

# Similar evolutionary trajectories for retrotransposon accumulation in mammals

Reuben M Buckley<sup>1</sup>, R Daniel Kortschak<sup>1</sup>, Joy M Raison<sup>1</sup>, David L Adelson<sup>1,\*</sup>

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

\* [david.adelson@adelaide.edu.au](mailto: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 conserved interaction between retrotransposon activity and conserved genome architecture. 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 the accumulation of retrotransposons. Through conservation of synteny, gene regulation and nuclear organisation, we have found that mammalian genomes 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 (Consortium et al. 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.

Mammalian genomes are hierarchically organised into compositionally distinct hetero- or

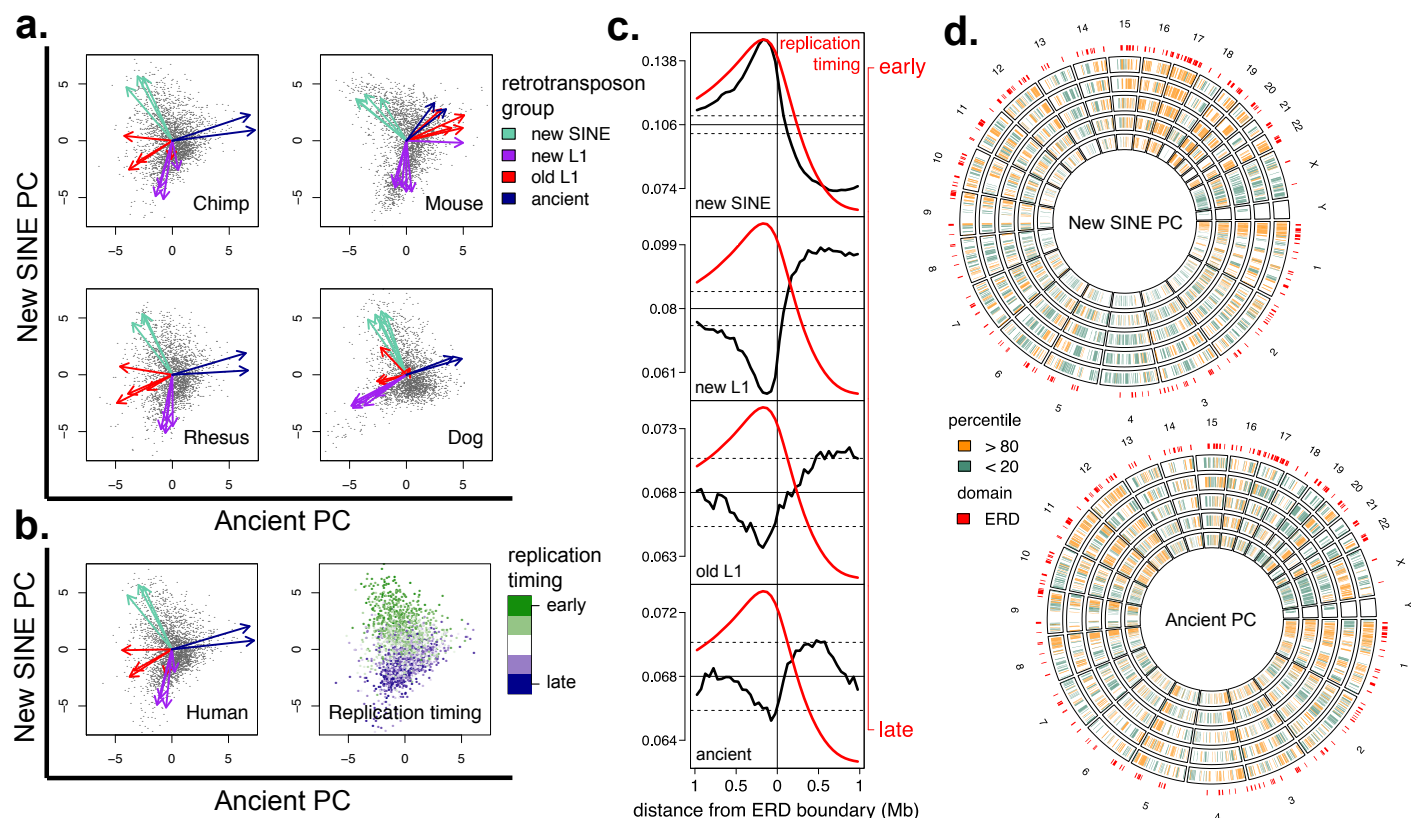
euchromatic large structural domains (Gibcus and Dekker 2013). These domains are largely composed of mobile self-replicating non-long terminal repeat (non-LTR) retrotransposons; with Long Interspersed Elements (LINEs) in heterochromatic regions and Short Interspersed Elements (SINEs) in euchromatic regions (Medstrand et al. 2002). The predominant LINE in most mammals is the ~6 kb long L1. This autonomously replicating element is responsible for the mobilisation of an associated non-autonomous SINE, usually ~300 bp long. Together, LINEs and SINEs occupy approximately 30% of the human genome (Lander et al. 2001), replicate via a well characterised RNA-mediated copy-and-paste mechanism (Cost et al. 2002) and co-evolve with host genomes (Kramerov and Vassetzky 2011; Chalopin et al. 2015; Furano et al. 2004).

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

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

To investigate the relationship between retrotransposon activity and genome evolution,

we began by characterising the distribution and accumulation of non-LTR retrotransposons 45  
within placental mammalian genomes. Next, we compared retrotransposon accumulation 46  
patterns in five separate evolutionary paths by humanising the repeat content (see methods) 47  
of the chimpanzee, rhesus macaque, mouse and dog genomes. Finally, we analysed human 48  
retrotransposon accumulation in large hetero- and euchromatic structural domains, focussing 49  
on regions surrounding genes, exons and regulatory elements. Our results suggested that 50  
accumulation of particular retrotransposon families follows from insertion into open chromatin 51  
found adjacent to regulatory elements and depends on local gene and regulatory element 52  
density. From this we propose a refined retrotransposon accumulation model in which 53  
random insertion of retrotransposons is primarily constrained by chromatin structure rather 54  
than local sequence composition. 55



**Figure 1. Large-scale genome distributions of retrotransposons are strongly associated with replication timing and conserved in distant mammalian species.** **a**, PCA of non-human genome retrotransposon content, each vector loading has been coloured according to the retrotransposon group it represents. PC1 and PC2 have been renamed according to the retrotransposon group whose variance they principally account for. **b**, PCA of human retrotransposon content and mean genome replication timing in HUVEC cells. **c**, Retrotransposon density per non-overlapping 50 kb intervals from a pooled set of ERD boundaries across all 16 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. **d**, 20% tails of New SINE and Ancient PC scores of humanised genomes plotted against human, large ERDs (> 2 Mb) from HUVEC cells are marked in red. Species from centre are human, chimpanzee, rhesus macaque, mouse and dog.

## Results

### Species selection and retrotransposon classification

We selected human, chimpanzee, rhesus macaque, mouse and dog as representative placental species because of their similar non-LTR retrotransposon composition (Fig. S1-S2) and phylogenetic relationships. Retrotransposon coordinates were obtained from the UCSC repeat

56

57

58

59

60

masker tables (Rosenbloom et al. 2015; Smit et al. 1996) and non-LTR retrotransposon families were grouped 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. S2). New L1s and new SINEs are retrotransposon families with high clade 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) and confirmed that our classification of retrotransposon groups agreed with ancestral and clade-specific periods of retrotransposon activity (Fig. S3).

## 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). From the PCA, we found that new SINEs and ancient elements strongly associated with the two major principal components (PC1 and PC2). Depending on this association we identified PC1 and PC2 as “New SINE PC” and “Ancient PC” respectively, or the converse (Fig. 1a). This showed that retrotransposon families from the same group accumulated in the same genomic regions. For all species examined, new SINEs were enriched in regions with few new L1s, and 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. 1a, S4). This mouse discordance has probably resulted from the increased genome turnover seen in the rodent lineage (Murphy et al. 2005) disrupting the distribution of ancestral retrotransposon families (Fig. S1-S2). As the relationship between mouse clade-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

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 cell line Repli-Seq data (Hansen et al. 2010a) from the UCSC genome browser. Repli-Seq measures the timing of genome replication during S-phase, where accessible euchromatic domains replicate early and inaccessible heterochromatic domains replicate late. Across our segmented human genome, we found a high degree of covariation between mean replication timing in HUVEC cells and New SINE PC scores (Fig. 1c), new SINEs associated with early replication and new L1s associated with late replication. This result is probably not specific to HUVEC cells alone, since early and late replicating regions from various independent cell lines exhibit a high degree of overlap (Fig. S5). In addition, by splitting L1s into old and new groups, we observed a strong association between replication timing and retrotransposon age that was not reported in previous analyses (Pope et al. 2014). To confirm these results, 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. 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. 1c). Because RD timing and genomic distributions of clade-specific retrotransposons are both largely conserved across human and mouse (Ryba 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 results showed that the genomic distribution of retrotransposons was similar across species (Fig 1a). To determine whether our observations resulted from retrotransposon insertion into orthologous regions, we used coordinate mappings between species to humanise retrotransposon family distributions and PC scores (see methods). From this, we found that retrotransposon families in different species that identified as the same group, accumulated

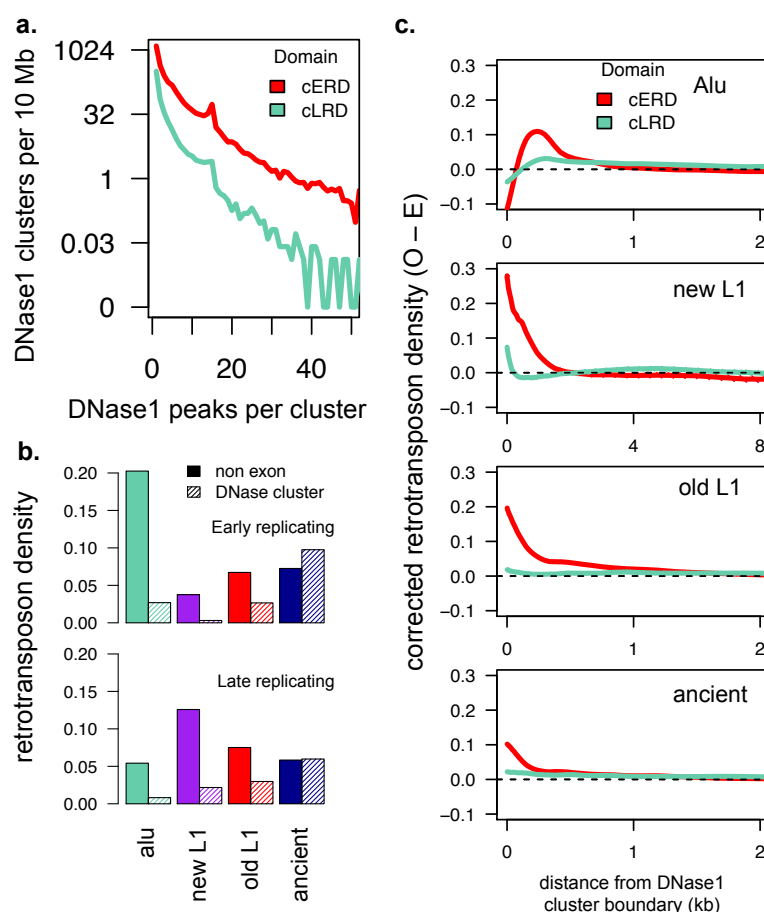
in regions with shared common ancestry (Fig. S6-S9). In addition, humanised genome segments from the 20% tails of the New SINE and Ancient PC score distributions showed high degrees of genomic overlap and associated with human RDs as described above (Fig. 1b). With regard to sequence conservation and retrotransposon accumulation, regions enriched for ancient elements shared the highest degree of pairwise similarity across our species (Fig. S10-S11). This demonstrates that regions enriched for ancient elements have likely been preserved throughout mammalian evolution (Adelson et al. 2009, 2010). Our results are consistent with retrotransposon accumulation overlying a conserved ancient 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 costly 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 (Consortium et al. 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. S12). 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. 2a). 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. 2b). Intriguingly, ancient element density in DNase1 clusters



remained relatively high, suggesting that some ancient elements may have been exapted. At 148  
 DNase1 cluster boundaries after removing interval size bias (Fig. S13-S14) (see methods), 149  
 retrotransposon density remained highly enriched in cERDs and close to expected levels in 150  
 cLRDs (Fig. 2c). This suggests that chromatin is likely to be open at highly active cluster 151  
 boundaries where insertion of retrotransposons is less likely to disrupt regulatory elements. 152  
 These results are consistent with an interaction between retrotransposon insertion, open 153  
 chromatin and regulatory activity, where insertions into open chromatin only persist if they 154  
 do not interrupt regulatory elements. 155

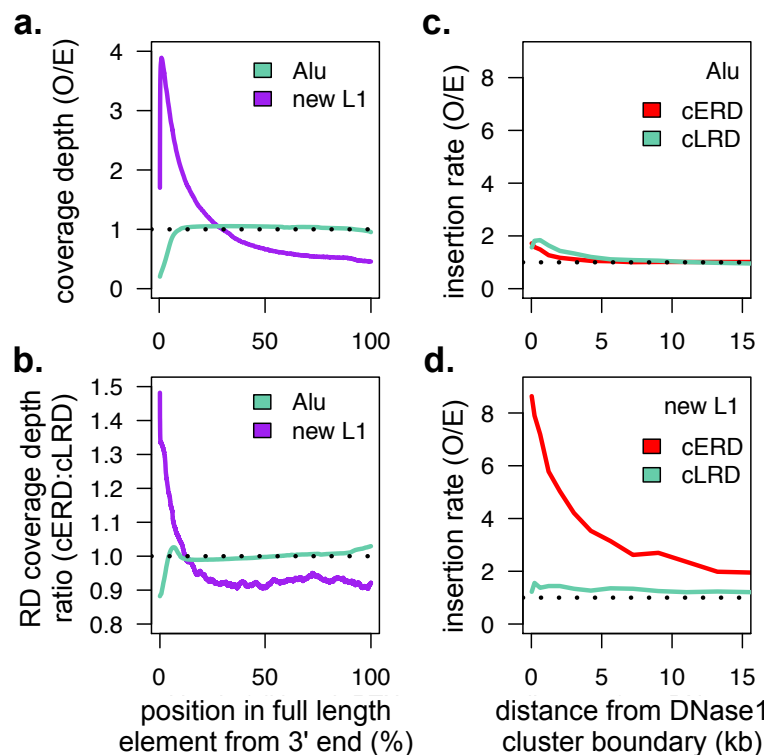


**Figure 2. Retrotransposon accumulation occurs in open chromatin near regulatory regions.** **a**, The activity of DNase1 clusters in cERDs and cLRDs. DNase1 clusters were identified by merging DNase1 hypersensitive sites across 15 tissues. Their activity levels were measured by the number of DNase1 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.

## 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 162  
and an over representation of new L1 3' ends, not seen with *Alu* elements (Fig. 3a). When we 163  
compared insertion size variation across cERDs and cLRDs we observed that smaller new L1s 164  
were enriched in cERDs and *Alu* elements showed no RD insertion size preference (Fig. 3b). 165  
The effect of insertion size on retrotransposon accumulation was estimated by comparing 166  
insertion rates of each retrotransposon group at DNase1 cluster boundaries in cERDs and 167  
cLRDs. We found that *Alu* insertion rates at DNase1 cluster boundaries were similarly 168  
above expected levels both in cERDs and cLRDs (Fig. 3c), whereas new L1 insertion rates 169  
at DNase1 cluster boundaries were further above expected levels in cERDs than cLRDs (Fig. 170  
3d). By comparing the insertion rate of new L1s — retrotransposons that exhibited RD 171  
specific insertion size variation — we found a negative correlation between element insertion 172  
size and gene/regulatory element density. Thus smaller elements, such as *Alu* elements, 173  
accumulate more in cERDs than do larger elements, such as new L1s, suggesting that smaller 174  
elements are more tolerated. 175

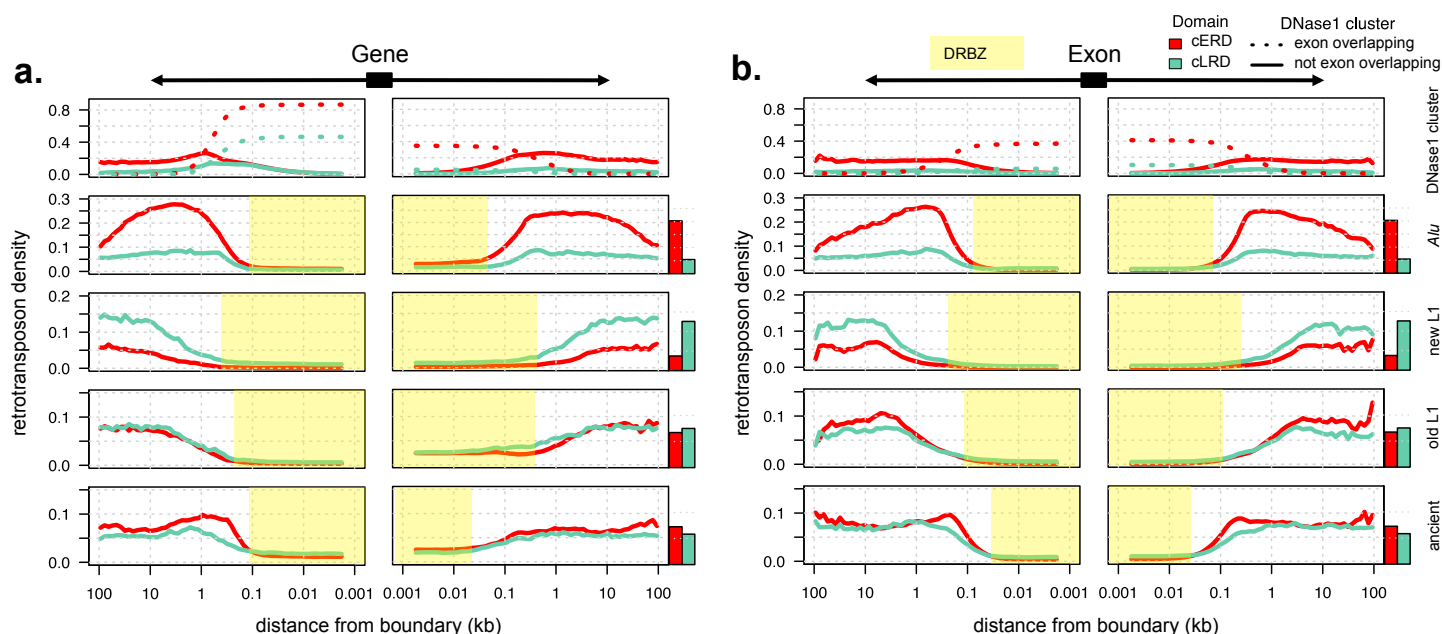


**Figure 3. Retrotransposon insertion size is inversely proportional to local regulatory element density.** **a**, Observed to expected ratio of retrotransposon position coverage depth measured from consensus 3' end. Expected retrotransposon position coverage depth was calculated as total retrotransposon 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.

## Retrotransposon insertion within gene and exon structures

Regulatory element organisation is largely shaped by gene and exon/intron structure which likely impacts the retrotransposon component of genome architecture. Therefore, we analysed retrotransposons and DNase1 clusters (exon overlapping and not exon overlapping) at the boundaries of genes and exons. Human RefSeq gene models were obtained from the UCSC genome browser and both intergenic and intronic regions were extracted (Table S4). At gene (Fig. 4a) and exon (Fig. 4b) boundaries, we found a high density of exon overlapping DNase1 clusters and depletion of retrotransposons. This created a depleted retrotransposon boundary zone (DRBZ) specific for each retrotransposon group, a region extending from the gene or exon boundary to the point where retrotransposon levels begin to increase. The size of each

DRBZ correlated with the average insertion size of each retrotransposon group, suggesting larger retrotransposons may have a greater capacity to disrupt important structural and regulatory genomic features. We also found that in cERDs the 5' gene boundary *Alu* DRBZ was larger than the 3' gene boundary *Alu* DRBZ. This difference was associated with increased exon overlapping DNase1 cluster density at 5' gene boundaries in cERDs (Fig. 4a), emphasising the importance of evolutionary constraints on promoter architecture. For ancient elements, their interval size corrected density approximately 1 kb from the 5' gene boundary was significantly higher than expected. This increase is consistent with exaptation of ancient elements into regulatory roles (Lowe et al. 2007) (Fig. S15-S18). Moreover, the density peak corresponding to uncorrected ancient elements also overlapped with that of not exon overlapping DNase1 clusters (Fig. 4a). Collectively, these results demonstrate the evolutionary importance of maintaining gene structure and regulation and how this in turn has canalised similar patterns of accumulation and distribution of retrotransposon families in different species over time.



**Figure 4. 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 not exon overlapping clusters. 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.

# Discussion

200

In our study, we compared several mammalian genomes and analysed chromatin structure at both small and large scales to better characterise retrotransposon accumulation. Our genome-wide comparisons across species were consistent with previous analyses that reported high levels of positional conservation for L1s and their associated SINEs (Chinwalla et al. 2002; Gibbs et al. 2004). Because new L1s and new SINEs underwent periods of activity after each of our sample species diverged from a common ancestor (Giordano et al. 2007), our observations are likely the result of a conserved interaction between retrotransposon activity and genome architecture. Previous analyses have attempted to capture this interaction through various retrotransposon accumulation models (Lander et al. 2001). Based on large-scale conservation of genome architecture and GC content (Chinwalla et al. 2002; Gibbs et al. 2004), the current model of retrotransposon accumulation suggests that random insertion of L1s and SINEs are similarly constrained by local sequence composition, where L1s are quickly purged from gene-rich regions via purifying selection at a higher rate than SINEs (Graham and Boissinot 2006; Gasior et al. 2007; Kvikstad and Makova 2010). However, this model fails to account for the demonstrated impact of chromatin structure on insertion site preference (Cost et al. 2001; Baillie et al. 2011).

201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215  
216

We used publicly available datasets to analyze the impact of chromatin architecture on retrotransposon accumulation. However, this approach is not without its limitations. For example, heritable retrotransposon insertions typically occur during embryogenesis or within the germline; developmental stages and tissue samples that were unavailable. To overcome such limitations we aggregated data from a range of biological contexts. Using this strategy, we increased the probability of capturing chromosomal domain structures and regulatory element sites present in embryonic and germline cell states.

217  
218  
219  
220  
221  
222  
223

From our analysis we found that 1) following preferential insertion into open chromatin domains, retrotransposons were tolerated adjacent to regulatory elements where they were less likely to cause harm; 2) element insertion size was a key factor affecting retrotransposon accumulation, where large elements accumulated in gene poor regions where they were less likely to perturb gene regulation; and 3) insertion patterns surrounding regulatory elements were persistent at the gene level. Based on these results, we propose a significant change to the current retrotransposon accumulation model; rather than random insertion constrained

224  
225  
226  
227  
228  
229  
230

by local sequence composition, we propose that insertion is instead primarily constrained  
by local chromatin structure. Following this, L1s and SINEs both preferentially insert into  
gene/regulatory element rich euchromatic domains, where L1s with their relatively high  
mutational burden are quickly eliminated via purifying selection at a much higher rate than  
SINEs. Over time this results in an enrichment of SINEs in euchromatic domains and an  
enrichment of L1s in heterochromatic domains.

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

## Methods

### Within species comparisons of retrotransposon genome distributions

Retrotransposon coordinates for each species were initially identified using RepeatMasker  
and obtained from 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, and canFam3 assemblies. Each species genome was segmented into 1  
Mb regions and the density of each retrotransposon family for each segment was calculated.  
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,  
for each retrotransposon family we calculated Pearson's correlation coefficient for each  
retrotransposon family across our genome segments.

### Across species comparisons of retrotransposon genome distributions

To compare genome distributions across species, we humanised a query species genome using  
mapping coordinates extracted from net AXT alignment files located on the UCSC genome  
browser (Table S1). First, genomes were filtered by discarding segments below a minimum

mapping fraction threshold, removing poorly represented regions (Fig S19a). Next, we used mapping coordinates to match fragments of query species segments to their corresponding human segments (Fig S19b). From this, the retrotransposon content and PC scores of the matched query segments were humanised following equation 1 (Fig S19c).

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

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 S19d). 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 S19e). 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 S19f). Spatial correlations and the P-values from measuring the effects of humanising and filtering were integrated into a heatmap (Fig S19g). The entire process was repeated several times at different minimum mapping fraction thresholds to optimally represent each retrotransposon families genomic distribution in a humanised genome (fig S20).

## Replication timing boundaries and constitutive replication timing domains

ERDs, LRDs, and timing transition regions (TTRs) for each dataset were previously identified using a deep neural network hidden Markov model (Table S2) (Liu et al. 2015). To determine RD boundary fluctuations of retrotransposon density, we defined ERD boundaries as the boundary of a TTR adjacent to an ERD. ERD boundaries from across each sample were pooled and retrotransposon density was calculated for 50 kb intervals from regions flanking each boundary 1 Mb upstream and downstream. Expected density and standard deviation for each retrotransposon group was derived from a background distribution generated by



calculating the mean of 500 randomly sampled 50 kb genomic bins within 2000 kb of each ERD boundary, replicated 10000 times. We also obtained Repli-Seq replication timing profiles from the UCSC genome browser as a wavelet signal (Table S2) (Hansen et al. 2010b). For each of our 50 kb intervals we calculated the mean replication timing from across each Repli-Seq sample. To identify cERDs and cLRDs, ERDs and LRDs classified by Liu *et al* (Liu et al. 2015) across each cell type were split into 1 kb intervals to find the intersection. If the classification of 12 out of 16 samples agreed at a certain region, we classified that region as belonging to a cERDs or a cLRDs, depending on that region's majority classification.

### **DNase1 cluster identification and activity**

DNase1 sites across 15 cell lines were found using DNase-seq and DNase-chip as part of the open chromatin synthesis dataset for ENCODE (Table S3) (Consortium et al. 2012). Regions where P-values of contiguous base pairs were below 0.05 were identified as significant DNase1 hypersensitive sites (Consortium et al. 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 S4). 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 S21). 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 grew as a function of distance from the boundary as position depth decreased. 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  $\sqrt{npq}$ , where  $n$  is the total number of bases at a given position and  $q$  is equal to  $1 - p$ .

## Interval size bias correction of retrotransposon densities

Interval boundary density is sensitive to retrotransposon insertion preferences into intervals of a certain size (Fig S22). 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 retrotransposon 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> . 345

## Additional Files 346

### Additional file 1 — Supplementary information 347

Figures S1–S22, Tables S1–S4. 348

### Competing interests 349

The authors declare that they have no competing interests. 350

### Author’s contributions 351

R.M.B., R.D.K., J.M.R., and D.L.A. designed research; R.M.B. performed research; and 352

R.M.B., R.D.K., and D.L.A. wrote the paper. 353

### Acknowledgements 354

For reviewing our manuscript and providing helpful advice we would like to thank the 355

following: Simon Baxter, Atma Ivancevic and Lu Zeng from the University of Adelaide; 356

Kirsty Kitto from Queensland University of Technology; and Udaya DeSilva from Oklahoma 357

State University. 358

### Availability of data and materials 359

All data was obtained from publicly available repositories, urls can be found in sup- 360

porting material (Table S1–S3). R scripts used to perform analyses can be found at 361

<https://github.com/AdelaideBioinfo/retrotransposonAccumulation>. 362

## References

Adelson, D., Raison, J., Garber, M., and Edgar, R. (2010). Interspersed repeats in the horse (equus caballus); spatial correlations highlight conserved chromosomal domains. *Animal genetics*, 41(s2):91–99.

- Adelson, D. L., Raison, J. M., and Edgar, R. C. (2009). Characterization and distribution of retrotransposons and simple sequence repeats in the bovine genome. *Proceedings of the National Academy of Sciences*, 106(31):12855–12860.
- 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.
- 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., Trocheset, 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., Fröncke, 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.
- Consortium, E. P. et al. (2012). An integrated encyclopedia of dna elements in the human genome. *Nature*, 489(7414):57–74.
- 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.
- 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. (2010a). Sequencing newly replicated dna reveals widespread plasticity in human replication timing. *Proceedings of the National Academy of Sciences*, 107(1):139–144.
- 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. (2010b). Sequencing newly replicated dna reveals widespread plasticity in human replication timing. *Proceedings of the National Academy of Sciences*, 107(1):139–144.
- Kramerov, D. and Vassetzky, N. (2011). Origin and evolution of sines in eukaryotic genomes. *Heredity*, 107(6):487–495.

- 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., Auvil, 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.
- 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.
- 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.
- Smit, A. F., Hubley, R., and Green, P. (1996). Repeatmasker open-3.0.