

Figure S1

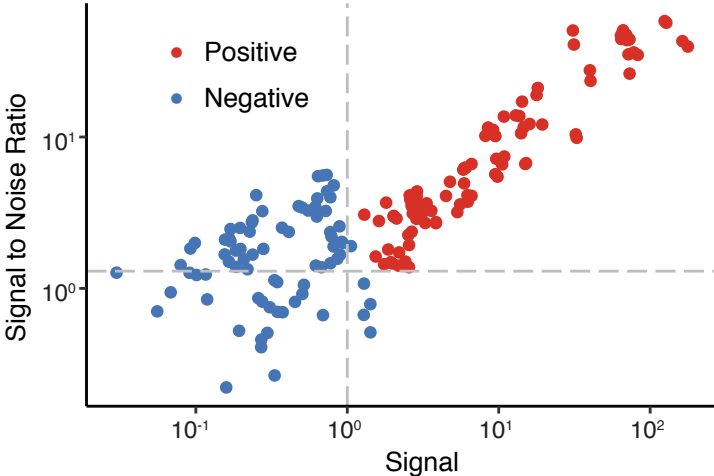


Figure S2

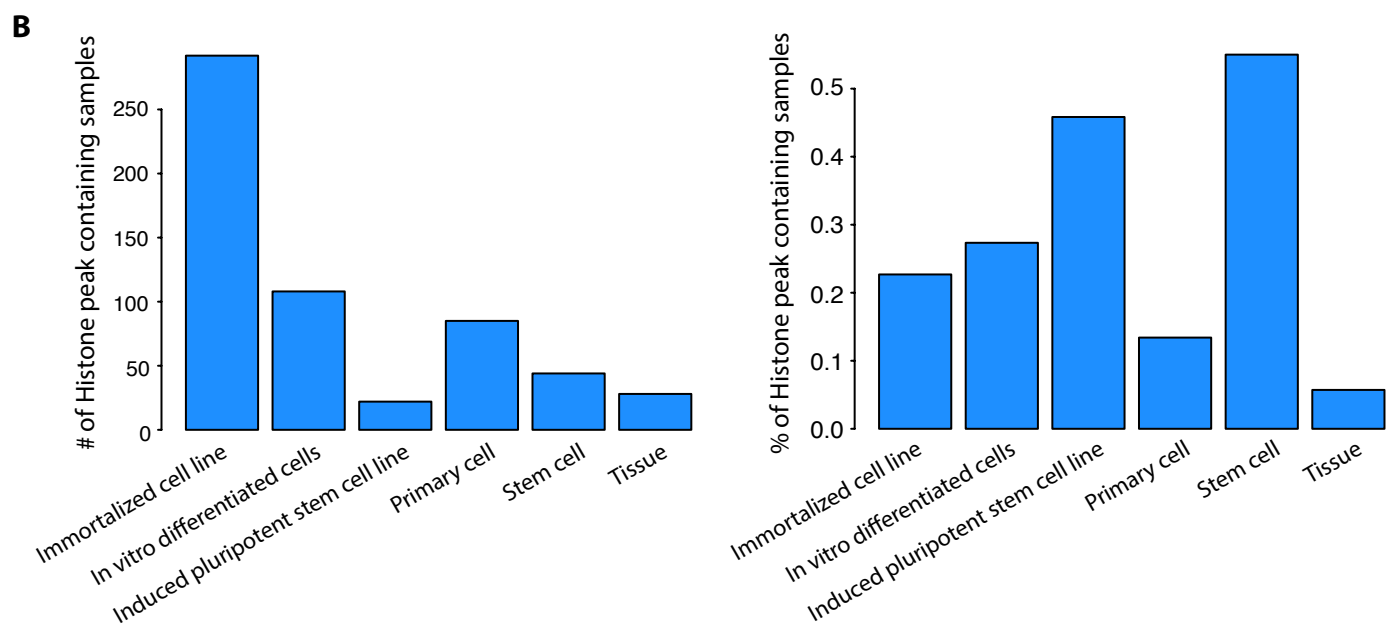
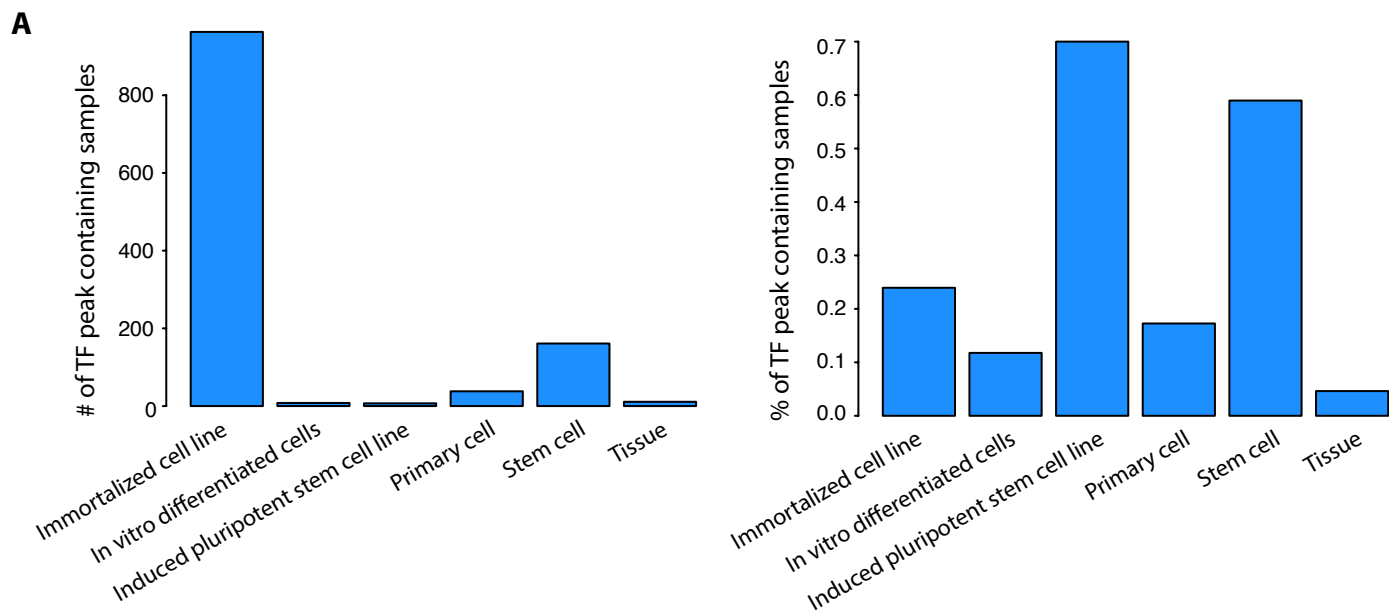


Figure S3

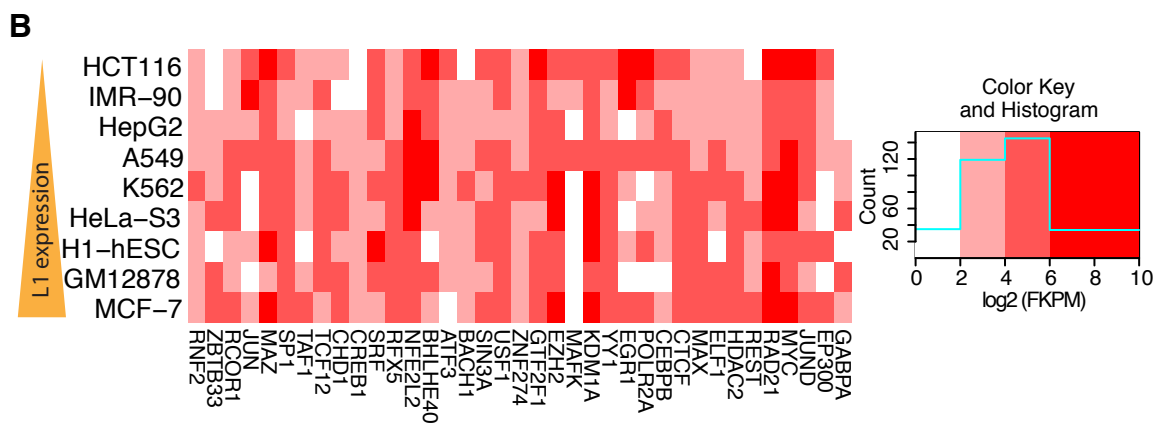
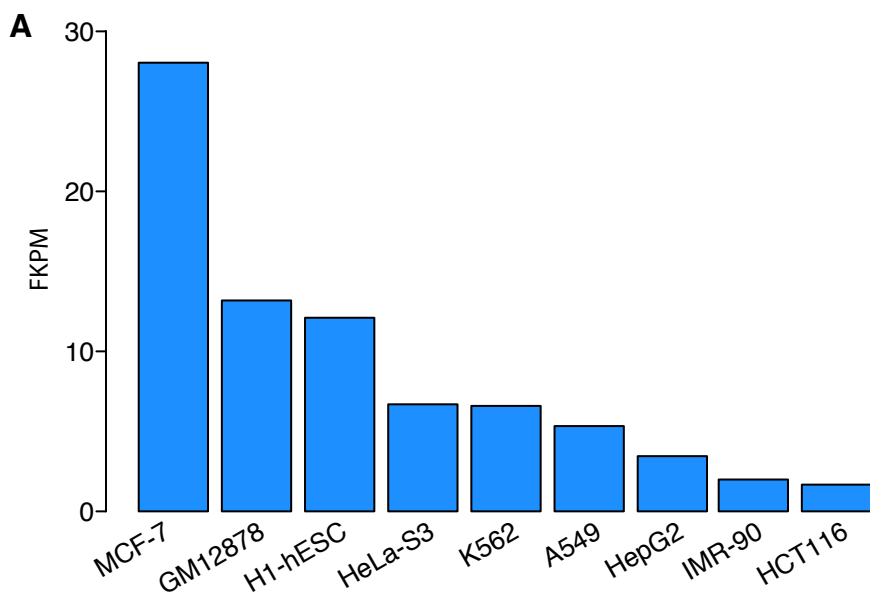


Figure S4

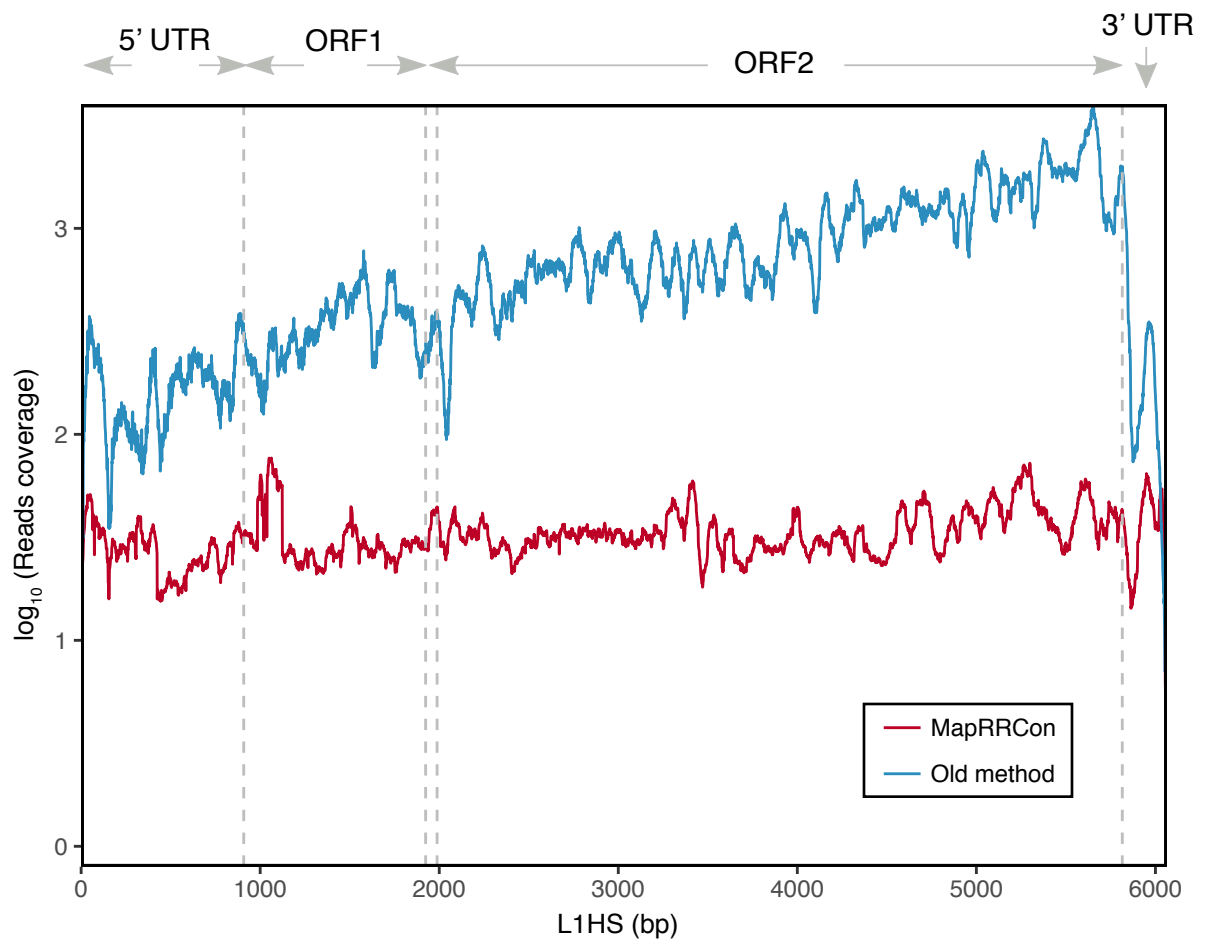


Figure S5

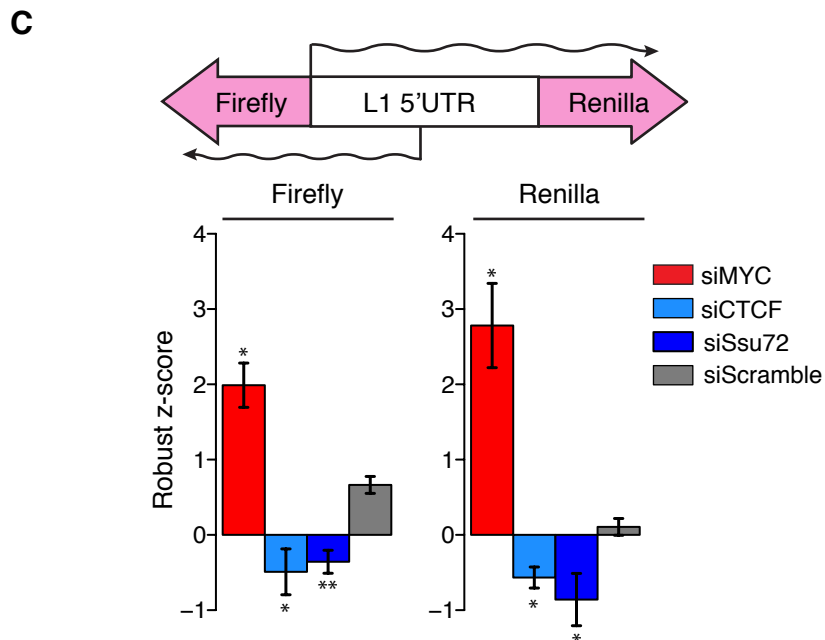
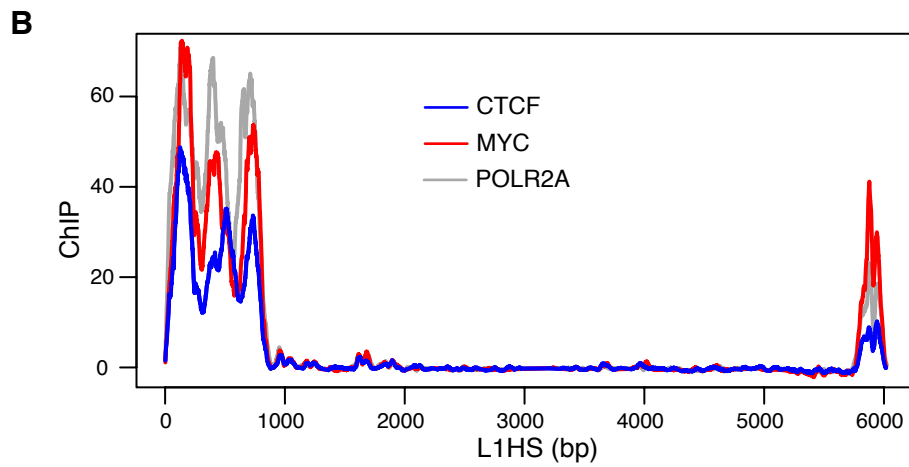
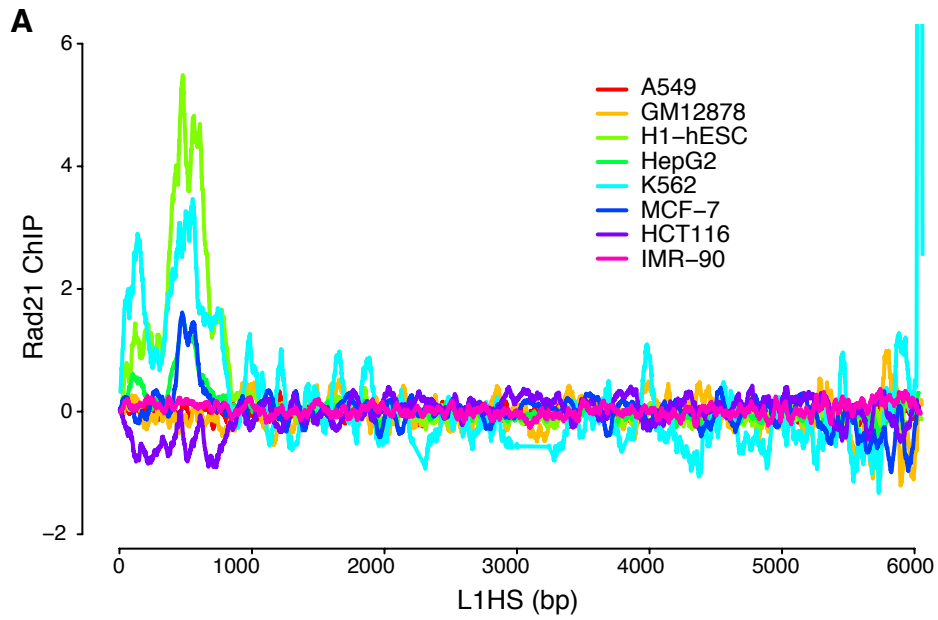


Figure S5

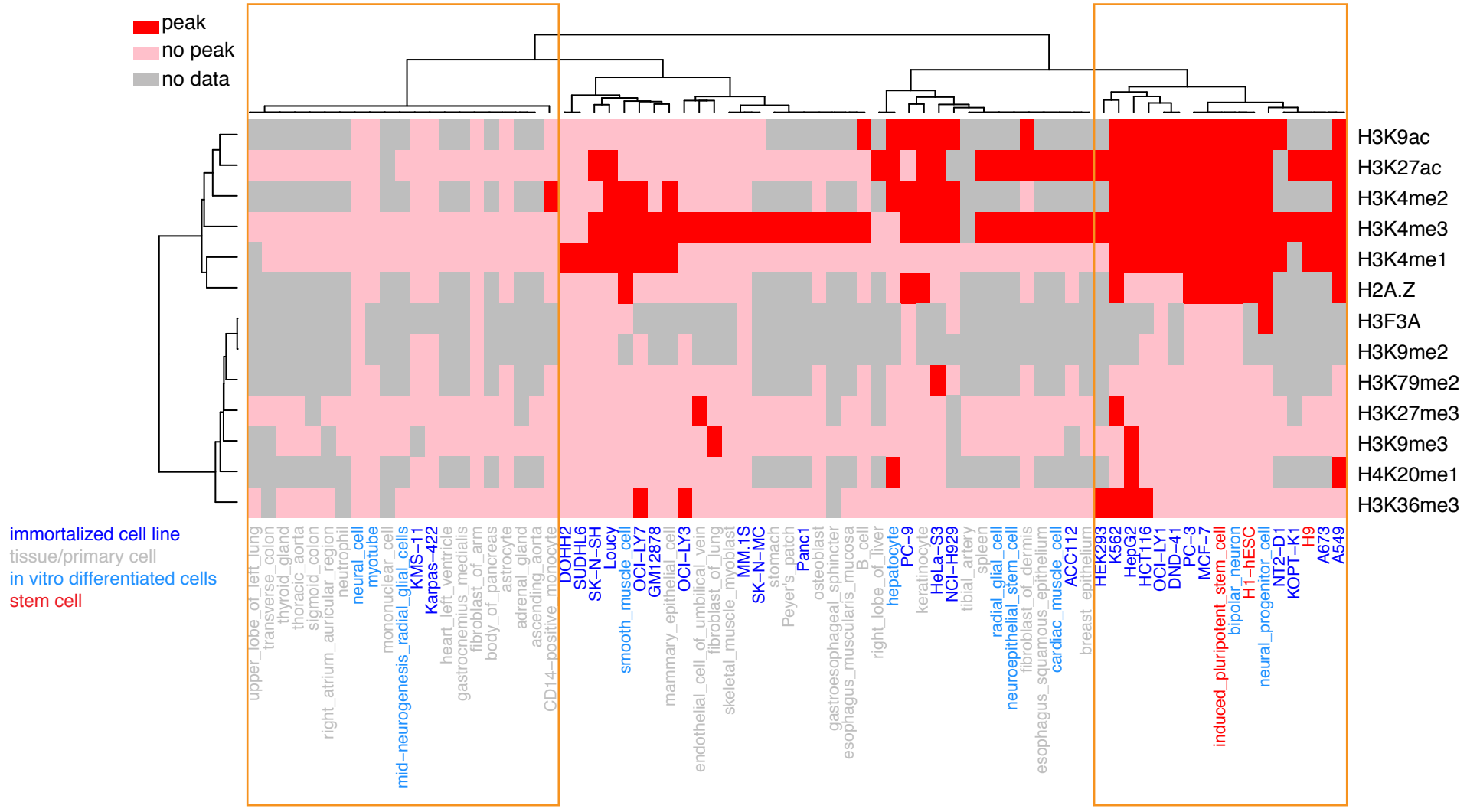


Figure S7

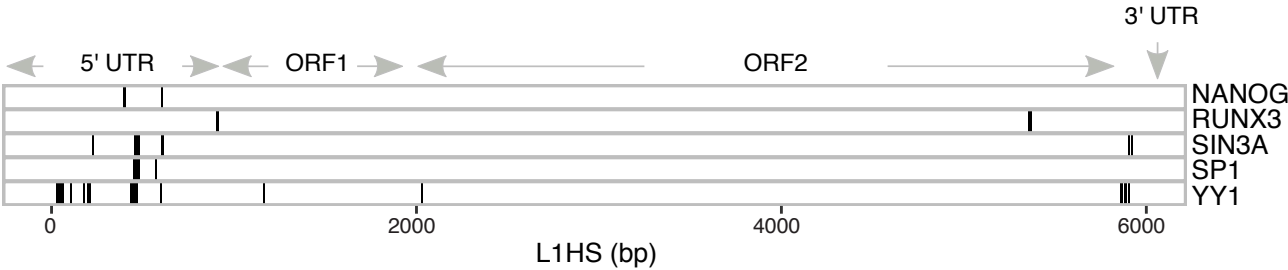
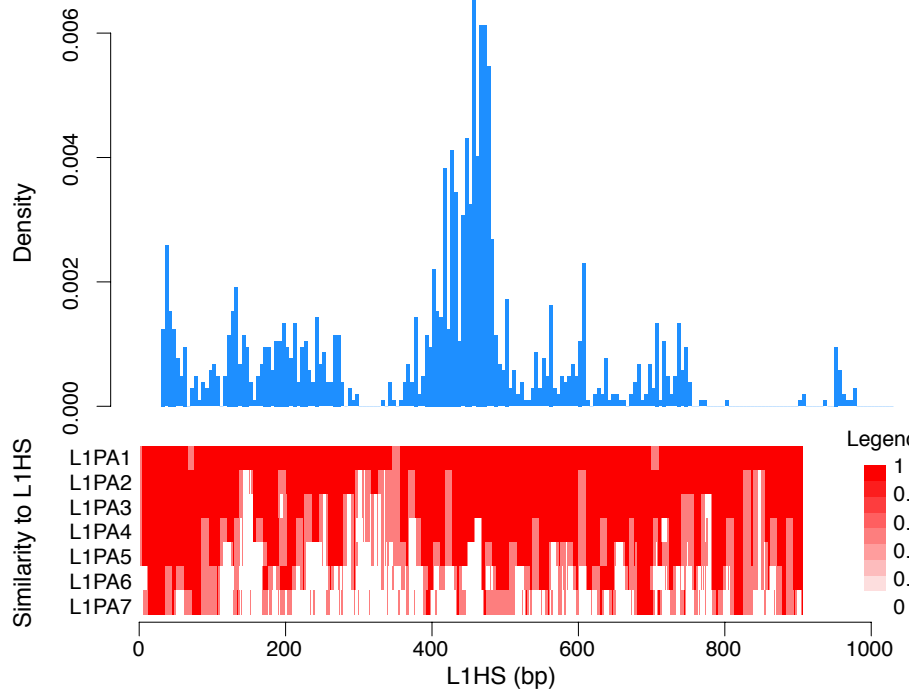


Figure S8

A



B

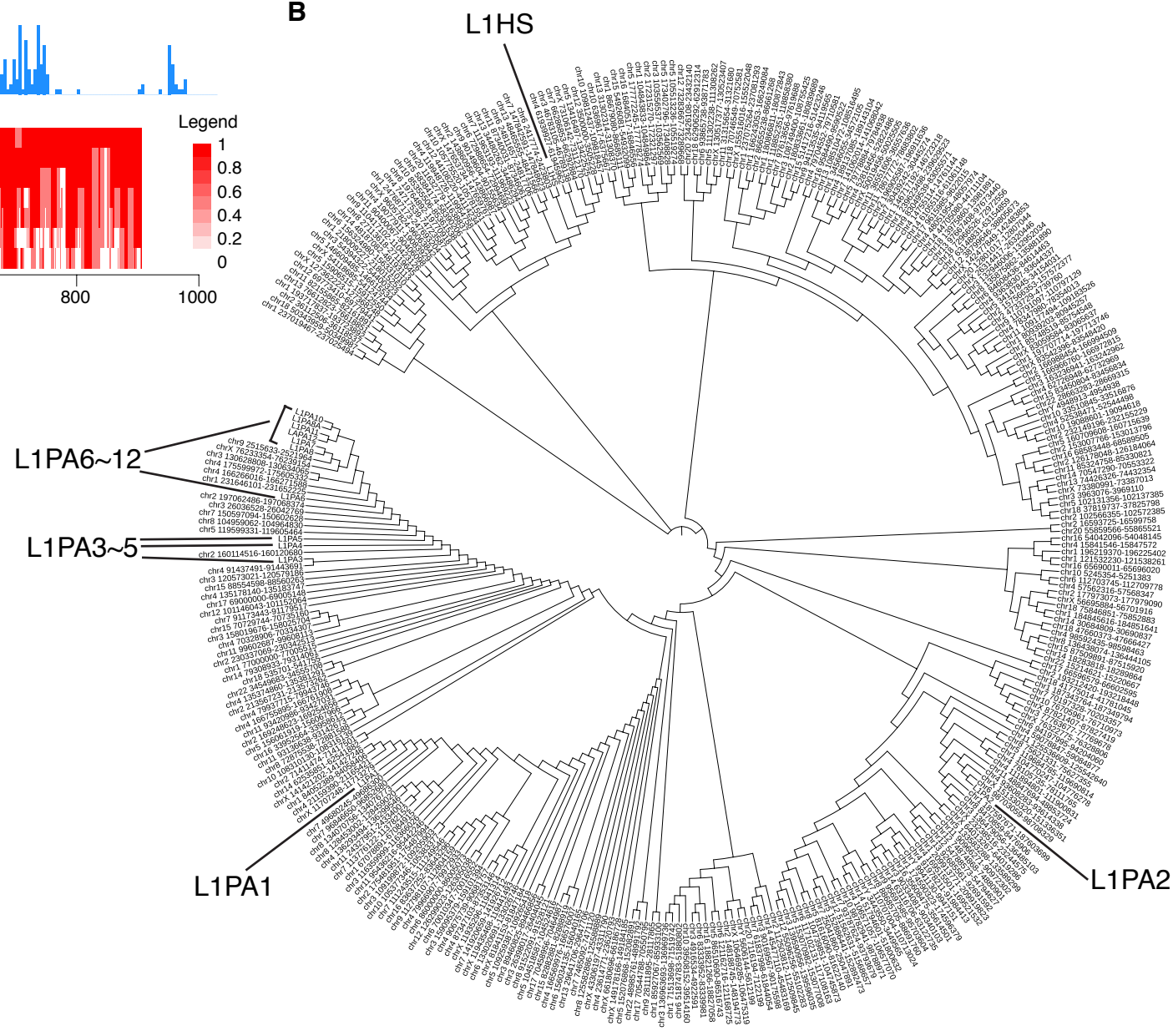
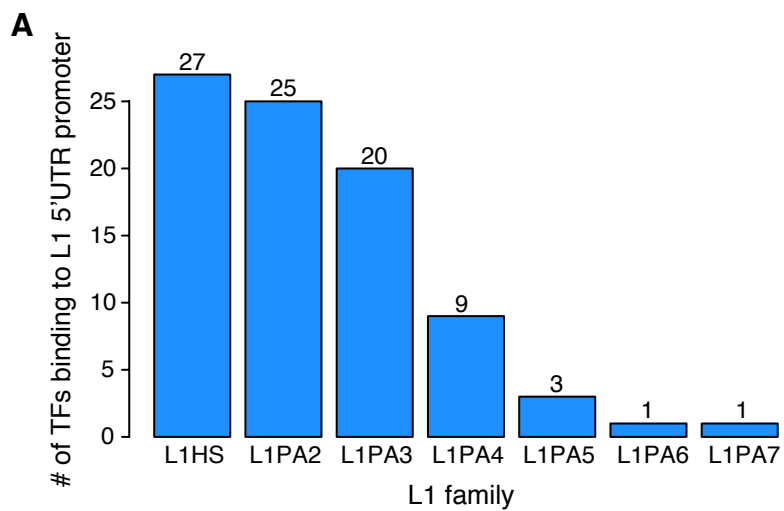
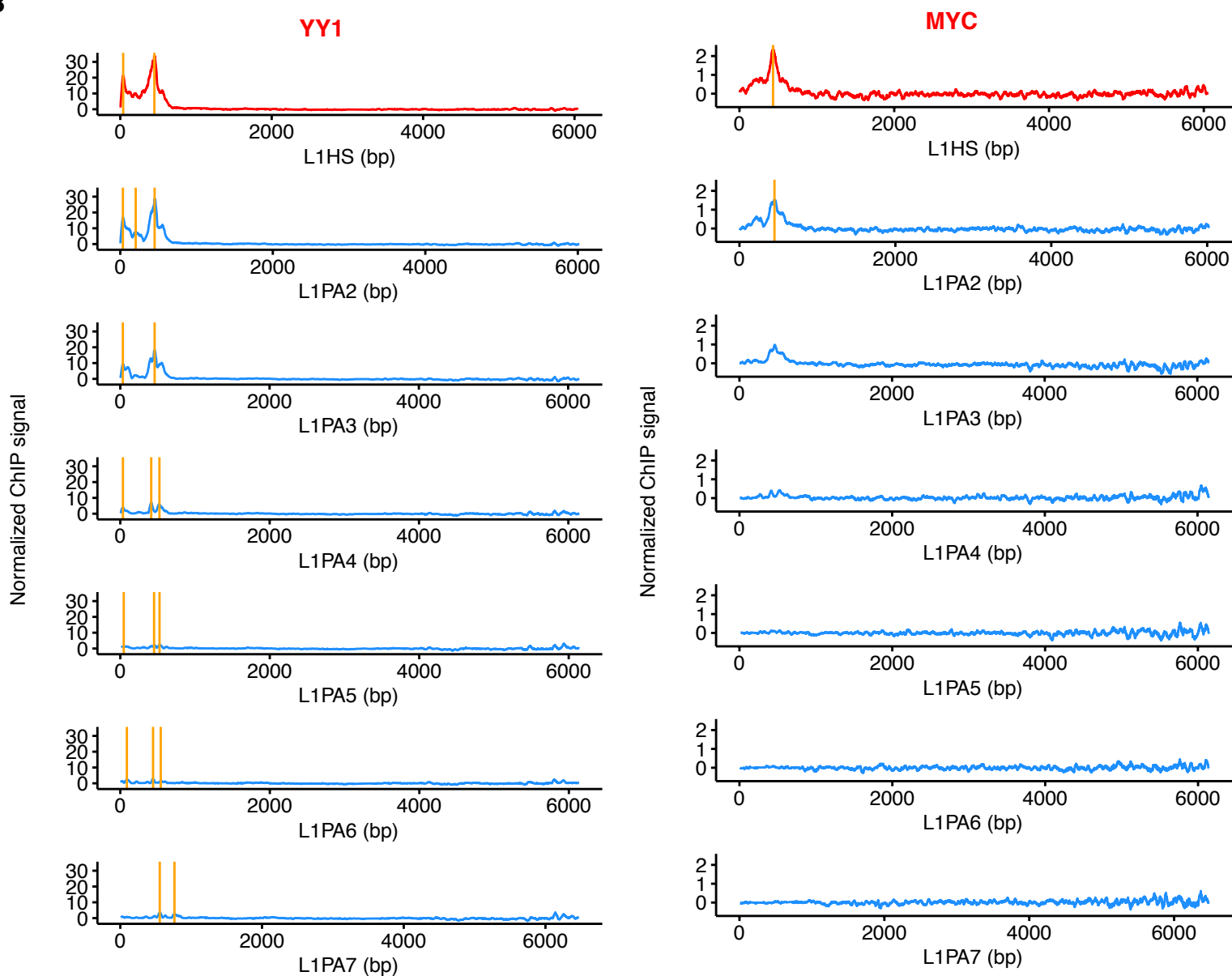


Figure S9



B



SI figure legends

Figure S1. Training datasets of the Peak Calling Algorithm

We optimized two parameters of the algorithm to fit our datasets. We used 91 positive and 77 negative datasets to train the algorithm. After plotting them based on intensity and signal-to-noise ratio, we found that these two datasets clustered into two populations, separated by two thresholds (normalized intensity >1; signal-to-noise ratio >1.3).

Figure S2. Overview of the TF/histone mark peaks in different biosamples.

AB. We plot the number of samples (without removing replicates) that contain peaks at L1. The right panel shows the percentage of peak-containing samples in each of the categorized groups.

Figure S3. Expression level of L1 and TFs in various cell lines.

A. L1 expression level in 9 ENCODE cell lines quantified by RNA-seq data. Three bars for each cell lines represents FPKM (Fragments Per Kilobase of transcript per Million mapped reads) values L1. **B.** Expression (log₂ FPKM) of a selected number of TFs (same set in Figure 2D) in 9 cell lines.

Figure S4. No 3' bias on L1HS is observed for RNA-seq data in 77 breast tumor samples.

The distributions of mapped reads are plotted along the L1HS, and each line represents one tumor sample. The boundaries of 5'UTR, ORF1, ORF2 and 3'UTR are shown with gray dashed lines. The lack of a 3' bias is consistent with most of the reads coming from full length transcripts and furthermore shows that the level of genomic DNA contamination in the RNA prep is very low.

Figure S5. CTCF may mediate L1 gene loops.

A. ChIP profiles of Rad21 in various cell lines on L1HS (replicates were averaged to form a single line). **B.** Colocalization of CTCF, c-Myc and RNA polymerase II ChIP-seq signals in MCF-7 cells. **C.** siRNA knockdown experiments of Ssu72. See Figure 4D legends for details.

Figure S6. Histone marks that enrich at L1 5'UTR in various samples.

We present the binding profiles of 13 histone marks in various samples with colors indicating whether peaks are detected at L1 5'UTR promoter. The samples are grouped into 4 categories shown by different colors on the x axis. Both dimensions are ranked by hierarchical clustering as described in Figure 1A. Two distinct clusters are marked with orange boxes: (1) The left box showed a cluster of samples that have no histone marks at L1 5'UTR promoter. (2) The right box showed a cluster of samples that have enriched active histone marks at L1 5'UTR promoter.

Figure S7. Peaks identified for known TFs.

Peak locations identified in this study of 5 published TFs, previously shown to bind and regulate L1 transcription, were plotted along L1HS.

Figure S8. Sequence similarity between L1 consensus sequences.

A. We calculated the pairwise sequence similarity between each of the L1PA1-xx consensus sequences to L1HS on a 10-bp sliding window with a step of 1-bp. For each window of the L1HS consensus sequence, we searched for sequences in L1PA with either a perfect match or 1-bp

mismatch. We then assigned the score 1 if there was a perfect matched sequence or as 0.5 if the best match had a 1-bp mismatch in each window. The score was set to 0 if the best matched sequence had more than 1-bp mismatches. The similarity scores are color coded (shown in red) and aligned with the peak density histogram (shown in blue). **B.** 332 Repeatmasker annotated L1HS insertion sequences are compared to L1 consensus sequences (L1HS and L1PAs) via phylogenetic analysis. Multiple sequence alignment was performed using Clustal Omega(1) and the tree was generated by iTOL (2).

Figure S9. TF binding profiles on L1PA families.

A. We ran MapRRCon on a subset of hESC datasets (182 datasets, 36 TFs: ATF3, BACH1, BRCA1, CHD1, CHD2, CTCF, EGR1, EP300, FOSL1, GABPA, GTF2F1, JUN, JUND, KDM1A, KDM6A, MAFK, MYC, NRF1, PHF8, POLR2A, RAD21, REST, RFX5, RXRA, SAP30, SIN3A, SIX5, SP1, SRF, TAF1, TAF7, TBP, USF1, USF2, YY1, ZNF143) using each one of the following consensus sequences and its corresponding insertion locations annotated by Repeatmasker. After combining the replicates, we plotted the number of TFs that show peaks at 5'UTR at each of the consensus sequences. **B.** Binding profiles of two TFs on each of the L1 consensus sequence are shown as examples. Orange lines indicate locations of peaks called by our peak-calling algorithm.

Reference

1. Sievers F, et al. (2011) Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol* 7(1):539–539.
2. Letunic I, Bork P (2016) Interactive tree of life (iTOL) v3: an online tool for the display and annotation of phylogenetic and other trees. *Nucl Acids Res* 44(W1):W242–5.