

1 **Evolution of transposable element-derived enhancer activity**

2

3 Alan Y. Du^{1,2}, Xiaoyu Zhuo^{1,2}, Vasavi Sundaram^{1,2}, Nicholas O. Jensen^{1,3,4}, Hemangi G.

4 Chaudhari^{1,2}, Nancy L. Saccone^{1,3}, Barak A. Cohen^{1,2}, and Ting Wang^{1,2}

5

6 ¹Department of Genetics, Washington University School of Medicine, St. Louis, MO, USA

7 ²The Edison Family Center for Genome Sciences and Systems Biology, Washington University

8 School of Medicine, St. Louis, MO, USA

9 ³Division of Biostatistics, Washington University School of Medicine, St. Louis, MO, USA

10 ⁴Department of Developmental Biology, Washington University School of Medicine, St. Louis,

11 MO, USA

12

13 **Abstract**

14 Many transposable elements (TEs) contain transcription factor binding sites and are implicated as
15 potential regulatory elements. However, TEs are rarely functionally tested for regulatory activity,
16 which in turn limits our understanding of how TE regulatory activity has evolved. We
17 systematically tested the human LTR18A subfamily for regulatory activity using massively
18 parallel reporter assay (MPRA) and found AP-1 and C/EBP-related binding motifs as drivers of
19 enhancer activity. Functional analysis of evolutionarily reconstructed ancestral sequences revealed
20 that LTR18A elements have generally lost regulatory activity over time through sequence changes,
21 with the largest effects occurring due to mutations in the AP-1 and C/EBP motifs. We observed
22 that the two motifs are conserved at higher rates than expected based on neutral evolution. Finally,
23 we identified LTR18A elements as potential enhancers in the human genome, primarily in
24 epithelial cells. Together, our results provide a model for the origin, evolution, and co-option of
25 TE-derived regulatory elements.

26

27 **Introduction**

28 Changes in gene regulation have long been implicated as crucial drivers in evolution¹. Since the
29 discovery of the SV40 enhancer element, enhancers have emerged as one of the major classes of
30 cis-regulatory sequences that can modulate gene expression^{2,3}. Due to several unique properties,
31 enhancers have emerged as excellent candidates upon which evolution can act. Enhancers are often
32 active depending on cellular context like cell type or response to stimuli. This modularity can
33 minimize functional trade-offs and allows selection to act more efficiently⁴. Furthermore,
34 redundant enhancers, or “shadow” enhancers, provide robustness in gene regulatory networks and
35 may allow for greater freedom to develop new functions^{5,6}.

36

37 The development of massively parallel reporter assays (MPRAs) has greatly accelerated our
38 understanding of enhancers by facilitating simultaneous testing of thousands of DNA sequences^{7–}
39 ¹⁰. MPRAs have been used to probe the enhancer potential of sequences underlying various
40 epigenetic marks¹¹, dissect enhancer logic through tiling and mutagenesis^{9,12,13}, and decipher the
41 effects of naturally occurring sequence variants^{8,14–16}. Several studies have also employed MPRA
42 to understand the evolution of fly and primate enhancers, revealing widespread enhancer
43 turnover^{17,18}.

44

45 Transposable elements (TEs) are repetitive DNA elements that represent a rich source of genetic
46 material for regulatory innovation¹⁹. In mammalian genomes, TEs have made substantial
47 contributions to the collection of transcription factor binding sites^{20–24}. These binding sites are
48 often enriched within certain TE subfamilies, groups of similar TE sequences that are derived from
49 a single ancestral origin. Individual copies of TE subfamilies can then be co-opted into gene

50 regulatory networks such as in pregnancy and innate immunity^{25,26}. Overall, TEs make up a quarter
51 of the regulatory epigenome in human²⁷, and by some estimates, the majority of primate-specific
52 regulatory sequences are derived from TEs^{28,29}. Despite these advances in the field, there remains
53 a gap in knowledge of how TEs obtain regulatory activity and how this activity changes over the
54 course of evolution.

55
56 As repetitive sequences, TEs offer a unique perspective into the evolution of cis-regulatory
57 elements. One intrinsic limitation for evolutionary studies is that each enhancer has one ortholog
58 per species barring duplication or deletion, which constrains the sample size for analysis. Within
59 a TE subfamily, each TE is descended from a common ancestor, with each copy evolving mostly
60 independently. This provides a large sample size to draw upon within even a single genome. To
61 serve as a representative subfamily, we selected LTR18A which we previously identified to be
62 enriched for MAFK transcription factor binding peaks and motifs²⁴.

63
64 Here, we aim to investigate the evolution of regulatory potential in the LTR18A subfamily using
65 MPRA. By using present day LTR18A sequences found across seven primate species, we
66 computationally reconstruct ancestral sequences during LTR18A evolution across a span of
67 roughly 75 million years. We apply tiling and motif-focused approaches to test reconstructed and
68 present day LTR18A sequences for enhancer activity. Using natural sequence variations between
69 LTR18A elements, we identify transcription factor binding sites that drive LTR18A enhancer
70 activity and validate them through mutagenesis. By annotating enhancer activity for the root and
71 intermediate ancestral LTR18A elements in our reconstructed phylogenetic tree, we investigate
72 the origin of enhancer activity for the LTR18A family as well as key mutations that have led to

73 changes in activity over time. Finally, we explore the influence of selection on LTR18A and the
74 possibility of co-option in the human epigenome.

75

76 **Results**

77 **Reconstruction of the LTR18A phylogenetic tree**

78 In order to reconstruct the evolutionary history of the LTR18A subfamily, we first identified high
79 confidence LTR18A elements in human and their orthologous elements in six other primate
80 species. The LTR18A subfamily is found in the Simiiformes taxa³⁰. From the Simiiformes, we
81 obtained RepeatMasker annotations for human (hg19), chimpanzee (panTro4), gorilla (gorGor3),
82 gibbon (nomLeu3), baboon (papAnu2), rhesus macaque (rheMac3), and marmoset (calJac3)
83 genomes. Due to the similarity of the LTR18A, LTR18B, and LTR18C consensus sequences, we
84 performed manual curation of hg19 LTR18A to select for LTR18A elements that are confidently
85 assigned to the subfamily. Briefly, we filtered out LTR18A elements that could be aligned to either
86 the LTR18B or LTR18C consensus, and we removed LTR18A elements that might be
87 misannotated using paired LTRs (Methods). Following these criteria, 181 out of 198 LTR18A
88 elements annotated by RepeatMasker are retained (Supplemental Table 1). Next, we found primate
89 orthologs for each hg19 LTR18A element by using synteny³¹. LTR18A elements that correspond
90 with multiple orthologs in the same genome, or vice versa, were excluded. Each hg19 LTR18A
91 element with its primate orthologs were considered an ortholog set. We further selected for
92 LTR18A pairs that have orthologs in chimpanzee, gorilla, and at least two of the four other
93 primates. In the end, 46 (consisting of 23 pairs) LTR18A ortholog sets were chosen for ancestral
94 reconstruction.

95

96 From our set of manually curated human LTR18A elements and their orthologs, we
97 computationally reconstructed the LTR18A phylogenetic tree using a two-step process. Based on
98 the unique characteristic of TEs to multiply by transposition and the presence of orthologous
99 copies in different primate genomes, we split our reconstruction of LTR18A evolution into two
100 phases corresponding to transposition and speciation (Figure 1A). For each of the 46 sets of
101 LTR18A orthologs, we aligned orthologs using MAFFT and then reconstructed ortholog ancestor
102 and intermediate sequences using PRANK^{32,33}. Then, using the ancestor sequences for the 46
103 LTR18A orthologs, we aligned and reconstructed the LTR18A subfamily ancestor as well as
104 intermediates predating speciation. PRANK was chosen for ancestral sequence and phylogenetic
105 tree reconstruction due to its ability to model insertions and deletions. However, PRANK tends to
106 be biased towards insertions in our reconstruction. Thus, we manually curated sequences following
107 PRANK reconstruction for both ortholog ancestors and subfamily ancestors (Methods).

108

109 Next, we evaluated our reconstructed LTR18A sequences to see if they are consistent with those
110 derived from other methods. TE consensus sequences are often used as a representation of the
111 ancestral state of the subfamily. Excluding insertions and deletions, our reconstructed LTR18A
112 subfamily ancestor has ~5.9% substitution rate relative to the LTR18A consensus sequence, which
113 is lower than the 16.1% subfamily average. This suggests that although we start from different
114 elements and use different methodologies, both our reconstruction and the RepBase consensus are
115 approaching each other. In addition to substitutions, our reconstructed ancestor also has ~8.0%
116 insertions compared to the consensus. The insertions appear to be caused by the consensus
117 dropping bases if the majority of elements do not have the base in the alignment, as well as
118 PRANK's tendency to include insertions when alignable sequence is present in more than one

119 element. The MAFK motif enriched in LTR18A was present in both our reconstructed subfamily
120 ancestor and the RepBase consensus. Overall, the topology of our reconstructed phylogenetic tree
121 resembles the tree generated from all hg19 LTR18A elements (Supplemental Figure 1). One
122 feature of note occurs in node 43, two nodes from the root of the tree (Figure 1B). Relative to the
123 subfamily ancestor, node 43 has a 27bp insertion that contains a C/EBP-related factor motif (Figure
124 1C). When we examined ortholog ancestor reconstructions for this insertion, three ortholog
125 ancestors have an alignable 27bp insert, and the insertion is present in all present-day primate
126 orthologs (Supplemental Figure 2). In hg19, 13/181 elements contain the insert. The insert-
127 containing elements are spread throughout most of the hg19 LTR18A phylogenetic tree, which is
128 consistent with a deep ancestral origin for the insert and occurrence in node 43 of our
129 reconstruction. Additionally, the C/EBP motif is also found in the LTR18A consensus and
130 enriched in the subfamily relative to genomic background. If the C/EBP motif is functionally
131 important, the insertion of a second C/EBP motif could be an ancestral gain of function mutation.
132 In conclusion, our reconstruction is able to generate a subfamily ancestor similar to the RepBase
133 consensus and reveals evolutionary events that would otherwise be missed.

134

135 **Identification of important TFBS motifs in LTR18A enhancers**

136 We designed our LTR18A MPRA library to assay elements at two resolutions (Figure 2). In one
137 half, we synthesized motif-focused regions for 1225 LTR18A elements found across seven primate
138 genomes, 280 ancestral reconstruction elements, and the RepBase consensus (Figure 2A).
139 Specifically, we took the sequence of each element aligning to the first 160bp of our reconstructed
140 ancestral node 43 (Methods). This allowed us to focus on the effects of sequence variation for both
141 the MAFK motif and the C/EBP motif. In the other half of the library, we synthesized 160bp tiles

142 at 10bp intervals of all pre-speciation ancestral reconstruction elements, ortholog ancestors and
143 their present-day hg19 elements, and the LTR18A consensus (Figure 2B). We cloned LTR18A
144 motif-focused regions and tiles upstream of a pGL4 vector with the hsp68 promoter (Figure 2C).

145

146 To understand cell type effects, we tested LTR18A for enhancer activity in HepG2 and K562 cell
147 lines. We calculated enrichment scores for each element by taking the log₂ of the RNA over DNA
148 ratio followed by normalization to the basal hsp68 promoter. Normalizing to the basal promoter
149 allowed us to have the same reference point between cell lines. Active elements were defined as
150 those with enrichment scores greater than 1, representing elements that increase transcription by
151 greater than twofold. When we compare the distribution of enrichment scores for HepG2 and
152 K562, we find that LTR18A elements are generally more active in HepG2 than K562 (Figure 3A).
153 Out of 1506 motif-focused sequences tested, 1004 were classified as active in HepG2 while only
154 52 were classified as active in K562. For genomic LTR18A, 786 (123 from hg19) were active in
155 HepG2 and 31 (4 from hg19) were active in K562. Enrichment scores are positively but poorly
156 correlated between HepG2 and K562 despite high correlations between biological replicates
157 ($p < 2.2e-16$, Figure 3B, Supplemental Figure 3), implying differential sequence features required
158 for enhancer activity between cell lines.

159

160 To identify important sequence features for enhancer activity, we took advantage of the natural
161 sequence variation within LTR18A elements. Using AME motif enrichment analysis³⁴, we asked
162 if active elements were enriched for motifs compared to the rest of elements as background.
163 Overall, 34.5% (20/58) motifs were enriched in active elements in both HepG2 and K562 (Figure
164 3C). Of the shared motifs, AP-1 (JUN, FOS, and ATF family) motifs were in the top 10 most

165 enriched for both cell lines. Top 10 most enriched motifs that were cell line specific include the
166 C/EBP family motifs and BATF3 for HepG2 and NRF1 in K562. As an orthologous method, we
167 investigated if individual nucleotide positions are associated with enhancer activity. As this is
168 analogous to genome-wide association studies (GWAS) but focused on sequence variation within
169 a TE subfamily, which we term TE-WAS, we adapted the GWAS tool PLINK to find significant
170 nucleotides^{35,36}. In HepG2, 6/11 JUN (AP-1) motif bases and 8/11 DBP (C/EBP family) motif
171 bases are significantly associated with increased enhancer activity (Figure 3D). In K562, after we
172 adjusted our cutoff for active elements to be an enrichment score of at least 0.5 to increase the
173 number of active elements from 52 to 239, 4/11 JUN motif bases and 0/11 DBP motif bases are
174 significantly associated with increased enhancer activity. In summary, both motif enrichment and
175 TE-WAS approaches implicate AP-1 motifs as important to both HepG2 and K562 LTR18A
176 enhancer activity while C/EBP-related motifs are HepG2-specific.

177

178 To validate the importance of C/EBP and AP-1 motifs to enhancer activity, we created targeted
179 mutations in the motif regions of LTR18A elements. We chose DBP to represent the C/EBP family
180 and JUN to represent the AP-1 family. We selected pairs of LTR18A orthologs of which one has
181 the motif and the other does not by FIMO motif scanning³⁷. For elements with the motif, we
182 mutated the motif bases to low information nucleotides based on the PWM. For elements without
183 the motif, we changed the motif aligned region to the consensus motif bases. To quantify the effect
184 of motif mutations on enhancer activity, we took the log₂ ratio of each motif mutated LTR18A
185 sequence to its native sequence (Figure 3E, 3F). On average, DBP mutation gain and loss lead to
186 a 2.07-fold increase and 2.36-fold decrease in enhancer activity respectively in HepG2. In contrast,
187 the same DBP mutations have little effect in K562. JUN gain and loss lead to 1.49-fold increase

188 and 1.68-fold decrease in HepG2 enhancer activity and 1.17-fold increase and 1.2-fold decrease
189 in K562 enhancer activity. Both DBP and JUN mutagenesis results are consistent with our previous
190 findings based on motif association.

191

192 **Evolution of LTR18A enhancer activity linked to sequence evolution**

193 One of our primary goals was to understand how enhancer activity of LTR18A as a subfamily
194 changed over time. To address this question, we synthesized 160bp tiles at 10bp intervals across
195 each LTR18A ancestral sequence, ortholog ancestor, and hg19 element used in reconstruction
196 (Figure 2B). After obtaining enrichment scores, we estimated nucleotide activity scores across
197 each element to infer their relative effects on enhancer activity using the SHARPR software for
198 MPRA tiling designs¹². Due to overall low activity in K562, we focus on HepG2 for evolutionary
199 analysis. When examining nucleotide activity scores across the length of our reconstructed
200 LTR18A subfamily ancestor, we observe regions of increased activity over basal. The C/EBP and
201 AP-1 motifs that we previously identified to be important for enhancer activity are embedded
202 within the largest active region located near the start of the sequence (Supplemental Figure 6).
203 Across LTR18A elements of our reconstructed phylogenetic tree, we were able to confirm that
204 regions of increased SHARPR nucleotide activity were enriched for C/EBP and AP-1 motifs. As
205 SHARPR nucleotide activity scores could discover the same biologically meaningful sequences
206 as our previous analyses, we took the sum of activity scores across each LTR18A element and
207 annotated them in our tree (Figure 4A). From a broad perspective, we were able to make several
208 observations. First, the most divergent (leftmost) lineage on the tree loses enhancer activity early,
209 and enhancer activity throughout the lineage remains low to the present day (Figure 4C). The low
210 regulatory activity of the lineage could be linked to its relatively low rate of expansion (27/181

211 LTR18A elements in the lineage) (Supplemental Figure 7). This low activity lineage contrasts with
212 the rest of the tree where evolutionary intermediates exhibit relatively high activity followed by
213 less active elements at ortholog ancestor and present-day elements. Indeed, the overall trend
214 appears to be that enhancer activity decreases over time, as shown by the decrease in mean
215 SHARPR sum with increasing divergence from the LTR18A subfamily ancestor (Figure 4B). On
216 the other hand, there is an increase in activity in the middle lineages, some of which persists to the
217 ortholog ancestors and present-day elements (Figure 4D). Finally, enhancer activity of present day
218 hg19 LTR18A elements and their corresponding ortholog ancestors are positively correlated with
219 mostly small differences in activity, implying that post-speciation evolution has had small effects
220 on regulatory potential overall (Supplementary Figure 8).

221
222 To further investigate why enhancer activity changes in our LTR18A tree, we looked at differences
223 in C/EBP and AP-1 motif presence using DBP and JUN as representatives. When elements are
224 categorized by the number of DBP and JUN motifs, the number of motifs is positively correlated
225 with SHARPR sum (Figure 4E). Furthermore, DBP or JUN loss correlates with a decrease in
226 SHARPR sum, with rare motif gains generally corresponding to increased SHARPR sums (Figure
227 4F). Due to the significance of the DBP motif, we evaluated ancestral node 43 as the sole
228 evolutionary intermediate that gained a second motif through an insertion event (Figure 1B). The
229 motif gain leads to an increase in SHARPR sum of ~39%, which is similar to the average effect
230 size of the DBP motif (~38%). This effect is validated by mutagenesis of our LTR18A subfamily
231 ancestor and consensus to have the same 27bp insertion (34% and 32% increase respectively) as
232 well as ablation of the second DBP motif in ancestral node 43 (41% decrease). In summary,
233 sequence evolution, especially at the C/EBP and AP-1 motifs, directly affects the ability of

234 LTR18A to act as regulatory elements, and most mutations have led to a decrease in regulatory
235 potential.

236

237 **Evidence of selection for enhancer associated C/EBP and AP-1 motifs**

238 Given that LTR18A has regulatory potential in certain cellular contexts like HepG2, we explored
239 the possibility of host exaptation through the lens of selection. We first asked if LTR18A elements
240 in chimpanzee, gorilla, gibbon, baboon, rhesus macaque, and marmoset have increased
241 substitution rates compared to their human orthologs with respect to the distance between
242 genomes. On average, LTR18A orthologs have slightly elevated substitution rates (12-32%) than
243 the corresponding genomes (Supplemental Table 2). The increased substitution rate holds true
244 even when only considering masked regions of the genome. Although it is possible that the
245 genomic background rate includes regions under selection, the LTR18A substitution rates across
246 primate species are overall inconsistent with purifying selection for the subfamily. Furthermore,
247 both PhyloP and PhastCons scores at LTR18A elements provide no evidence of selection at the
248 subfamily level across 30 mammals, including 27 primates^{38,39} (Supplemental Figure 9).

249

250 While there is no evidence that LTR18A as a whole is under selection, it is possible that certain
251 regions within LTR18A are. We aligned LTR18A elements in each of our seven primate species
252 to the LTR18A consensus and tested sliding 10bp windows for increased conservation compared
253 to the average window. Overall, 29% (707/2429) of all 10bp windows are significantly more
254 conserved than the average window. The majority (84%) of conserved 10bp sliding windows are
255 shared across all seven primates for a total of 24.5% (85/347) possible 10bp windows covering
256 58% of the LTR18A consensus (208/357bp) being classified as conserved. Shared, conserved

257 regions defined by our sliding window analysis contain transcription factor motifs, including AP-1
258 and C/EBP (Figure 5A).

259

260 Since C/EBP and AP-1 motifs are critical for enhancer activity, we hypothesized that the motifs
261 provided by LTR18A have been under selection and consequently exhibit higher conservation than
262 expected under a neutral model of evolution. To obtain the background motif conservation rates,
263 we adapted a method previously used in yeast⁴⁰. Briefly, we take the sum of probabilities for all
264 sequences that match a motif PWM, with each sequence probability calculated starting from the
265 LTR18A consensus and the observed transition and transversion rate of the LTR18A subfamily.
266 As in previous analyses, we chose DBP and JUN to represent C/EBP and AP-1. Expected
267 conservation rates for DBP and JUN are consistent across species, ranging from 38.7% in
268 marmoset to 44.8% in human for DBP and 34.1% in marmoset to 39.3% in human for JUN (Table
269 1). Meanwhile, observed DBP and JUN conservation rates are on average 69.3% and 59.3%,
270 respectively, which is 26.4% and 21.6% higher than expected. This indicates that C/EBP and AP-1
271 motifs from the ancestral LTR18A sequence are being retained and may be under selection.
272 Measuring conservation from the LTR18A consensus includes the transposition phase of TE
273 evolution, which could select for C/EBP and AP-1 motifs due to enhancing transcription of the
274 ERV. To address conservation specifically during primate evolution, we recalculated conservation
275 rates by comparing human LTR18A elements to their primate orthologs. Generally, DBP and JUN
276 motifs are significantly more conserved than expected (Table 2). The one exception is JUN for the
277 human-chimpanzee comparison, which might be due to low human-chimpanzee divergence. We
278 also confirmed higher motif conservation rates during transposition+speciation and speciation
279 phases using simulations based on observed transition and transversion rates (Figure 5B, 5C).

280 Together, our analysis suggests that C/EBP and AP-1 motifs contributed by LTR18A have been
281 under selection in primates both before and after speciation.

282

283 **Human LTR18A has epigenetic signatures of active regulatory elements**

284 Our MPRA reveals that LTR18A elements have the sequence features to be activating regulatory
285 elements depending on cellular context. To explore the relationship between regulatory potential
286 from MPRA and enhancer function in the genome, we examined epigenetic marks in HepG2 and
287 K562 using ENCODE data⁴¹. We first profiled LTR18A elements overlapping ATAC peaks for
288 open chromatin, which is a common epigenetic feature for active regulatory elements. In HepG2,
289 LTR18A is not enriched for ATAC peaks, with only 5 LTR18A elements overlapping with peaks.
290 On the other hand, K562 has 11 overlapping LTR18A elements. This contrasts with the high
291 MPRA activity in HepG2 relative to K562. Additionally, H3K27ac and H3K4me1, histone marks
292 commonly associated with active enhancers, are also low across LTR18A in HepG2 and K562
293 (Supplemental Figure 10). We hypothesized that epigenetic repression of LTR18A may be the
294 cause for the lack of active enhancer marks in HepG2. Consistent with this hypothesis, repressive
295 histone mark H3K9me3 is enriched over LTR18A compared to the surrounding genomic region
296 (Supplemental Figure 10). These results suggest that although LTR18A elements possess the
297 sequence features necessary for enhancer activity, they can be epigenetically silenced.

298

299 While most of the LTR18A subfamily is unlikely to be active in HepG2 and K562, we sought to
300 ascertain the contribution of LTR18A to the regulatory genome across human cell types and
301 tissues. To get a global perspective, we overlapped LTR18A elements with candidate cis-
302 regulatory elements (cCREs) as defined by ENCODE Registry V2 across 839 cell/tissue types⁴¹.

303 Despite the limited number of cell/tissue types (25) that have full classification of cCREs, 69 of
304 198 (34.8%) LTR18A elements overlap with a cCRE, most of which (87%) have enhancer-like
305 signatures (ELS) in at least one cell/tissue type. This represents 29.3% of all LTR18A bases which
306 is about 3.1x enriched over the genomic background ($p < 3.5e-10$, BEDTools fisher). Among fully
307 classified cell/tissue types, keratinocytes have the highest number of LTR18A elements associated
308 with ELS, followed by PC-3 and PC-9 cell lines (Figure 6A). LTR18A is not restricted to a single
309 cell/tissue type, as some LTR18A elements are associated with cCREs in multiple cell/tissue types
310 (Figure 6B). Across all 839 cell/tissue types, cell types with the most LTR18A overlapping cCREs
311 largely consist of epithelial cells, such as MCF10A, mammary epithelial cells, esophagus epithelial
312 cells, and foreskin keratinocytes (Figure 6C). To corroborate cCRE results which are based on
313 DNase hypersensitivity, H3K27ac, H3K4me3, and CTCF ChIP-seq, LTR18A elements were
314 intersected with ENCODE ATAC-seq peaks across 46 cell/tissue types. Similar to cCREs,
315 LTR18A is especially enriched for ATAC peaks in epithelial cells/tissues foreskin keratinocytes
316 and esophagus mucosa (11.4x and 16.1x enrichment over background respectively, BEDTools
317 fisher). While certainly not comprehensive, the available epigenetic data supports an active
318 enhancer-like state for LTR18A with the highest enrichment in epithelial cells.

319

320 As LTR18A enhancer potential is influenced by sequence variation especially at transcription
321 factor binding sites, we sought to understand whether transcription factor motifs are associated
322 with active epigenetic states. Without considering cell/tissue type, we found no transcription factor
323 motif to be significantly associated with LTR18A overlapping cCREs relative to other LTR18A.
324 Due to the cell type specific nature of most enhancers, we identified motifs enriched in cCRE
325 associated LTR18A in the top cell/tissue types (Figure 6D). Many of the most common motifs are

326 of AP-1 transcription factors. Another common motif is NFIC, which is consistent with an
327 activating role previously described in cancer and could serve a similar role in activating LTR18A
328 elements⁴². Of note, the C/EBP-related factor HLF is enriched only in the MCF10A cell line. Using
329 ATAC data, we confirmed AP-1 and NFIC motifs as enriched in LTR18A elements associated
330 with active epigenetic states in foreskin keratinocytes and esophagus mucosa. Altogether, these
331 results suggest that LTR18A elements become epigenetically activated in epithelial cells primarily
332 through AP-1 transcription factors and NFIC.

333

334 **Discussion**

335 Since Britten and Davidson first hypothesized how repetitive elements could influence the
336 development of gene regulatory networks, a growing number of studies have shown the
337 contribution of TEs as regulatory modules⁴³. Using LTR18A as a representative subfamily, we
338 performed the first systematic functional testing of regulatory potential for a TE subfamily using
339 MPRA. By taking advantage of the natural sequence variation across elements, we identify AP-1
340 and C/EBP-related motifs as important drivers of LTR18A regulatory activity. This regulatory
341 activity is highly dependent on cell context, with LTR18A displaying much higher activity in
342 HepG2 than in K562. However, the sequence potential for regulatory activity does not necessarily
343 reflect activity in the genome, as shown by LTR18A elements rarely associating with active
344 epigenetic marks in HepG2. Due to general repression of TEs, we believe that similarly silenced
345 TEs with the potential for enhancer activity may be common. These inactive TEs may be latent
346 under epigenetic control, but there remains the possibility that a changing epigenome such as
347 during tumorigenesis can reactivate them⁴⁴.

348

349 Another unique aspect of this study is leveraging the phylogenetic relationship between LTR18A
350 elements within human and across primate species to investigate the origin and evolution of
351 regulatory activity in the subfamily. Previous research has implicated two evolutionary paths
352 through which TE sequence can contribute to the spread of regulatory modules. The first case is
353 when the ancestral TE originally possesses the driving regulatory features, such as the p53 binding
354 site in LTR10 and MER61 or the STAT1 binding site in MER41B^{20,26}. A second possibility exists
355 where the ancestral TE gains the regulatory module in one lineage through mutation before
356 amplification, such as the 10bp deletion in ISX relative to ISY in *D. miranda* that recruits the
357 MSL-complex⁴⁵. In the LTR18A family, we observe both scenarios. Both C/EBP and AP-1 motifs
358 are found in the LTR18A consensus and our reconstructed subfamily ancestor, and many elements
359 retain the motifs to the present day. Divergence from the ancestor over time, especially at the two
360 motifs, is correlated with a decrease in regulatory activity. In addition to the two consensus motifs,
361 a second C/EBP motif is gained through an insertion at an early evolutionary timepoint. This
362 second C/EBP motif further increases the regulatory potential of LTR18A. Ultimately, however,
363 few present-day elements have maintained the second motif. This could be explained by negative
364 selection or a deletion bias from the sequence similarity of the insertion with the upstream
365 sequence. It also plausible that our evolutionary reconstruction makes an incorrect assumption
366 about the timing of the second C/EBP motif, and each one occurred independently rather than
367 through a common ancestor. If this scenario is true, recurrent insertions in TEs may be more
368 common than previously thought.

369

370 An intriguing possibility is the relationship between TE regulatory potential and genomic
371 expansion. In our reconstructed LTR18A phylogenetic tree, we observe loss of enhancer activity

372 in the leftmost lineage going as far back as its lineage ancestor. This low enhancer activity lineage
373 corresponds to the earliest diverging branch in the human LTR18A subfamily phylogenetic tree
374 and composes only $\sim 1/6$ of all elements. On the other hand, the major lineage of LTR18A has
375 enhancer activity throughout transposition. The stark contrast between the two lineages in
376 enhancer activity and abundance leads us to speculate that the regulatory potential of LTR18A was
377 directly related to its ability to expand in the genome. This is perhaps unsurprising, as transcription
378 is typically the first step of transposition and provides the substrate for integration of
379 retrotransposons. However, one important consequence is that transcription factor binding sites
380 that contribute to TE regulatory potential could be enriched within a subfamily due to biased
381 lineage amplification. This appears to have been the case for the recently reclassified LTR7
382 subfamilies, each of which possess a unique set of transcription factor motifs and underwent a
383 wave of genomic expansion to fill different early embryonic niches⁴⁶. It will be important for future
384 studies to distinguish between selection and passive enrichment of transcription factor binding
385 sites through lineage amplification.

386

387 To compare ancestral and present day LTR18A elements, we tested all elements within the same
388 cell line. This assumes that HepG2 and K562 cells provide the same *trans* environment as the
389 equivalent primate and ancestral cell types. Previous studies suggest that transcription factor
390 binding and subsequent activation of transcription are deeply conserved from humans to flies^{47,48}.
391 Klein et al. make a similar assumption in their study of liver enhancer evolution in primates and
392 find the same general trend that present-day elements have lost enhancer activity relative to the
393 ancestral state¹⁸.

394

395 Most TEs are thought be under neutral evolution and do not significantly impact phenotype. We
396 find that LTR18A elements as a whole have higher mutation rates than genomic average and do
397 not exhibit signs of selection based on phyloP and phastCons scores. Despite the lack of evidence
398 for selection at the element level, AP-1 and C/EBP binding motifs found within LTR18A are more
399 conserved than expected under the neutral model of evolution. This suggests that selection does
400 not need to apply to entire TEs and instead acts on functional units found within each element.
401 Indeed, we find that at least a third of LTR18A elements have enhancer associated epigenetic
402 marks, and in some cell/tissue types, the active elements are enriched for the conserved AP-1
403 motif. Although the C/EBP motif is not significantly enriched with active elements outside of
404 MCF10A, we suspect that the motif is important in other cell/tissue types that have yet to be
405 profiled.

406

407 **Methods:**

408 **LTR18A manual curation for ancestral reconstruction**

409 We downloaded RepeatMasker 4.0.5 (Repeat Library 20140131) annotations for human (hg19),
410 chimpanzee (panTro4), gorilla (gorGor3), gibbon (nomLeu3), rhesus macaque (rheMac3), and
411 marmoset (calJac3) genomes⁴⁹. For baboon (papAnu2) which is not available on
412 www.repeatmasker.org, we ran RepeatMasker 4.1.0 using the RepBase RepeatMasker library
413 20170127. Since LTR18A consensus sequences are 98% similar between the two repeat libraries,
414 we believe that most if not all LTR18A elements will be identified in papAnu2 in the same way
415 as the other primate genomes. For the closest two subfamilies, LTR18B and LTR18C consensus
416 sequences are ~75% and 67% similar to the LTR18A consensus respectively.

417 For manual curation, we examined the alignment of each annotated LTR18A element and removed
418 the element if it satisfied any of our filtering criteria (Supplemental Table 1). First, we exclude
419 LTR18A elements that have significant alignments to LTR18B or LTR18C. RepeatMasker outputs
420 alignment scores for each repetitive element, some of which have multiple significant alignment
421 scores for different subfamily consensus sequences. RepeatMasker then chooses the subfamily
422 with the highest alignment score to annotate elements with the same ID. A consequence of this
423 method is that fragmented elements can be annotated for the same subfamily even when the highest
424 scoring alignment differs for each fragment. Since LTR18B and LTR18C consensus sequences are
425 ~75% and 67% similar to LTR18A respectively, some LTR18A elements have significant
426 alignments to LTR18B and/or LTR18C. Thus, we discard these elements with multiple possible
427 alignments to avoid ambiguity from subfamily assignment. Second, we use paired LTR
428 information to remove LTR18A elements that have discordant annotations. Due to the mechanism
429 of ERV retrotransposition, we expect non-solo LTR18A elements to exist as same orientation pairs

430 that are separated by the ERV internal region. Using this logic, we reasoned that paired LTRs that
431 are assigned to different subfamilies have uncertain annotation.

432 To find LTR18A ortholog sets for ancestral reconstruction, we searched for LTR18A element pairs
433 that fulfilled several requirements. First, the hg19 LTR18A elements must have orthologs in
434 chimpanzee and gorilla. Second, elements must have orthologs in at least two of the other primate
435 species: gibbon, baboon, rhesus macaque, and marmoset. Third, hg19 LTR18A elements must be
436 >250bp (>70% of consensus) in length. Finally, both elements of a pair need to pass all
437 requirements to be selected for ancestral reconstruction. Orthologs were defined using the chain
438 files from UCSC to find LTR18A elements within the same syntenic blocks³¹.

439 Ancestral reconstruction of both ortholog ancestors and subfamily ancestors used MAFFT and
440 PRANK followed by manual curation^{32,33}. To generate ortholog ancestors, we aligned ortholog
441 sets (e.g. human, chimpanzee, gorilla, gibbon, baboon orthologs) using MAFFT multiple sequence
442 alignment. We used the alignments to produce ancestral and intermediate sequences as well as the
443 phylogenetic tree using PRANK. The PRANK phylogenetic trees typically reflected the expected
444 evolutionary relationship between the seven primate species. Next, we manually adjusted ortholog
445 ancestors to remove unlikely insertions. We focused on insertions rather than deletions due to the
446 possibility of insertions propagating up the tree. We determined insertion sites by examining the
447 multiple sequence alignment of ortholog ancestors and finding gaps in the alignment created by
448 insertions in only a few ortholog ancestors. Generally, we used parsimony when deciding to keep
449 or remove an insertion. For example, if the insertion is present in only one primate lineage, then it
450 is less likely for the insertion to have existed in the ortholog ancestor. Our reasoning is that an
451 insertion in the ortholog ancestor and subsequent deletion in the other lineages requires at least
452 two mutation events, whereas a single insertion in one primate lineage requires only one mutation

453 event. After manual curation of ortholog ancestors, we used MAFFT and PRANK to reconstruct
454 the phylogenetic tree and sequences of LTR18A subfamily ancestral sequences. We again applied
455 parsimony to manually adjust the LTR18A subfamily ancestor.

456

457 **LTR18A MPRA library construction**

458 The MPRA library was designed to consist of a motif-focused half and a tiling half. To design the
459 motif-focused half of our MPRA library, we took advantage of the relatedness of TEs within the
460 same subfamily. Similar to RepeatMasker, we can align all LTR18A elements to a reference
461 sequence. Instead of using the subfamily consensus sequence, we used our reconstructed ancestral
462 node 43 to perform pairwise global alignments to all present-day and reconstructed elements.
463 Then, we took the sequence of each element aligned to the first 160bp of ancestral node 43. We
464 filtered out elements that have fewer than 70bp due to deletions and elements that have more than
465 160bp due to insertions. We also removed elements that contain a restriction site that we used for
466 cloning. In total, 1255/1387 RepeatMasker annotated LTR18A elements across seven primate
467 genomes and all 280 reconstructed elements were included. For the tiling half of the library, we
468 selected all pre-speciation ancestral reconstruction elements, ortholog ancestors and their present-
469 day hg19 elements, and the LTR18A consensus. We then synthesized 160bp tiles at 10bp intervals
470 spanning each selected element. In addition to motif-focused and tiled sequences, we selected 456
471 elements for reverse complements, 37 pairs of elements for JUN mutagenesis, and 46 pairs of
472 elements for DBP mutagenesis. Elements for mutagenesis were chosen based on the closest
473 primate ortholog with/without the motif. JUN motifs were mutated to TCACCAATGGT and DBP
474 motifs were mutated to TCCCACAGCAT. Non-motif containing elements were mutated to
475 GCTGAGTCATG for JUN and ATTATGTAACC for DBP. For positive and negative controls,

476 we selected 223 regions from a previous study by Ernst et al.¹². 30 dinucleotide shuffled LTR18A
477 RepBase consensus sequences were included as a second set of negative controls⁵⁰. Each
478 synthesized sequence was tagged with 10 unique barcodes. To control for differences in overall
479 library activity between cell lines, we included a set of sequences that would leave only the basal
480 hsp68 promoter tagged with 300 barcodes. Oligos were ordered from Agilent and structured as
481 follows: 5' priming sequence containing NheI site (CGGTATCTAAGAgctagcGT)/CRE/EcoRI
482 site/Filler (if necessary)/BglII site/BamHI site/constant 'G'/barcode/constant 'A'/AgeI site/3'
483 priming site (ATTAGCATGTCGTG)¹¹. Total length of oligos was 230bp. In total, 5918 elements
484 were synthesized using 59470 unique barcodes.

485 The MPRA library was constructed as previously described with some adjustments. An AgeI site
486 was introduced upstream of the SV40 polyA signal and the BamHI site downstream of the SV40
487 polyA signal was deleted using the QuikChange Lightning site-directed mutagenesis kit (Agilent).
488 Synthesized oligos were amplified with 0.05pmol of template per 50µL PCR reaction for seven
489 cycles using MPRA library amplification primers. A total of 32 reactions were performed.
490 Following amplification and gel purification, oligos were cloned into a pGL backbone with the
491 AgeI insert using NheI and AgeI sites. Multiple ligations were pooled, purified by PCR cleanup
492 (Nucleospin), and transformed into 5-alpha electrocompetent *E. coli* (NEB). The hsp68 promoter
493 driving dsRed reporter was cloned using EcoRI and BamHI sites. The MPRA library with the
494 hsp68 promoter and dsRed reporter was purified and transformed into *E. coli* before plasmid
495 extraction. The final library was concentrated by ethanol precipitation.

496

497 **Cell culture and library transfection**

498 Cell culture and library transfections were performed as previously described¹¹. K562 cells were
499 grown in RPMI 1640 with L-glutamine (Gibco) + 10% Fetal Bovine Serum (FBS) + 1%
500 penicillin/streptomycin. HepG2 cells were grown in Dulbecco's Modified Eagle Medium with
501 high glucose, L-glutamine, and without sodium pyruvate + 10% FBS + 1%
502 penicillin/streptomycin. For each of three replicates, 5 μ g of library was transfected into 1.2 million
503 cells using Neon electroporation (Life Technologies). For K562, electroporation parameters were
504 three 10 millisecond pulses at 1450V. For HepG2, electroporation parameters were three 20
505 millisecond pulses at 1230V. As a transfection control, 0.5 μ g of pmaxGFP (Lonza) was used.

506

507 **Measurement of library expression**

508 RNA extraction was performed 24 hours after transfection using PureLink RNA Mini Kit with on-
509 column DNase treatment (Life Technologies) followed by DNase I treatment using TURBO DNA-
510 free kit (Invitrogen). Samples were prepared for RNA-seq as previously described¹¹. First strand
511 cDNA synthesis was performed using Superscript III Reverse Transcriptase (Life Technologies).
512 Barcodes were amplified from cDNA from three transfections and three technical replicates of
513 DNA from the plasmid library. Amplified barcodes were digested with KpnI and EcoRI and ligated
514 to Illumina adapters. Ligation products were further amplified, after which replicates and plasmid
515 library DNA input were pooled for sequencing. We obtained over 1000x average coverage for
516 each transfection replicate and the DNA input. For each tested element, we added up read counts
517 for all of its barcodes and filtered out those with fewer than 5 total counts in any transfection
518 replicate or DNA input. Reads were then normalized to counts per million (CPM). Expression of
519 an element was calculated as RNA CPM/DNA CPM. Expression was normalized to the average
520 of Basal construct transfection replicates. Finally, enrichment score was calculated as the log₂ of

521 normalized expression. Enrichment scores of elements were highly reproducible across
522 transfection replicates in HepG2 (average $R^2=0.904$) while moderately reproducible in K562
523 (average $R^2=0.666$) (Supplemental Figure 3). We confirmed that orientation does not have large
524 effects on enrichment score in both HepG2 and K562 (Supplemental Figure 4). We also found that
525 selected control sequences from Ernst et al. follow expected trends for both their original
526 annotations as well as redefined annotations based on expression values from Ernst et al. MPRA
527 results (Supplemental Figure 5). Enrichment scores of elements are provided in Supplemental
528 Data.

529

530 **TE-WAS analysis of nucleotides and motifs**

531 LTR18A sequences were first globally aligned pairwise to the ancestral node 43 sequence as
532 reference⁵¹. Individual pairwise alignments were then combined based on the common reference.
533 Positions that had bases (not gaps) in less than 20% of all LTR18A sequences were removed. This
534 filter retained all consensus base positions.

535 GWAS analysis tool PLINK was used to identify nucleotides significantly associated with the
536 phenotype, such as MPRA activity/inactivity or ATAC peak³⁶. We limited tested nucleotides at
537 each position to the most common nucleotide at the position across LTR18A sequences to give us
538 greater confidence based on sample size. We ran PLINK association analysis using the above-
539 described alignment and MPRA active/inactive annotations for each element based on enrichment
540 score. Nucleotides were deemed significant if $p\text{-value} < 5e-5$.

541 From the list of significant nucleotides in TE-WAS, we identified transcription factor motifs that
542 are overrepresented based on information content. Information content at each significant
543 nucleotide was calculated from each motif's position frequency matrix with the background

544 nucleotide frequencies of 0.25. The information content of significant nucleotides within each
545 motif was then compared to a background expectation derived from 1000 random shuffles of
546 significant nucleotides for the phenotype. Motifs were identified if they had higher information
547 content from significant nucleotides than background using t-test and more than significant
548 nucleotide within the motif.

549

550 **Evolutionary analysis using SHARPR**

551 From tiled MPRA, we calculated regulatory activity for full length elements using SHARPR with
552 a few adjustments¹². For each tile of an element, the previously calculated enrichment score was
553 used as input for SHARPR infer with the default varpriors of 1 and 50. Each inferred 10bp step
554 was then normalized to the mean inferred value for randomly shuffled Basal elements as
555 background. SHARPR combine and interpolate commands were used to generate the SHARPR
556 nucleotide activity scores. Finally, full length element activities were calculated as the sum of
557 nucleotide scores across each element.

558 To validate the SHARPR approach, we identified motifs that were enriched in peaks, or regions
559 of high nucleotide activity. Peaks were defined as regions with nucleotide activity scores greater
560 than three standard deviations above the Basal mean. Enriched motifs were then identified in peak
561 regions using AME using shuffled sequence as background³⁴.

562

563 **Transcription factor motif conservation**

564 For sliding window conservation analysis, we aligned all present-day genomic LTR18A elements
565 to the RepBase consensus sequence using the previously defined method. Conservation, defined
566 as percent match to the consensus, was calculated for each 10bp window for each element in each

567 species. Windows with gaps or degenerate bases in at least half of the total window length (≥ 5)
568 were excluded. The mean conservation was then calculated for each 10bp window separately for
569 each species. Windows were determined to be significantly conserved using t-test comparing
570 conservation across elements in the window against conservation across all windows, with a p-
571 value threshold of 0.05 after Bonferroni correction. Only windows that were conserved in all seven
572 primate species were kept for further analysis. Motif scanning by FIMO was performed to find
573 transcription factor motifs fully within conserved windows³⁷.

574 For JUN and DBP transcription factor motif conservation analysis, transition and transversion rates
575 in the LTR18A subfamily were calculated for each species. The neutral expectation for motif
576 conservation was calculated as previously described⁴⁰. We identified all kmers of the motif length
577 which are found by FIMO³⁷. The total motif conservation probability was calculated as the sum of
578 the probabilities for each motif kmer. We used the RepBase consensus sequence as the ancestral
579 LTR18A state. To represent post-speciation conservation, we used hg19 orthologs as the reference
580 to compare to other primate LTR18A elements. The observed motif conservation rate was
581 calculated for each species based on the percentage of elements that retain the motif. Elements
582 with gaps in the alignment to its reference were excluded. Statistical significance was determined
583 by one sample test of proportions and a p-value threshold of 0.05. We also simulated transcription
584 factor motif conservation rates for each primate species. Each simulation consisted of randomly
585 mutating nucleotides in the motif region of each LTR18A element based on the observed transition
586 and transversion rates. 1000 simulations were performed for each motif.

587

588 **Overlap of LTR18A with genomic annotations**

589 The cCRE genome annotations and various epigenetic datasets such as ATAC-seq, histone ChIP-
590 seq, and WGBS were downloaded from ENCODE⁴¹. The phyloP and phastCons scores were
591 downloaded from ENCODE and converted to bedGraph³¹. Overlaps with LTR18A elements were
592 obtained by BEDTools intersect with the criteria of at least 50% LTR18A length overlapping with
593 a cCRE or epigenetic mark peak⁵². Enrichment of LTR18A in cCREs and ATAC peaks was
594 obtained by BEDTools fisher using the same criteria. Heatmaps at and around LTR18A were
595 generated using deepTools⁵³.

596

597 **Identification of motifs associated with cCRE overlapping LTR18A**

598 Fisher's exact test was used to determine if transcription factor binding motifs in LTR18A
599 elements are associated with cCRE overlap. Motifs that had adjusted p-values below 0.05 were
600 considered significant. The top six cell/tissue types were selected for analysis as they provided the
601 greatest number of LTR18A elements overlapping cCREs.

602

603 **Subfamily age estimate**

604 The average divergence, weighted by copy length, was calculated for the LTR18A subfamily using
605 the RepeatMasker output for hg19. The age was obtained by using the average divergence and the
606 average mammalian genome mutation rate of 2.2×10^{-9} per base per year⁵⁴.

607

- 608 1. King, M. C. & Wilson, A. C. Evolution at two levels in humans and chimpanzees. *Science*
609 (80-). **188**, 107–116 (1975).
- 610 2. Banerji, J., Rusconi, S. & Schaffner, W. Expression of a beta-globin gene is enhanced by
611 remote SV40 DNA sequences. *Cell* **27**, 299–308 (1981).
- 612 3. Moreau, P. *et al.* The SV40 72 base repair repeat has a striking effect on gene expression
613 both in SV40 and other chimeric recombinants. *Nucleic Acids Res.* **9**, 6047–6068 (1981).
- 614 4. Wray, G. A. The evolutionary significance of cis-regulatory mutations. *Nat. Rev. Genet.*
615 *2007* **8**, 206–216 (2007).
- 616 5. Hong, J. W., Hendrix, D. A. & Levine, M. S. Shadow enhancers as a source of evolutionary
617 novelty. *Science* vol. 321 1314 (2008).
- 618 6. Cannavò, E. *et al.* Shadow Enhancers Are Pervasive Features of Developmental Regulatory
619 Networks. *Curr. Biol.* **26**, 38–51 (2016).
- 620 7. Patwardhan, R. P. *et al.* High-resolution analysis of DNA regulatory elements by synthetic
621 saturation mutagenesis. *Nat. Biotechnol.* **27**, 1173–5 (2009).
- 622 8. Patwardhan, R. P. *et al.* Massively parallel functional dissection of mammalian enhancers
623 in vivo. *Nat. Biotechnol.* **30**, 265–70 (2012).
- 624 9. Melnikov, A. *et al.* Systematic dissection and optimization of inducible enhancers in human
625 cells using a massively parallel reporter assay. *Nat. Biotechnol.* **30**, 271–277 (2012).
- 626 10. Kwasnieski, J. C., Mogno, I., Myers, C. A., Corbo, J. C. & Cohen, B. A. Complex effects
627 of nucleotide variants in a mammalian cis-regulatory element. *Proc. Natl. Acad. Sci. U. S.*
628 *A.* **109**, 19498–503 (2012).
- 629 11. Kwasnieski, J. C., Fiore, C., Chaudhari, H. G. & Cohen, B. A. High-throughput functional
630 testing of ENCODE segmentation predictions. *Genome Res.* **24**, 1595–602 (2014).

- 631 12. Ernst, J. *et al.* Genome-scale high-resolution mapping of activating and repressive
632 nucleotides in regulatory regions. *Nat. Biotechnol.* **34**, 1180–1190 (2016).
- 633 13. Chaudhari, H. G. & Cohen, B. A. Local sequence features that influence AP-1 cis-regulatory
634 activity. *Genome Res.* **28**, 171 (2018).
- 635 14. Vockley, C. M. *et al.* Massively parallel quantification of the regulatory effects of
636 noncoding genetic variation in a human cohort. *Genome Res.* **25**, 1206–1214 (2015).
- 637 15. Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants
638 using a Multiplexed Reporter Assay. *Cell* **165**, 1519–1529 (2016).
- 639 16. Ulirsch, J. C. *et al.* Systematic Functional Dissection of Common Genetic Variation
640 Affecting Red Blood Cell Traits. *Cell* **165**, 1530–1545 (2016).
- 641 17. Arnold, C. D. *et al.* Quantitative genome-wide enhancer activity maps for five *Drosophila*
642 species show functional enhancer conservation and turnover during cis-regulatory
643 evolution. *Nat. Genet.* **46**, 685–692 (2014).
- 644 18. Klein, J. C., Keith, A., Agarwal, V., Durham, T. & Shendure, J. Functional characterization
645 of enhancer evolution in the primate lineage. *Genome Biol.* **19**, 99 (2018).
- 646 19. Feschotte, C. Transposable elements and the evolution of regulatory networks. *Nat. Rev.*
647 *Genet.* **9**, 397–405 (2008).
- 648 20. Wang, T. *et al.* Species-specific endogenous retroviruses shape the transcriptional network
649 of the human tumor suppressor protein p53. *Proc. Natl. Acad. Sci. U. S. A.* **104**, 18613–8
650 (2007).
- 651 21. Bourque, G. *et al.* Evolution of the mammalian transcription factor binding repertoire via
652 transposable elements. *Genome Res.* **18**, 1752–62 (2008).
- 653 22. Kunarso, G. *et al.* Transposable elements have rewired the core regulatory network of

- 654 human embryonic stem cells. *Nat. Genet.* **42**, 631–634 (2010).
- 655 23. Schmidt, D. *et al.* Waves of retrotransposon expansion remodel genome organization and
656 CTCF binding in multiple mammalian lineages. *Cell* **148**, 335–348 (2012).
- 657 24. Sundaram, V. *et al.* Widespread contribution of transposable elements to the innovation of
658 gene regulatory networks. *Genome Res.* **24**, 1963–76 (2014).
- 659 25. Lynch, V. J., Leclerc, R. D., May, G. & Wagner, G. P. Transposon-mediated rewiring of
660 gene regulatory networks contributed to the evolution of pregnancy in mammals. *Nat Genet*
661 **43**, 1154–1159 (2011).
- 662 26. Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory evolution of innate immunity through
663 co-option of endogenous retroviruses. *Science* **351**, 1083–7 (2016).
- 664 27. Pehrsson, E. C., Choudhary, M. N. K., Sundaram, V. & Wang, T. The epigenomic landscape
665 of transposable elements across normal human development and anatomy. *Nat. Commun.*
666 *2019 101* **10**, 1–16 (2019).
- 667 28. Jacques, P. É., Jeyakani, J. & Bourque, G. The Majority of Primate-Specific Regulatory
668 Sequences Are Derived from Transposable Elements. *PLoS Genet.* **9**, 1003504 (2013).
- 669 29. Trizzino, M. *et al.* Transposable elements are the primary source of novelty in primate gene
670 regulation. *Genome Res.* **27**, 1623–1633 (2017).
- 671 30. Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community resource
672 of transposable element families, sequence models, and genome annotations. *Mob. DNA* **12**,
673 1–14 (2021).
- 674 31. Kuhn, R. M., Haussler, D. & James Kent, W. The UCSC genome browser and associated
675 tools. *Brief. Bioinform.* **14**, 144–161 (2013).
- 676 32. Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid multiple

- 677 sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–66
678 (2002).
- 679 33. Löytynoja, A. Phylogeny-aware alignment with PRANK. *Methods Mol. Biol.* **1079**, 155–
680 170 (2014).
- 681 34. McLeay, R. C. & Bailey, T. L. Motif Enrichment Analysis: A unified framework and an
682 evaluation on ChIP data. *BMC Bioinformatics* **11**, 1–11 (2010).
- 683 35. Purcell, S. *et al.* PLINK: A Tool Set for Whole-Genome Association and Population-Based
684 Linkage Analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 685 36. Chang, C. C. *et al.* Second-generation PLINK: Rising to the challenge of larger and richer
686 datasets. *Gigascience* **4**, 7 (2015).
- 687 37. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.
688 *Bioinformatics* **27**, 1017–1018 (2011).
- 689 38. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
690 genomes. *Genome Res.* **15**, 1034–1050 (2005).
- 691 39. Pollard, K. S., Hubisz, M. J., Rosenbloom, K. R. & Siepel, A. Detection of nonneutral
692 substitution rates on mammalian phylogenies. *Genome Res.* **20**, 110 (2010).
- 693 40. Doniger, S. W., Huh, J. & Fay, J. C. Identification of functional transcription factor binding
694 sites using closely related *Saccharomyces* species. *Genome Res.* **15**, 701 (2005).
- 695 41. Moore, J. E. *et al.* Expanded encyclopaedias of DNA elements in the human and mouse
696 genomes. *Nature* **583**, 699–710 (2020).
- 697 42. Fane, M., Harris, L., Smith, A. G. & Piper, M. Nuclear factor one transcription factors as
698 epigenetic regulators in cancer. *Int. J. Cancer* **140**, 2634–2641 (2017).
- 699 43. Britten, R. J. & Davidson, E. H. Repetitive and non-repetitive DNA sequences and a

- 700 speculation on the origins of evolutionary novelty. *Q. Rev. Biol.* **46**, 111–138 (1971).
- 701 44. Jang, H. S. *et al.* Transposable elements drive widespread expression of oncogenes in
702 human cancers. *Nat. Genet.* **51**, 611–617 (2019).
- 703 45. Ellison, C. & Bachtrog, D. Dosage Compensation via Transposable Element Mediated
704 Rewiring of a Regulatory Network. *Science (80-.)*. **342**, 846–850 (2013).
- 705 46. Carter, T. A. *et al.* Mosaic cis-regulatory evolution drives transcriptional partitioning of
706 HERVH endogenous retrovirus in the human embryo. *bioRxiv* 2021.07.08.451617 (2021)
707 doi:10.1101/2021.07.08.451617.
- 708 47. Nitta, K. R. *et al.* Conservation of transcription factor binding specificities across 600
709 million years of bilateria evolution. *Elife* **2015**, (2015).
- 710 48. Stampfel, G. *et al.* Transcriptional regulators form diverse groups with context-dependent
711 regulatory functions. *Nature* **528**, 147 (2015).
- 712 49. Smit, A., Hubley, R. & Green, P. RepeatMasker Open-4.0. 2013-2015
713 <<http://www.repeatmasker.org>>.
- 714 50. Bailey, T. L., Johnson, J., Grant, C. E. & Noble, W. S. The MEME Suite. *Nucleic Acids*
715 *Res.* **43**, W39–W49 (2015).
- 716 51. Needleman, S. B. & Wunsch, C. D. A general method applicable to the search for
717 similarities in the amino acid sequence of two proteins. *J. Mol. Biol.* **48**, 443–453 (1970).
- 718 52. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
719 features. *Bioinformatics* **26**, 841–842 (2010).
- 720 53. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data
721 analysis. *Nucleic Acids Res.* **44**, W160 (2016).
- 722 54. Kumar, S. & Subramanian, S. Mutation rates in mammalian genomes. *Proc. Natl. Acad.*

723 *Sci.* **99**, 803–808 (2002).

724

725

726 **Acknowledgements**

727 We would like to thank J. Hoisington-López and M.L. Jaeger from The Edison Family Center for
728 Genome Sciences & Systems Biology (CGSSB) for assistance with sequencing. This work was
729 funded by NIH grant numbers R01HG007175, U01CA200060, U01HG009391, U41HG010972,
730 and U24HG012070. A.Y.D. was supported by NHGRI training grant T32 HG000045.

731

732 **Contributions**

733 A.Y.D., V.S., and T.W. designed the study. X.Z. contributed to evolutionary analysis. N.O.J. and
734 N.L.S. contributed to TE-WAS analysis. A.Y.D. performed the MPRA with contributions by
735 H.G.C. and B.A.C. in design and analysis. The manuscript was prepared by A.Y.D. and T.W. with
736 input from authors.

RepBase consensus, ancestral node 45 (#45#, subfamily ancestor), and ancestral node 43. Motifs in the sequences are boxed. DBP is shown to represent C/EBP-related motifs.

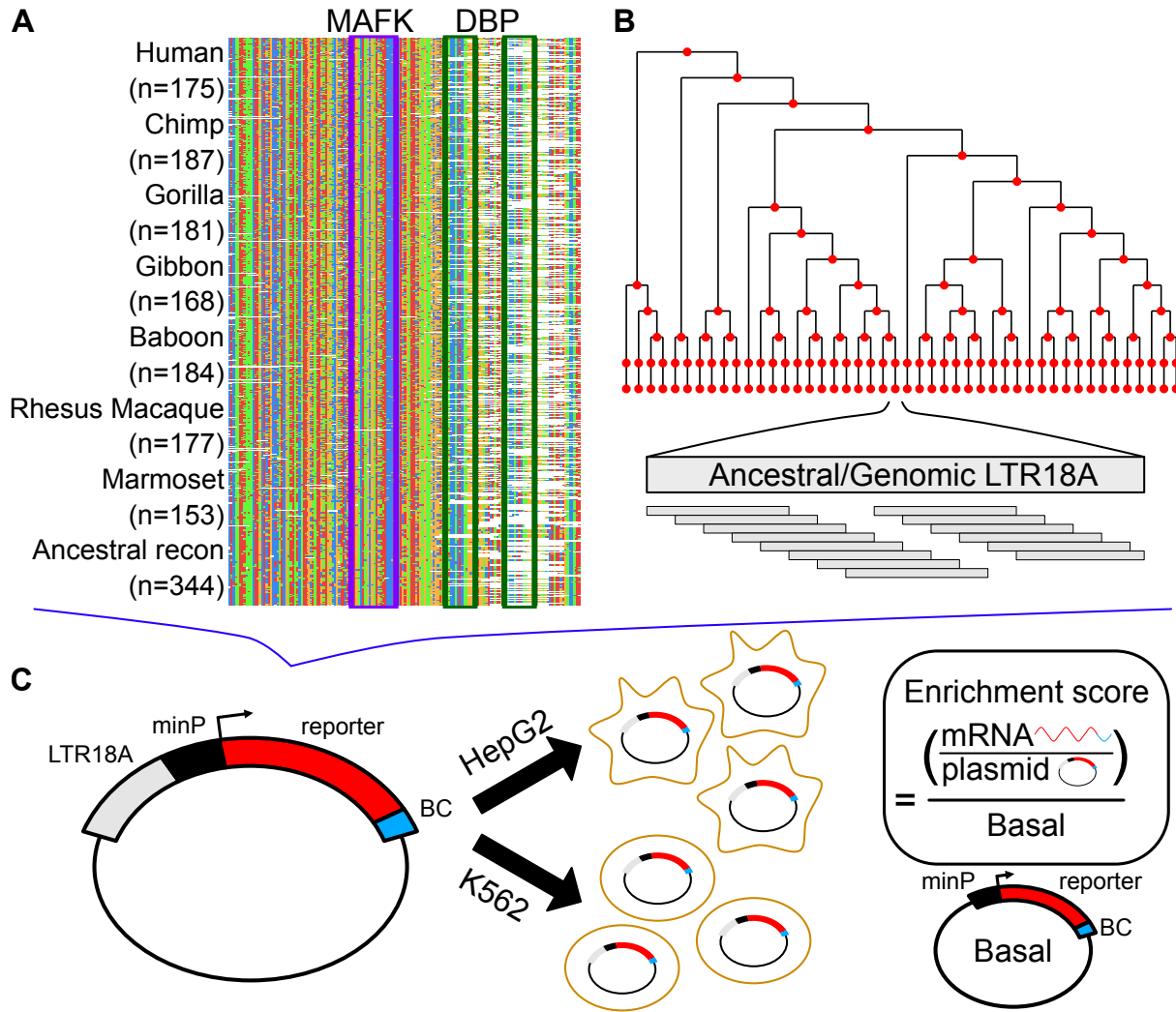


Figure 2: Schematic of MPRA. A) Sequence alignment of motif-focused regions to test primate and ancestral reconstructed LTR18A elements. MAFK and DBP motif regions are boxed. B) Tiling of ancestral and hg19 genomic LTR18A elements in reconstructed phylogenetic tree. All elements were tiled with 160bp tiles at 10bp intervals. C) Plasmid construct and enrichment score calculation. Each LTR18A fragment was integrated upstream of a minimal promoter (minP) and tagged with 10 unique barcodes (BC). The MPRA library was transfected into HepG2 and K562 cells. Enrichment scores are \log_2 ratios of RNA/DNA normalized to Basal.

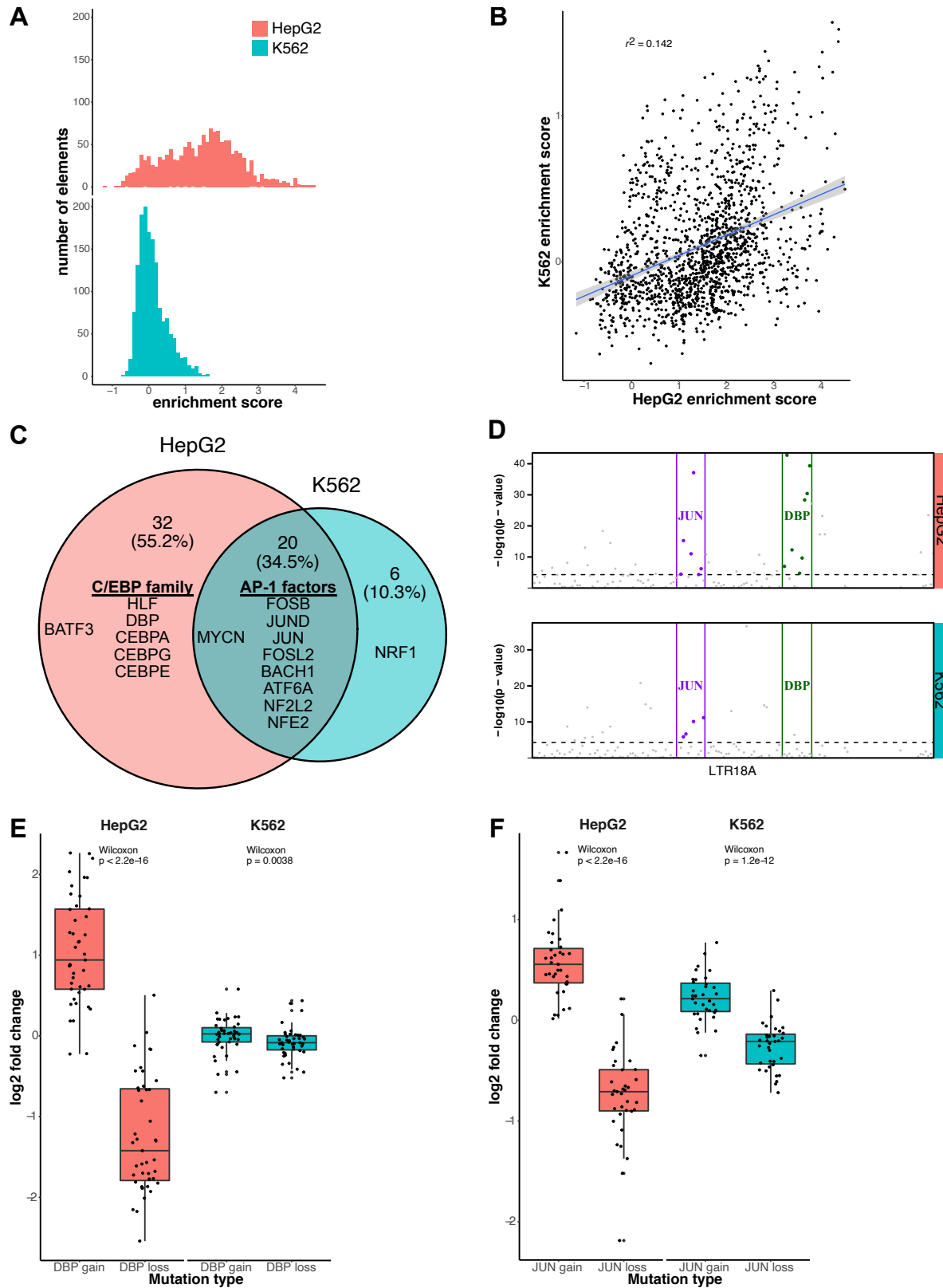


Figure 3: AP-1 motifs drive enhancer activity in HepG2 and K562 while C/EBP motifs are HepG2 specific. A) Distribution of enrichment scores of LTR18A motif focused regions in HepG2 and K562. B) Correlation of enrichment scores between HepG2 and K562. C) Overlap of

motifs significantly associated with active LTR18A. Top 10 transcription factor motifs for both cell lines are displayed. AP-1 and C/EBP-related transcription factors are grouped. D) TEWAS significant nucleotides associated with active LTR18A. JUN and DBP motifs representing AP-1 and C/EBP-related motifs are boxed. Significant positions ($p < 5e-5$, above dotted line) within the two motifs that are associated with active elements are highlighted. E) DBP mutagenesis effects on enhancer activity. F) JUN mutagenesis effects on enhancer activity. *P* values were derived from two-tailed Mann-Whitney *U* tests.

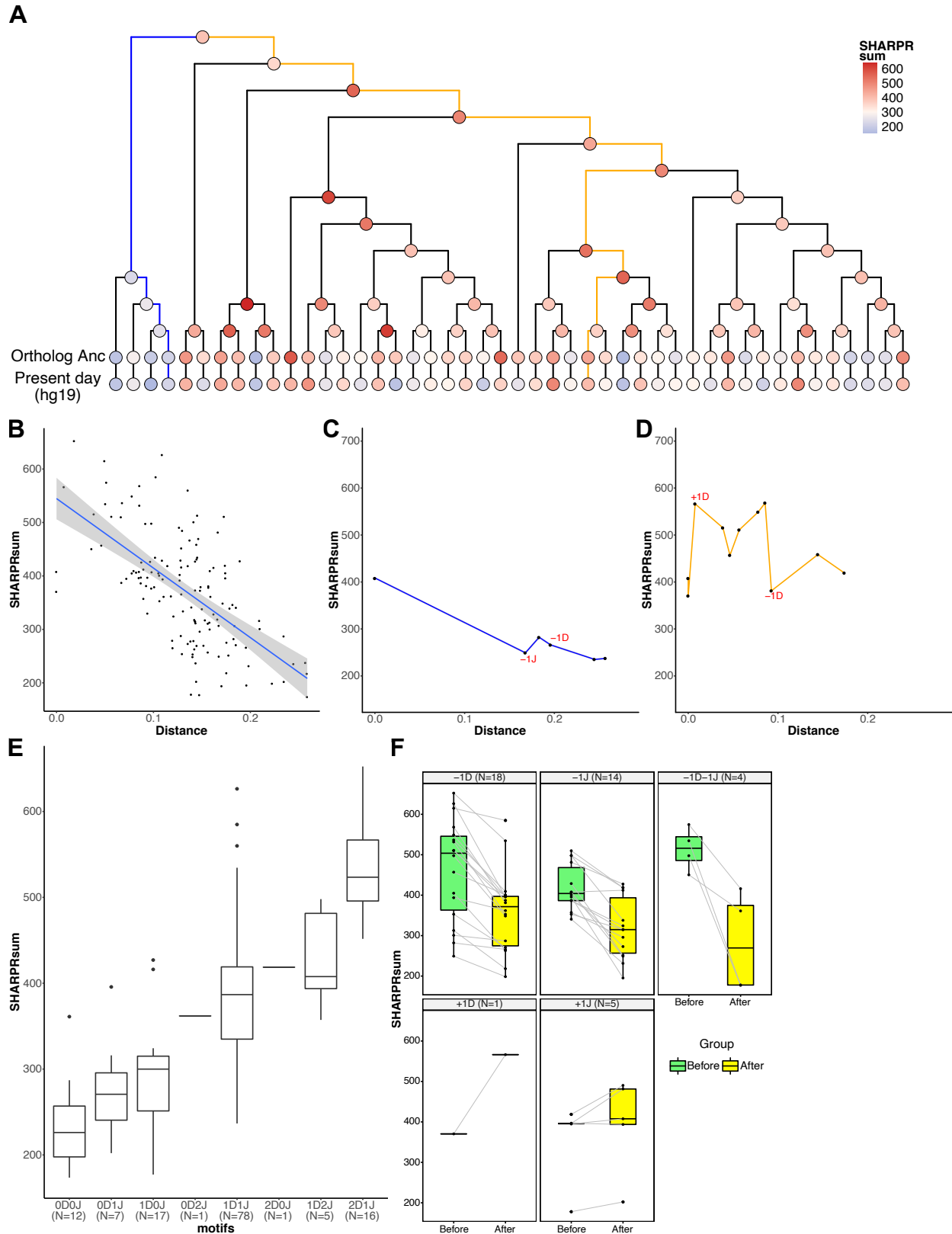


Figure 4: Evolution of regulatory activity in LTR18A in HepG2. A) Phylogenetic tree of reconstructed ancestral LTR18A annotated at each node/element with the sum of SHARPR nucleotide activity scores. B) Correlation of SHARPR sum and distance (substitution rate) from

subfamily ancestor for each LTR18A in the phylogenetic tree. C) Example of regulatory activity evolution along the blue path in A. Motif changes are labeled in red (D = DBP, J = JUN). D) Same as C, but for the orange path in A. E) Distribution of SHARPR sums for phylogenetic tree elements separated by DBP and JUN motif content. F) Motif associated changes in SHARPR sum. Each motif change in the phylogenetic tree is shown with the before and after motif change SHARPR sums connected by a line.

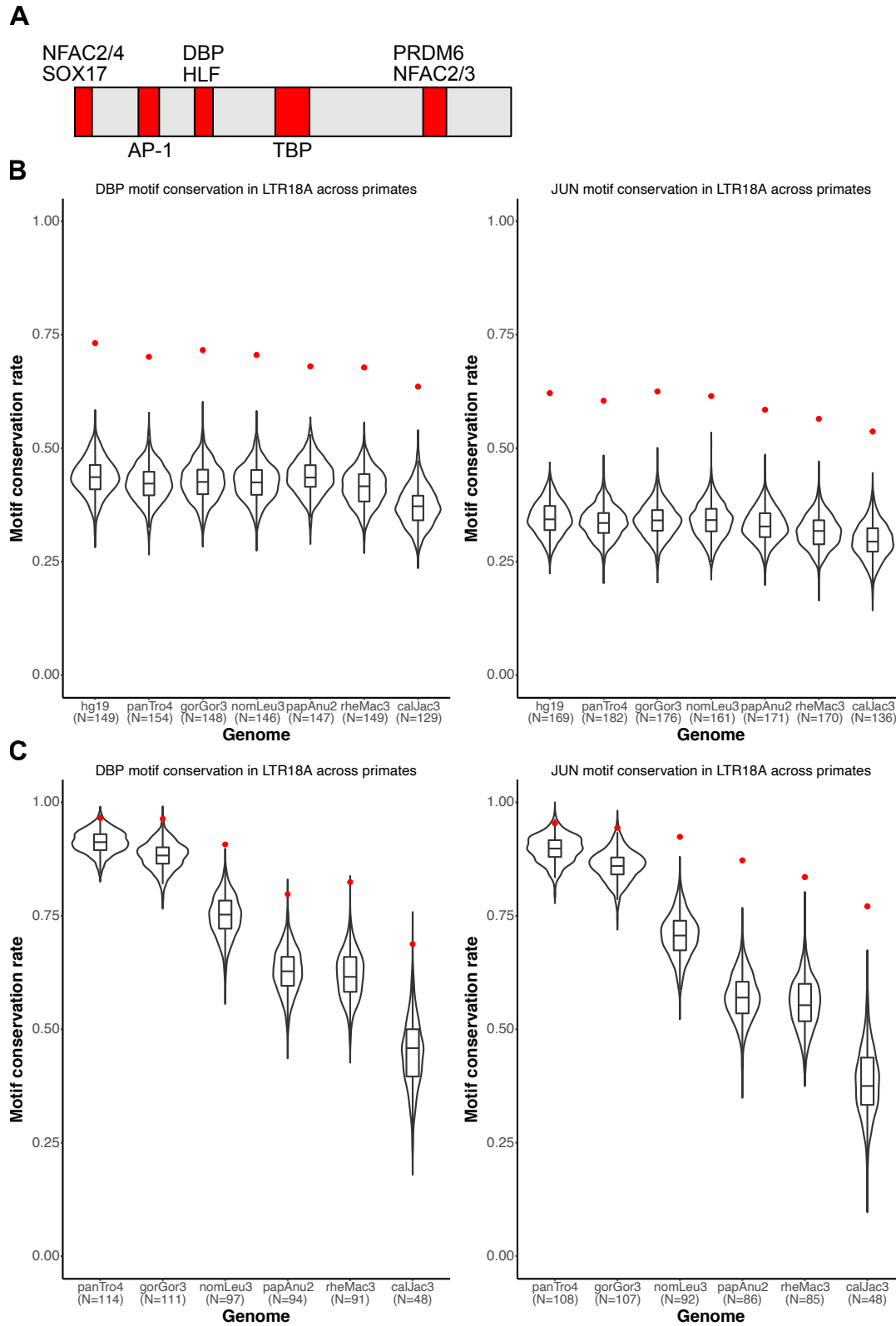


Figure 5: DBP and JUN motifs are more conserved than expected. A) Motifs that are fully encompassed within shared, conserved 10bp sliding windows across seven primate species. Motif locations in red are relative to the LTR18A RepBase consensus sequence. B) Distribution

of expected neutral DBP and JUN motif conservation rates from the consensus motif. 1000 simulations are displayed for each species. The observed conservation rate is shown by the red point. C) Same as B, but for conservation rates from the hg19 ortholog as reference.

Table 1: DBP and JUN motif conservation from RepBase consensus (ancestral), neutral evolution expectation vs. observed

Motif: DBP_HUMAN.H11MO.0.B						
Species	Total possible elements	Expected conserved probability	Expected conserved number	Observed conserved number	Observed conserved proportion	p-value
hg19	149	44.77%	66.71	109	73.15%	1.61E-12
panTro4	154	43.70%	67.30	108	70.13%	1.89E-11
gorGor3	148	43.85%	64.90	106	71.62%	4.96E-12
nomLeu3	146	44.10%	64.39	103	70.55%	6.12E-11
papAnu2	147	42.94%	63.12	100	68.03%	3.97E-10
rheMac3	149	42.17%	62.84	101	67.79%	1.22E-10
calJac3	129	38.71%	49.93	82	63.57%	3.39E-09
Motif: JUN_HUMAN.H11MO.0.A						
Species	Total possible elements	Expected conserved probability	Expected conserved number	Observed conserved number	Observed conserved proportion	p-value
hg19	169	39.34%	66.49	105	62.13%	6.63E-10
panTro4	182	38.54%	70.14	110	60.44%	6.33E-10
gorGor3	176	38.65%	68.02	110	62.50%	4.05E-11
nomLeu3	161	38.61%	62.16	99	61.49%	1.23E-09
papAnu2	171	37.58%	64.27	100	58.48%	8.41E-09
rheMac3	170	37.01%	62.92	96	56.47%	7.43E-08
calJac3	136	34.07%	46.33	73	53.68%	7.01E-07

Table 2: DBP and JUN motif conservation from hg19 ortholog as reference, neutral evolution expectation vs. observed

Motif: DBP_HUMAN.H11MO.0.B						
Species	Total possible elements	Expected conserved probability	Expected conserved number	Observed conserved number	Observed conserved proportion	p-value
panTro4	114	92.33%	105.26	110	96.49%	0.0476
gorGor3	111	89.42%	99.25	107	96.40%	0.0084
nomLeu3	97	76.83%	74.53	88	90.72%	0.0006
papAnu2	94	65.84%	61.89	75	79.79%	0.0022
rheMac3	91	64.71%	58.89	75	82.42%	0.0002
calJac3	48	47.71%	22.90	33	68.75%	0.0018
Motif: JUN_HUMAN.H11MO.0.A						
Species	Total possible elements	Expected conserved probability	Expected conserved number	Observed conserved number	Observed conserved proportion	p-value
panTro4	108	91.08%	98.37	103	95.37%	0.0590
gorGor3	107	87.70%	93.84	101	94.39%	0.0175
nomLeu3	92	73.86%	67.95	85	92.39%	2.62E-05
papAnu2	86	62.02%	53.33	75	87.21%	7.41E-07
rheMac3	85	60.87%	51.74	71	83.53%	9.29E-06
calJac3	48	44.93%	21.57	37	77.08%	3.77E-06

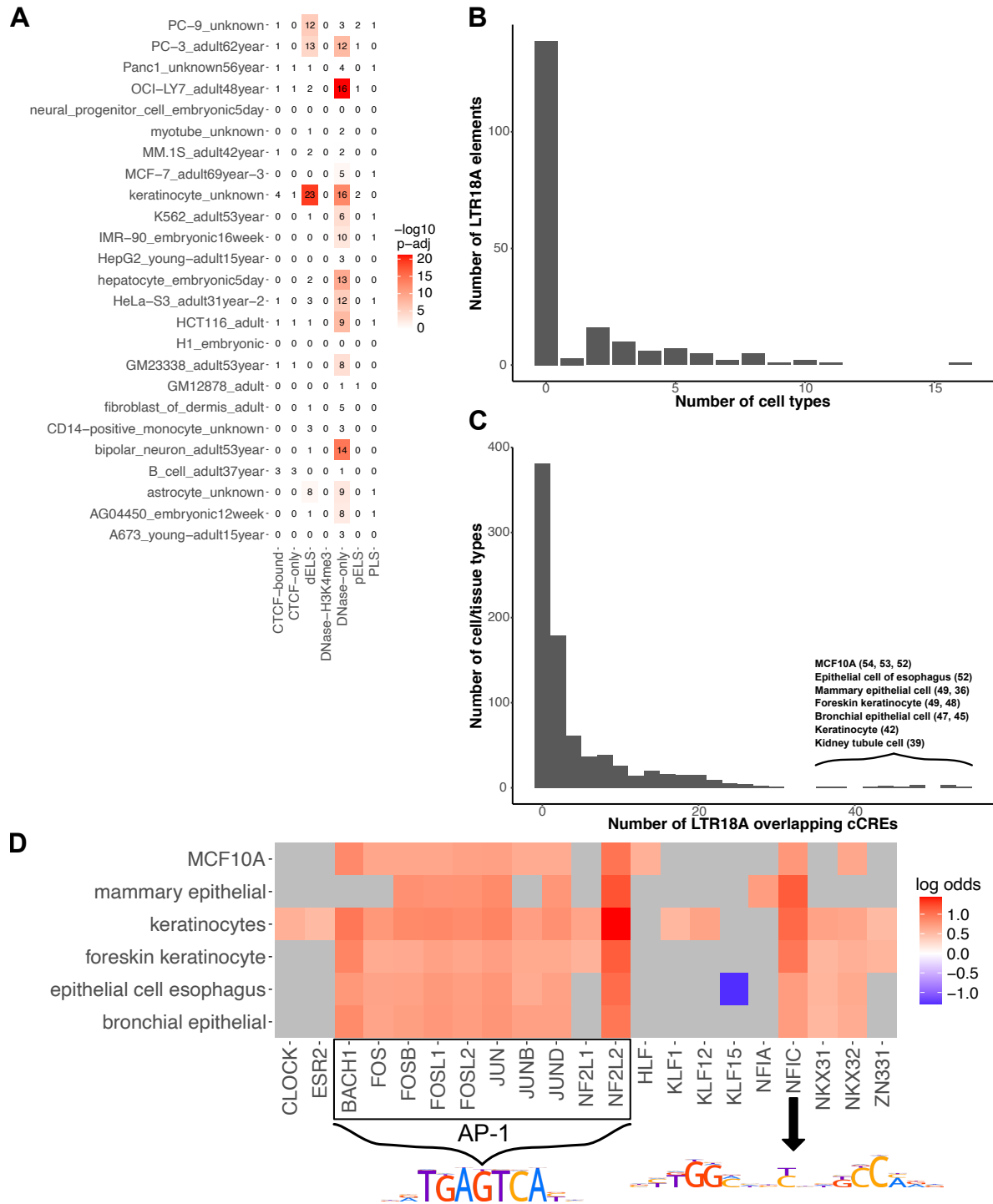


Figure 6: LTR18A elements are associated with enhancer epigenetic marks in human. A) Overlap of LTR18A with ENCODE cCREs across 25 full classification cell/tissue types (dELS, distal enhancer-like signature; pELS, proximal enhancer-like signature; PLS, promoter-like signature). The number of elements that overlap with cCREs are shown as well as their $-\log_{10}$ adjusted p-value by bedtools fisher. B) Distribution of LTR18A elements overlapping cCREs

across multiple full classification cell/tissue types. C) Distribution of cell/tissue types overlapping LTR18A elements. The top cell/tissue types are displayed with the number of LTR18A elements that overlap with a cCRE. D) Motifs associated with the cCRE-overlapping LTR18A elements from the top cell/tissue types in C. Grey indicates non-significance at adjusted p-value threshold of 0.05. PWMs for JUN (AP-1 related factors) and NFIC are shown.