

Similar evolutionary trajectories for retrotransposon accumulation in mammals

Reuben M Buckley¹, R Daniel Kortschak¹, Joy M Raison¹, David L Adelson^{1,*}

1 Department of Genetics and Evolution, The University of Adelaide, North Tce, 5005, Adelaide, Australia

* david.adelson@adelaide.edu.au

Keywords: Transposable element, Genome Evolution, Genome Architecture

Running title: Similar evolutionary trajectories in mammals

Abstract

The factors guiding retrotransposon insertion site preference are not well understood. Different types of retrotransposons share common replication machinery and yet occupy distinct genomic domains. Autonomous long interspersed elements accumulate in gene-poor domains and their non-autonomous short interspersed elements accumulate in gene-rich domains. To determine genomic factors that contribute to this discrepancy we analysed the distribution of retrotransposons within the framework of chromosomal domains and regulatory elements. Using comparative genomics, we identified large-scale conserved patterns of retrotransposon accumulation across several mammalian genomes. Importantly, retrotransposons that were active after our sample-species diverged accumulated in orthologous regions. This suggested a similar evolutionary interaction between retrotransposon activity and conserved genome architecture across our species. In addition, we found that retrotransposons accumulated at regulatory element boundaries in open chromatin, where accumulation of particular retrotransposon types depended on insertion size and local regulatory element density. From our results, we propose a model where density and distribution of genes and regulatory elements canalise retrotransposon accumulation. Through conservation of synteny, gene regulation and nuclear organisation, mammalian genomes with dissimilar retrotransposons follow similar evolutionary trajectories.

Introduction

An understanding of the dynamics of evolutionary changes in mammalian genomes is critical for understanding the diversity of mammalian biology. Most work on mammalian molecular evolution is on protein coding genes, based on the assumed centrality of their roles and because of the lack of appropriate methods to identify the evolutionary conservation of apparently non-conserved, non-coding sequences. Consequently, this approach addresses only a tiny fraction (less than 2%) of a species' genome, leaving significant gaps in our understanding of evolutionary processes (ENCODE Project Consortium 2012; Lander et al. 2001). In this report we describe how large scale positional conservation of non-coding, repetitive DNA sheds light on the possible conservation of mechanisms of genome evolution, particularly with respect to the acquisition of new DNA sequences.

1
2
3
4
5
6
7
8
9
10
11

Mammalian genomes are hierarchically organised into compositionally distinct hetero- or 12
euchromatic large structural domains (Gibcus and Dekker 2013). These domains are largely 13
composed of mobile self-replicating non-long terminal repeat (non-LTR) retrotransposons; 14
with Long INterspersed Elements (LINEs) in heterochromatic regions and Short INterspersed 15
Elements (SINEs) in euchromatic regions (Medstrand et al. 2002). The predominant LINE 16
in most mammals is the ~6 kb long L1. In many mammal genomes, this autonomously 17
replicating element is responsible for the mobilisation of an associated non-autonomous 18
SINE, usually ~300 bp long. Together, LINEs and SINEs occupy approximately 30% of 19
the human genome (Lander et al. 2001), replicate via a well characterised RNA-mediated 20
copy-and-paste mechanism (Cost et al. 2002) and co-evolve with host genomes (Kramerov 21
and Vassetzky 2011; Chalopin et al. 2015; Furano et al. 2004). 22

The accumulation of L1s and their associated SINEs into distinct genomic regions depends 23
on at least one of two factors. 1) Each element's insertion preference for particular genomic 24
regions and 2) the ability of particular genomic regions to tolerate insertions. According to 25
the current retrotransposon accumulation model, both L1s and SINEs likely share the same 26
insertion patterns constrained by local sequence composition. Therefore, their accumulation 27
in distinct genomic regions is a result of region specific tolerance to insertions. Because L1s 28
are believed to have a greater capacity than SINEs to disrupt gene regulatory structures, 29
they are evolutionarily purged from gene-rich euchromatic domains at a higher rate than 30
SINEs. Consequently, this selection asymmetry in euchromatic gene-rich regions causes L1s 31
to become enriched in gene-poor heterochromatic domains (Lander et al. 2001; Graham and 32
Boissinot 2006; Gasior et al. 2007; Kvikstad and Makova 2010). 33

An important genomic feature, not explored in the accumulation model, is the chro- 34
matin structure that surrounds potential retrotransposon insertion sites. Retrotransposons 35
preferentially insert into open chromatin (Cost et al. 2001; Baillie et al. 2011), which is 36
usually found overlapping gene regulatory elements. As disruption of regulatory elements 37
can often be harmful, this creates a fundamental evolutionary conflict for retrotransposons; 38
their immediate replication may be costly to the overall fitness of the genome in which they 39
reside. Therefore, rather than local sequence composition or tolerance to insertion alone, 40
retrotransposon accumulation is more likely to be constrained by an interaction between 41
retrotransposon expression, openness of chromatin, susceptibility of a particular site to alter 42
gene regulation, and the capacity of an insertion to impact on fitness. 43

To investigate the relationship between retrotransposon activity and genome evolution, we began by characterising the distribution and accumulation of non-LTR retrotransposons within placental mammalian genomes. Next, we compared retrotransposon accumulation patterns in eight separate evolutionary paths by ‘humanising’ the repeat content (see methods) of the chimpanzee, rhesus macaque, mouse, rabbit, dog, horse and cow genomes. Finally, we analysed human retrotransposon accumulation in large hetero- and euchromatic structural domains, focusing on regions surrounding genes, exons and regulatory elements. Our results suggest that accumulation of particular retrotransposon families follows from insertion into open chromatin found adjacent to regulatory elements and depends on local gene and regulatory element density. From this we propose a refined retrotransposon accumulation model in which random insertion of retrotransposons is primarily constrained by chromatin structure rather than local sequence composition.

Materials and Methods

Within species comparisons of retrotransposon genome distributions

Retrotransposon coordinates for each species were initially identified using RepeatMasker and obtained from either the RepeatMasker website or UCSC genome browser (Table S1) (Smit et al. 1996; Rosenbloom et al. 2015). We grouped retrotransposon elements based on repeat IDs used in Giordano *et al* (Giordano et al. 2007). Retrotransposon coordinates were extracted from hg19, mm9, panTro4, rheMac3, oruCun2, equCab2, susScr2, and canFam3 assemblies. Each species genome was segmented into 1 Mb regions and the density of each retrotransposon family for each segment was calculated. Retrotransposon density of a given genome segment is equal to a segments total number of retrotransposon nucleotides divided by that segments total number of mapped nucleotides (non-N nucleotides). From this, each species was organised into an n -by- p data matrix of n genomic segments and p retrotransposon families. Genome distributions of retrotransposons were then analysed using principle component analysis (PCA) and correlation analysis. For correlation analysis, we used our genome segments to calculate Pearson’s correlation coefficient between each pair-wise combination of retrotransposon families within a species.

Across species comparisons of retrotransposon genome distributions

72

To compare genome distributions across species, we humanised a segmented query species genome using mapping coordinates extracted from net AXT alignment files located on the UCSC genome browser (Table S1). First, poorly represented regions were removed by filtering out genome segments that fell below a minimum mapping fraction threshold (Fig. 1a). Next, we used mapping coordinates to match fragments of query species segments to their corresponding human segments (Fig. 1b). From this, the retrotransposon content and PC scores of the matched query segments were humanised following equation 1 (Fig. 1c).

73

74

75

76

77

78

79

$$c_i^* = \frac{\sum_j c_{ij} l_j^Q / q_j}{\sum_j l_j^R / r}, \quad (1)$$

where c_{ij} is the density of retrotransposon family i in query segment j , l_j^Q is the total length of the matched fragments between query segment j and the reference segment, l_j^R is the total length of the reference segment fragments that match query segment j , q_j is the total length of the query segment j , and r is the total length of the reference segment. The result c_i^* is the humanised coverage fraction of retrotransposon family i that can now be compared to a specific reference segment. Once genomes were humanised, Pearson's correlation coefficient was used to determine the conservation between retrotransposon genomic distributions (Fig. 1d). Using the Kolmogorov-Smirnov test, we measured the effect of humanising by comparing the humanised query retrotransposon density distribution to the query filtered retrotransposon density distribution (Fig. 1e). The same was done to measure the effect of filtering by comparing the segmented human retrotransposon density distribution to the human filtered retrotransposon density distribution (Fig. 1f). Our Pearson's correlation coefficients and P-values from measuring the effects of humanising and filtering were integrated into a heatmap (Fig. 1g). This entire process was repeated at different minimum mapping fraction thresholds to optimally represent each retrotransposon families genomic distribution in a humanised genome (fig S1).

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

Replication timing profiles, boundaries and constitutive domains

96

Genome-wide replication timing data for human and mouse were initially generated as part of the ENCODE project and were obtained from UCSC genome browser (Table S2-S3) (Yue et al.

97

98

2014; ENCODE Project Consortium 2012). For human genome-wide replication timing we 99
used Repli-Seq smoothed wavelet signals generated by the UW ENCODE group (ENCODE 100
Project Consortium 2012), in each cell-line we calculated the mean replication timing per 101
1Mb genome segment. For mouse genome-wide replication timing we used Repli-Chip wave 102
signals generated by the FSU ENCODE group (Yue et al. 2014). Since two replicates were 103
performed on each cell-line, we first calculated each cell-lines mean genome-wide replication 104
timing and then used this value to calculate the mean replication timing per 1Mb genome 105
segment. By calculating mean replication timing per 1 Mb segment we were able to easily 106
compare large-scale genome-wide replication timing patterns across cell-lines. We obtained 107
early replication domains (ERDs), late replication domains (LRDs) and timing transition 108
regions (TTRs) from the gene expression omnibus (accession ID GSE53984) (Table S2). 109
Replication domains for each dataset were identified using a deep neural network hidden 110
Markov model (Liu et al. 2015). To determine RD boundary fluctuations of retrotransposon 111
density, we defined ERD boundaries as the boundary of a TTR adjacent to an ERD. ERD 112
boundaries from across each sample were pooled and retrotransposon density was calculated 113
for 50 kb intervals from regions flanking each boundary 1 Mb upstream and downstream. 114
Expected density and standard deviation for each retrotransposon group was derived from a 115
background distribution generated by calculating the mean of 500 randomly sampled 50 kb 116
genomic bins within 2000 kb of each ERD boundary, replicated 10000 times. To generate 117
replication timing profiles for our ERD boundaries, we also calculated the mean replication 118
timing per 50 kb intervals from across each human Repli-Seq sample. To identify constitutive 119
ERDs and LRDs (cERDs and cLRDs), ERDs and LRDs classified by Liu *et al* (Liu et al. 120
2015) across each cell type were evenly split into 1 kb intervals. If the classification of 12
of 16 samples agreed across a certain 1 kb interval, we classified that region as belonging to 122
a cERDs or cLRDs, depending the region's majority classification of the 1 kb interval. 123

DNase1 cluster identification and activity

 124

DNase1 sites across 15 cell lines were found using DNase-seq and DNase-chip as part of the 125
open chromatin synthesis dataset for ENCODE generated by Duke University's Institute for 126
Genome Sciences & Policy, University of North Carolina at Chapel Hill, University of Texas 127
at Austin, European Bioinformatics Institute and University of Cambridge, Department 128
of Oncology and CR-UK Cambridge Research Institute (Table S4) (ENCODE Project 129

Consortium 2012). Regions where P-values of contiguous base pairs were below 0.05 were identified as significant DNase1 hypersensitive sites (ENCODE Project Consortium 2012). From this we extracted significant DNase1 hypersensitive sites from each sample and pooled them. DNase1 hypersensitive sites were then merged into DNase1 clusters. Cluster activity was calculated as the number of total overlapping pooled DNase1 hypersensitive sites. We also extracted intervals between adjacent DNase1 clusters to look for enrichment of retrotransposons at DNase1 cluster boundaries.

Extraction of intergenic and intron intervals

hg19 RefSeq gene annotations obtained from UCSC genome browser were used to extract a set of introns and intergenic intervals (Table S5). RefSeq gene annotations were merged and intergenic regions were classified as regions between the start and end of merged gene models. We used the strandedness of gene model boundaries to classify adjacent intergenic region boundaries as upstream or downstream. We discarded intergenic intervals adjacent to gene models where gene boundaries were annotated as both + and - strand. Regions between adjacent RefSeq exons within a single gene model were classified as introns. Introns interrupted by exons in alternatively spliced transcripts and introns overlapped by other gene models were excluded. Upstream and downstream intron boundaries were then annotated depending on the strandedness of the gene they were extracted from.

Interval boundary density of retrotransposons

Intervals were split in half and positions were reckoned relative to the feature adjacent boundary, where the feature was either a gene, exon, or DNase1 cluster (Fig. S2). To calculate the retrotransposon density at each position, we measured the fraction of bases at each position annotated as a retrotransposon. Next, we smoothed retrotransposon densities by calculating the mean and standard deviation of retrotransposon densities within an expanding window, where window size grows as a linear function of distance from the boundary. This made it possible to accurately compare the retrotransposon density at positions where retrotransposon insertions were sparse and density levels at each position fluctuated drastically. At positions with a high base pair density a small window was used and at positions with a low base pair density a large window was used. Expected retrotransposon density p was calculated as the total proportion of bases covered by retrotransposons across

all intervals. Standard deviation at each position was calculated as $\sqrt{np(1-p)}$, where n is
the total number of bases at a given position.

Interval size bias correction of retrotransposon densities

Interval boundary density is sensitive to retrotransposon insertion preferences into intervals
of a certain size (Fig. S3). To determine interval size retrotransposon density bias, we
grouped intervals according to size and measured the retrotransposon density of each interval
size group. Retrotransposon density bias was calculated as the observed retrotransposon
density of an interval size group divided by the expected retrotransposon density, where the
expected retrotransposon density is the total retrotransposon density across all intervals.
Next, using the intervals that contribute to the position depth at each position adjacent
to feature boundaries, we calculated the mean interval size. From this we corrected retro-
transposon density at each position by dividing the observed retrotransposon density by the
retrotransposon density bias that corresponded with that position's mean interval size.

Software and data analysis

All statistical analyses were performed using R (R Core Team 2015) with the packages
GenomicRanges (Lawrence et al. 2013) and rtracklayer (Lawrence et al. 2009). R scripts
used to perform analyses can be found at:
<https://github.com/AdelaideBioinfo/retrotransposonAccumulation> .

Results

Species selection and retrotransposon classification

We selected human, chimpanzee, rhesus macaque, mouse, rabbit, dog, horse and pig as
representative placental species because of their similar non-LTR retrotransposon composition
(Fig. S4-S5) and phylogenetic relationships. Retrotransposon coordinates were obtained
from UCSC repeat masker tables and the online repeat masker database (Rosenbloom et al.
2015; Smit et al. 1996). We grouped non-LTR retrotransposon families according to repeat
type and period of activity as determined by genome-wide defragmentation (Giordano et al.
2007). Retrotransposons were placed into the following groups; new L1s, old L1s, new SINEs

and ancient elements (for families in each group see Fig. S5). New L1s and new SINEs are retrotransposon families with high lineage specificity and activity, while old L1s and ancient elements (SINE MIRs and LINE L2s) are retrotransposon families shared across taxa. We measured sequence similarity within retrotransposon families as percentage mismatch from family consensus sequences (Bao et al. 2015). We found that more recent lineage-specific retrotransposon families had accumulated a lower percentage of substitutions per element than older families (Fig. S6-S13). This confirmed that our classification of retrotransposon groups agreed with ancestral and lineage-specific periods of retrotransposon activity.

Genomic distributions of retrotransposons

To analyse the large scale distribution of retrotransposons, we segmented each species genome into adjacent 1 Mb regions, tallied retrotransposon distributions, performed principal component analysis (PCA) and pairwise correlation analysis (see methods). For PCA, our results showed that retrotransposon families from the same group tended to accumulate in the same genomic regions. We found that each individual retrotransposon group was usually highly weighted in one of the two major principal components (PC1 and PC2) (Fig. 2). Depending on associations between PCs and particular retrotransposon families we identified PC1 and PC2 as either the “lineage-specific PC” or the “ancestral PC”. Along the lineage-specific PC, new SINEs and new L1s were highly weighted, where in all species new SINEs were enriched in regions with few new L1s. Alternatively, along the ancestral PC, old L1s and ancient elements were highly weighted, where in all species except mouse — where ancient elements and old L1s were co-located — ancient elements were enriched in regions with few old L1s (Fig. 2-3a,S14). The discordance observed in mouse probably resulted from the increased genome turnover and rearrangement seen in the rodent lineage potentially disrupting the distribution of ancestral retrotransposon families (Murphy et al. 2005; Capilla et al. 2016). In addition, the genome-wide density of ancestral retrotransposons in mouse was particularly low compared to our other species (Fig. S4-S5). However, as the relationship between mouse lineage-specific new retrotransposons is maintained, this discordance does not impact on downstream analyses. These results show that most genomic context associations between retrotransposon families are conserved across our sample species.

Retrotransposon accumulation and chromatin environment 216

In human and mouse, LINES and SINEs differentially associate with distinct chromatin environments (Ashida et al. 2012). To determine how our retrotransposon groups associate with chromatin accessibility, we obtained ENCODE generated human cell line Repli-Seq data and mouse cell line Repli-ChIP data from the UCSC genome browser (ENCODE Project Consortium 2012; Yue et al. 2014). Repli-Seq and Repli-ChIP both measure the timing of genome replication during S-phase, where accessible euchromatic domains replicate early and inaccessible heterochromatic domains replicate late. Across our segmented genomes, we found a high degree of covariation between genome-wide mean replication timing and lineage-specific PC scores (Fig. 3a), new SINEs associated with early replication and new L1s associated with late replication. In addition, by splitting L1s into old and new groups, we showed a strong association between replication timing and retrotransposon age that was not reported in previous analyses (Pope et al. 2014). These results are probably not specific to a particular cell line, since genome-wide replication timing patterns are mostly highly correlated across cell lines from either species (Table S6). Moreover, early and late replicating domains from various human cell lines exhibit a high degree of overlap (Fig. S15). To confirm that lineage-specific retrotransposon accumulation associates with replication timing, we analysed retrotransposon accumulation at the boundaries of previously identified replication domains (RDs) (Liu et al. 2015). We focused primarily on early replicating domain (ERD) boundaries rather than late replicating domain (LRD) boundaries because ERD boundaries mark the transition from open chromatin states to closed chromatin states and overlap with topologically associated domain (TAD) boundaries (Pope et al. 2014). Consistent with our earlier results, significant density fluctuations at ERD boundaries were only observed for new L1s and new SINEs (Fig. 3b). Because RD timing and genomic distributions of clade-specific retrotransposons are both largely conserved across human and mouse (Ryba et al. 2010; Yaffe et al. 2010), these results suggest that the relationship between retrotransposon accumulation and RD timing may be conserved across mammals.

The genomic distribution of retrotransposons is conserved across species

Our earlier results showed that the genomic distribution of retrotransposons is similar across species (Fig. 2). To determine whether our observations resulted from retrotransposon insertion into orthologous regions, we humanised segmented genomes of non-human species. Humanisation, began with a segmented human genome, a segmented non-human mammalian genome, and a set of pairwise alignments between both species. Using the pairwise alignments we calculated the percentage of nucleotides from each human segment that aligned to a specific non-human segment and vice-versa. This made it possible to remodel the retrotransposon content of each non-human genome segment within the human genome and essentially humanise non-human mammalian genomes (Fig. 1) (see methods). To test the precision of our humanisation process, we used the Kolmogorov-Smirnov test to compare the humanised retrotransposon density distribution of a specific retrotransposon family, to the non-humanised retrotransposon density distribution of that same retrotransposon family (Fig. S1). If the Kolmogorov-Smirnov test returned a low P-value, this suggested that the humanisation process for a given retrotransposon family had a low level of precision. Therefore, to increase our precision we used a minimum mapping fraction threshold to discard genomic segments that had only had a small amount of aligning regions between each genome. The motivation behind this was that genomic segments with a small amount of aligning sequence were the ones most likely to inaccurately represent non-human retrotransposon genomic distributions when humanised. However, it is important to note that our increase in precision requires a trade-off in accuracy. By discarding genomic segments below a certain threshold we sometimes removed a significant fraction of our non-human genomes from the analysis. In addition, this approach disproportionately affected retrotransposons such as new L1s, as they were most enriched in segments with a small amount of aligning regions between each genome (Fig.S16-S17). To overcome this, we humanised each non-human genome at minimum mapping fraction thresholds of 0, 10, 20, 30, 40 and 50 percent and recorded the percentage of the genome that remained. We found that most retrotransposon families were precisely humanised at a minimum mapping fraction threshold of 10%. In non-human species where humanisation was most precise, a minimum mapping fraction threshold of 10% resulted in greater than 90% of the human and non-human genome remaining in the

analysis (Fig. 4,S18-S24). After humanising each non-human genome, we performed pairwise correlation analysis (see methods) between the genomic distributions of each humanised and human retrotransposon family. Our results showed that retrotransposon families in different species that were identified as the same group showed relatively strong correlations, suggesting that they accumulated in regions with shared common ancestry (Fig. 4,S18-S24). Next, we assessed the level of conservation of retrotransposon accumulation patterns across all of our species. For each retrotransposon group in each humanised genome, we identified the top 10% retrotransposon dense genome segments. We found that when these segments were compared with the human genome, there was a relatively high degree of overlap (Fig. 5a-b). These results suggest that lineage-specific retrotransposon accumulation may follow an ancient conserved mammalian genome architecture.

Retrotransposon insertion in open chromatin surrounding regulatory elements

Retrotransposons preferentially insert into open chromatin, yet open chromatin usually overlaps gene regulatory elements. As stated above, this creates a fundamental evolutionary conflict for retrotransposons; their immediate replication may be detrimental to the overall fitness of the genome in which they reside. To investigate retrotransposon insertion/accumulation dynamics at open chromatin regions, we analysed DNase1 hypersensitive activity across 15 cell lines in both ERDs and LRDs. DNase1 hypersensitive sites obtained from the UCSC genome browser (ENCODE Project Consortium 2012) were merged into DNase1 clusters and DNase1 clusters overlapping exons were excluded. As replication is sometimes cell type-specific we also constructed a set of constitutive ERDs and LRDs (cERDs and cLRDs) (see methods). Based on previous analyses, cERDs and cLRDs likely capture RD states present during developmental periods of heritable retrotransposition (Rivera-Mulia et al. 2015). Our cERDs and cLRDs capture approximately 50% of the genome and contain regions representative of genome-wide intron and intergenic genome structure (Fig. S25). In both cERDs and cLRDs, we measured DNase1 cluster activity by counting the number of DNase1 peaks that overlapped each cluster. We found that DNase1 clusters in cERDs were much more active than DNase1 clusters in cLRDs (Fig. 6a). Next, we analysed retrotransposon accumulation both within and at the boundaries of DNase1 clusters. Consistent

with disruption of gene regulation by retrotransposon insertion, non-ancient retrotransposon groups were depleted from DNase1 clusters (Fig. 6b). Intriguingly, ancient element density in DNase1 clusters remained relatively high, suggesting that some ancient elements may have been exapted. At DNase1 cluster boundaries after removing interval size bias (Fig. S26-S27) (see methods), retrotransposon density remained highly enriched in cERDs and close to expected levels in cLRDs (Fig. 6c). This suggests that chromatin is likely to be open at highly active cluster boundaries where insertion of retrotransposons is less likely to disrupt regulatory elements. These results are consistent with an interaction between retrotransposon insertion, open chromatin and regulatory activity, where insertions into open chromatin only persist if they do not interrupt regulatory elements.

Retrotransposon insertion size and regulatory element density

L1s and their associated SINEs differ in size by an order of magnitude, retrotranspose via the L1-encoded chromatin-sensitive L1ORF2P and accumulate in compositionally distinct genomic domains (Cost et al. 2001; Baillie et al. 2011). This suggests that retrotransposon insertion size determines observed accumulation patterns. L1 and *Alu* insertions occur via target-primed reverse transcription which is initiated at the 3' end of each element. With L1 insertion, this process often results in 5' truncation, causing extensive insertion size variation and an over representation of new L1 3' ends, not seen with *Alu* elements (Fig. 7a). When we compared insertion size variation across cERDs and cLRDs we observed that smaller new L1s were enriched in cERDs and *Alu* elements showed no RD insertion size preference (Fig. 7b). The effect of insertion size on retrotransposon accumulation was estimated by comparing insertion rates of each retrotransposon group at DNase1 cluster boundaries in cERDs and cLRDs. We found that *Alu* insertion rates at DNase1 cluster boundaries were similarly above expected levels both in cERDs and cLRDs (Fig. 7c), whereas new L1 insertion rates at DNase1 cluster boundaries were further above expected levels in cERDs than cLRDs (Fig. 7d). By comparing the insertion rate of new L1s — retrotransposons that exhibited RD specific insertion size variation — we observed a negative correlation between element insertion size and gene/regulatory element density. Thus smaller elements, such as *Alu* elements, accumulate more in cERDs than do larger elements, such as new L1s, suggesting that smaller elements are more tolerated.

Retrotransposon insertion within gene and exon structures 334

Regulatory element organisation is largely shaped by gene and exon/intron structure which 335
likely impacts the retrotransposon component of genome architecture. Therefore, we analysed 336
retrotransposons and DNase1 clusters (exon-overlapping and exon non-overlapping) at the 337
boundaries of genes and exons. Human RefSeq gene models were obtained from the UCSC 338
genome browser and both intergenic and intronic regions were extracted (Table S5). At 339
gene (Fig. 8a) and exon (Fig. 8b) boundaries, we found a high density of exon overlapping 340
DNase1 clusters and depletion of retrotransposons. This created a depleted retrotransposon 341
boundary zone (DRBZ) specific for each retrotransposon group, a region extending from 342
the gene or exon boundary to the point where retrotransposon levels begin to increase. 343
The size of each DRBZ correlated with the average insertion size of each retrotransposon 344
group, consistent with larger retrotransposons having a greater capacity to disrupt important 345
structural and regulatory genomic features. We also found that in cERDs the 5' gene 346
boundary *Alu* DRBZ was larger than the 3' gene boundary *Alu* DRBZ. This difference was 347
associated with increased exon overlapping DNase1 cluster density at 5' gene boundaries 348
in cERDs (Fig. 8a), emphasising the importance of evolutionary constraints on promoter 349
architecture. For ancient elements, their retrotransposon density at approximately 1 kb from 350
the 5' gene boundary, when corrected for interval size bias, was significantly higher than 351
expected. This increase is consistent with exaptation of ancient elements into regulatory roles 352
(Lowe et al. 2007) (Fig. S28-S31). Moreover, the density peak corresponding to uncorrected 353
ancient elements also overlapped with that of exon non-overlapping DNase1 clusters (Fig. 354
8a). Collectively, these results demonstrate the evolutionary importance of maintaining gene 355
structure and regulation and how this in turn has canalised similar patterns of accumulation 356
and distribution of retrotransposon families in different species over time. 357

Discussion 358

A conserved architectural framework shapes the genomic distribu- 359 tion of ancestral retrotransposons 360

The majority of divergence between our sample species has taken place over the last 100 361
million years. Throughout this time period many genomic rearrangements have occurred, 362

causing a great deal of karyotypic variation. However, we found that the genomic distributions 363
of ancestral elements remained conserved. The evolutionary forces preserving the ancestral 364
genomic distributions of these elements remain unclear. 365

One suggestion is that ancestral elements play essential roles in mammalian organisms. 366
Our results in Fig. 6b and 8a suggest that ancient elements have been exapted. Their 367
accumulation within open chromatin sites is consistent with their roles as *cis*-regulatory 368
element, such as MIR elements that perform as TFBSs and enhancers (Bourque et al. 2008; 369
Jjingo et al. 2014). Similarly, L1s also carry binding motifs for DNA-binding proteins. L1 370
elements that were active prior to the boreoeutherian ancestor bind a wide variety of KRAB 371
zinc-finger proteins (KZFPs), most of which have unknown functions (Imbeault et al. 2017). 372
In terms of genome structural roles, some human MIR elements have been identified as 373
insulators, separating open chromatin regions from closed chromatin regions (Wang et al. 374
2015). While these MIR insulators function independently of CTCF binding, their mechanism 375
of action remains largely unknown. Despite this, when a human MIR insulator was inserted 376
into the zebrafish genome it was able to maintain function (Wang et al. 2015). This suggests 377
that MIR insulators recruit a highly conserved insulator complex and maintain insulator 378
function across the mammalian lineage. Collectively, these findings identified a number 379
of examples where ancestral elements are associated with important biological roles. This 380
may suggest that genomic distributions of ancestral elements are conserved across mammals 381
because they play conserved biological roles across mammals. However, it is necessary to 382
draw a distinction between evolutionary conservation of an ancient functional element and 383
evolutionary conservation of large-scale genomic distributions of retrotransposons. This 384
is important because for most of our sample species, ancient elements and old L1s each 385
occupy approximately 7% of each of their genomes (Fig. S4). Compared to the 0.04% of the 386
human genome that is comprised of transposable elements under purifying selection (Lowe 387
et al. 2007), this suggests that the vast majority of ancestral elements may not actually play 388
conserved roles in mammalian biology. 389

Rather than ancestral elements playing a conserved role in genome maintenance, their 390
genomic distributions may instead remain conserved as a consequence of evolutionary 391
dynamics occurring at higher order levels of genome architecture. TADs have been identified 392
as a fundamental unit of genome structure, they are approximately 900 kb in length and 393
contain highly self interacting regions of chromatin (Dixon et al. 2012). Despite large-scale 394

genomic rearrangements, the boundaries between TADs have remained conserved across 395
mammals (Dixon et al. 2012). An analysis involving rhesus macaque, dog, mouse and 396
rabbit, identified TAD boundaries at the edge of conserved syntenic regions associating 397
with evolutionary breakpoints between genomic rearrangements (Rudan et al. 2015). This 398
suggests that genome rearrangements occur primarily along TAD boundaries leaving TADs 399
themselves largely intact. Similarly, TAD architecture could also be the driving force behind 400
the observed frequent reuse of evolutionary breakpoints throughout mammalian genome 401
evolution (Murphy et al. 2005). Together these findings suggest that TADs form part of a 402
conserved evolutionary framework whose boundaries are sensitive to genomic rearrangements. 403
Therefore, the current observed genomic distributions of ancestral retrotransposons reflects 404
mostly ancestral retrotransposons that inserted within TADs rather than at their boundaries. 405
This is because elements that accumulated near TAD boundaries were most likely lost 406
through recurrent genomic rearrangements and genome turnover. 407

Another example supporting the idea that conserved genomic distributions are shaped 408
by a conserved architectural evolutionary framework can be found in the rodent lineage. 409
Rodents have experienced rates of genome reshuffling two orders of magnitude greater than 410
other mammalian lineages (Capilla et al. 2016). This has caused rodent genomes to contain a 411
higher number of evolutionary breakpoints, many of which are rodent-specific (Capilla et al. 412
2016). From our analysis we found that old L1s and ancient elements each occupied only 1% 413
of the mouse genome (Fig. S4), with similar levels of ancient elements within the rat genome 414
(Gibbs et al. 2004). Compared to our other species where the genomes are approximately 415
7% ancient elements and old L1s each (S4), rodent genomes are significantly depleted of 416
ancestral elements. Together, these findings show a negative correlation between ancestral 417
retrotransposon content and rate of genome rearrangements, suggesting that increased 418
rates of genome rearrangements can strongly impact the genomic distributions of ancestral 419
retrotransposons. In addition, the large number of rodent specific evolutionary breakpoints 420
may explain why the genomic distribution of ancestral elements in mouse is discordant with 421
our other species. Specifically, ancient elements and old L1s in mouse accumulated in similar 422
regions, whereas in each of our other species ancient elements and old L1s accumulated in 423
almost opposite regions as defined by PC1 (Fig. 2,3a). 424

Conserved genome architecture drives the accumulation patterns of lineage-retrotransposons

Across mammals, lineage-specific retrotransposons are responsible for the vast majority of lineage-specific DNA gain (Kapusta et al. 2017). Throughout our sample-species we found that new SINEs and new L1s independently accumulated in similar regions in different species. These results suggest there is a high degree of conservation surrounding their insertion mechanisms and genomic environments. Since, L1 conservation in mammals is well documented in the literature and our new SINE families all replicate using L1 machinery, mainly we spend this section discussing the role of conserved genome architecture (Ivancevic et al. 2016; Vassetzky and Kramerov 2013).

Earlier, we discussed the importance of TADs and how they form a fundamental component of conserved genome architecture. This same architectural framework may also shape the accumulation pattern of lineage specific retrotransposons. TAD boundaries separate the genome into regions comprised of genes that are largely regulated by a restricted set of nearby enhancers. Moreover, TADs are subject to large-scale changes in chromatin structure, where individual TADs are known to switch between open and closed chromatin states in a cell type-specific manner (Dixon et al. 2012). One method of capturing shifts in chromatin state between TADs is to measure genome-wide replication timing (Pope et al. 2014). This is because replication timing associates with the genomes accessibility to replication machinery. Accessible regions that comprise an open chromatin structure replicate early while inaccessible regions with a closed chromatin structure replicate late. Genome-wide replication timing follows a domain-like organisation, where large contiguous regions either replicate at earlier or later stages of mitosis. Importantly, ERD boundaries directly overlap TAD boundaries, supporting the notion that TADs are also fundamental units of large-scale chromatin state organisation (Pope et al. 2014). Previously, LINE and SINE accumulation patterns were associated with TAD and RD genome architecture, where LINEs were enriched in LRDs and SINEs were enriched in ERDs (Hansen et al. 2010; Rivera-Mulia et al. 2015; Pope et al. 2014; Ashida et al. 2012). Unlike our analysis, these earlier studies decided not to separate LINEs into ancestral and lineage-specific families. Despite this difference, Fig. 3 shows that our results are consistent with earlier analyses, except for our observation that only lineage-specific retrotransposon families are associated with replication timing. Therefore,

by separating L1s and SINEs according to period of activity, we observed much stronger 456
associations between replication timing and retrotransposon accumulation than previously 457
reported (Pope et al. 2014; Ashida et al. 2012). Since replication timing and boundaries 458
between TADs and RDs are conserved across mammalian species (Ryba et al. 2010; Yaffe 459
et al. 2010; Pope et al. 2014; Dixon et al. 2012), our results suggest that domain-level 460
genome architecture likely plays a role in shaping conserved lineage-specific retrotransposon 461
accumulation patterns. 462

While our species genomes are conserved at a structural level, conserved patterns of 463
lineage-specific retrotransposon accumulation can have significant evolutionary impacts. 464
new SINEs accumulate in ERDs which tend to be highly active gene-rich genomic regions. 465
However, despite the fact that all of our new SINE families follow L1 mediated replication, 466
they stem from unique origins. For example, Primate-specific *Alu* elements are derived from 467
7SL RNA and carnivora-specific SINEC elements are derived from tRNA (Quentin 1994; 468
Coltman and Wright 1994). Due to their large-scale accumulation patterns this means that 469
new SINEs in mammalian genomes simultaneously drive convergence in genome architecture 470
and divergence in genome sequence composition. This is especially important because SINEs 471
are also a large source of evolutionary innovation for gene regulation. In human, various 472
individual *Alu* elements have been identified as bona fide enhancers with many more believed 473
to be proto-enhancers serving as a repertoire for birth of new enhancers (Su et al. 2014). 474
Similarly, in dog, mouse and opossum, lineage specific SINEs carry CTCF binding sites and 475
have driven the expansion of species-specific CTCF binding patterns (Schmidt et al. 2012). 476

Like new SINEs, new L1s also accumulate in similar regions in different species. However, 477
unlike new SINEs, lineage-specific mammalian L1 elements most likely stem from a common 478
ancestor (Furano et al. 2004). This means that individual new L1 elements in different 479
species are more likely than species-specific SINEs to share similar sequence composition 480
(Ivancevic et al. 2016). Therefore, LRDs, which are enriched for new L1s, may show higher 481
levels of similarity for genome sequence composition than ERDs, which are enriched for 482
new SINEs. Considering results from genome-wide alignments between mammals, this may 483
be counter intuitive, mainly because the surrounding sequence in new L1 enriched regions 484
exhibits poor sequence conservation (Fig. S16-S17). However, it is important to realise that 485
similar sequence composition is not the same as sequence conservation itself, especially at 486
the level of mammalian genome architecture. Sequence composition refers to the kinds of 487

sequences in a particular region rather than the entire sequence of the region itself. For 488
example, binding sites for the same transcription factor in different species are sometimes 489
located in similar regions yet differ in position relative to their target genes (Kunarso et al. 490
2010). So while genome-wide alignments may suggest low levels of genome conservation or 491
high levels of turnover, sequence composition within these regions remains similar and can 492
still be indicative of conserved function. Therefore with the accumulation of new L1s after 493
species divergence, it is likely that sequence conservation decreases at a much faster rate 494
than compositional similarity. For new L1s enriched in similar regions in different species, 495
this may have important functional consequences. Recently, highly conserved ancient KZFPs 496
were discovered to bind to members of both old and new L1 families in human (Imbeault 497
et al. 2017). This suggests that new L1s in humans may be interchangeable with old L1s 498
and play important roles in highly conserved gene regulatory networks. Therefore, because 499
new L1s in different species share similar sequences and their accumulation patterns are 500
also conserved, new L1s may actively preserve ancient gene regulatory networks across the 501
mammalian lineage. 502

A chromatin based model of retrotransposon accumulation 503

Analysis of repetitive elements in mammalian genome sequencing projects has consistently 504
revealed that L1s accumulate in GC-poor regions and their mobilised SINEs accumulate in 505
GC-rich regions (Lander et al. 2001; Gibbs et al. 2004; Chinwalla et al. 2002). Our results 506
were consistent with this and showed that accumulation patterns of new SINEs and new 507
L1s were conserved across species and corresponded with distinct genomic environments. 508
Since these elements both replicate via the same machinery, their accumulation patterns 509
are most likely shaped by how insertion of each element type interacts with its immediate 510
genomic environment. The current model of retrotransposon accumulation begins with 511
random insertion, constrained by local sequence composition, followed by immediate selection 512
against harmful insertions (Graham and Boissinot 2006; Gasior et al. 2007; Kvikstad and 513
Makova 2010). During early embryogenesis or in the germline, it is believed retrotransposons 514
in individual cells randomly insert into genomic loci that contain a suitable insertion motif. 515
Because this process is assumed to be random, new insertions can occasionally interrupt 516
essential genes or gene regulatory structures. These insertions are usually harmful, causing 517

the individual cell carrying them to be quickly removed from the population. This process of 518
purifying selection prevents harmful insertions from being passed down to the next generation 519
and plays a large role in shaping retrotransposon accumulation patterns. According to this 520
model, because of their size difference L1s are considered to have a more harmful impact on 521
nearby genes and gene regulatory structures than SINEs. New L1 insertion into GC-rich 522
regions, which are also gene-rich, are more likely to cause harm than if new SINEs inserted 523
into those same regions. Therefore, new L1s are evolutionary purged from GC rich regions 524
causing them to become enriched in gene-poor AT-rich regions. While this model is simple, it 525
fails to take into account the impact of chromatin structure that constrains retrotransposon 526
insertion preference. Therefore, we decided to analyse retrotransposon accumulation at the 527
level of large-scale chromosomal domains and fine-scale open chromatin sites. 528

Our results showed that lineage-specific retrotransposons accumulated at the boundaries of 529
open chromatin sites. This was particularly striking as it appeared to reconcile insertion into 530
open chromatin with the risk of disrupting regulatory elements. Single cell analysis has shown 531
somatic retrotransposition events correlate with preferable insertion into open chromatin 532
sites or within actively expressed genes (Klawitter et al. 2016; Upton et al. 2015; Baillie 533
et al. 2011). However, because open chromatin usually surrounds regulatory elements these 534
kinds of insertions can be a major cause of genetic disease (Wimmer et al. 2011). Therefore, 535
retrotransposons accumulate in open chromatin regions where their insertion is less likely 536
to disrupt regulatory elements. We further demonstrated the impact of retrotransposon 537
insertion by considering element insertion size. Our results showed that shorter L1s were 538
much more likely to insert close to open chromatin sites surrounding regulatory elements 539
than larger L1s. This suggested that L1 insertions were much more likely than *Alu* insertions 540
to impact on gene regulatory structures due to their larger insertion size. At this point, it 541
should be noted that chromatin state can be highly dynamic, switching between open and 542
closed states depending on cell type (ENCODE Project Consortium 2012). Importantly, 543
heritable retrotransposon insertions typically occur during embryogenesis or within the 544
germline. However, chromatin state data for these developmental stages and tissue samples 545
was unavailable. To overcome this limitation we aggregated data from a range of biological 546
contexts. The underlying assumption behind this strategy was that open chromatin sites 547
found in at least one cell likely contain regulatory elements that may be reused in another 548
cell type. Therefore, by using this strategy, we increased the probability of capturing 549

chromosomal domain structures and regulatory element sites present in embryonic and 550
germline cell states. 551

An important aspect of both our refined model and the current model of retrotransposon 552
accumulation is the immediate evolutionary impact of retrotransposon insertions. Specifically, 553
at what rate do embryonic and germline retrotransposition events occur and what proportion 554
of these events escape purifying selection? Answering this question is a challenging task 555
primarily limited by the availability of samples at the correct developmental time periods. 556
Ideally we would require genome sequencing data from a large population of germline or 557
embryonic cells derived from a similar genetic background. Given that data, we could identify 558
new insertions before they have undergone selection and compare their retrotransposition rates 559
to retrotransposition rates inferred from population data. Alternatively, retrotransposition 560
rates have been measured in somatic cells and stem-cell lines. In hippocampal neurons 561
and glia, L1 retrotransposition occurs at rates of 13.7 and 6.5 events per cell, where in 562
human induced pluripotent stem cells retrotransposition rates are approximately 1 event per 563
cell (Klawitter et al. 2016; Upton et al. 2015). In neurons, L1s insertions were enriched in 564
neuronally expressed genes and in human induced pluripotent stem cells, L1s were found 565
to insert near transcription start sites, disrupting the expression of some genes (Klawitter 566
et al. 2016; Upton et al. 2015; Baillie et al. 2011). This suggests L1s are particularly active 567
in humans, able to induce a large amount of variation and disrupt gene regulation and 568
function. It is also important to note that the estimated L1 heritable retrotransposition rate 569
is approximately one event per 95 to 270 births (Ewing and Kazazian 2010), suggesting that 570
many insertions are removed from the germline cell population. For *Alu* elements this rate is 571
much greater, *Alu* elements are estimated to undergo heritable retrotransposition at a rate 572
of one event per 20 births (Cordaux et al. 2006). These findings support the notion that 573
the majority of retrotransposon insertions are likely to be evolutionarily purged from the 574
genome. 575

In summary, by analysing open chromatin sites, we found that 1) following preferential 576
insertion into open chromatin domains, retrotransposons were tolerated adjacent to regu- 577
latory elements where they were less likely to cause harm; 2) element insertion size was 578
a key factor affecting retrotransposon accumulation, where large elements accumulated in 579
gene poor regions where they were less likely to perturb gene regulation; and 3) insertion 580
patterns surrounding regulatory elements were persistent at the gene level. From this we 581

propose a significant change to the current retrotransposon accumulation model; rather 582
than random insertion constrained by local sequence composition, we propose that insertion 583
is instead primarily constrained by local chromatin structure. Therefore, L1s and SINEs 584
both preferentially insert into gene/regulatory element rich euchromatic domains, where L1s 585
with their relatively high mutational burden are quickly eliminated via purifying selection 586
at a much higher rate than SINEs. Over time this results in an enrichment of SINEs in 587
euchromatic domains and an enrichment of L1s in heterochromatic domains. 588

Conclusion 589

In conjunction with large scale conservation of synteny (Chowdhary et al. 1998), gene 590
regulation (Chan et al. 2009) and the structure of RDs/TADs (Dixon et al. 2012; Ryba et al. 591
2010), our findings suggest that large scale positional conservation of old and new non-LTR 592
retrotransposons results from their association with the regulatory activity of large genomic 593
domains. Therefore we propose that similar constraints on insertion and accumulation of 594
clade specific retrotransposons in different species can define common trajectories for genome 595
evolution. 596

Additional Files 597

Additional file 1 — Supplementary information 598

Figures S1–S31, Tables S1–S6. 599

Competing interests 600

The authors declare that they have no competing interests. 601

Author's contributions 602

R.M.B., R.D.K., J.M.R., and D.L.A. designed research; R.M.B. performed research; and 603
R.M.B., R.D.K., and D.L.A. wrote the paper. 604

Acknowledgements

605

For reviewing our manuscript and providing helpful advice we would like to thank the
following: Simon Baxter, Atma Ivancevic and Lu Zeng from the University of Adelaide;
Kirsty Kitto from Queensland University of Technology; and Udaya DeSilva from Oklahoma
State University.

606

607

608

609

Availability of data and materials

610

All data was obtained from publicly available repositories, urls can be found in supporting
material (Table S1–S4). R scripts used to perform analyses can be found at
<https://github.com/AdelaideBioinfo/retrotransposonAccumulation>.

611

612

613

References

- Ashida, H., Asai, K., and Hamada, M. (2012). Shape-based alignment of genomic landscapes in multi-scale resolution. *Nucleic acids research*, 40(14):6435–6448.
- Baillie, J. K., Barnett, M. W., Upton, K. R., Gerhardt, D. J., Richmond, T. A., De Sapio, F., Brennan, P. M., Rizzu, P., Smith, S., Fell, M., et al. (2011). Somatic retrotransposition alters the genetic landscape of the human brain. *Nature*, 479(7374):534–537.
- Bao, W., Kojima, K. K., and Kohany, O. (2015). Repbase update, a database of repetitive elements in eukaryotic genomes. *Mobile DNA*, 6(1):1.
- Bourque, G., Leong, B., Vega, V. B., Chen, X., Lee, Y. L., Srinivasan, K. G., Chew, J.-L., Ruan, Y., Wei, C.-L., Ng, H. H., et al. (2008). Evolution of the mammalian transcription factor binding repertoire via transposable elements. *Genome research*, 18(11):1752–1762.
- Capilla, L., Sánchez-Guillén, R. A., Farre, M., Paytuví-Gallart, A., Malinverni, R., Ventura, J., Larkin, D. M., and Ruiz-Herrera, A. (2016). Mammalian comparative genomics reveals genetic and epigenetic features associated with genome reshuffling in rodentia. *Genome Biology and Evolution*, page evw276.
- Chalopin, D., Naville, M., Plard, F., Galiana, D., and Volff, J.-N. (2015). Comparative analysis of transposable elements highlights mobilome diversity and evolution in vertebrates. *Genome biology and evolution*, 7(2):567–580.

- Chan, E. T., Quon, G. T., Chua, G., Babak, T., Trochesset, M., Zirngibl, R. A., Aubin, J., Ratcliffe, M. J., Wilde, A., Brudno, M., et al. (2009). Conservation of core gene expression in vertebrate tissues. *Journal of biology*, 8(3):1.
- Chinwalla, A. T., Cook, L. L., Delehaunty, K. D., Fewell, G. A., Fulton, L. A., Fulton, R. S., Graves, T. A., Hillier, L. W., Mardis, E. R., McPherson, J. D., et al. (2002). Initial sequencing and comparative analysis of the mouse genome. *Nature*, 420(6915):520–562.
- Chowdhary, B. P., Raudsepp, T., Frönicke, L., and Scherthan, H. (1998). Emerging patterns of comparative genome organization in some mammalian species as revealed by zoo-fish. *Genome research*, 8(6):577–589.
- Coltman, D. W. and Wright, J. M. (1994). Can sines: a family of trna-derived retroposons specific to the superfamily canoidea. *Nucleic acids research*, 22(14):2726–2730.
- Cordaux, R., Hedges, D. J., Herke, S. W., and Batzer, M. A. (2006). Estimating the retrotransposition rate of human alu elements. *Gene*, 373:134–137.
- Cost, G. J., Feng, Q., Jacquier, A., and Boeke, J. D. (2002). Human l1 element target-primed reverse transcription in vitro. *The EMBO Journal*, 21(21):5899–5910.
- Cost, G. J., Golding, A., Schlissel, M. S., and Boeke, J. D. (2001). Target dna chromatinization modulates nicking by l1 endonuclease. *Nucleic acids research*, 29(2):573–577.
- Dixon, J. R., Selvaraj, S., Yue, F., Kim, A., Li, Y., Shen, Y., Hu, M., Liu, J. S., and Ren, B. (2012). Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*, 485(7398):376–380.
- ENCODE Project Consortium (2012). An integrated encyclopedia of dna elements in the human genome. *Nature*, 489(7414):57–74.
- Ewing, A. D. and Kazazian, H. H. (2010). High-throughput sequencing reveals extensive variation in human-specific l1 content in individual human genomes. *Genome research*, 20(9):1262–1270.
- Furano, A. V., Duvernell, D. D., and Boissinot, S. (2004). L1 (line-1) retrotransposon diversity differs dramatically between mammals and fish. *Trends in Genetics*, 20(1):9–14.

- Gasior, S. L., Preston, G., Hedges, D. J., Gilbert, N., Moran, J. V., and Deininger, P. L. (2007). Characterization of pre-insertion loci of de novo l1 insertions. *Gene*, 390(1):190–198.
- Gibbs, R. A., Weinstock, G. M., Metzker, M. L., Muzny, D. M., Sodergren, E. J., Scherer, S., Scott, G., Steffen, D., Worley, K. C., Burch, P. E., et al. (2004). Genome sequence of the brown norway rat yields insights into mammalian evolution. *Nature*, 428(6982):493–521.
- Gibcus, J. H. and Dekker, J. (2013). The hierarchy of the 3d genome. *Molecular cell*, 49(5):773–782.
- Giordano, J., Ge, Y., Gelfand, Y., Abrusán, G., Benson, G., and Warburton, P. E. (2007). Evolutionary history of mammalian transposons determined by genome-wide defragmentation. *PLoS Comput Biol*, 3(7):e137.
- Graham, T. and Boissinot, S. (2006). The genomic distribution of l1 elements: the role of insertion bias and natural selection. *BioMed Research International*, 2006.
- Hansen, R. S., Thomas, S., Sandstrom, R., Canfield, T. K., Thurman, R. E., Weaver, M., Dorschner, M. O., Gartler, S. M., and Stamatoyannopoulos, J. A. (2010). Sequencing newly replicated dna reveals widespread plasticity in human replication timing. *Proceedings of the National Academy of Sciences*, 107(1):139–144.
- Imbeault, M., Hellebood, P.-Y., and Trono, D. (2017). Krab zinc-finger proteins contribute to the evolution of gene regulatory networks. *Nature*, 543(7646):550–554.
- Ivancevic, A. M., Kortschak, R. D., Bertozzi, T., and Adelson, D. L. (2016). Lines between species: Evolutionary dynamics of line-1 retrotransposons across the eukaryotic tree of life. *Genome Biology and Evolution*, 8(11):3301–3322.
- Jjingo, D., Conley, A. B., Wang, J., Mariño-Ramírez, L., Lunyak, V. V., and Jordan, I. K. (2014). Mammalian-wide interspersed repeat (mir)-derived enhancers and the regulation of human gene expression. *Mobile DNA*, 5(1):1.
- Kapusta, A., Suh, A., and Feschotte, C. (2017). Dynamics of genome size evolution in birds and mammals. *Proceedings of the National Academy of Sciences*, 114(8):E1460–E1469.
- Klawitter, S., Fuchs, N. V., Upton, K. R., Muñoz-Lopez, M., Shukla, R., Wang, J., Garcia-Cañadas, M., Lopez-Ruiz, C., Gerhardt, D. J., Sebe, A., et al. (2016). Reprogramming

- triggers endogenous l1 and alu retrotransposition in human induced pluripotent stem cells. *Nature communications*, 7.
- Kramerov, D. and Vassetzky, N. (2011). Origin and evolution of sines in eukaryotic genomes. *Heredity*, 107(6):487–495.
- Kunarso, G., Chia, N.-Y., Jeyakani, J., Hwang, C., Lu, X., Chan, Y.-S., Ng, H.-H., and Bourque, G. (2010). Transposable elements have rewired the core regulatory network of human embryonic stem cells. *Nature genetics*, 42(7):631–634.
- Kvikstad, E. M. and Makova, K. D. (2010). The (r) evolution of sine versus line distributions in primate genomes: sex chromosomes are important. *Genome research*, 20(5):600–613.
- Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., et al. (2001). Initial sequencing and analysis of the human genome. *Nature*, 409(6822):860–921.
- Lawrence, M., Gentleman, R., and Carey, V. (2009). rtracklayer: an r package for interfacing with genome browsers. *Bioinformatics*, 25:1841–1842.
- Lawrence, M., Huber, W., Pagès, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M., and Carey, V. (2013). Software for computing and annotating genomic ranges. *PLoS Computational Biology*, 9.
- Liu, F., Ren, C., Li, H., Zhou, P., Bo, X., and Shu, W. (2015). De novo identification of replication-timing domains in the human genome by deep learning. *Bioinformatics*, page btv643.
- Lowe, C. B., Bejerano, G., and Haussler, D. (2007). Thousands of human mobile element fragments undergo strong purifying selection near developmental genes. *Proceedings of the National Academy of Sciences*, 104(19):8005–8010.
- Medstrand, P., Van De Lagemaat, L. N., and Mager, D. L. (2002). Retroelement distributions in the human genome: variations associated with age and proximity to genes. *Genome research*, 12(10):1483–1495.
- Murphy, W. J., Larkin, D. M., Everts-van der Wind, A., Bourque, G., Tesler, G., Auvin, L., Beever, J. E., Chowdhary, B. P., Galibert, F., Gatzke, L., et al. (2005). Dynamics of

- mammalian chromosome evolution inferred from multispecies comparative maps. *Science*, 309(5734):613–617.
- Pope, B. D., Ryba, T., Dileep, V., Yue, F., Wu, W., Denas, O., Vera, D. L., Wang, Y., Hansen, R. S., Canfield, T. K., et al. (2014). Topologically associating domains are stable units of replication-timing regulation. *Nature*, 515(7527):402–405.
- Quentin, Y. (1994). Emergence of master sequences in families of retroposons derived from 7sl rna. *Genetica*, 93(1-3):203–215.
- R Core Team (2015). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rivera-Mulia, J. C., Buckley, Q., Sasaki, T., Zimmerman, J., Didier, R. A., Nazor, K., Loring, J. F., Lian, Z., Weissman, S., Robins, A. J., et al. (2015). Dynamic changes in replication timing and gene expression during lineage specification of human pluripotent stem cells. *Genome research*.
- Rosenbloom, K. R., Armstrong, J., Barber, G. P., Casper, J., Clawson, H., Diekhans, M., Dreszer, T. R., Fujita, P. A., Guruvadoo, L., Haeussler, M., et al. (2015). The ucsc genome browser database: 2015 update. *Nucleic acids research*, 43(D1):D670–D681.
- Rudan, M. V., Barrington, C., Henderson, S., Ernst, C., Odom, D. T., Tanay, A., and Hadjur, S. (2015). Comparative hi-c reveals that ctcf underlies evolution of chromosomal domain architecture. *Cell reports*, 10(8):1297–1309.
- Ryba, T., Hiratani, I., Lu, J., Itoh, M., Kulik, M., Zhang, J., Schulz, T. C., Robins, A. J., Dalton, S., and Gilbert, D. M. (2010). Evolutionarily conserved replication timing profiles predict long-range chromatin interactions and distinguish closely related cell types. *Genome research*, 20(6):761–770.
- Schmidt, D., Schwalie, P. C., Wilson, M. D., Ballester, B., Gonçalves, Â., Kutter, C., Brown, G. D., Marshall, A., Flicek, P., and Odom, D. T. (2012). Waves of retrotransposon expansion remodel genome organization and ctcf binding in multiple mammalian lineages. *Cell*, 148(1):335–348.
- Smit, A. F., Hubley, R., and Green, P. (1996). Repeatmasker open-3.0.

- Su, M., Han, D., Boyd-Kirkup, J., Yu, X., and Han, J.-D. J. (2014). Evolution of alu elements toward enhancers. *Cell reports*, 7(2):376–385.
- Upton, K. R., Gerhardt, D. J., Jesuadian, J. S., Richardson, S. R., Sánchez-Luque, F. J., Bodea, G. O., Ewing, A. D., Salvador-Palomeque, C., van der Knaap, M. S., Brennan, P. M., et al. (2015). Ubiquitous l1 mosaicism in hippocampal neurons. *Cell*, 161(2):228–239.
- Vassetzky, N. S. and Kramerov, D. A. (2013). Sinebase: a database and tool for sine analysis. *Nucleic acids research*, 41(D1):D83–D89.
- Wang, J., Vicente-García, C., Seruggia, D., Moltó, E., Fernandez-Miñán, A., Neto, A., Lee, E., Gómez-Skarmeta, J. L., Montoliu, L., Lunyak, V. V., et al. (2015). Mir retrotransposon sequences provide insulators to the human genome. *Proceedings of the National Academy of Sciences*, 112(32):E4428–E4437.
- Wimmer, K., Callens, T., Wernstedt, A., and Messiaen, L. (2011). The nf1 gene contains hotspots for l1 endonuclease-dependent de novo insertion. *PLoS Genet*, 7(11):e1002371.
- Yaffe, E., Farkash-Amar, S., Polten, A., Yakhini, Z., Tanay, A., and Simon, I. (2010). Comparative analysis of dna replication timing reveals conserved large-scale chromosomal architecture. *PLoS Genet*, 6(7):e1001011.
- Yue, F., Cheng, Y., Breschi, A., Vierstra, J., Wu, W., Ryba, T., Sandstrom, R., Ma, Z., Davis, C., Pope, B. D., et al. (2014). A comparative encyclopedia of dna elements in the mouse genome. *Nature*, 515(7527):355–364.

Figures

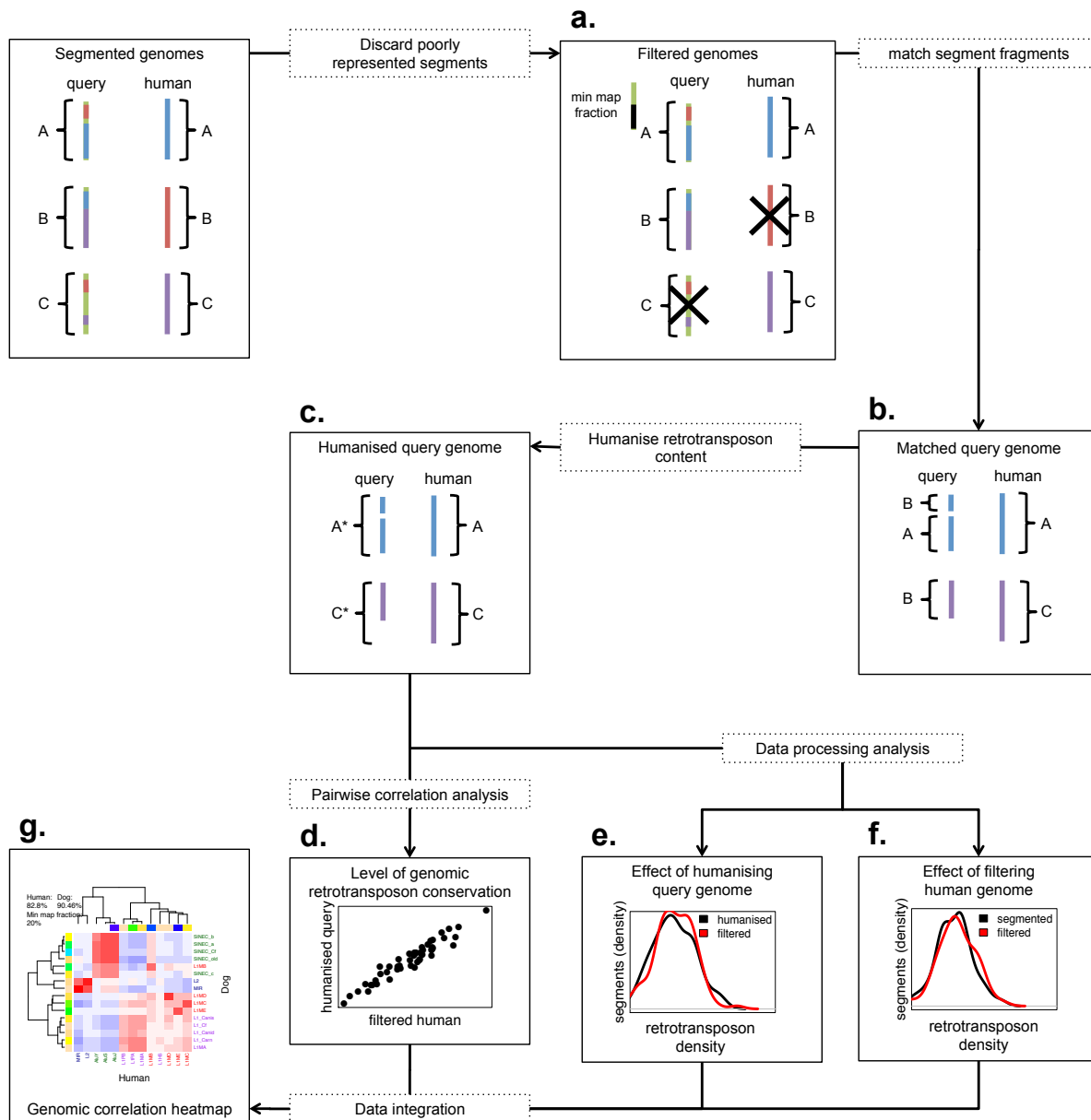


Figure 1. Overview of humanising retrotransposon distributions. **a**, Genomes are segmented and filtered according to a minimum mapping fraction threshold, removing poorly represented segments from both species. The black X shows which segments were not able to reach the minimum mapping fraction threshold. **b**, Fragments of query species' genome segments are matched to their corresponding human genome segments using genome alignments. **c**, Query species genomes are humanised following equation 1. **d**, Pairwise genomic correlations are measured between each humanised retrotransposon family and each human retrotransposon family. **e**, The effect of humanising on retrotransposon density distributions is measured by performing a Kolmogorov-Smirnov test between the humanised query retrotransposon density distribution and the filtered query retrotransposon density distribution. **f**, The effect of filtering on retrotransposon density distributions is measured by performing a Kolmogorov-Smirnov test between the segmented human retrotransposon density distribution and the filtered human retrotransposon density distribution. **g**, The pairwise correlation analysis results and the P-values from the Kolmogorov-Smirnov tests are integrated into heatmaps (Fig. 4,S18-S22) that compare the genomic relationships of retrotransposons between species.

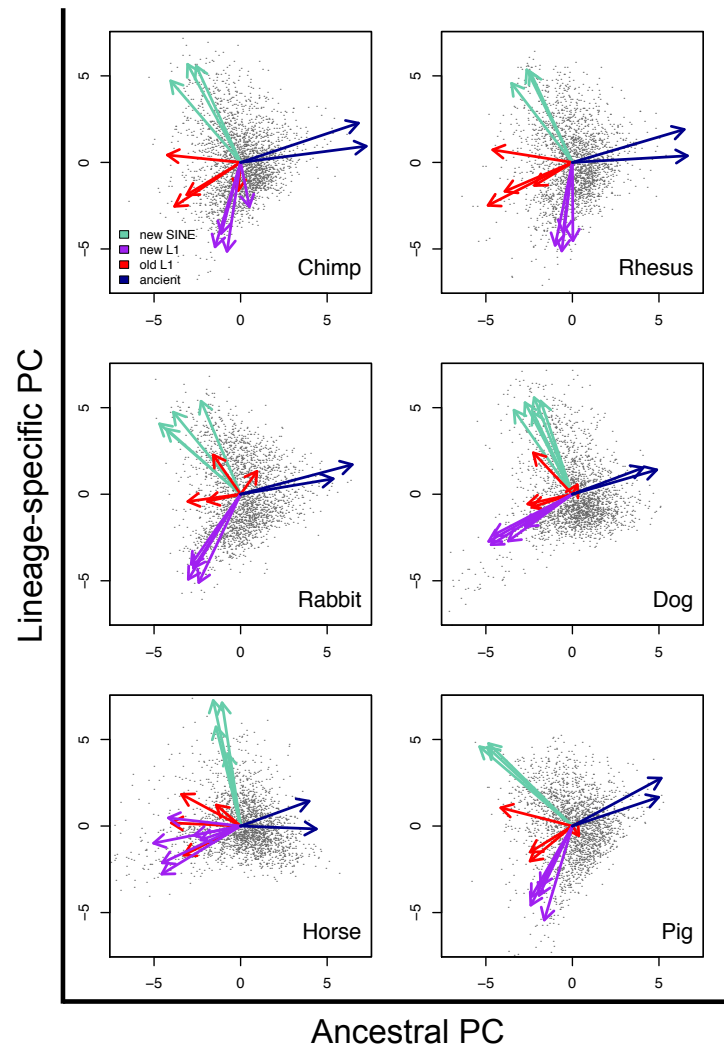


Figure 2. Similar genomic distributions of retrotransposons across mammals. Principal Component 1 and Principal Component 2 of non-human and non-mouse genome retrotransposon content, each vector loading has been coloured according to the retrotransposon group it represents. Principal components have been renamed according to the retrotransposon group whose variance they principally account for.

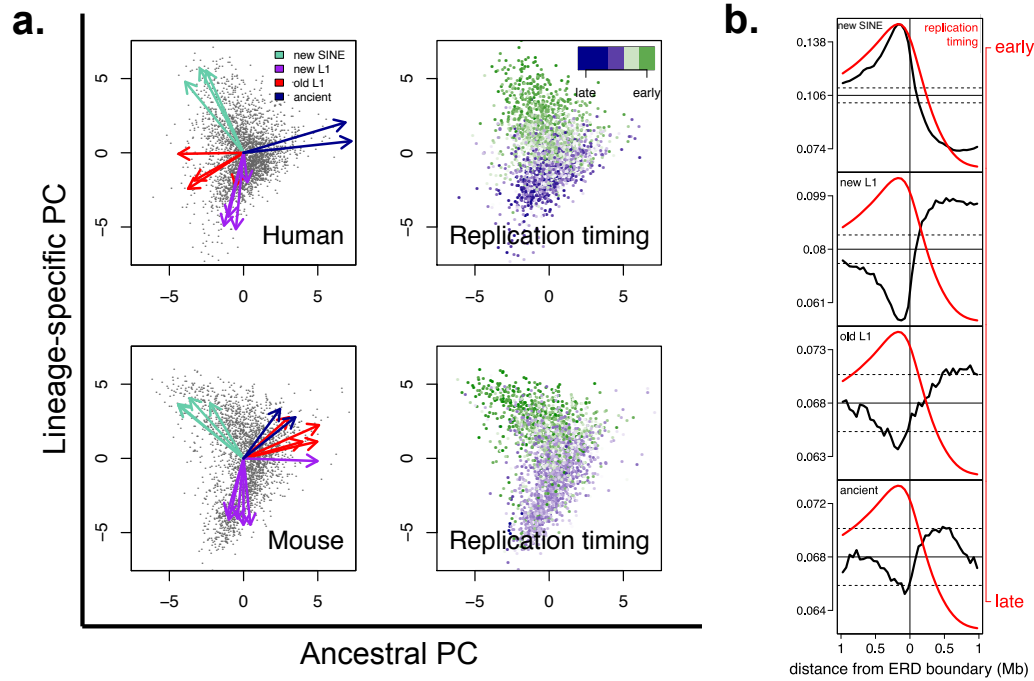


Figure 3. a, PCA of human and mouse retrotransposon content and mean genome replication timing in human HUVEC cells and mouse EpiSC-5 cells. b, Retrotransposon density per non-overlapping 50 kb intervals from a pooled set of ERD boundaries across all 16 human cell lines. Black dashed lines indicate 2 standard deviations from the mean (solid horizontal black line). Red line indicates mean replication timing across all samples.

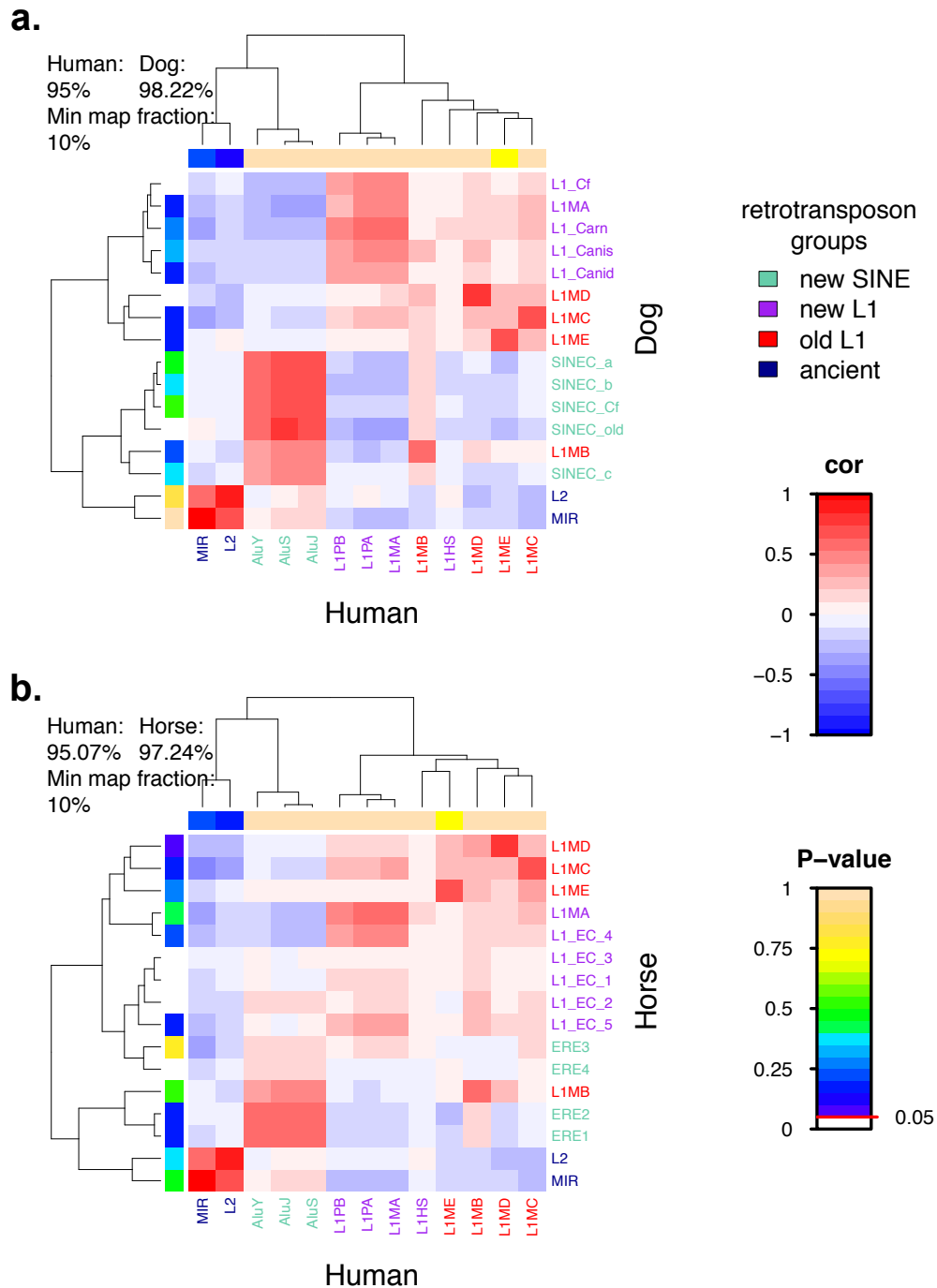


Figure 4. Genome-wide spatial correlations of humanised retrotransposon families. Heatmap colours represent Pearson's correlation coefficient for genomic distributions between humanised **a**, dog and human retrotransposon families, and humanised **b**, horse and human retrotransposon families. Values at the top left of each heatmap reflect the proportion of each genome analysed after filtering at a 10% minimum mapping fraction threshold (Fig. 1a). Dog and horse P-values represent the effect of humanising on filtered non-human retrotransposon density distributions (Fig. 1e). Human P-values represent the effect of filtering on the human retrotransposon density distributions (Fig. 1f).

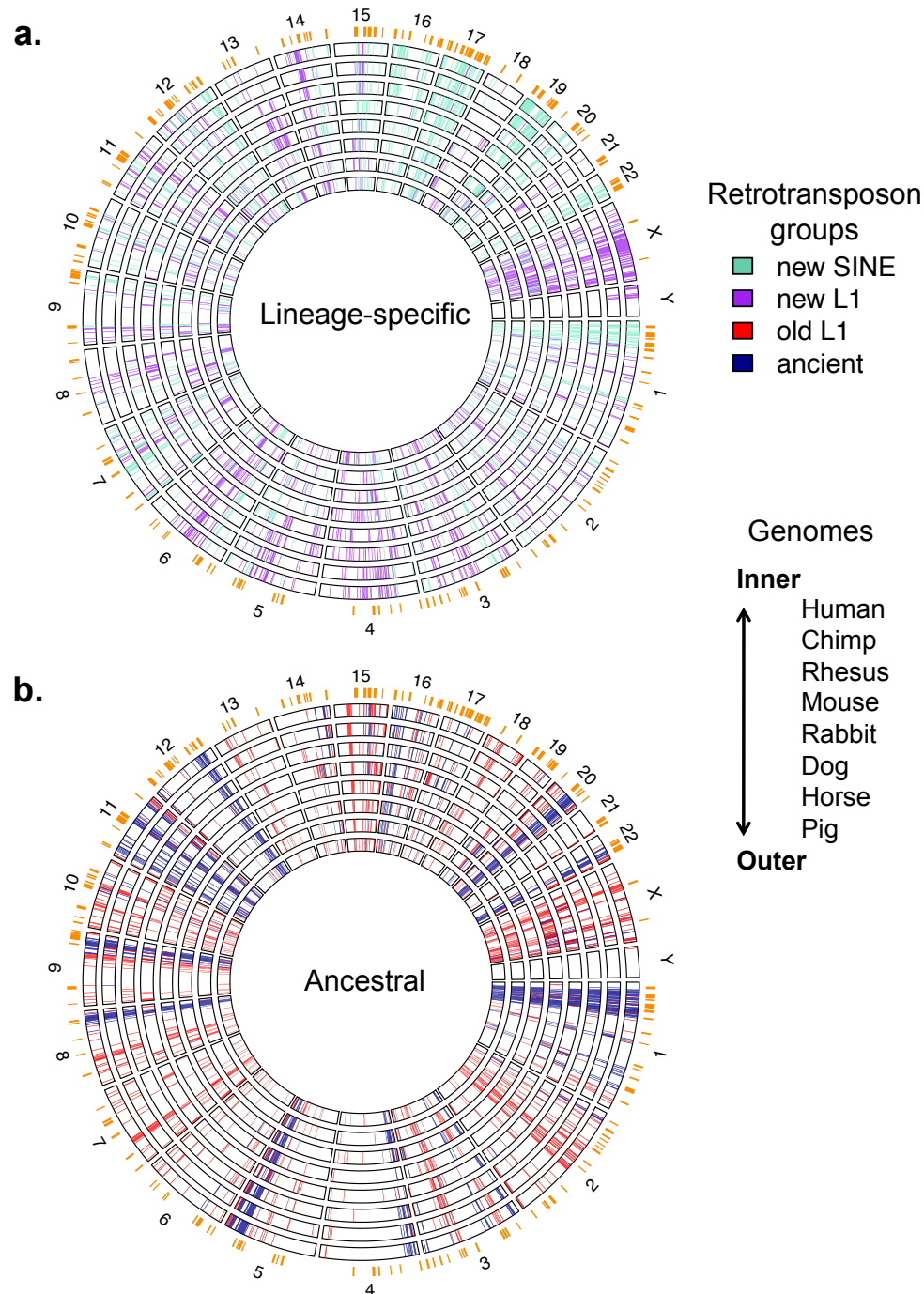


Figure 5. Retrotransposon accumulation patterns are conserved across mammals. **a**, Top 10% of genome segments based on retrotransposon density of new SINEs and new L1s. **b**, Top 10% of genome segments based on retrotransposon density of ancient elements and old L1s. In both **a** and **b**, segments for non-human genomes were ranked according to their humanised values. Large ERDs (> 2 Mb) from HUVEC cells are marked in orange.

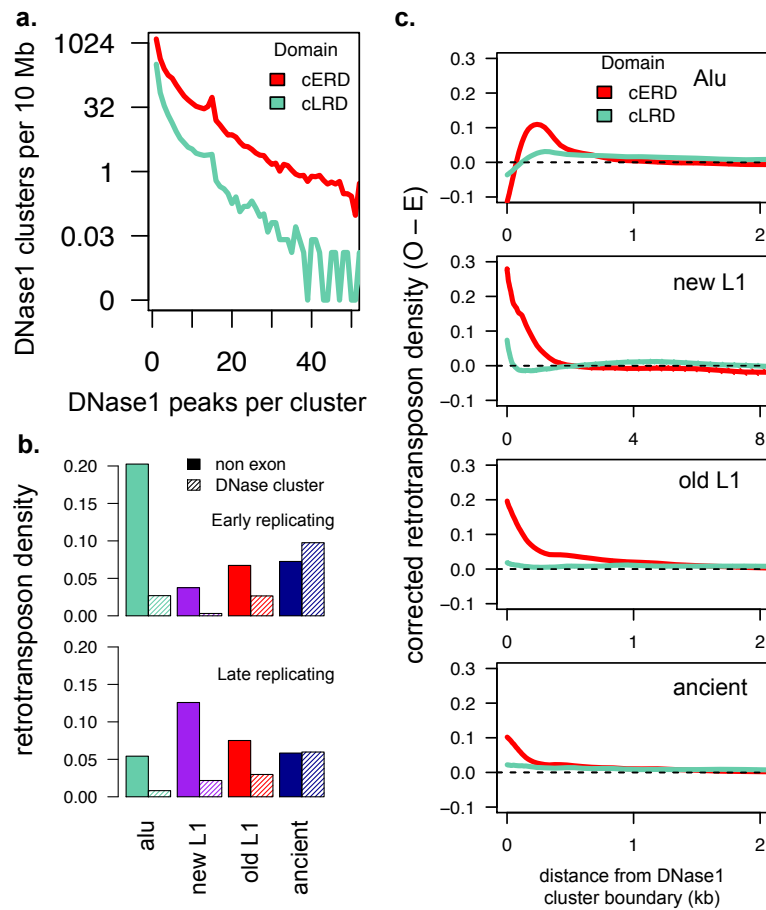


Figure 6. Retrotransposon accumulation occurs in open chromatin near regulatory regions. **a**, The activity of DNase1 clusters in cERDs and cLRDs. DNase1 clusters were identified by merging DNase1 hypersensitive sites across 15 tissues. Their activity levels were measured by the number of DNase1 hypersensitive sites overlapping each DNase1 cluster. **b**, Retrotransposon density of non-exonic regions and DNase1 clusters in cERDs and cLRDs. **c**, Observed minus expected retrotransposon density at the boundary of DNase1 clusters corrected for interval size bias (see methods). Expected retrotransposon density was calculated as each group's non-exonic total retrotransposon density across cERDs and cLRDs. A confidence interval of 3 standard deviations from expected retrotransposon density was also calculated, however the level of variation was negligible.

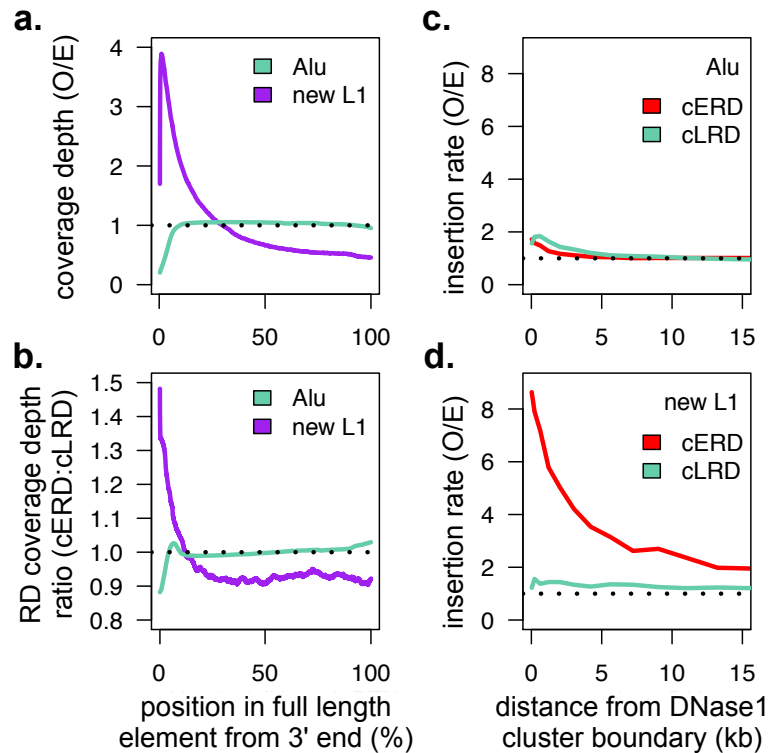


Figure 7. Retrotransposon insertion size is inversely proportional to local regulatory element density. **a**, Observed to expected ratio of retrotransposon position coverage depth measured from consensus 3' end. Expected retrotransposon position coverage depth was calculated as total retrotransposon coverage over consensus element length. We used 6 kb as the consensus new L1 length and 300 bp as the consensus *Alu* length. **b**, New L1 and *Alu* position density ratio (cERDs:cLRDs). **c**, *Alu* and **d**, new L1 observed over expected retrotransposon insertion rates at DNase1 cluster boundaries in cERDs and cLRDs. Insertion rates were measured by prevalence of 3' ends and expected levels were calculated as the per Mb insertion rate across cERDs and cLRDs.

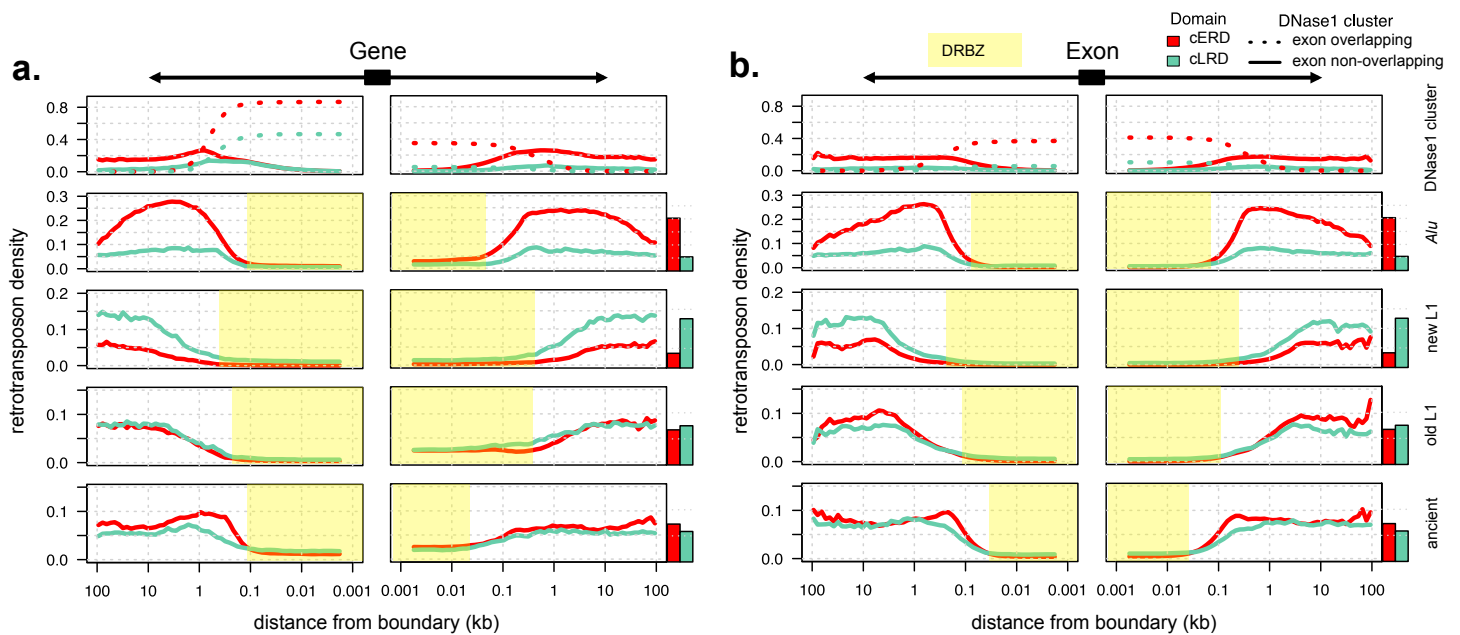


Figure 8. Retrotransposon accumulation within intergenic and intronic regions correlates with the distribution of DNase1 clusters. Density of DNase1 clusters and retrotransposons at each position upstream and downstream of genes and exons in **a**, intergenic and **b**, intronic regions. For DNase1 clusters, dotted lines represent exon overlapping clusters and solid lines represent clusters that do not overlap exons. For retrotransposons, solid lines represent the uncorrected retrotransposon density at exon and gene boundaries. Bar plots show expected retrotransposon density across cERDs and cLRDs. Highlighted regions outline DRBZs, regions extending from the gene or exon boundary to the point where retrotransposon levels begin to increase.