR7 insertions (data not shown). The others were dropped from the

676    analysis.

677    ChIP-seq reads were aligned to the hg19 human reference genome using the Bowtie2.

678    All reads with phred score less than 33 and PCR duplicates were removed using

679    bowtie2 and Picard tools respectively. ChIP-seq peaks were called by MACS2 with the

680    parameters in "narrow" mode for TFs and "broad" mode for histone modifications,

681    keeping FDR < 1%. ENCODE-defined blacklisted regions were excluded from called

682    peaks. For phyloregulatory analysis (Fig. 2), we then converted hg19 to hg38 (no alt)

683    coordinates via UCSC *liftover* (100% of coordinates lifted) and intersected these peak

684    with the loci from LTR7 subfamilies using bedtools with any overlap. For ChIP-seq

685    binding enrichment on a subset of marks following motif analysis (Fig 5), 70% overlap of

686    peak and LTR was required. Enrichment of a given TF within LTR7 subfamilies was

687    calculated using enrichR package in R, using the customized in-house codes (see the

688    codes on GitHub for the detailed analysis pipelines and calculation of enrichment

689    score).

690    Phyloregulatory analysis:

691    Peaks from external ChIP-seq datasets were intersected with LTR7 insertions (Quinlan

692    and Hall, 2010). LTR7 insertions that intersected with >1bp of peaks were counted as

693    positive for the respective mark. We repeated this analysis with a range of overlap

694    requirements from extending the LTR 500bp into unique DNA to 70% overlap and found

695    few differential calls (data not shown). The phylogenetic tree rooted on 7b (ggtree) was

696    combined with these binary data (ggheat).

697    "Highly transcribed" (fpkm >2) and "chimeric" HERVH from H1 cells (GSE54726) (Wang

698    et al., 2014) were intersected with LTR7 similarly to ChIP-seq data. Those which

699    intersected LTR7 were marked as "RNA-seq" or "chimeric" respectively. GRO-seq

700    profiles from H1 cells (Estaras et al.) (GSE64758) were created for windows 10bp

701    upstream and 8kb downstream of 5' and solo LTR7 (Ramírez et al., 2016). The most

702    visible signal was confined to the top 7th of insertions (Figure supplement 2). All LTR7

703    were subdivided into septiles, due to visible signal being confined to the top 7th of

704    insertions; those of the top septile were labeled "GRO-seq".

705    Peak proportion heatmap generation and statistical analysis:

706   Tables with the proportion of solo and 5' LTRs from a given subfamily positive for select

707   marks (phyloregulatory analysis) were used to generate heatmaps with the R package

708   ggplot (ggheat) (Ginestet, 2011). Those with padj<0.05 (Chi-square Bonferroni

709   correction n=147 tests for a total of 21 marks examined) were considered significantly

710   enriched in 7up1. Enrichment for non-LTR7up subfamilies was not tested. While not all

711   tested marks are displayed in the main text, statistical analysis was performed with all

712   tested marks (n=147) (supp. file 8). For comparing transcribed 7up to untranscribed

713   7up, 18 pairwise comparisons were made (supp. file 9).

714   Aggregate signal heatmap generation:

715   GRO-seq (H1 cells - GSE64758), whole-genome bisulfite sequencing (WGBS-seq – H1

716   cells), and H3K9me3 ChIP-seq (H1 – primed - GSE78099) bams were retrieved from

717   (Estarás et al., 2015), (Dunham et al., 2012), and (Theunissen et al., 2016) respectively.

718   Deeptools (Ramírez et al., 2016) was used to visualize these marks by LTR7 subfamily

719   division in windows 10bp upstream and 8kb downstream of the most 5' position in the

720   LTR (Figure supplement 2).

721   Orthologous insertion aging:

722   Human coordinates for 7b, 7c, and 7y and LTR7 used in alignments and tree generation

723   were lifted over (Kent et al., 2002; Raney et al., 2014) from GRCh38/hg38 (Miga et al.,

724   2014) to Clint_PTRv2/panTro6 (Waterson et al., 2005), Kamilah_GGO_v0/gorGor6

725   (Scally et al., 2012), Susie_PABv2/ponAbe3 (Locke et al., 2011), GGSC

726   Nleu3.0/nomLeu3 (Carbone et al., 2014), or Mmul_10/rheMac10 (Gibbs et al., 2007).

727   Those that were successfully lifted over from human to non-human primate were then

728   lifted over back to human. Only those that survived both liftovers (1:1 orthologous) were

729   counted as present in non-human primates. The proportion of those orthologous to

730   human and total number of orthologous was plotted with ggplot2.

731   Terminal branch length aging:

732   Terminal branch lengths from the LTR7 phylogenetic tree (Fig. 1B) were extracted and

733   plotted with ggplot2. Similarly aged subfamilies were inferred from means here and from

734   orthologous insertion aging for statistical testing. Three total groups were tested for

34

735     differences in means (7up1/7up2/7u2 vs. 7d1/7d2/7u1 vs. 7bc/o) via Wilcox rank-sum

736     test with Bonferroni multiple testing correction.

737     <u>Identification of recombination breakpoints and consensus parsimony tree generation:</u>

738     Major recombination breakpoints were identified by eye from the consensus sequence

739     MSA, where SNPs and structural rearrangements seemed to have different

740     relationships between blocks. Putative block recombination events were identified by

741     looking for shared shapes in the block consensus MSA (Fig. 4A). To test if these were

742     truly recombination events and could not be explained by evolution by common

743     descent, inter-block sequence relationship differences were tested by generating

744     parsimony trees and comparing these to the overall phylogenetic structure from Fig. 1A.

745     Parsimony trees were generated in SEAVIEW, treating all gaps as unknown states

746     (except in the case of 2b, where the entire sequence is gaps and gaps were not treated

747     differently than other sequence), bootstrapped 5000 times with the option "more

748     thorough tree search". Differences in block parsimony trees and the overall phylogeny

749     that had bootstrap support were marked in red and included in Fig. 4D,7.

750     <u>7up consensus block 2a 2b alignment and parsimony tree:</u>

751     LTR7up blocks 2a and 2b (Fig. 4) appeared to share sequence. To determine if block

752     2b was the result of a duplication of 2a, we extracted these sequences from the

753     LTR7up1 consensus and aligned them with blastn (NCBI web version) with default

754     settings. To determine the relationship of all HERVH LTR 2a and 2b blocks, we

755     performed a muscle alignment (default settings) of all 2a and 2b from all HERVH LTR

756     consensus sequences and then generated a parsimony tree with 5000 bootstraps with

757     SEAVIEW with the option "more thorough tree search".

758     <u>New LTR7B/C/Y consensus generation and remasking of human genome:</u>

759     Consensus sequences for LTR7 subfamilies were generated using the tree from figure

760     1b (see above). For LTR7b/c/y, we used the alignment and tree comprising all HERVH

761     LTR (Figure supplement 5). To do this, we identified nodes with >0.95 ultrafast

762     bootstrap support that were comprised of predominately (>90%) of previously annotated

763     LTR7b, LTR7c, or LTR7y. These sequences were used to generate majority-rule

764 consensus sequences for their respective subfamily. We generated 2 mutually-

765 exclusive LTR7c consensus sequences (LTR7C1 and LTR7C2) due to the high

766 sequence divergence of LTR7C. Both of these subfamilies were merged into "LTR7C"

767 after remasking.

768 Parsing previously annotated LTR7 into 8 subfamilies and evidence of recurrent

769 recombination events caused concern that HERVH LTRs may be misannotated in the

770 repeat masker annotations. To compensate, we remasked (Smit et al., 2013)

771 GRCh38/hg38 (excluding alt chromosomes) with a custom library consisting of the new

772 consensus sequences for LTR7 subfamilies, new consensus sequences for 7b, 7y, and

773 7c (see above) based on the HERVH LTR tree from Fig. 4, and HERVH-int (dfam). We

774 also included annotated consensus sequences from dfam for MER48, MER39, AluYk3,

775 and MST1N2, who we found a HERVH only library also masked to a limited degree

776 (data not shown). With this library, we ran RepeatMasker with crossmatch and

777 "sensitive" settings: -e crossmatch -a -s -no_is. Changes in annotations can be found in

778 (HERVH_LTRremasking.xlsx)

779 Embryonic HERVH subfamily expression analysis:

780 We downloaded the raw single-cell RNA-seq datasets from early human embryos and

781 embryonic stem cells (GSE36552) and the EPI, PE, TE cells (GSE66507) in sra format.

782 Following the conversion of raw files into fastq format, the quality was determined by

783 using the FastQC. We removed two nucleotides from the ends as their quality scores

784 were highly variable compared with the rest of the sequences in RNA-seq reads. Prior

785 to aligning the resulting reads, we first curated the reference genome annotations using

786 the LTR7 classification, shown in the manuscript. We extracted the genes (genecode

787 V19), and LTR7 subfamilies (see figure 5) genomic sequences and combined them to

788 generate a reference transcriptome. These sequences were then appended, comprising

789 the coding-sequences plus UTRs of genes and locus-level LTR7 subfamilies sequences

790 in fasta format. We then annotated every fasta sequences with their respective genes or

791 LTR7 subfamilies IDs. To guide the transcriptome assembly, we also appended the

792 each of the resulting contigs and modelled them in gtf format that we utilized for the

793 expression quantification. Next, we indexed the concatenated genes and LTR7

794    subfamilies transcriptome and genome reference sequences using 'salmon' (Patro et

795    al., 2017). Finally, we aligned the trimmed sequencing reads against the curated

796    reference genome. The 'salmon' tool quantified the counts and normalized expression

797    (Transcripts per million (TPM)) for each single cell RNAseq sample. Overall, this

798    approach enabled us to simultaneously calculate LTR7 subfamilies and protein-coding

799    gene expression using expected maximization algorithms. Data integration of obtained

800    count matrix, normalization at logarithmic scale, and scaling were performed as per the

801    "Seurat V.3.7" (http://satijalab.org/seurat/) guidelines. The annotations of cell-types were

802    taken as it was classified in original studies. We calculated differential expression and

803    tested their significance level using Kruskal–Wallis test by comparing cell-types of

804    interest with the rest of the cells. The obtained p-values were further adjusted by the

805    Benjamini-Hochberg method to calculate the False Discovery Rate (FDR). All the

806    statistics and visualization of RNA-seq were performed on R (https://www.r-project.org/).

807    Motif Enrichment:

808    For each subfamily of LTR7 elements, all re-annotated elements were aligned against

809    the subfamily consensus sequence using MUSCLE (Edgar, 2004). These multiple-

810    sequence alignments were then split based on the recombination block positions in the

811    consensus sequence. The consensus sequence was then removed. Binding motif

812    position-weight matricies were downloaded from HOMER (Heinz et al., 2010) and were

813    used to perform pairwise motif enrichment using the command 'homer2 find'. For

814    LTR7up1 enrichment (Fig. 6A - testing which motifs were enriched in LTR7up1

815    compared to other subfamilies), enrichment was only calculated for LTR7up1 and the

816    motifs with a -log(p-value) cutoff of 1x10-5 were kept. For enrichment in all subfamilies

817    (supp. files 3,4) – testing all subfamilies against all others), every pairwise subfamily

818    combination within each block was tested and all results are displayed.

819

820    SOX2 ChIP-seq signal on LTR7:

821    SOX2 ChIP-seq and whole-cell extract datasets from primed hESCs were downloaded

822    in fastq format from GEO ID GSE125553 (Bayerl et al., 2021). Fastq reads were

823    mapped against the hg19 reference genome with the bowtie2 parameters: –*very-*

824     *sensitive-local*. All unmapped reads with Phred score < 33 and putative PCR duplicates

825     were removed using *Picard* and *samtools*. All the ChIP-seq narrow peaks were called

826     by MACS2 (FDR <  0.01). To generate a set of unique peaks, we merged ChIP-seq

827     peaks within 50 bp of one another using the *mergeBed* function from bedtools. We then

828     intersected these peak sets with LTR7 subgroups from hg19 repeat-masked

829     coordinates using bedtools *intersectBed* with 50% overlap. LTR7up1 and LTR7up2

830     were harboring the highest number of peaks compared with the rest of the subgroups.

831     To illustrate the enrichment over the LTR7 subgroups, we first extended 500 basepairs

832     from upstream and downstream coordinates from the left boundary of each

833     LTR7subgroups. These 1KB windows were further divided into 10 bps bins. The

834     normalized ChIP-seq signal over the local lambda (piled up bedGraph outputs from

835     MACS2) was counted in each bin. These counts were then normalized by the total

836     number of mappable reads per million in given samples and presented as signal per

837     million per 10 bps. Finally, these values were averaged across the loci for each bin to

838     illustrate the subfamilies' level of ChIP-seq enrichment. Replicates were merged prior to

839     plotting. Note: Pearson's correlation coefficient between replicates across the bins was

840     found to be r > 0.90.

841

842     Luciferase reporter assay:

843     The inserts (LTR7 variants or EF1a promoter) with restriction enzyme overhangs were

844     ordered from Genewiz and cloned into pGL3-basic plasmid upstream of the firefly

845     reporter gene (E1751, Promega). Minipreps were prepared with QIAprep Spin Miniprep

846     kit (Qiagen). Plasmids were sequenced to ensure the correct sequence and

847     directionality of the insert. 24 h before transfection, human iPSC WTC-11 (Coriell

848     Institute) cells were plated on Vitronectin (Thermo Fisher Scientific) coated 12-well

849     plates in Essential 8 Flex medium (Thermo Fisher Scientific) with E8 supplement

850     (Thermo Fisher Scientific), Rock inhibitor and 2.5% penicillin-streptomycin. Cells were

851     co-transfected with 800 ng of plasmid of interest and 150 ng plasmid containing EF1a

852     upstream of GFP for normalization with Lipofectamine Stem transfection reagent

853     (Thermo Fisher scientific) according to manufacturer's instructions. 48 h after

854     transfection, cell pellet was harvested and luciferase activity was measured with

855     Luciferase Reporter Assay kit (Promega) on Glomax (Promega) according to

856     instructions. Transfection efficiency and cell count was normalized with GFP.

857     1. 7down:

858     GCTAGCTGTCAGGCCTCTGAGCCCAAGCTAAGCCATCATATCCCCTGTGACCTGC

859     ACGTACACATCCAGATGGCCGGTTCCTGCCTTAACTGATGACATTCCACCACAAAA

860     GAAGTGAAAATGGCCTGTTCCTGCCTTAACTGATGACATTATCTTGTGAAATTCCTT

861     CTCCTGGCTCATCCTGGCTCAAAAGCTCCCCTACTGAGCACCTTGTGACCCCCACT

862     CCTGCCCGCCAGAGAACAACCCCCCTTTGACTGTAATTTTCCTTTACCTACCCAAA

863     TCCTATAAAACGGCCCCACCCCTATCTCCCTTCGCTGACTCTCTTTTCGGACTCAG

864     CCCGCCTGCACCCAGGTGAAATAAACAGCTTTATTGCTCACACAAAGCCTGTTTGG

865     TGGTCTCTTCACACGGACGCGCATGCTCGAG

866     2. LTR7upcons:

867     GCTAGCTGTCAGGCCTCTGAGCCCAAGCCAAGCCATCGCATCCCCTGTGACTTGC

868     ACGTATACGCCCAGATGGCCTGAAGTAACTGAAGAATCACAAAAGAAGTGAATATG

869     CCCTGCCCCACCTTAACTGATGACATTCCACCACAAAAGAAGTGTAAATGGCCGGT

870     CCTTGCCTTAAGTGATGACATTACCTTGTGAAAGTCCTTTTCCTGGCTCATCCTGGC

871     TCAAAAAGCACCCCCACTGAGCACCTTGCGACCCCCACTCCTGCCCGCCAGAGAA

872     CAAACCCCCTTTGACTGTAATTTTCCTTTACCTACCCAAATCCTATAAAACGGCCCC

873     ACCCTTATCTCCCTTCGCTGACTCTCTTTTCGGACTCAGCCCGCCTGCACCCAGGT

874     GAAATAAACAGCCATGTTGCTCACACAAAGCCTGTTTGGTGGTCTCTTCACACGGA

875     CGCGCATGCTCGAG

876     5. LTR7upcons_AAAGAAG_deletion:

877     GCTAGCTGTCAGGCCTCTGAGCCCAAGCCAAGCCATCGCATCCCCTGTGACTTGC

878     ACGTATACGCCCAGATGGCCTGAAGTAACTGAAGAATCACAAAAGAAGTGAATATG

879     CCCTGCCCCACCTTAACTGATGACATTCCACCATTGTAAATGGCCGGTCCTTGCCT

880     TAAGTGATGACATTACCTTGTGAAAGTCCTTTTCCTGGCTCATCCTGGCTCAAAAG

881    CACCCCCACTGAGCACCTTGCGACCCCCACTCCTGCCCGCCAGAGAACAAACCCC

882    CTTTGACTGTAATTTTCCTTTACCTACCCAAATCCTATAAAACGGCCCCACCCTTAT

883    CTCCCTTCGCTGACTCTCTTTTCGGACTCAGCCCGCCTGCACCCAGGTGAAATAAA

884    CAGCCATGTTGCTCACACAAAGCCTGTTTGGTGGTCTCTTCACACGGACGCGCAT

885    GCTCGAG

886    5'NheI highlighted in Yellow

887    3'XhoI highlighted in Cyan

888    Babaian A, Mager DL. 2016. Endogenous retroviral promoter exaptation in human cancer. *Mobile DNA*
889    **7**:24. doi:10.1186/s13100-016-0080-x

890    Bailly-Bechet M, Haudry A, Lerat E. 2014. "One code to find them all": a perl tool to conveniently parse
891    RepeatMasker output files. *Mobile DNA* **5**:13. doi:10.1186/1759-8753-5-13

892    Bandelt HJ, Forster P, Röhl A. 1999. Median-joining networks for inferring intraspecific phylogenies.
893    *Molecular Biology and Evolution* **16**:37–48. doi:10.1093/oxfordjournals.molbev.a026036

894    Bannert N, Kurth R. 2004. Retroelements and the human genome: New perspectives on an old relation.
895    *PNAS* **101**:14572–14579. doi:10.1073/pnas.0404838101

896    Bao W, Kojima KK, Kohany O. 2015. Repbase Update, a database of repetitive elements in eukaryotic
897    genomes. *Mobile DNA* **6**:11. doi:10.1186/s13100-015-0041-9

898    Bayerl J, Ayyash M, Shani T, Manor YS, Gafni O, Massarwa R, Kalma Y, Aguilera-Castrejon A, Zerbib M,
899    Amir H, Sheban D, Geula S, Mor N, Weinberger L, Naveh Tassa S, Krupalnik V, Oldak B, Livnat N, Tarazi S,
900    Tawil S, Wildschutz E, Ashouokhi S, Lasman L, Rotter V, Hanna S, Ben-Yosef D, Novershtern N, Viukov S,
901    Hanna JH. 2021. Principles of signaling pathway modulation for enhancing human naive pluripotency
902    induction. *Cell Stem Cell* S1934-5909(21)00158–2. doi:10.1016/j.stem.2021.04.001

903    Bergsland M, Ramsköld D, Zaouter C, Klum S, Sandberg R, Muhr J. 2011. Sequentially acting Sox
904    transcription factors in neural lineage development. *Genes Dev* **25**:2453–2464.
905    doi:10.1101/gad.176008.111

906    Blakeley P, Fogarty NME, del Valle I, Wamaitha SE, Hu TX, Elder K, Snell P, Christie L, Robson P, Niakan
907    KK. 2015. Defining the three cell lineages of the human blastocyst by single-cell RNA-seq. *Development*
908    **142**:3151–3165. doi:10.1242/dev.123547

909    Boyer LA, Lee TI, Cole MF, Johnstone SE, Levine SS, Zucker JP, Guenther MG, Kumar RM, Murray HL,
910    Jenner RG, Gifford DK, Melton DA, Jaenisch R, Young RA. 2005. Core Transcriptional Regulatory Circuitry
911    in Human Embryonic Stem Cells. *Cell* **122**:947–956. doi:10.1016/j.cell.2005.08.020

912    Bruno M, Mahgoub M, Macfarlan TS. 2019. The Arms Race Between KRAB–Zinc Finger Proteins and
913    Endogenous Retroelements and Its Impact on Mammals. *Annual Review of Genetics* **53**:393–416.
914    doi:10.1146/annurev-genet-112618-043717

915    Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. 2009. trimAl: a tool for automated alignment
916    trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**:1972–1973.
917    doi:10.1093/bioinformatics/btp348

918    Carbone L, Alan Harris R, Gnerre S, Veeramah KR, Lorente-Galdos B, Huddleston J, Meyer TJ, Herrero J,
919    Roos C, Aken B, Anaclerio F, Archidiacono N, Baker C, Barrell D, Batzer MA, Beal K, Blancher A, Bohrson
920    CL, Brameier M, Campbell MS, Capozzi O, Casola C, Chiatante G, Cree A, Damert A, de Jong PJ, Dumas L,
921    Fernandez-Callejo M, Flicek P, Fuchs NV, Gut I, Gut M, Hahn MW, Hernandez-Rodriguez J, Hillier LW,
922    Hubley R, Ianc B, Izsvák Z, Jablonski NG, Johnstone LM, Karimpour-Fard A, Konkel MK, Kostka D, Lazar
923    NH, Lee SL, Lewis LR, Liu Y, Locke DP, Mallick S, Mendez FL, Muffato M, Nazareth LV, Nevonen KA,
924    O'Bleness M, Ochis C, Odom DT, Pollard KS, Quilez J, Reich D, Rocchi M, Schumann GG, Searle S, Sikela
925    JM, Skollar G, Smit A, Sonmez K, Hallers B ten, Terhune E, Thomas GWC, Ullmer B, Ventura M, Walker JA,

926 Wall JD, Walter L, Ward MC, Wheelan SJ, Whelan CW, White S, Wilhelm LJ, Woerner AE, Yandell M, Zhu
927 B, Hammer MF, Marques-Bonet T, Eichler EE, Fulton L, Fronick C, Muzny DM, Warren WC, Worley KC,
928 Rogers J, Wilson RK, Gibbs RA. 2014. Gibbon genome and the fast karyotype evolution of small apes.
929 *Nature* **513**:195–201. doi:10.1038/nature13679

930 Chambers I, Smith A. 2004. Self-renewal of teratocarcinoma and embryonic stem cells. *Oncogene*
931 **23**:7150–7160. doi:10.1038/sj.onc.1207930

932 Chang N-C, Rovira Q, Wells JN, Feschotte C, Vaquerizas JM. 2021. A genomic portrait of zebrafish
933 transposable elements and their spatiotemporal embryonic expression. *bioRxiv* 2021.04.08.439009.
934 doi:10.1101/2021.04.08.439009

935 Charlesworth B, Langley CH. 1986. THE EVOLUTION OF SELF-REGULATED TRANSPOSITION OF
936 TRANSPOSABLE ELEMENTS. *Genetics* **112**:359–383.

937 Chernomor O, von Haeseler A, Minh BQ. 2016. Terrace Aware Data Structure for Phylogenomic
938 Inference from Supermatrices. *Systematic Biology* **65**:997–1008. doi:10.1093/sysbio/syw037

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

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

943 Cordaux R, Hedges DJ, Batzer MA. 2004. Retrotransposition of Alu elements: how many sources? *Trends*
944 *in Genetics* **20**:464–467. doi:10.1016/j.tig.2004.07.012

945 Corsinotti A, Wong FC, Tatar T, Szczerbinska I, Halbritter F, Colby D, Gogolok S, Pantier R, Liggat K,
946 Mirfazeli ES, Hall-Ponsele E, Mullin NP, Wilson V, Chambers I. 2017. Distinct SoxB1 networks are
947 required for naïve and primed pluripotency. *eLife* **6**:e27746. doi:10.7554/eLife.27746

948 Cosby RL, Chang N-C, Feschotte C. 2019. Host–transposon interactions: conflict, cooperation, and
949 cooption. *Genes Dev* **33**:1098–1116. doi:10.1101/gad.327312.119

950 Deniz Ö, Ahmed M, Todd CD, Rio-Machin A, Dawson MA, Branco MR. 2020. Endogenous retroviruses are
951 a source of enhancers with oncogenic potential in acute myeloid leukaemia. *Nat Commun* **11**:3506.
952 doi:10.1038/s41467-020-17206-4

953 Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis CA, Doyle F, Epstein CB, Frietze S, Harrow J, Kaul R,
954 Khatun J, Lajoie BR, Landt SG, Lee B-K, Pauli F, Rosenbloom KR, Sabo P, Safi A, Sanyal A, Shoresh N,
955 Simon JM, Song L, Trinklein ND, Altshuler RC, Birney E, Brown JB, Cheng C, Djebali S, Dong X, Dunham I,
956 Ernst J, Furey TS, Gerstein M, Giardine B, Greven M, Hardison RC, Harris RS, Herrero J, Hoffman MM,
957 Iyer S, Kellis M, Khatun J, Kheradpour P, Kundaje A, Lassmann T, Li Q, Lin X, Marinov GK, Merkel A,
958 Mortazavi A, Parker SCJ, Reddy TE, Rozowsky J, Schlesinger F, Thurman RE, Wang J, Ward LD, Whitfield
959 TW, Wilder SP, Wu W, Xi HS, Yip KY, Zhuang J, Bernstein BE, Birney E, Dunham I, Green ED, Gunter C,
960 Snyder M, Pazin MJ, Lowdon RF, Dillon LAL, Adams LB, Kelly CJ, Zhang J, Wexler JR, Green ED, Good PJ,
961 Feingold EA, Bernstein BE, Birney E, Crawford GE, Dekker J, Elnitski L, Farnham PJ, Gerstein M, Giddings
962 MC, Gingeras TR, Green ED, Guigó R, Hardison RC, Hubbard TJ, Kellis M, Kent WJ, Lieb JD, Margulies EH,
963 Myers RM, Snyder M, Stamatoyannopoulos JA, Tenenbaum SA, Weng Z, White KP, Wold B, Khatun J, Yu

964   Y, Wrobel J, Risk BA, Gunawardena HP, Kuiper HC, Maier CW, Xie L, Chen X, Giddings MC, Bernstein BE,
965   Epstein CB, Shoresh N, Ernst J, Kheradpour P, Mikkelsen TS, Gillespie S, Goren A, Ram O, Zhang X, Wang
966   L, Issner R, Coyne MJ, Durham T, Ku M, Truong T, Ward LD, Altshuler RC, Eaton ML, Kellis M, Djebali S,
967   Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W, Schlesinger F, Xue C,
968   Marinov GK, Khatun J, Williams BA, Zaleski C, Rozowsky J, Röder M, Kokocinski F, Abdelhamid RF, Alioto
969   T, Antoshechkin I, Baer MT, Batut P, Bell I, Bell K, Chakrabortty S, Chen X, Chrast J, Curado J, Derrien T,
970   Drenkow J, Dumais E, Dumais J, Duttagupta R, Fastuca M, Fejes-Toth K, Ferreira P, Foissac S, Fullwood
971   MJ, Gao H, Gonzalez D, Gordon A, Gunawardena HP, Howald C, Jha S, Johnson R, Kapranov P, King B,
972   Kingswood C, Li G, Luo OJ, Park E, Preall JB, Presaud K, Ribeca P, Risk BA, Robyr D, Ruan X, Sammeth M,
973   Sandhu KS, Schaeffer L, See L-H, Shahab A, Skancke J, Suzuki AM, Takahashi H, Tilgner H, Trout D,
974   Walters N, Wang H, Wrobel J, Yu Y, Hayashizaki Y, Harrow J, Gerstein M, Hubbard TJ, Reymond A,
975   Antonarakis SE, Hannon GJ, Giddings MC, Ruan Y, Wold B, Carninci P, Guigó R, Gingeras TR, Rosenbloom
976   KR, Sloan CA, Learned K, Malladi VS, Wong MC, Barber GP, Cline MS, Dreszer TR, Heitner SG, Karolchik D,
977   Kent WJ, Kirkup VM, Meyer LR, Long JC, Maddren M, Raney BJ, Furey TS, Song L, Grasfeder LL, Giresi PG,
978   Lee B-K, Battenhouse A, Sheffield NC, Simon JM, Showers KA, Safi A, London D, Bhinge AA, Shestak C,
979   Schaner MR, Ki Kim S, Zhang ZZ, Mieczkowski PA, Mieczkowska JO, Liu Z, McDaniell RM, Ni Y, Rashid NU,
980   Kim MJ, Adar S, Zhang Z, Wang T, Winter D, Keefe D, Birney E, Iyer VR, Lieb JD, Crawford GE, Li G,
981   Sandhu KS, Zheng M, Wang P, Luo OJ, Shahab A, Fullwood MJ, Ruan X, Ruan Y, Myers RM, Pauli F,
982   Williams BA, Gertz J, Marinov GK, Reddy TE, Vielmetter J, Partridge E, Trout D, Varley KE, Gasper C, The
983   ENCODE Project Consortium, Overall coordination (data analysis coordination), Data production leads
984   (data production), Lead analysts (data analysis), Writing group, NHGRI project management (scientific
985   management), Principal investigators (steering committee), Boise State University and University of
986   North Carolina at Chapel Hill Proteomics groups (data production and analysis), Broad Institute Group
987   (data production and analysis), Cold Spring Harbor U of G Center for Genomic Regulation, Barcelona,
988   RIKEN, Sanger Institute, University of Lausanne, Genome Institute of Singapore group (data production
989   and analysis), Data coordination center at UC Santa Cruz (production data coordination), Duke
990   University E University of Texas, Austin, University of North Carolina-Chapel Hill group (data production
991   and analysis), Genome Institute of Singapore group (data production and analysis), HudsonAlpha
992   Institute C UC Irvine, Stanford group (data production and analysis). 2012. An integrated encyclopedia of
993   DNA elements in the human genome. *Nature* **489**:57–74. doi:10.1038/nature11247

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

996   Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
997   *Acids Res* **32**:1792–1797. doi:10.1093/nar/gkh340

998   Eickbush TH, Malik HS. 2002. Origins and Evolution of Retrotransposons. *Mobile DNA II* 1111–1144.
999   doi:10.1128/9781555817954.ch49

1000  Estarás C, Benner C, Jones KA. 2015. SMADs and YAP Compete to Control Elongation of β-Catenin:LEF-1-
1001  Recruited RNAPII during hESC Differentiation. *Molecular Cell* **58**:780–793.
1002  doi:10.1016/j.molcel.2015.04.001

1003  Faulkner GJ, Kimura Y, Daub CO, Wani S, Plessy C, Irvine KM, Schroder K, Cloonan N, Steptoe AL,
1004  Lassmann T, Waki K, Hornig N, Arakawa T, Takahashi H, Kawai J, Forrest ARR, Suzuki H, Hayashizaki Y,

43

1005    Hume DA, Orlando V, Grimmond SM, Carninci P. 2009. The regulated retrotransposon transcriptome of
1006    mammalian cells. *Nat Genet* **41**:563–571. doi:10.1038/ng.368

1007    Fernandes JD, Zamudio-Hurtado A, Clawson H, Kent WJ, Haussler D, Salama SR, Haeussler M. 2020. The
1008    UCSC repeat browser allows discovery and visualization of evolutionary conflict across repeat families.
1009    *Mobile DNA* **11**:13. doi:10.1186/s13100-020-00208-w

1010    Feschotte C. 2008. Transposable elements and the evolution of regulatory networks. *Nat Rev Genet*
1011    **9**:397–405. doi:10.1038/nrg2337

1012    Fort A, Hashimoto K, Yamada D, Salimullah M, Keya CA, Saxena A, Bonetti A, Voineagu I, Bertin N, Kratz
1013    A, Noro Y, Wong C-H, de Hoon M, Andersson R, Sandelin A, Suzuki H, Wei C-L, Koseki H, Hasegawa Y,
1014    Forrest ARR, Carninci P. 2014. Deep transcriptome profiling of mammalian stem cells supports a
1015    regulatory role for retrotransposons in pluripotency maintenance. *Nature Genetics* **46**:558–566.
1016    doi:10.1038/ng.2965

1017    Gafni O, Weinberger L, Mansour AA, Manor YS, Chomsky E, Ben-Yosef D, Kalma Y, Viukov S, Maza I,
1018    Zviran A, Rais Y, Shipony Z, Mukamel Z, Krupalnik V, Zerbib M, Geula S, Caspi I, Schneir D, Shwartz T,
1019    Gilad S, Amann-Zalcenstein D, Benjamin S, Amit I, Tanay A, Massarwa R, Novershtern N, Hanna JH. 2013.
1020    Derivation of novel human ground state naive pluripotent stem cells. *Nature* **504**:282–286.
1021    doi:10.1038/nature12745

1022    Gemmell P, Hein J, Katzourakis A. 2019. The Exaptation of HERV-H: Evolutionary Analyses Reveal the
1023    Genomic Features of Highly Transcribed Elements. *Front Immunol* **10**. doi:10.3389/fimmu.2019.01339

1024    Gemmell P, Hein J, Katzourakis A. 2015. Orthologous endogenous retroviruses exhibit directional
1025    selection since the chimp-human split. *Retrovirology* **12**. doi:10.1186/s12977-015-0172-6

1026    Gibbs RA, Rogers J, Katze MG, Bumgarner R, Weinstock GM, Mardis ER, Remington KA, Strausberg RL,
1027    Venter JC, Wilson RK, Batzer MA, Bustamante CD, Eichler EE, Hahn MW, Hardison RC, Makova KD, Miller
1028    W, Milosavljevic A, Palermo RE, Siepel A, Sikela JM, Attaway T, Bell S, Bernard KE, Buhay CJ,
1029    Chandrabose MN, Dao M, Davis C, Delehaunty KD, Ding Y, Dinh HH, Dugan-Rocha S, Fulton LA, Gabisi RA,
1030    Garner TT, Godfrey J, Hawes AC, Hernandez J, Hines S, Holder M, Hume J, Jhangiani SN, Joshi V, Khan
1031    ZM, Kirkness EF, Cree A, Fowler RG, Lee S, Lewis LR, Li Z, Liu Yih-shin, Moore SM, Muzny D, Nazareth LV,
1032    Ngo DN, Okwuonu GO, Pai G, Parker D, Paul HA, Pfannkoch C, Pohl CS, Rogers Y-H, Ruiz SJ, Sabo A,
1033    Santibanez J, Schneider BW, Smith SM, Sodergren E, Svatek AF, Utterback TR, Vattathil S, Warren W,
1034    White CS, Chinwalla AT, Feng Y, Halpern AL, Hillier LW, Huang X, Minx P, Nelson JO, Pepin KH, Qin X,
1035    Sutton GG, Venter E, Walenz BP, Wallis JW, Worley KC, Yang S-P, Jones SM, Marra MA, Rocchi M, Schein
1036    JE, Baertsch R, Clarke L, Csürös M, Glasscock J, Harris RA, Havlak P, Jackson AR, Jiang H, Liu Yue, Messina
1037    DN, Shen Y, Song HX-Z, Wylie T, Zhang L, Birney E, Han K, Konkel MK, Lee J, Smit AFA, Ullmer B, Wang H,
1038    Xing J, Burhans R, Cheng Z, Karro JE, Ma J, Raney B, She X, Cox MJ, Demuth JP, Dumas LJ, Han S-G,
1039    Hopkins J, Karimpour-Fard A, Kim YH, Pollack JR, Vinar T, Addo-Quaye C, Degenhardt J, Denby A, Hubisz
1040    MJ, Indap A, Kosiol C, Lahn BT, Lawson HA, Marklein A, Nielsen R, Vallender EJ, Clark AG, Ferguson B,
1041    Hernandez RD, Hirani K, Kehrer-Sawatzki H, Kolb J, Patil S, Pu L-L, Ren Y, Smith DG, Wheeler DA, Schenck
1042    I, Ball EV, Chen R, Cooper DN, Giardine B, Hsu F, Kent WJ, Lesk A, Nelson DL, O'Brien WE, Prüfer K,
1043    Stenson PD, Wallace JC, Ke H, Liu X-M, Wang P, Xiang AP, Yang F, Barber GP, Haussler D, Karolchik D,

1044    Kern AD, Kuhn RM, Smith KE, Zwieg AS. 2007. Evolutionary and Biomedical Insights from the Rhesus
1045    Macaque Genome. *Science* **316**:222–234. doi:10.1126/science.1139247

1046    Ginestet C. 2011. ggplot2: Elegant Graphics for Data Analysis. *Journal of the Royal Statistical Society:*
1047    *Series A (Statistics in Society)* **174**:245–246. doi:https://doi.org/10.1111/j.1467-985X.2010.00676_9.x

1048    Glinsky GV. 2015. Transposable Elements and DNA Methylation Create in Embryonic Stem Cells Human-
1049    Specific Regulatory Sequences Associated with Distal Enhancers and Noncoding RNAs. *Genome Biol Evol*
1050    **7**:1432–1454. doi:10.1093/gbe/evv081

1051    Göke J, Lu X, Chan Y-S, Ng H-H, Ly L-H, Sachs F, Szczerbinska I. 2015. Dynamic Transcription of Distinct
1052    Classes of Endogenous Retroviral Elements Marks Specific Populations of Early Human Embryonic Cells.
1053    *Cell Stem Cell* **16**:135–141. doi:10.1016/j.stem.2015.01.005

1054    Goodchild NL, Wilkinson DA, Mager DL. 1993. Recent Evolutionary Expansion of a Subfamily of RTVL-H
1055    Human Endogenous Retrovirus-like Elements. *Virology* **196**:778–788. doi:10.1006/viro.1993.1535

1056    Gouy M, Guindon S, Gascuel O. 2010. SeaView Version 4: A Multiplatform Graphical User Interface for
1057    Sequence Alignment and Phylogenetic Tree Building. *Molecular Biology and Evolution* **27**:221–224.
1058    doi:10.1093/molbev/msp259

1059    Haig D. 2016. Transposable elements: Self-seekers of the germline, team-players of the soma. *BioEssays*
1060    **38**:1158–1166. doi:10.1002/bies.201600125

1061    Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, Cheng JX, Murre C, Singh H, Glass CK. 2010.
1062    Simple combinations of lineage-determining transcription factors prime cis-regulatory elements
1063    required for macrophage and B cell identities. *Mol Cell* **38**:576–589. doi:10.1016/j.molcel.2010.05.004

1064    Hermant C, Torres-Padilla M-E. 2021. TFs for TEs: the transcription factor repertoire of mammalian
1065    transposable elements. *Genes Dev* **35**:22–39. doi:10.1101/gad.344473.120

1066    Imbeault M, Helleboid P-Y, Trono D. 2017. KRAB zinc-finger proteins contribute to the evolution of gene
1067    regulatory networks. *Nature* **543**:550–554. doi:10.1038/nature21683

1068    Ito J, Sugimoto R, Nakaoka H, Yamada S, Kimura T, Hayano T, Inoue I. 2017. Systematic identification and
1069    characterization of regulatory elements derived from human endogenous retroviruses. *PLOS Genetics*
1070    **13**:e1006883. doi:10.1371/journal.pgen.1006883

1071    Izsvák Z, Wang J, Singh M, Mager DL, Hurst LD. 2016. Pluripotency and the endogenous retrovirus
1072    HERVH: Conflict or serendipity? *BioEssays* **38**:109–117. doi:10.1002/bies.201500096

1073    Jacobs FM, Greenberg D, Nguyen N, Haeussler M, Ewing AD, Katzman S, Paten B, Salama SR, Haussler D.
1074    2014. An evolutionary arms race between KRAB zinc finger genes 91/93 and SVA/L1 retrotransposons.
1075    *Nature* **516**:242–245. doi:10.1038/nature13760

1076    Jacques P-É, Jeyakani J, Bourque G. 2013. The Majority of Primate-Specific Regulatory Sequences Are
1077    Derived from Transposable Elements. *PLOS Genetics* **9**:e1003504. doi:10.1371/journal.pgen.1003504

1078    Jern P, Sperber GO, Ahlsén G, Blomberg J. 2005. Sequence Variability, Gene Structure, and Expression of
1079    Full-Length Human Endogenous Retrovirus H. *Journal of Virology* **79**.

1080　Jern P, Sperber GO, Blomberg J. 2004. Definition and variation of human endogenous retrovirus H.
1081　*Virology* **327**:93–110. doi:10.1016/j.virol.2004.06.023

1082　Jetzt AE, Yu H, Klarmann GJ, Ron Y, Preston BD, Dougherty JP. 2000. High rate of recombination
1083　throughout the human immunodeficiency virus type 1 genome. *J Virol* **74**:1234–1240.
1084　doi:10.1128/jvi.74.3.1234-1240.2000

1085　Johnson WE. 2019. Origins and evolutionary consequences of ancient endogenous retroviruses. *Nat Rev*
1086　*Microbiol* **17**:355–370. doi:10.1038/s41579-019-0189-2

1087　Kalkan T, Smith A. 2014. Mapping the route from naive pluripotency to lineage specification.
1088　*Philosophical Transactions of the Royal Society B: Biological Sciences* **369**:20130540.
1089　doi:10.1098/rstb.2013.0540

1090　Kelley D, Rinn J. 2012. Transposable elements reveal a stem cell-specific class of long noncoding RNAs.
1091　*Genome Biology* **13**:R107. doi:10.1186/gb-2012-13-11-r107

1092　Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler and D. 2002. The Human
1093　Genome Browser at UCSC. *Genome Res* **12**:996–1006. doi:10.1101/gr.229102

1094　Khan H, Smit A, Boissinot S. 2006. Molecular evolution and tempo of amplification of human LINE-1
1095　retrotransposons since the origin of primates. *Genome Res* **16**:78–87. doi:10.1101/gr.4001406

1096　Kinoshita M, Barber M, Mansfield W, Cui Y, Spindlow D, Stirparo GG, Dietmann S, Nichols J, Smith A.
1097　2021. Capture of Mouse and Human Stem Cells with Features of Formative Pluripotency. *Cell Stem Cell*
1098　**28**:453-471.e8. doi:10.1016/j.stem.2020.11.005

1099　Kojima KK. 2018. Human transposable elements in Repbase: genomic footprints from fish to humans.
1100　*Mobile DNA* **9**:2. doi:10.1186/s13100-017-0107-y

1101　Kong Y, Rose CM, Cass AA, Williams AG, Darwish M, Lianoglou S, Haverty PM, Tong A-J, Blanchette C,
1102　Albert ML, Mellman I, Bourgon R, Greally J, Jhunjhunwala S, Chen-Harris H. 2019. Transposable element
1103　expression in tumors is associated with immune infiltration and increased antigenicity. *Nat Commun*
1104　**10**:5228. doi:10.1038/s41467-019-13035-2

1105　Krönung SK, Beyer U, Chiaramonte ML, Dolfini D, Mantovani R, Dobbelstein M. 2016. LTR12 promoter
1106　activation in a broad range of human tumor cells by HDAC inhibition. *Oncotarget* **7**:33484–33497.
1107　doi:10.18632/oncotarget.9255

1108　Kunarso G, Chia N-Y, Jeyakani J, Hwang C, Lu X, Chan Y-S, Ng H-H, Bourque G. 2010. Transposable
1109　elements have rewired the core regulatory network of human embryonic stem cells. *Nature Genetics*
1110　**42**:631–634. doi:10.1038/ng.600

1111　Lai MM. 1992. RNA recombination in animal and plant viruses. *Microbiol Rev* **56**:61–79.

1112　Lanciano S, Cristofari G. 2020. Measuring and interpreting transposable element expression. *Nat Rev*
1113　*Genet* **21**:721–736. doi:10.1038/s41576-020-0251-y

1114　Leigh JW, Bryant D. 2015. popart: full-feature software for haplotype network construction. *Methods in*
1115　*Ecology and Evolution* **6**:1110–1116. doi:https://doi.org/10.1111/2041-210X.12410

1116   Locke DP, Hillier LW, Warren WC, Worley KC, Nazareth LV, Muzny DM, Yang S-P, Wang Z, Chinwalla AT,
1117   Minx P, Mitreva M, Cook L, Delehaunty KD, Fronick C, Schmidt H, Fulton LA, Fulton RS, Nelson JO,
1118   Magrini V, Pohl C, Graves TA, Markovic C, Cree A, Dinh HH, Hume J, Kovar CL, Fowler GR, Lunter G,
1119   Meader S, Heger A, Ponting CP, Marques-Bonet T, Alkan C, Chen L, Cheng Z, Kidd JM, Eichler EE, White S,
1120   Searle S, Vilella AJ, Chen Y, Flicek P, Ma J, Raney B, Suh B, Burhans R, Herrero J, Haussler D, Faria R,
1121   Fernando O, Darré F, Farré D, Gazave E, Oliva M, Navarro A, Roberto R, Capozzi O, Archidiacono N, Valle
1122   GD, Purgato S, Rocchi M, Konkel MK, Walker JA, Ullmer B, Batzer MA, Smit AFA, Hubley R, Casola C,
1123   Schrider DR, Hahn MW, Quesada V, Puente XS, Ordoñez GR, López-Otín C, Vinar T, Brejova B, Ratan A,
1124   Harris RS, Miller W, Kosiol C, Lawson HA, Taliwal V, Martins AL, Siepel A, RoyChoudhury A, Ma X,
1125   Degenhardt J, Bustamante CD, Gutenkunst RN, Mailund T, Dutheil JY, Hobolth A, Schierup MH, Ryder
1126   OA, Yoshinaga Y, de Jong PJ, Weinstock GM, Rogers J, Mardis ER, Gibbs RA, Wilson RK. 2011.
1127   Comparative and demographic analysis of orang-utan genomes. *Nature* **469**:529–533.
1128   doi:10.1038/nature09687

1129   Loewer S, Cabili MN, Guttman M, Loh Y-H, Thomas K, Park IH, Garber M, Curran M, Onder T, Agarwal S,
1130   Manos PD, Datta S, Lander ES, Schlaeger TM, Daley GQ, Rinn JL. 2010. Large intergenic non-coding RNA-
1131   RoR modulates reprogramming of human induced pluripotent stem cells. *Nat Genet* **42**:1113–1117.
1132   doi:10.1038/ng.710

1133   Löytynoja A, Goldman N. 2010. webPRANK: a phylogeny-aware multiple sequence aligner with
1134   interactive alignment browser. *BMC Bioinformatics* **11**:579–579. doi:10.1186/1471-2105-11-579

1135   Lu X, Sachs F, Ramsay L, Jacques P-É, Göke J, Bourque G, Ng H-H. 2014. The retrovirus HERVH is a long
1136   noncoding RNA required for human embryonic stem cell identity. *Nature Structural & Molecular Biology*
1137   **21**:423–425. doi:10.1038/nsmb.2799

1138   Macfarlan TS, Gifford WD, Driscoll S, Lettieri K, Rowe HM, Bonanomi D, Firth A, Singer O, Trono D, Pfaff
1139   SL. 2012. Embryonic stem cell potency fluctuates with endogenous retrovirus activity. *Nature* **487**:57–
1140   63. doi:10.1038/nature11244

1141   Mager DL, Freeman JD. 1995. HERV-H Endogenous Retroviruses: Presence in the New World Branch but
1142   Amplification in the Old World Primate Lineage. *Virology* **213**:395–404. doi:10.1006/viro.1995.0012

1143   Mager DL, Freeman JD. 1987. Human endogenous retroviruslike genome with type C pol sequences and
1144   gag sequences related to human T-cell lymphotropic viruses. *J Virol* **61**:4060–4066.
1145   doi:10.1128/jvi.61.12.4060-4066.1987

1146   Matsuda E, Garfinkel DJ. 2009. Posttranslational interference of Ty1 retrotransposition by antisense
1147   RNAs. *Proc Natl Acad Sci U S A* **106**:15657–15662. doi:10.1073/pnas.0908305106

1148   Miao B, Fu S, Lyu C, Gontarz P, Wang T, Zhang B. 2020. Tissue-specific usage of transposable element-
1149   derived promoters in mouse development. *Genome Biol* **21**:1–25. doi:10.1186/s13059-020-02164-3

1150   Miga KH, Newton Y, Jain M, Altemose N, Willard HF, Kent WJ. 2014. Centromere reference models for
1151   human chromosomes X and Y satellite arrays. *Genome Res* **24**:697–707. doi:10.1101/gr.159624.113

1152   Nakamura T, Yamada KD, Tomii K, Katoh K. 2018. Parallelization of MAFFT for large-scale multiple
1153   sequence alignments. *Bioinformatics* **34**:2490–2492. doi:10.1093/bioinformatics/bty121

1154   Nichols J, Smith A. 2009. Naive and Primed Pluripotent States. *Cell Stem Cell* **4**:487–492.
1155   doi:10.1016/j.stem.2009.05.015

1156   Niwa H. 2007. How is pluripotency determined and maintained? *Development* **134**:635–646.
1157   doi:10.1242/dev.02787

1158   Niwa H, Nakamura A, Urata M, Shirae-Kurabayashi M, Kuraku S, Russell S, Ohtsuka S. 2016. The
1159   evolutionally-conserved function of group B1 Sox family members confers the unique role of Sox2 in
1160   mouse ES cells. *BMC Evolutionary Biology* **16**:173. doi:10.1186/s12862-016-0755-4

1161   Ohnuki M, Tanabe K, Sutou K, Teramoto I, Sawamura Y, Narita M, Nakamura Michiko, Tokunaga Y,
1162   Nakamura Masahiro, Watanabe A, Yamanaka S, Takahashi K. 2014. Dynamic regulation of human
1163   endogenous retroviruses mediates factor-induced reprogramming and differentiation potential. *PNAS*
1164   **111**:12426–12431. doi:10.1073/pnas.1413299111

1165   Patro R, Duggal G, Love MI, Irizarry RA, Kingsford C. 2017. Salmon provides fast and bias-aware
1166   quantification of transcript expression. *Nat Methods* **14**:417–419. doi:10.1038/nmeth.4197

1167   Peaston AE, Evsikov AV, Graber JH, Vries WN de, Holbrook AE, Solter D, Knowles BB. 2004.
1168   Retrotransposons Regulate Host Genes in Mouse Oocytes and Preimplantation Embryos. *Developmental*
1169   *Cell* **7**:597–606. doi:10.1016/j.devcel.2004.09.004

1170   Pérez-Losada M, Arenas M, Galán JC, Palero F, González-Candelas F. 2015. Recombination in viruses:
1171   Mechanisms, methods of study, and evolutionary consequences. *Infect Genet Evol* **30**:296–307.
1172   doi:10.1016/j.meegid.2014.12.022

1173   Pontis J, Planet E, Offner S, Turelli P, Duc J, Coudray A, Theunissen TW, Jaenisch R, Trono D. 2019.
1174   Hominoid-Specific Transposable Elements and KZFPs Facilitate Human Embryonic Genome Activation
1175   and Control Transcription in Naive Human ESCs. *Cell Stem Cell* **24**:724-735.e5.
1176   doi:10.1016/j.stem.2019.03.012

1177   Posada D, Crandall KA. 2001. Intraspecific gene genealogies: trees grafting into networks. *Trends in*
1178   *Ecology & Evolution* **16**:37–45. doi:10.1016/S0169-5347(00)02026-7

1179   Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
1180   *Bioinformatics* **26**:841–842. doi:10.1093/bioinformatics/btq033

1181   Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S, Dündar F, Manke T. 2016.
1182   deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Research*
1183   **44**:W160–W165. doi:10.1093/nar/gkw257

1184   Raney BJ, Dreszer TR, Barber GP, Clawson H, Fujita PA, Wang T, Nguyen N, Paten B, Zweig AS, Karolchik
1185   D, Kent WJ. 2014. Track data hubs enable visualization of user-defined genome-wide annotations on the
1186   UCSC Genome Browser. *Bioinformatics* **30**:1003–1005. doi:10.1093/bioinformatics/btt637

1187   Rebollo R, Romanish MT, Mager DL. 2012. Transposable elements: an abundant and natural source of
1188   regulatory sequences for host genes. *Annu Rev Genet* **46**:21–42. doi:10.1146/annurev-genet-110711-
1189   155621

1190   Römer C, Singh M, Hurst LD, Izsvák Z. 2017. How to tame an endogenous retrovirus: HERVH and the
1191   evolution of human pluripotency. *Current Opinion in Virology*, Animal models for viral diseases •
1192   Paleovirology **25**:49–58. doi:10.1016/j.coviro.2017.07.001

1193   Rossant J, Tam PPL. 2017. New Insights into Early Human Development: Lessons for Stem Cell Derivation
1194   and Differentiation. *Cell Stem Cell* **20**:18–28. doi:10.1016/j.stem.2016.12.004

1195   Santoni FA, Guerra J, Luban J. 2012. HERV-H RNA is abundant in human embryonic stem cells and a
1196   precise marker for pluripotency. *Retrovirology* **9**:111. doi:10.1186/1742-4690-9-111

1197   Scally A, Dutheil JY, Hillier LW, Jordan GE, Goodhead I, Herrero J, Hobolth A, Lappalainen T, Mailund T,
1198   Marques-Bonet T, McCarthy S, Montgomery SH, Schwalie PC, Tang YA, Ward MC, Xue Y, Yngvadottir B,
1199   Alkan C, Andersen LN, Ayub Q, Ball EV, Beal K, Bradley BJ, Chen Y, Clee CM, Fitzgerald S, Graves TA, Gu Y,
1200   Heath P, Heger A, Karakoc E, Kolb-Kokocinski A, Laird GK, Lunter G, Meader S, Mort M, Mullikin JC,
1201   Munch K, O'Connor TD, Phillips AD, Prado-Martinez J, Rogers AS, Sajjadian S, Schmidt D, Shaw K,
1202   Simpson JT, Stenson PD, Turner DJ, Vigilant L, Vilella AJ, Whitener W, Zhu B, Cooper DN, de Jong P,
1203   Dermitzakis ET, Eichler EE, Flicek P, Goldman N, Mundy NI, Ning Z, Odom DT, Ponting CP, Quail MA,
1204   Ryder OA, Searle SM, Warren WC, Wilson RK, Schierup MH, Rogers J, Tyler-Smith C, Durbin R. 2012.
1205   Insights into hominid evolution from the gorilla genome sequence. *Nature* **483**:169–175.
1206   doi:10.1038/nature10842

1207   Schön U, Seifarth W, Baust C, Hohenadl C, Erfle V, Leib-Mösch C. 2001. Cell Type-Specific Expression and
1208   Promoter Activity of Human Endogenous Retroviral Long Terminal Repeats. *Virology* **279**:280–291.
1209   doi:10.1006/viro.2000.0712

1210   Simon-Loriere E, Holmes EC. 2011. Why do RNA viruses recombine? *Nat Rev Microbiol* **9**:617–626.
1211   doi:10.1038/nrmicro2614

1212   Smit AF, Hubley R, Green P. 2013. RepeatMasker Open-4.0.

1213   Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. 2021. The Dfam community resource of transposable
1214   element families, sequence models, and genome annotations. *Mobile DNA* **12**:2. doi:10.1186/s13100-
1215   020-00230-y

1216   Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, Snyder MP, Wang T. 2014. Widespread contribution of
1217   transposable elements to the innovation of gene regulatory networks. *Genome Res* **24**:1963–1976.
1218   doi:10.1101/gr.168872.113

1219   Sundaram V, Wysocka J. 2020. Transposable elements as a potent source of diverse cis-regulatory
1220   sequences in mammalian genomes. *Philosophical Transactions of the Royal Society B: Biological Sciences*
1221   **375**:20190347. doi:10.1098/rstb.2019.0347

1222   Svoboda P, Stein P, Anger M, Bernstein E, Hannon GJ, Schultz RM. 2004. RNAi and expression of
1223   retrotransposons MuERV-L and IAP in preimplantation mouse embryos. *Developmental Biology*
1224   **269**:276–285. doi:10.1016/j.ydbio.2004.01.028

1225   Takahashi K, Nakamura M, Okubo C, Kliesmete Z, Ohnuki M, Narita M, Watanabe A, Ueda M, Takashima
1226   Y, Hellmann I, Yamanaka S. 2021. The pluripotent stem cell-specific transcript ESRG is dispensable for
1227   human pluripotency. *PLOS Genetics* **17**:e1009587. doi:10.1371/journal.pgen.1009587

1228    Tang F, Barbacioru C, Bao S, Lee C, Nordman E, Wang X, Lao K, Surani MA. 2010. Tracing the Derivation
1229    of Embryonic Stem Cells from the Inner Cell Mass by Single-Cell RNA-Seq Analysis. *Cell Stem Cell* **6**:468–
1230    478. doi:10.1016/j.stem.2010.03.015

1231    Theunissen TW, Friedli M, He Y, Planet E, O'Neil RC, Markoulaki S, Pontis J, Wang H, Iouranova A,
1232    Imbeault M, Duc J, Cohen MA, Wert KJ, Castanon R, Zhang Z, Huang Y, Nery JR, Drotar J, Lungjangwa T,
1233    Trono D, Ecker JR, Jaenisch R. 2016. Molecular Criteria for Defining the Naive Human Pluripotent State.
1234    *Cell Stem Cell* **19**:502–515. doi:10.1016/j.stem.2016.06.011

1235    Thomas J, Perron H, Feschotte C. 2018. Variation in proviral content among human genomes mediated
1236    by LTR recombination. *Mobile DNA* **9**:36. doi:10.1186/s13100-018-0142-3

1237    Thompson PJ, Macfarlan TS, Lorincz MC. 2016. Long Terminal Repeats: From Parasitic Elements to
1238    Building Blocks of the Transcriptional Regulatory Repertoire. *Molecular Cell* **62**:766–776.
1239    doi:10.1016/j.molcel.2016.03.029

1240    Urusov FA, Nefedova LN, Kim AI. 2011. Analysis of the tissue- and stage-specific transportation of the
1241    Drosophila melanogaster gypsy retrotransposon. *Russ J Genet Appl Res* **1**:507–510.
1242    doi:10.1134/S2079059711060104

1243    Vargiu L, Rodriguez-Tomé P, Sperber GO, Cadeddu M, Grandi N, Blikstad V, Tramontano E, Blomberg J.
1244    2016. Classification and characterization of human endogenous retroviruses; mosaic forms are common.
1245    *Retrovirology* **13**:7. doi:10.1186/s12977-015-0232-y

1246    Wang J, Xie G, Singh M, Ghanbarian AT, Raskó T, Szvetnik A, Cai H, Besser D, Prigione A, Fuchs NV,
1247    Schumann GG, Chen W, Lorincz MC, Ivics Z, Hurst LD, Izsvák Z. 2014. Primate-specific endogenous
1248    retrovirus-driven transcription defines naive-like stem cells. *Nature* **516**:405–409.
1249    doi:10.1038/nature13804

1250    Wang Z, Oron E, Nelson B, Razis S, Ivanova N. 2012. Distinct lineage specification roles for NANOG,
1251    OCT4, and SOX2 in human embryonic stem cells. *Cell Stem Cell* **10**:440–454.
1252    doi:10.1016/j.stem.2012.02.016

1253    Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2—a multiple
1254    sequence alignment editor and analysis workbench. *Bioinformatics* **25**:1189–1191.
1255    doi:10.1093/bioinformatics/btp033

1256    Waterson RH, Lander ES, Wilson RK, The Chimpanzee Sequencing and Analysis Consortium. 2005. Initial
1257    sequence of the chimpanzee genome and comparison with the human genome. *Nature* **437**:69–87.
1258    doi:10.1038/nature04072

1259    Wolf G, de Iaco A, Sun M-A, Bruno M, Tinkham M, Hoang D, Mitra A, Ralls S, Trono D, Macfarlan TS.
1260    2020. KRAB-zinc finger protein gene expansion in response to active retrotransposons in the murine
1261    lineage. *eLife* **9**:e56337. doi:10.7554/eLife.56337

1262    Yang P, Wang Y, Macfarlan TS. 2017. The role of KRAB-ZFPs in transposable element repression and
1263    mammalian evolution. *Trends Genet* **33**:871–881. doi:10.1016/j.tig.2017.08.006

1264    Yu H-L, Zhao Z-K, Zhu F. 2013. The role of human endogenous retroviral long terminal repeat sequences
1265    in human cancer (Review). *International Journal of Molecular Medicine* **32**:755–762.
1266    doi:10.3892/ijmm.2013.1460

1267    Zhang Y, Li T, Preissl S, Amaral ML, Grinstein JD, Farah EN, Destici E, Qiu Y, Hu R, Lee AY, Chee S, Ma K, Ye
1268    Z, Zhu Q, Huang H, Fang R, Yu L, Izpisua Belmonte JC, Wu J, Evans SM, Chi NC, Ren B. 2019.
1269    Transcriptionally active HERV-H retrotransposons demarcate topologically associating domains in human
1270    pluripotent stem cells. *Nature Genetics* **51**:1380–1388. doi:10.1038/s41588-019-0479-7

1271