

RESEARCH

Open Access



Co-opted transposons help perpetuate conserved higher-order chromosomal structures

Mayank NK Choudhary, Ryan Z. Friedman, Julia T. Wang, Hyo Sik Jang, Xiaoyu Zhuo and Ting Wang*

Abstract

Background: Transposable elements (TEs) make up half of mammalian genomes and shape genome regulation by harboring binding sites for regulatory factors. These include binding sites for architectural proteins, such as CTCF, RAD21, and SMC3, that are involved in tethering chromatin loops and marking domain boundaries. The 3D organization of the mammalian genome is intimately linked to its function and is remarkably conserved. However, the mechanisms by which these structural intricacies emerge and evolve have not been thoroughly probed.

Results: Here, we show that TEs contribute extensively to both the formation of species-specific loops in humans and mice through deposition of novel anchoring motifs, as well as to the maintenance of conserved loops across both species through CTCF binding site turnover. The latter function demonstrates the ability of TEs to contribute to genome plasticity and reinforce conserved genome architecture as redundant loop anchors. Deleting such candidate TEs in human cells leads to the collapse of conserved loop and domain structures. These TEs are also marked by reduced DNA methylation and bear mutational signatures of hypomethylation through evolutionary time.

Conclusions: TEs have long been considered a source of genetic innovation. By examining their contribution to genome topology, we show that TEs can contribute to regulatory plasticity by inducing redundancy and potentiating genetic drift locally while conserving genome architecture globally, revealing a paradigm for defining regulatory conservation in the noncoding genome beyond classic sequence-level conservation.

Keywords: 3D genome, Loops, Evolution, Conservation, Transposable elements, Binding site turnover

Background

The 3D organization of various genomes has been mapped at high resolution using a variety of methods [1–5]. While genome folding is largely conserved in mammals [1, 4], the genetic forces shaping its emergence and evolution remain poorly understood. Two distinct yet mutually non-exclusive models [6] have recently gained much traction: that of phase separation [7] and of loop extrusion [8, 9] by factors such as Cohesin that colocalizes extensively with CTCF throughout the genome. In relation to the latter, TEs are known to contain and disseminate functional regulatory sequences [10–13] including that of CTCF. In contrast to relying on point mutations to evolve a functional CTCF binding

site, TE transposition presents an attractive model for rapid regulatory sequence dissemination and regime building [14–17]. Hence, we hypothesized that TEs have been a rich source of sequence for the assembly and tinkering of higher-order chromosomal structures. We studied the influence of all repetitive elements (REs) in establishing higher-order chromosomal structures and, more specifically, the role of TEs in the evolution of these higher-order chromosomal structures in humans and mice.

Results

We examined REs' contribution to loop anchor CTCF sites using published genome-wide chromosome conformation capture data from assays including ChIA-PET [2] and Hi-C in human (GM12878, HeLa, HMEC, IMR90, K562, NHEK) and mouse (ESCs, NSCs, CH12-LX) cell lines [1, 18]. We determined that 398 out of

* Correspondence: twang@genetics.wustl.edu

The Edison Family Center for Genome Sciences & Systems Biology, Department of Genetics, Washington University, 4515 McKinley Avenue, Campus Box 8510, St. Louis, MO 63110, USA



3159 (12.6%) unique loop anchor CTCF sites were derived from REs in the mouse lymphoblastoid cell line. These RE-derived CTCF sites help establish 451 out of 2718 (16.6%) loops with discernible, unique CTCF loop anchors (Fig. 1a, b). In the corresponding human lymphoblastoid cell line, REs contributed 935 out of 8324 (11.2%) unique loop anchor CTCF sites that help

establish 1244 out of 8007 (15.6%) loops. Overall, REs contributed 9–15% of the anchor CTCF sites that result in 12–18% loops in humans and 12–23% of the anchor CTCF sites that result in 15–27% loops in mouse, across a variety of cell lines (Fig. 1a, b). Interestingly, the proportion of RE-derived loops and anchor CTCF sites in mouse ESCs is significantly higher than NSCs and

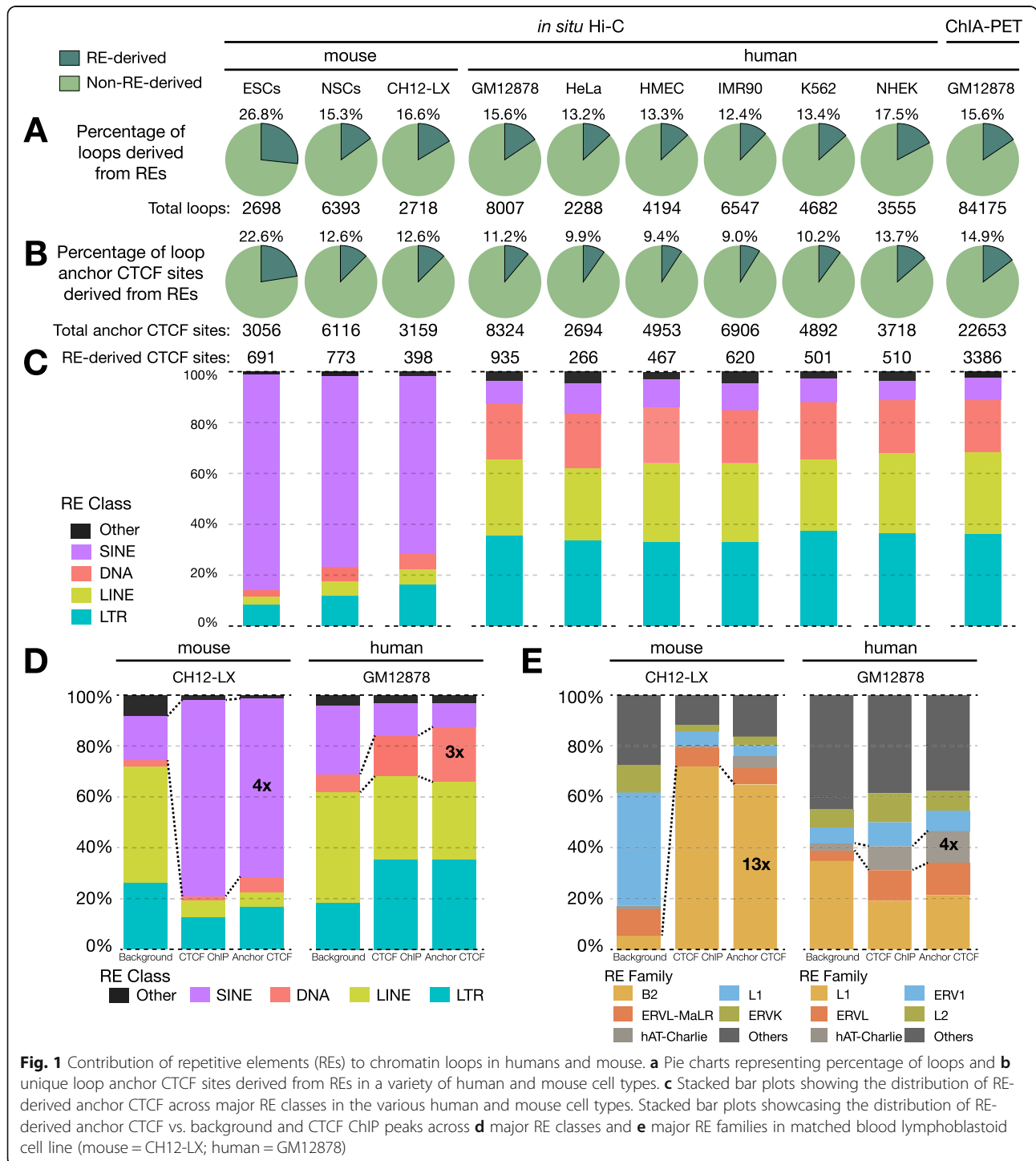


Fig. 1 Contribution of repetitive elements (REs) to chromatin loops in humans and mouse. **a** Pie charts representing percentage of loops and **b** unique loop anchor CTCF sites derived from REs in a variety of human and mouse cell types. **c** Stacked bar plots showing the distribution of RE-derived anchor CTCF across major RE classes in the various human and mouse cell types. Stacked bar plots showcasing the distribution of RE-derived anchor CTCF vs. background and CTCF ChIP peaks across **d** major RE classes and **e** major RE families in matched blood lymphoblastoid cell line (mouse = CH12-LX; human = GM12878)

CH12-LX cells. This observation could potentially be driven by both: genome-wide demethylation of transposable elements in ESCs leading to a higher proportion of TE-derived CTCF motifs accessible for CTCF binding as well as fewer chromatin loops observed in mouse ESCs [18].

In both species, RE-derived loop anchor CTCF sites were largely derived from TEs (> 95%) and their class of origin (SINE, LINE, LTR, DNA) showed a species-biased distribution (Fig. 1c). Using the highest resolution in-situ Hi-C maps in matched lymphoblastoid cell types in mice (CH12-LX) and humans (GM12878), we compared the composition of the RE-derived loop anchor CTCF sites. While the mouse lineage was profoundly shaped by the SINEs (70%, 4× enrichment over background), the human lineage was overrepresented by retroviral LTR elements and DNA transposons (36% and 22%, 2× and 3× enriched over the background respectively) (Fig. 1d). At the family level, the B2 SINEs in mice were 13-fold enriched over background and contributed 65% of TE-

derived loop anchor CTCF sites. In humans, the hAT-Charlie family of DNA transposons contributed 13% of TE-derived loop anchor CTCF sites, a 4-fold enrichment over background (Fig. 1e). These contributions are underestimates as we have yet to (i) uniquely identify all loop anchor CTCF sites (especially in repetitive regions) and (ii) annotate all repetitive elements, especially ancient TEs that have diverged far from their identity [19]. Further, we looked at the cell-type specificity of these loop anchor CTCF sites in humans and see that 1334 out of 2017 (66%) RE-derived loop anchor CTCF sites were found in only one cell type (Additional file 1: Figure S1A). However, we did not find any specific TE family that enriches for cell-type specific loop anchor CTCF sites in the cell lines profiled (Additional file 1: Figure S1B).

To study the evolution of chromatin loops, we compared their conservation (Fig. 2a, Methods) in matched human and mouse cell types. Briefly, we used the lift-Over tool [20] to compare loops across species and required exactly one reciprocal match (reciprocal best hit)

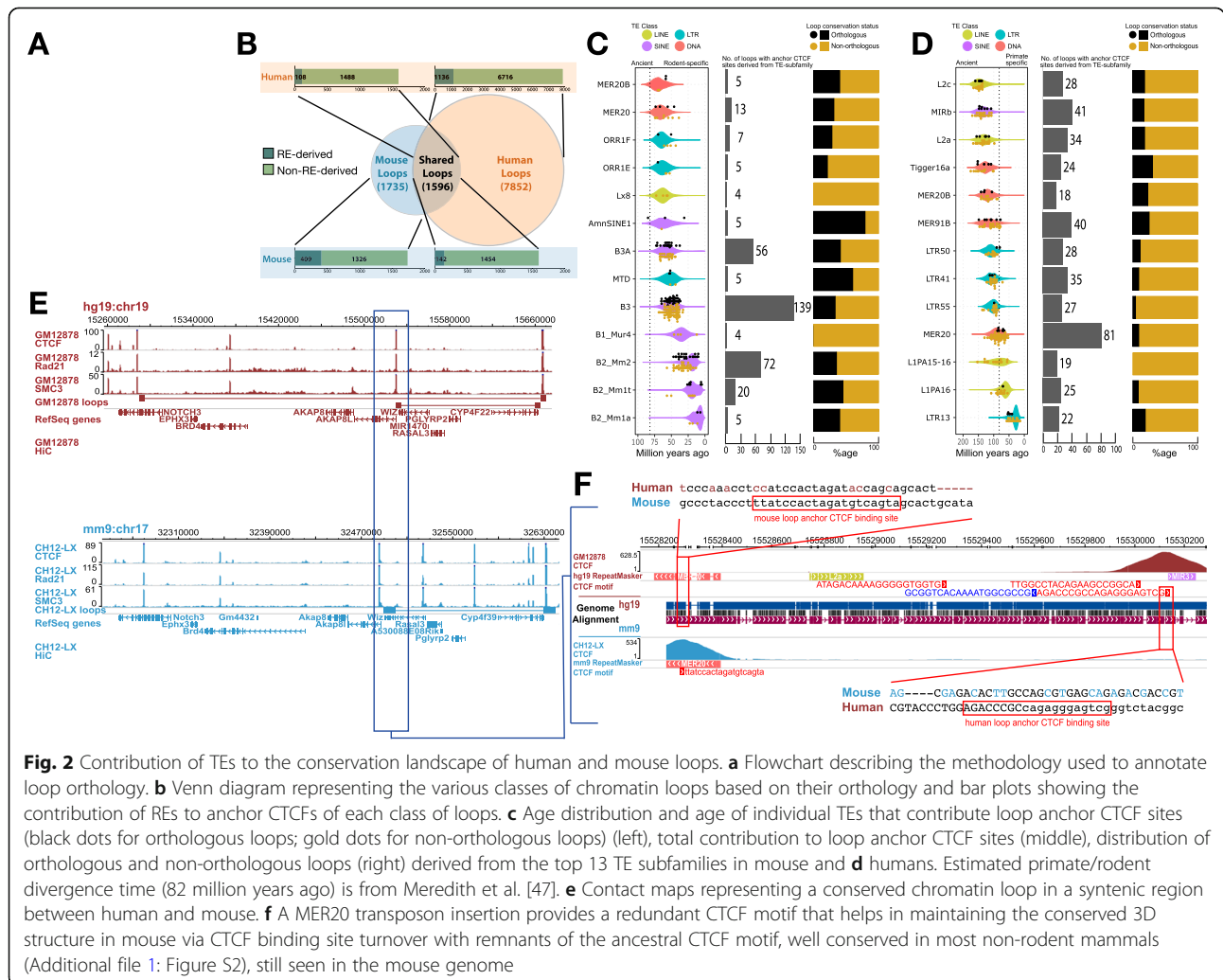


Fig. 2 Contribution of TEs to the conservation landscape of human and mouse loops. **a** Flowchart describing the methodology used to annotate loop orthology. **b** Venn diagram representing the various classes of chromatin loops based on their orthology and bar plots showing the contribution of REs to anchor CTCFs of each class of loops. **c** Age distribution and age of individual TEs that contribute loop anchor CTCF sites (black dots for orthologous loops; gold dots for non-orthologous loops) (left), total contribution to loop anchor CTCF sites (middle), distribution of orthologous and non-orthologous loops (right) derived from the top 13 TE subfamilies in mouse and **d** humans. Estimated primate/rodent divergence time (82 million years ago) is from Meredith et al. [47]. **e** Contact maps representing a conserved chromatin loop in a syntenic region between human and mouse. **f** A MER20 transposon insertion provides a redundant CTCF motif that helps in maintaining the conserved 3D structure in mouse via CTCF binding site turnover with remnants of the ancestral CTCF motif, well conserved in most non-rodent mammals (Additional file 1: Figure S2), still seen in the mouse genome

to designate conserved loops. We found that 48% of all mouse loops (1596 out of 3331) had a loop call in the corresponding syntenic region in humans (Additional file 2: Table S1.1). Our observation is in close agreement with prior studies [1, 4] that show about half of all higher-order chromosomal structures to be conserved. We then sought to characterize the contribution of TEs to various classes of loops based on their orthology.

We compared the origin of loop anchor CTCF sites of orthologous loops in mouse and human. We found that out of 1596 orthologous loops, 142 (8.9%) in mouse and 108 (6.7%) in human had at least one TE-derived loop anchor CTCF site (Fig. 2b). In addition to orthologous loops, TE-derived loop anchor CTCF sites also gave rise to 24% (409 out of 1735) and 15% (1136 out of 7852) non-orthologous (species-specific) loops in mouse and humans, respectively (Fig. 2b), consistent with the appreciable role of TEs in genome innovation [14–16, 21, 22]. Overall, the majority of TE-derived loop anchors in mouse were established by a handful of young TE subfamilies (B3, B2_Mm2, B3A, B2_Mm1t) that expanded in the rodent lineage [23] (Fig. 2c). In contrast, multiple TE subfamilies of varying evolutionary ages contributed diffusely to CTCF loop anchors in humans (Fig. 2d). Altogether, TEs in humans contributed to fewer orthologous loops and distributed over more TE subfamilies than in mouse.

Intriguingly, 123/142 (87%) TE-derived orthologous loops in mouse were discordant for TEs in humans (Additional file 2: Table S1.2). In the sense, while the loops in humans were anchored at the putative ancestral CTCF binding sites, the syntenic ancestral CTCF motifs were largely degraded or deleted in mouse and the loops were now anchored at CTCF sites derived from nearby, co-opted TEs instead. One such example is an orthologous loop at the 5' end of the *Akap8l* gene (Fig. 2e) maintained in mouse by a MER20 element transposed ~1.5 kb upstream of the degraded ancestral motif which was well conserved in most non-rodent mammals (Additional file 1: Figure S2). The degradation of the ancestral CTCF motif derived from an ancient MIR3 element that is over 147 million years old (see “Methods”) incapacitates CTCF binding as evidenced by the CTCF ChIP track (Fig. 2f). In contrast, the younger MER20 element that inserted ~90 million years ago harbored strong CTCF binding, providing an anchor site to maintain the conserved loop in mouse. Similarly, we find that 89/108 (82%) TE-derived orthologous loops in human GM12878 cells were discordant for TEs in mouse (Additional file 2: Table S1.3). We hypothesized that TEs provide redundant CTCF sites and mediated binding site turnover for CTCF contributing to conserved genome folding events between human and mouse.

Moreover, the 123 turned-over loops in mouse represent 127 turnover events (4 loops had both loop anchors turned-over) mediated by 124 unique loop anchors (3 turned-over loop anchors tethered 2 loops each). Out of the 124 unique loop anchors, 61 events represent turnover of the left loop anchor and 63 events represent turnover of the right loop. In terms of CTCF motif orientation—for the 61 left loop anchor turnover events, 53 were positive and 8 were negative, and for the 63 right loop anchor turnover events, 45 were negative and 18 were positive (chi-square test, p value = 5.3×10^{-11}). Similarly, in humans the 89 turned-over loops represent 93 turnover events (4 loops had both loop anchors turned-over) were mediated by 84 unique loop anchors (1 turned-over loop anchor tethered 3 loops, and 7 loop anchors tethered 2 loops each). Out of the 84 unique loop anchors, 43 events represent turnover of the left loop anchor (43 positive orientation CTCF motif and 0 negative orientation CTCF motif), and 41 events represent turnover of the right loop (40 positive orientation CTCF motif and 1 negative orientation CTCF motif) (chi-square test, p value = 3.6×10^{-19}). These results further lend credence to the loop extrusion model [8] and suggest that TE exaptation is more likely when the orientation of the inserted TE (and the underlying CTCF motif provided) is compatible with the local loop structure.

mm9 CH12-LX ($n = 124$)	Left loop anchor	Right loop anchor
Positive CTCF motif	53	18
Negative CTCF motif	8	45
hg19 GM12878 ($n = 84$)	Left loop anchor	Right loop anchor
Positive CTCF motif	43	1
Negative CTCF motif	0	40

Since the mouse genome is replete with repeat-derived CTCF sites [23] that could interfere with the targeted study of specific TE candidates, we decided to validate these hypotheses in human cell lines.

Here we examine two candidate TEs that maintain conserved higher-order chromosomal structures in humans: one belonging to the L1M3f subfamily of LINES, and the other belonging to the LTR41 subfamily of endogenous-retrovirus-derived long terminal repeat (LTR). The former TE replaces the function of a lost ancestral CTCF site (Additional file 1: Figure S3), while the latter is functionally redundant for an ancestral CTCF site still present in humans (Additional file 1: Figure S4). These two TEs were specifically chosen as they could be unambiguously attributed to the genome folding function (no other CTCF/Cohesin binding site in the

vicinity). Using CRISPR-Cas9, we obtained clones of GM12878 cells bearing homozygous deletions of the L1M3f and LTR41 elements, respectively (Additional file 1: Figure S5, Additional file 3: Table S2.4). We then performed HYbrid-Capture on the in situ Hi-C library (Hi-C²) to examine the effect of the TE deletion on the local 3D structure [8] (Additional file 3: Table S2.1, S2.2, S2.3).

The L1M3f-derived CTCF site was positioned at a conserved domain border and anchored three chromatin loops (Additional file 1: Figure S3). Upon deletion of this L1M3f, the conserved local chromosomal structure collapsed as evidenced by (i) the loss of focal enrichment in the homozygous TE knockout (KO) contact map in comparison to the wild-type (WT) contact map and (ii) the fusion of two neighboring domains (Hi-C² results: Fig. 3a, Hi-C results: Additional file 1: Figure S6). The Virtual 4C plot anchored at the region surrounding the L1M3f element showed three distinct peaks (corresponding to the three loops in the WT cell line), which were lost in the KO (Δ L1M3f) cell line. We also found that cross-domain interactions significantly increased from 8% in WT to 19% in KO cell lines ($\sim 2.4\times$, Welch's *t*-test *p* value $< 1.5 \times 10^{-16}$, Additional file 3: Table S2.5) across the L1M3f-established domain boundary, a change specific to the targeted domain and not seen in a control domain from a nearby region (Fig. 3c). Thus, the L1M3f element is necessary for maintaining the conserved loops and domain boundary in humans. It represents a novel class of binding site turnover [24–27] for CTCF leading to conservation in terms of function via establishment of long-range interactions and potentially the underlying gene regulation, but not in local primary sequence.

Our second candidate was a species-specific LTR41-derived CTCF site (“c” in Fig. 3d, e) that replaced an ancestral CTCF site derived from a much older TE (“d” in Fig. 3d, e) of the MER82 subfamily that is conserved in humans and mouse. The ancestral MER82-derived CTCF site was “decommissioned” as the LTR41 insertion (after the primate-rodent split) provided a negative orientation CTCF motif upstream of the MER82 element. Based on the loop extrusion model, the LTR41-derived CTCF motif would be encountered before the MER82-derived CTCF site and hence the ancestral site is mostly decommissioned in present-day human genome as evidenced by the drastically reduced CTCF binding (Additional file 1: Figure S4B). In the WT contact map, we observed a bright focal enrichment corresponding to CTCF binding sites a–c suggesting a looping interaction. In contrast, there was little focal enrichment corresponding to a–d (Fig. 3d, top row). Additionally, in the WT Virtual 4C track anchored on “a,” we observed a clear peak corresponding to LTR41 (“c”) suggesting an

a–c loop (Fig. 3e). Upon deletion of LTR41, the conserved loop's anchor is offset to the MER82-derived CTCF site (“d”) downstream of the LTR41 as evidenced by the shift in the focal enrichment in the KO contact map (Fig. 3d, bottom row) and an increase in the KO Virtual 4C peak corresponding to the MER82-derived CTCF site (i.e., a–d loop) (Fig. 3e, Additional file 1: Figure S7). Upon anchoring the Virtual 4C on a 5-kb window containing LTR41 (c), we observed a peak loss at “a” corresponding to the loss of the a–c loop in the KO, an interaction that existed in the WT cells (Fig. 3f). With the ~ 39 kb shift of the anchor site, the half-megabase scale chromosomal structure around the anchor region remained largely preserved (Additional file 1: Figure S4C). Upon deletion of this TE candidate, the local sequence configuration probably resembled that of the pre TE-insertion, ancestral genome. This example therefore illustrates a potential path by which the local 3D genome evolved upon insertion of the LTR41 element as well as the plasticity TEs, like LTR41 and MER82 in this case, can encode in their host genomes by providing redundant CTCF binding sites.

These results support the hypothesis that TEs are able to contribute regulatory robustness and strengthen conserved regulatory architecture as redundant or “shadow” loop anchors. The mouse genome that underwent a lineage-specific expansion of SINE B2s [23], which carry a CTCF binding motif, is saturated with such events.

TEs are typically silenced by host repressive machineries including DNA and histone methylation [28–30]. However, a small fraction of TEs escape epigenetic silencing and provide functional regulatory elements for the host in a process termed exaptation [31–34]. Since CTCF is a methylation-sensitive chromatin factor and only binds to unmethylated DNA [35, 36], we examined the DNA methylation levels of loop anchor CTCF sites of orthologous loops (“Methods”). We found that TE-derived CTCF sites were marked by reduced DNA methylation, similar to their non-TE derived genomic counterparts (Fig. 4a). To understand the DNA methylation dynamics through evolution, we took advantage of the differential mutation rate of 5-methylcytosine (5mC) to thymine (T) [37]. Unmethylated cytosines (C) mutate to T at a lower rate than 5mC; thus, methylated DNA exhibits higher frequency of C to T mutations [38]. We found that TEs involved in turnover events had a significantly lower frequency of methylation-associated C-to-T and G-to-A mutations compared to an identically sampled background of TEs not involved in looping (1000 simulations), but no difference in all other combined substitutions (summarized human results: Fig. 4b; full human and mouse results: Additional file 1: Figure S8, Figure S9, Additional file 4:

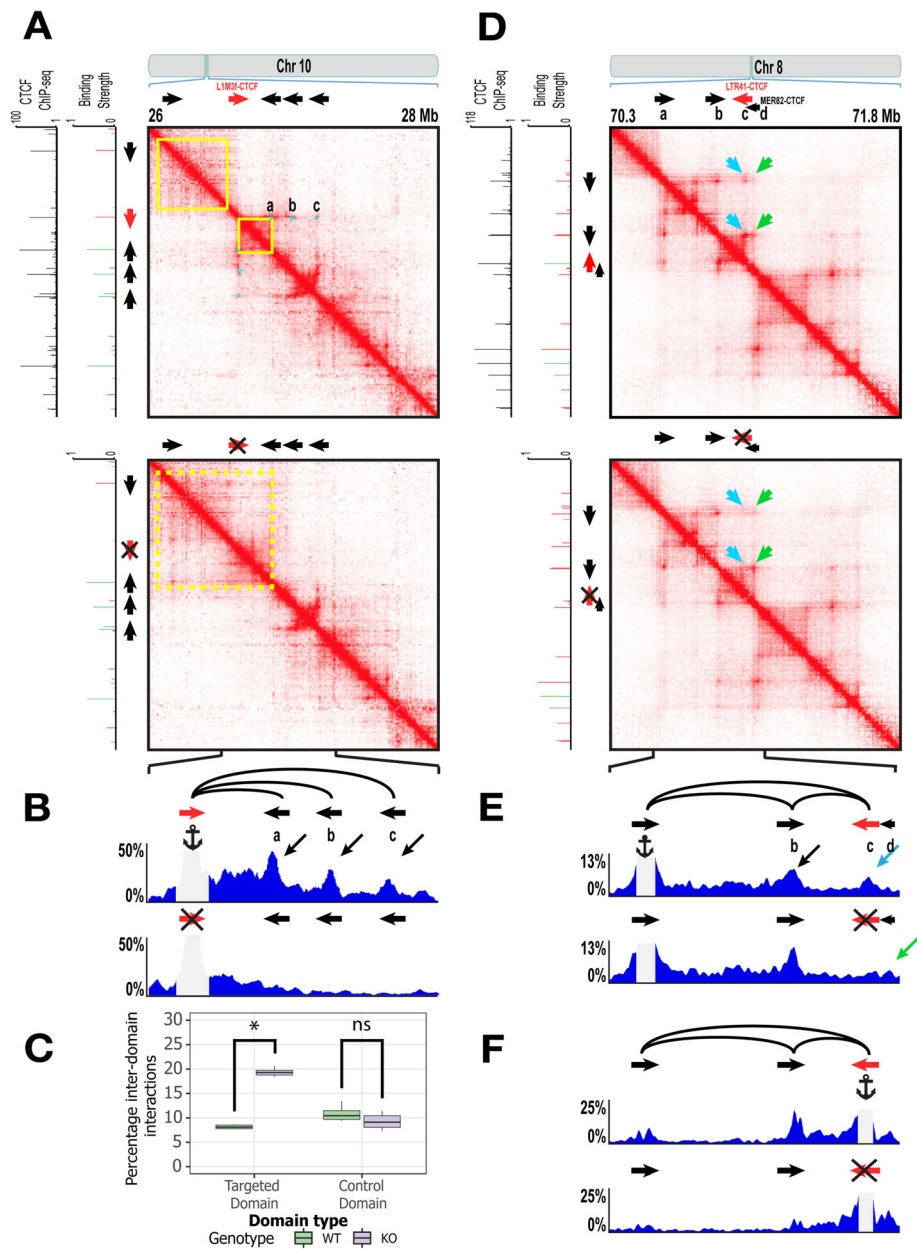


Fig. 3 TEs are necessary for maintaining conserved higher-order chromosomal structures in humans. **a** Results of a CRISPR/Cas9-based deletion of an L1M3f element at chr10:26–28 Mb in GM187278 cells. Mega-contact maps (details in “Methods”) generated using Hi-C² technology for the (top) WT locus and (bottom) KO (Δ L1M3f) locus. **b** Virtual 4C plot displaying total percent interactions emanating from an anchor on a 5-kb window containing the L1M3f element. **c** Boxplot measuring the percent inter-domain interactions (Additional file 3: Table S2.5) across the targeted domain and a control domain (boundaries unaffected by CRISPR edits) using subsampled contact maps (details in “Methods”). **d** Results of CRISPR/Cas9-based deletion of an LTR41 element at chr8:70.3–71.8 Mb in GM12878 cells. Mega-contact maps generated in Hi-C² experiments for the (top) WT locus and (bottom) KO (Δ LTR41) locus. **e** Virtual 4C plot displaying total percent interactions emanating from an anchor on a 5-kb window containing the left anchor CTCF of the conserved loop, and **f** the LTR41 element

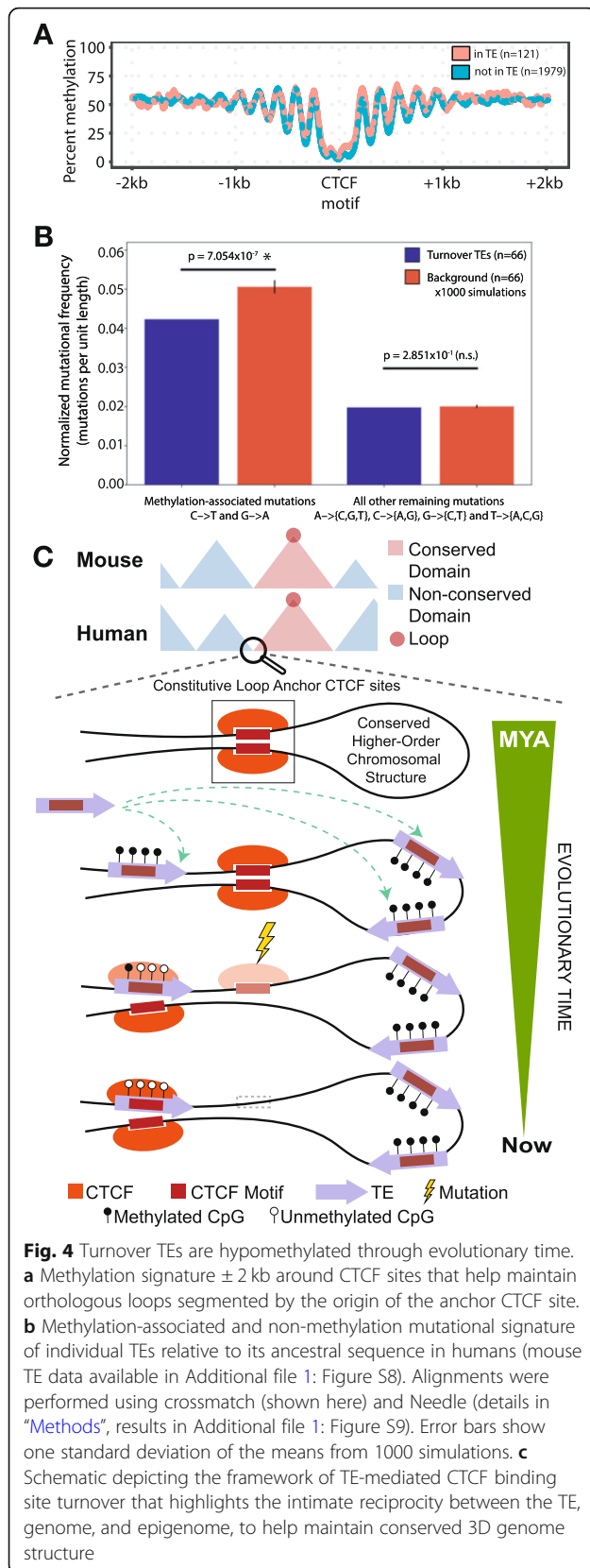


Table S3). These results suggest that TEs providing CTCF binding site turnover were hypomethylated over evolutionary time to maintain their functional role, compared to other TE copies (Fig. 4c).

Discussion

TEs have substantially contributed to higher order chromatin structures by serving as chromatin loop anchors—a large fraction of which were found to be species-specific, confirming TEs’ role in genome innovation. Pioneering work in the last decade has extensively outlined this contribution of TEs in shaping gene regulatory networks by depositing new TF binding sites in host genomes, leading to the origins of novel phenotypes like innate immunity and pregnancy in mammals. Herein lies the catch: research to date showcases the role of TEs in bringing novelty and new regulatory functions to the host genome. Hence, TEs have long been considered a source of genetic innovation. However, by comparing topologies instead of raw DNA sequences in this study, for the first time, we have been able to reveal the role of TEs in 3D genome conservation. This seemingly counter-intuitive role of species-specific parasitic sequences in helping maintain ancestral genome architecture is fundamentally different from all current and previous work regarding TEs’ role in gene regulation. This role is mediated by a long-postulated, classic genetic phenomena of binding site turnover—for CTCF in this case. Redundant TE-derived CTCF sites in the vicinity of conserved chromatin anchor/boundary can sometimes take over from the conserved anchor/boundary element, thus slightly shifting the anchor/boundary site while largely maintaining the 3D structure. Certain TE subfamilies like mouse SINE B2s contain pre-existing CTCF motifs within them, while others like mouse RLTR30 provide sequence fodder which upon a couple of specific point mutations can acquire CTCF binding and potentiate this binding site turnover.

In this study, 123 turnover events were observed in mouse on the basis of 3331 annotated loops (3.7%) whereas in humans 89 turnover events were observed out of 9448 loops (0.94%). This four-fold higher rate of turnover events in mouse highlights differences in between species and the turnover phenomenon being investigated. The higher rate of loop anchor CTCF turnover in the mouse genome was amplified by the arrival of CTCF-motif containing B2 SINEs. The genome is replete with such events and we have for the first time functionally dissected and validated them in the context of 3D genome conservation, opening up the doors for such investigations in the field for enhancer or promoter turnover events.

The *fons et origo* of CTCF motifs in B2 SINEs has been extensively researched. B2 SINEs are derived from

tRNA genes. Mouse tRNA genes have been shown to possess classical insulator activity and the potential to function as boundary elements [39]. Moreover, CTCF-binding enrichment in B2 SINEs and repeat-driven dispersal of CTCF binding has been shown to be a fundamental, ancient, and still highly active mechanism of genome evolution in mammalian lineages [23].

Similarly, the role CTCF motifs in viral genome regulation has been a topic of tremendous interest and investigation. In EBV, this control involves direct binding of CTCF across the viral genome and the formation of three-dimensional loops between virus promoters and enhancers [40]. CTCF is important in the regulation of gene expression of a number of human DNA viruses [41]. It also plays a critical role in epigenetic regulation of viral gene expression to establish and/or maintain a form of latent infection that can reactivate efficiently [42]. Recent evidence has also shown that HTLV-1 inserts an ectopic CTCF binding site forming loops between the provirus and host genome, altering expression of proviral and host gene [43]. CTCF has also been shown to promote HSV-1 lytic transcription by facilitating the elongation of RNA Pol II and preventing silenced chromatin on the viral genome [44]. Moreover, one can speculate that having a CTCF motif can not only help in maintaining viral genome confirmation but can also help insulate the chromatin activity of the neighborhood wherein the virus inserts into the host genome. It may also increase the chances of long-range interactions taking place which can sometimes bring in other TFs and/or polymerase, leading to enhanced transcription at the site of viral integration.

Our in-depth analysis of 3D genome structures upon genetic manipulation of candidate TEs revealed principles of how 3D genome evolves. In one example, a human TE provided a conserved chromatin boundary and loop anchor, whereas the ancestral CTCF site had decayed. Upon deletion, the chromatin domains collapsed, and loops eliminated, underscoring the importance of the TE in maintaining the local 3D genome structure.

In another case where a human TE provided a similarly conserved boundary and loop anchor, the ancestral CTCF site was still recognizable but was decommissioned. Deletion of the TE resulted in reinstallation of the ancestral CTCF site to form a slightly shifted boundary and loop anchor, and the local chromatin domains were largely preserved. In this second case that we validated, we potentially undid the events that took place during the course of (tens of millions of years) evolution by removing a young TE (LTR41) and having the ancestral “decommissioned” TE (MER82) re-uptake its function, thereby “reversing” the path of evolution in a dish (in days). Thus,

experimentally demonstrating the evolutionary impact of TE-derived CTCF sites. Moreover, the concept of such shadow loop anchors residing in TEs that can be activated upon escape from epigenetic silencing is extremely crucial to take into account for studies pertaining to diseases of the epigenome like certain cancers, their treatment and therapy. This study also underscores the redundancy that exists in the genome when it comes to CTCF binding sites and can potentially explain why we may not always see a drastic change in 3D genome structure upon deleting CTCF binding sites.

It is important to remember that the contribution of TEs outlined in this manuscript are underestimates as we have yet to (i) uniquely identify all loop anchor CTCF sites (especially in highly repetitive regions), (ii) annotate all repetitive elements, especially ancient TEs that have diverged far from their identity [19], and (iii) identify other architectural proteins and expand this framework beyond just CTCF-derived loop anchors.

While most studies highlight TEs’ role in innovating new functions by providing novel regulatory elements such as enhancers and promoters, we implicate the role of TEs in functional conservation inviting us to reexamine this unconventional role—perhaps many novel regulatory elements derived from TEs are not creating new functions, but rather providing redundant genetic material thus contributing to the robustness of gene regulatory networks. These findings will undoubtedly stimulate investigations to explore the multitude modes of regulatory evolution mediated by TEs. Recently, TAD boundaries have been shown to frequently harbor clusters of CTCF sites that contribute to cohesin stabilization and are critical for the functional stability of higher-order chromatin structure [45]. Indeed, some of the CTCF sites in these clusters are TE-derived, further appreciating the role of TEs in the maintenance of higher-order chromosomal structures in mammals. The transcriptional activation of retrotransposons has also been linked to the restructuring of genome architecture during human cardiomyocyte development [46].

A caveat of the analysis presented in this study is that the *in situ* Hi-C maps (re-analyzed in this study) of the 9 cell lines were sequenced to varying depths and thus differ in their resolution and “completeness” of loop annotations. Hence, due to this limitation of publicly available high-resolution Hi-C data, our findings likely represent a lower bound of TEs’ involvement in shaping both the conserved and species-specific 3D genome. These analyses need to be revisited as and when higher-resolution datasets are available.

Lastly, our study opens the doors for population-scale genetic variation studies that identify polymorphic TE insertions to be reconciled with population-scale 3D genome and regulatory variation. These future

explorations will present yet another vignette of TEs and their very many roles in accelerating adaptive evolution.

Conclusions

Taken together, our findings reveal a formerly uncharacterized role that TEs have played in the evolution of higher-order chromosomal structures in mammals. TEs have contributed a substantial number of loop anchors in mouse and human 3D genomes, a fraction of which were co-opted to help maintain conserved higher-order chromosomal structures. TE transposition provides redundant CTCF motifs and a novel method for CTCF binding site turnover to maintain regulatory conservation (defined here as the preservation of long-range chromosomal interactions, loop, and boundary formation), by compensating for the loss of local primary sequence—local sequence that would have otherwise allowed the assessment of purifying selection. Deletion of these TEs in human cell lines eliminated the chromatin loops that they anchor and resulted in collapse of conserved chromatin structure, as expected by our hypothesis. More strikingly, we demonstrate that in another case the loop anchor shifted to an alternative TE-derived CTCF site nearby, resulting in largely unchanged chromatin structure, underscoring the dynamic nature and robustness of the 3D genome upon TE infiltration. These TEs that maintain conserved chromatin loops via turnover are hypomethylated through deep time, an observation that highlights the intimate interplay between genome, epigenome, and 3D genome in evolution. This research provides a foundation to study the impact of TEs and expand our understanding of chromosomal folding—its emergence, maintenance, and transformation—in the context of evolving genomes. Ultimately, our study reveals how selfish genetic elements, regardless of their origins, can be repurposed to provide redundant TF motifs and maintain latent genome sanctity and regulatory fidelity by conserving 3D genome structure.

Methods

Dataset GEO accession numbers

The genomic data analyzed in this study were obtained from publicly available datasets. Hi-C datasets were obtained from GSE63525 (mouse: CH12; humans: GM12878, HeLa, HMEC, IMR90, K562, NHEK). GM12878 ChIA-PET dataset was obtained from GSE72816. GM12878 CTCF ChIP-seq datasets were obtained from ENCODE (ENCSR000AKB and ENCSR000DZN). CH12-LX CTCF ChIP-seq datasets were obtained from Mouse ENCODE (ENCSR000ERM and ENCSR000DIU). WGBS methylation dataset for GM12878 was also obtained from ENCODE, GEO: GSE86765 (ENCSR890UQO). Mouse ESC and NSC Hi-C data was obtained from PMID: 30414923.

Loop anchor CTCF-RE intersection

We generated a list of unique anchor CTCF sites using the HiCCUPS output [1] for various mentioned cell lines. We then overlapped loop anchor CTCF motifs identified using HiCCUPS [1] with *RepeatMasker* (RMSK v4.0.7, for hg19 and mm9) and required at least 10 bp of the core CTCF motif to intersect with a repetitive element (RE) to call it a RE-derived loop anchor CTCF site. Further, only loops with (i) at least one known RE-derived anchor CTCF site or (ii) two non-RE derived anchor CTCF sites were taken into consideration for analysis of RE-derived loop counts, because we can definitively say whether the loops and their loop anchor CTCF sites were derived from REs or not. Loops with both unidentified loop anchor CTCF sites, or one unidentified and one non-RE derived anchor CTCF site were not considered as there is the possibility of having at least one of the other anchor CTCF sites derived from a RE. We followed the same methodology when considering ChIA-PET loops.

TE class and family distribution

We ran *RepeatMasker* v4.0.7 with the *-s* slow search parameter on the hg19 and mm9 genomes to obtain a comprehensive list of REs in the genome and their corresponding subfamily, family, and class annotations. We used RE counts (generated as previously outlined) to characterize their distribution to loop anchor CTCF sites. For characterizing RE-derived CTCF binding peaks, we repurposed a previously used strategy [10]. Briefly, we required that the centers of the MACS-called peaks of ENCODE-generated CTCF ChIP datasets overlapped with RE fragments. We used the length distribution of various RE family and classes in the entire genome as the background distribution.

Loop orthology check

We used *liftOver* [20] to convert CH12-LX loop annotations from mm9 mouse genome coordinates to hg19 human genome coordinates. We used various sequence match rates (*minMatch* = 0.05, 0.1..., 1) to convert CH12-LX mouse peaks from mm9 genome coordinates to hg19 genome coordinates. To optimize for the *minMatch* parameter, we generated ten shuffled (randomized) peak annotations by using *bedtools shuffle -chrom* command to permute their location on the chromosome of origin. *minMatch* parameter of 0.1 was chosen for *liftOver* analyses henceforth, as it resulted in the greatest number of features being lifted over (on average) and lower coefficient of variation across the 10 simulated sets. We lifted over 3245 out of 3331 mouse peaks from mm9 to hg19, using the *minMatch* 0.1, to facilitate cross-species peak annotation comparison. To call a mouse feature conserved in humans, we required that

the loop anchor pairs individually lie within a min (half of loop length, vicinity threshold) window of an existing loop anchor pair. The vicinity threshold was put in place to account for cross-species liftOver errors and facilitate comparison of higher-order chromosomal features that vary from 120 kb to 125 Mb in length (in mouse). We tested multiple vicinity thresholds ranging from 500 bp to 100 Mb and identified false discovery rates using simulated sets of mouse features and comparing them to the orthology observed between the real CH12-LX (mouse) and GM12878 (human) features. We decided to use 50 kb as the vicinity threshold as it corresponded to a false discovery less than 0.1. We found that 1688 CH12-LX mouse peaks overlapped at least one corresponding peak in GM12878 human lymphoblastoid cells. We performed a similar analysis to compare “muritized” human features (liftOver from GM12878) to actual mouse features (CH12-LX). We found that 1900 GM12878 human peaks overlapped at least one corresponding peak in CH12-LX mouse lymphoblastoid cells. We then filtered for features that displayed reciprocal matches (reciprocal best hits) in the two comparisons (mouse-to-human and human-to-mouse) as stated above. Finally, we curated the list by considering genic, epigenomic, and transcriptomic synteny to pick exactly one orthologous human loop to a corresponding mouse loop, to enlist 1596 high-confidence orthologous peak calls (Additional file 2: Table S1.1). A flowchart of the pipeline is shown in Fig. 2a.

TE age estimation

Species divergence times were based on [47]. Repeat ages were estimated by dividing the percent divergence of extant copies from the consensus sequence by the species neutral substitution rate. Substitution rates (mutations/year) used were as follows: humans: 2.2×10^{-9} ; mouse: 4.5×10^{-9} , from [48]. Jukes-Cantor and Kimura distances were calculated by aligning each TE to its consensus sequence and counting all possible mutations (see below). Single nucleotide substitution counts were normalized by the length of the genomic TE minus the number of insertions (gaps in the consensus). These mutation rates were then used to calculate the Jukes-Cantor and Kimura distances for each genomic TE.

Candidate selection and filtering

After curating the list of conserved loops, we looked for TE-derived orthologous loops in humans that were discordant for TEs in mouse. After identifying the list of TE-derived CTCF turnover events in humans, we comprehensively surveyed the local CTCF binding landscape (CTCF ChIP-seq peaks) to ensure (i) there were no other CTCF binding sites in the vicinity that could function as loop anchors in humans (in the first case); and

(ii) there was only one other unique CTCF binding site, i.e., the ancestral CTCF motif (in the second case). We also ensured that the TE insertion from which the loop anchor CTCF site was derived was human-specific and not present in mouse (Additional file 2: Table S1.2). We repeated this analysis to identify TE-mediated turnover in mouse as well (Additional file 2: Table S1.3). We also identified events wherein TEs mediated turnover events both in mouse and human (Additional file 2: Table S1.4). One possible explanation for this observation is that similar selective pressures (like the need to maintain higher-order chromosomal structure) led to the convergent co-option of species-specific TEs at syntenic locus, independently in both the genomes.

Cell culture methods

GM12878 cell lines were grown between 200K and 800K cells/ml in 10-ml cultures in T-25 flasks, in a humidified incubator with 95% CO₂ at 37 °C in RPMI1640 media (Gibco, 1187-085) supplemented with 15% fetal bovine serum (Corning, 35-011-CV) and 100 U/ml penicillin-streptomycin (Gibco, 15140-122) as per the ENCODE standards.

CRISPR-Cas9 mediated genome engineering

Our CRISPR workflow consisted of the following steps: We identified turned over chromatin loops that are maintained by TEs, with unique, convergently oriented TE-derived CTCF motifs within loop anchors [1]. We used two independent CRISPR sgRNA design engines CRISPOR [49] and CRISPRScan [50] to rationally design multiple pairs of sgRNAs that have high cutting efficiency and minimum off-target effects. We used pU6-(BbsI)_CBh-Cas9-T2A-BFP plasmid (Addgene, 64323) and pU6-(BbsI)_CBh-Cas9-T2A-mCherry plasmid (Addgene, 64324) as the CRISPR delivery vectors. For each sgRNA, we designed and annealed two single-stranded oligos with compatible overhangs that can be cloned into BbsI-digested BFP and mCherry CRISPR vectors through standard ligation techniques. For every pair of sgRNAs, we constructed BFP-CRISPR vectors and mCherry-CRISPR vectors that express sgRNAs targeting upstream and downstream of the candidate TEs, respectively. BFP-CRISPR vectors and mCherry-CRISPR vectors each were co-transfected into GM12878 cells in antibiotic-free media using the Neon transfection system. After 24 h of incubation, the transfected cells were analyzed by flow cytometry (Beckman Coulter MoFlo) for BFP-positive and mCherry-positive subpopulations. Transfection efficiencies were usually between 3 and 5%. We single-cell sorted these double-positive fluorescent cells into 96-well plates for clone expansion and allowed to grow for 21–28 days. After that, 20–48 clones were screened per transfection. Genomic DNA from CRISPR

clones was extracted using *Quick-DNA* Miniprep kit for genotyping and validated with Sanger sequencing. Details of sequences used to generate clones used in this study are listed in Additional file 3: Table S2.4. We then performed in situ Hi-C on the select edited cell lines and performed hybrid selection on the in situ Hi-C libraries for a region around the targeted TE to generate Hi-C² libraries that can easily and cheaply be sequenced to read off the effects of our TE deletions on local genome folding.

Hi-C² probe design

To design probes targeting the two regions for Hybrid Capture Hi-C (Hi-C²), we followed a similar approach as [8]. In short, we (i) identified all MboI restriction sites within the target region, (ii) designed our bait probe sequences to target sequences within a certain distance of the MboI restriction sites as Hi-C ligation junctions occur between them, and (iii) followed a similar three-pass probe design strategy sequentially increasing various parameters like the distance of the probe from the MboI restriction site, the number of repetitive bases, the GC content, and probe density in gaps with relaxed probe design quality filters. We then removed overlapping probes or probes with identical sequences. After all three passes, we identified 2741 unique probes covering region 1 (chr10:26-28 Mb; 1.37 probes/kb) and 1856 probes covering region 2 (chr8:70.3-71.8 Mb; 1.24 probes/kb). Fifteen-base pair primer sequences (unique for each region, details in Additional file 3: Table S2.3) were then appended to both ends of the 120-bp probe sequence to facilitate single oligo pool synthesis and subsequent amplification of region-specific sub-pools. Probe construction and hybrid selection was then followed with sequences specific to this study using the same strategy detailed in [8].

Hi-C experiments

The Hi-C datasets used in our analyses were generated using the in situ Hi-C protocol standardized by the 4DN consortia. In brief, the in situ Hi-C protocol involves crosslinking cells with 1% formaldehyde for 10 min, permeabilizing them with nuclei intact, digesting the DNA with MboI (4-cutter restriction enzyme), filling the 5'-overhangs while incorporating biotin-14-dATP (a biotinylated nucleotide), followed by ligating the resulting blunt-end fragments, shearing the DNA to a 400–700-bp fragment size, capturing the biotinylated ligation junctions with streptavidin beads, building an Illumina library with 10–12 rounds of PCR amplification, and finally analyzing the resulting fragments with paired-end sequencing. The resulting library was always shallow sequenced to 500 K–4 M reads to check for library build quality looking at key statistics such as complexity,

number of Hi-C contacts, inter vs. intrachromosomal interactions, and long-range vs/ short-range intrachromosomal interactions. Libraries that passed the quality check were either sequenced deeper and/or used as pools for subsequent Hi-C² experiments.

For our genome engineering experiments, we generated 14 in situ Hi-C libraries (Additional file 3: Table S2.1) from GM12878 cells. We also generated 18 in situ Hi-C² libraries from various genome-engineered GM12878 cell lines on which hybrid selection was performed. All in situ Hi-C libraries generated as part of this study are detailed in Additional file 3: Table S2.2. All the Hi-C data was processed using the computational pipeline described in full detail in [1]. Hi-C libraries were sequenced to a depth of between 624K and 333M reads (on average, 63.8M reads). Hi-C² libraries were sequenced to a depth of between 6.7M and 168M reads (on average, 35.8M reads). All data was initially processed using the pipeline published in [1] and visualized on the desktop and web version of Juicebox. We combined Hi-C and Hi-C² contact maps corresponding to the same genotype and the same locus using the Juicer's mega.sh script as these are in essence “biological” replicates, to generate higher resolution megamaps.

Analysis of cross-domain interactions

We subsampled the Hi-C² corresponding to the R1-WT megamap (containing 46 M reads) and R1-KO (containing 56 M reads) for 5 M reads, 10 times to create 10 independent R1-WT and R1-KO mini-maps. For each of these Hi-C maps, we used the Juicer Tools dump command to extract the VC_sqrt normalized contact matrix. Intradomain interactions were defined as interaction that (i) originate and terminate in domain 1 or (ii) originate and terminate in domain 2. Interdomain interactions were defined as interactions that originate in domain 1 and terminate in domain 2. We then calculate percentage of cross-domain interactions for each of the mini-maps using the formula: $(\text{number of interdomain-interactions}) \times 100 / ((\text{number of intradomain-interactions}) + (\text{number of interdomain-interactions}))$. The percentage of cross-domain interactions were calculated for the target domain as well as a control domain. The distribution of cross-domain interactions across the targeted domain was found to be significantly different in the KO vs. the WT (*t*-test: two-sample assuming unequal variances, *p* value = 1.40668×10^{-16}). The distribution of cross-domain interactions across a nearby control domain however was not found to be significantly different in the KO vs. the WT (*t*-test: two-sample assuming unequal variances, *p* value = 0.013254165). Raw simulation data and statistics are provided in Additional file 3: Table S2.5.

For Additional file 1: Figure S3C, we used the Hi-C megamap corresponding to R2-WT and R2-KO to

retrieve raw interaction counts at a 100 kb resolution. Percent cross-domain interactions were calculated using the formula stated above. We calculated the enrichment of cross-domain interactions in the LTR41-DKO w.r.t. the WT across the targeted domain as well as a nearby control domain.

DNA methylation analysis

We generated a methylation metaplot representing the mean CpG methylation value from WGBS data (ENCODE dataset: ENCFF835NTC) of 20-bp sliding windows, centered on CTCF motifs (and ± 2 kb around it) segmented by their origin/TE derivation status.

Analysis of TE mutational profile

TE consensus construction

For most of the TE subfamilies, we retrieved the consensus sequences from the RepBase library (RepBase 22.02, RepeatMaskerEdition20170127) [51]. However, LINE elements are fragmented to 5' end, ORF2, and 3' end regions in RepBase library. To reconstruct full-length LINE consensus, we identified TE fragments in human and mouse genome using RepeatMasker and compared the standard output (.out file) with the alignment output (.align file) from the same RepeatMasker run [52]. For each LINE element in the standard output, we summarized which 5' end, ORF2, and 3' end fragments have been used most to construct the full-length element. Then we use EMBOSS Water local alignment algorithm to align the three pieces together and generated the full-length LINE consensus sequences [53].

Crossmatch alignments

We ran RepeatMasker 4.0.7 on the mm9 and hg19 genomes using crossmatch as the search engine. We then parsed the alignment file to determine the substitution rates between the ancestral sequence and the genomic element. For each genomic element, we counted the number of A-to-C, A-to-G, A-to-T, C-to-A, C-to-G, C-to-T, G-to-A, G-to-C, G-to-T, T-to-A, T-to-C, and T-to-G substitutions (single nucleotide substitutions), where the first nucleotide indicates the ancestral sequence and the second nucleotide indicates the genomic sequence. We ignored any substitutions that involved ambiguous nucleotides. We also counted the number of insertions and deletions. All substitution frequencies were normalized by the length of the genomic sequence to estimate the substitution rates in each TE. Any genomic TE with a length less than 20% of the ancestral sequence was filtered out. For each single nucleotide substitution, we calculated the average substitution rate in two subsets of TEs (details below). We also calculated the combined C-to-T and G-to-A substitution rate (methylation-associated substitutions) and the combined

rate of all other substitutions (non-methylation-associated substitutions) to compare the rate of DNA methylation-induced mutations to other mutations. The methylation substitution rate was computed by taking the average of the C-to-T and G-to-A rates for each TE and then averaging over turnover events. The non-methylation substitution rate was computed by taking the average of all other (ten) single nucleotide substitutions for each TE and then averaging over turnover events.

We generated a background distribution by repeating this analysis on 1000 permutations of all genomic TEs. We first calculated the frequency of each TE subfamily in the set of turnover events. For each permutation, we randomly selected genomic TEs (not involved in anchoring loops) from each subfamily to reflect their frequency in turnover events. The single nucleotide substitution rate, methylation-associated substitution rate, and non-methylation-associated substitution rate were calculated as described above. The distribution of all substitution rates from the permutations follows a normal distribution (KS test, $p > 0.0036$, Bonferroni correction $\alpha = 0.05$ for $N = 14$ hypotheses, Additional file 4: Table S3.1). The background distribution was then used to perform a left-tailed z -test. We did not compute a two-tailed p value because our null hypothesis is that the observed mutation rates are greater than or equal to the background distribution mean. For the 12 single nucleotide substitutions, we used Bonferroni correction to account for multiple hypotheses.

Needle realignments

RepeatMasker performs post-processing after running crossmatch, so coordinates and TE subfamily assignments in the .out file do not always reflect the contents of the .align file. To improve our estimates of mutation rates, we realigned each TE to its matched consensus sequences. We extracted the genomic and subfamily consensus sequence using the coordinates reported in the .out file. We then performed a global alignment using EMBOSS Needle v6.6.0.0 using a gap open penalty of 10, a gap extension penalty of 0.5, and the EDNAFULL scoring matrix. We used the alignment to recompute single nucleotide substitutions for each TE and then repeated the same analysis we used for crossmatch alignments. We did not filter out TEs with a length less than 20% of the ancestral sequence because this filter was originally put in place to account for discrepancies between the .align and .out files. As before, the distribution of all substitution rates from the permutations follows a normal distribution (KS test, $p > 0.0036$, Bonferroni correction $\alpha = 0.05$ for $N = 14$ hypotheses, Additional file 4: Table S3.2).

Supplementary information

The online version of this article (<https://doi.org/10.1186/s13059-019-1916-8>) contains supplementary material, which is available to authorized users.

Additional file 1: Figure S1. Contribution of repetitive element (RE) families to cell-type specific loop anchor CTCF sites. **Figure S2.** Phylogenetic reconstruction of CTCF binding motifs that anchor the conserved chromatin loop near AKAP8L gene. **Figure S3.** Landscape of the higher-order chromosomal structures conserved in human and mouse maintained by an L1M3f-derived CTCF binding site turnover event in humans. **Figure S4.** Landscape of the higher-order chromosomal structures conserved in human and mouse maintained by an LTR41-derived CTCF binding site turnover event in humans. **Figure S5.** Experimental design and results of CRISPR-Cas9 mediated deletions of L1M3f and LTR41. **Figure S6.** L1M3f-derived CTCF site is required for maintaining conserved boundary and fidelity of long-range interactions. **Figure S7.** CTCF binding site turnover from an LTR41-derived CTCF motif to MER82-derived CTCF site. **Figure S8.** TEs that potentiate loop anchor turnover in human and mouse show mutational signatures of hypomethylation through evolutionary time (related to Fig. 4B). **Figure S9.** TEs that potentiate loop anchor turnover in human and mouse show mutational signatures of hypomethylation through evolutionary time (related to Fig. 4B).

Additional file 2: Table S1. This file contains 4 supplementary tables for this paper. Table S1.1: List of orthologous mouse-human loops mapped with details. Table S1.2: List of orthologous mouse-human loops that are turnover events in mouse mapped with TE details. Table S1.3: List of orthologous mouse-human loops that are turnover events in human mapped with TE details. Table S1.4: List of orthologous mouse-human loops whose anchors are derived from TEs in both humans and mouse, mapped with TE details.

Additional file 3: Table S2. This file contains 5 supplementary tables for this paper. Table S2.1: Hi-C Experiments. Table S2.2: Hi-C² Experiments. Table S2.3: Hi-C² Region/Probe Information. Table S2.4: Targeted TEs and genotype information for all mutant clones. Table S2.5: Crossdomain interaction stats (related to Fig. 3C).

Additional file 4: Table S3.1. KS-test statistic to test distribution of various nucleotide substitutions (crossmatch). Table S3.2. KS-test statistic to test distribution of various nucleotide substitutions (needle).

Additional file 5: Review History. This file contains reviewers comments and author responses.

Acknowledgements

We thank members of the Wang Lab for helpful discussions related to the project; Jessica Hoisington-López and Maria Lynn Jaeger from The Edison Family Center for Genome Sciences and Systems Biology for assistance with sequencing; Matthew Patana and Daniel Schweppe from the Siteman Flow Cytometry core for FACS expertise.

Third party data

Mouse ESC and NSC HiC data was obtained from PMID: 30414923. Other Hi-C datasets were obtained from GSE63525 (mouse: CH12; humans: GM12878, HeLa, HMEC, IMR90, K562, NHEK). GM12878 ChIA-PET dataset was obtained from GSE72816. GM12878 CTCF ChIP-seq datasets were obtained from ENCODE (ENCSR000AKB and ENCSR000DZN). CH12-LX CTCF ChIP-seq datasets were obtained from Mouse ENCODE (ENCSR000ERM and ENCSR000DIU). WGBS methylation dataset for GM12878 was also obtained from ENCODE, GEO: GSE86765 (ENCSR890UQO).

Peer review information

Yixin Yao was the primary editor on this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

One-sentence summary

Co-option of transposable elements maintains conserved 3D genome structures via CTCF binding site turnover in human and mouse.

Review history

The review history is available as Additional file 5.

Funding

M.N.K.C. was partly supported by the Precision Medicine Pathway, Washington University; H.S.J. was partly supported by NIH grant T32 GM007067; X.Z. was partly supported by R25DA027995; T.W. is supported by R01HG007175, U01CA200060, U24ES026699, U01HG009391, U41HG010972, and American Cancer Society RSG-14-049-01-DMC.

Availability of data and materials

Datasets generated and analyzed in this current study are available on Gene Expression Omnibus under the accession number GSE141550 [54] (<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi>).

Authors' contributions

MNKC and TW conceived and designed this study; MNKC analyzed the data, performed the experiments, generated the sequencing libraries, and wrote the manuscript with inputs from TW; RZF, JTW, and XZ contributed text and revised the manuscript; HSJ contributed the reagents and resources; RZF performed the mutation frequency simulations along with MNKC; XZ generated the TE ancestral sequences and TE alignments; TW supervised the project. All authors subsequently edited and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Received: 5 July 2019 Accepted: 8 December 2019

Published online: 24 January 2020

References

- Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, et al. A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*. 2014;159(7):1665–80.
- Tang Z, Luo OJ, Li X, Zheng M, Zhu JJ, Szalaj P, et al. CTCF-mediated human 3D genome architecture reveals chromatin topology for transcription. *Cell*. 2015;163(7):1611–27.
- Sexton T, Yaffe E, Kenigsberg E, Bantignies F, Leblanc B, Hoichman M, et al. Three-dimensional folding and functional organization principles of the *Drosophila* genome. *Cell*. 2012;148(3):458–72.
- Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, et al. Topological domains in mammalian genomes identified by analysis of chromatin interactions. *Nature*. 2012;485(7398):376–80.
- Gorkin DU, Leung D, Ren B. The 3D genome in transcriptional regulation and pluripotency. *Cell Stem Cell*. 2014;14:771–5.
- Schwarzer W, Abdennur N, Goloborodko A, Pekowska A, Fudenberg G, Loe-Mie Y, et al. Two independent modes of chromatin organization revealed by cohesin removal. *Nature*. 2017;551(7678):51–6.
- Strom AR, Emelyanov AV, Mir M, Fyodorov DV, Darzacq X, Karpen GH. Phase separation drives heterochromatin domain formation. *Nature*. 2017; 547(7662):241–5.
- Sanborn AL, Rao SSP, Huang S-C, Durand NC, Huntley MH, Jewett AI, et al. Chromatin extrusion explains key features of loop and domain formation in wild-type and engineered genomes. *Proc Natl Acad Sci*. 2015;112(47): E6456–65.
- Fudenberg G, Imakaev M, Lu C, Goloborodko A, Abdennur N, Mirny LA. Formation of Chromosomal Domains by Loop Extrusion. *Cell Rep*. 2016; 15(9):2038–49.
- Sundaram V, Cheng Y, Ma Z, Li D, Xing X, Edge P, et al. Widespread contribution of transposable elements to the innovation of gene regulatory networks. *Genome Res*. 2014;24(12):1963–76.
- Bourque G, Leong B, Vega VB, Chen X, Yen LL, Srinivasan KG, et al. Evolution of the mammalian transcription factor binding repertoire via transposable elements. *Genome Res*. 2008;18(11):1752–62.

12. Kunarso G, Chia N-Y, Jeyakani J, Hwang C, Lu X, Chan Y-S, et al. Transposable elements have rewired the core regulatory network of human embryonic stem cells. *Nat Genet.* 2010;42(7):631–4.
13. Jacques PÉ, Jeyakani J, Bourque G. The Majority of Primate-Specific Regulatory Sequences Are Derived from Transposable Elements. *PLoS Genet.* 2013;9(5):e1003504.
14. Wang T, Zeng J, Lowe CB, Sellers RG, Salama SR, Yang M, et al. Species-specific endogenous retroviruses shape the transcriptional network of the human tumor suppressor protein p53. *Proc Natl Acad Sci U S A.* 2007; 104(47):18613–8.
15. Chuong EB, Elde NC, Feschotte C. Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Science.* 2016;351(6277): 1083–7.
16. Lynch VJ, Leclerc RD, May G, Wagner GP. Transposon-mediated rewiring of gene regulatory networks contributed to the evolution of pregnancy in mammals. *Nat Genet.* 2011;43(11):1154–9.
17. Britten RJ, Davidson EH. Repetitive and non-repetitive DNA sequences and a speculation on the origins of evolutionary novelty. *Q Rev Biol.* 1971;46(2): 111–38.
18. Pękowska A, Klaus B, Xiang W, Severino J, Daigle N, Klein FA, et al. Gain of CTCF-Anchored Chromatin Loops Marks the Exit from Naive Pluripotency. *Cell Syst.* 2018;7(5):482–95.
19. de Koning APJ, Gu W, Castoe TA, Batzer MA, Pollock DD. Repetitive elements may comprise over two-thirds of the human genome. *PLoS Genet.* 2011;7(12):e1002384.
20. Hinrichs AS. The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res.* 2006;34:D590–8.
21. Feschotte C, Pritham EJ. DNA Transposons and the Evolution of Eukaryotic Genomes. *Annu Rev Genet.* 2007;41:331–68.
22. Sundaram V, Wang T. Transposable element mediated innovation in gene regulatory landscapes of cells: re-visiting the “gene-battery” model. *BioEssays.* 2018;40:1700155.
23. Schmidt D, Schwalie PC, Wilson MD, Ballester B, Gonçalves A, Kutter C, et al. Waves of retrotransposon expansion remodel genome organization and CTCF binding in multiple mammalian lineages. *Cell.* 2012;148(1–2):335–48.
24. Ludwig MZ, Patel NH, Kreitman M. Functional analysis of eve stripe 2 enhancer evolution in *Drosophila*: rules governing conservation and change. *Development.* 1998;125(5):949–58.
25. Moses AM, Pollard DA, Nix DA, Iyer VN, Li XY, Biggin MD, et al. Large-scale turnover of functional transcription factor binding sites in *Drosophila*. *PLoS Comput Biol.* 2006;2(10):e130.
26. Venkataram S, Fay JC. Is transcription factor binding site turnover a sufficient explanation for cis-regulatory sequence divergence? *Genome Biol Evol.* 2010;2:851–8.
27. Villar D, Flicek P, Odom DT. Evolution of transcription factor binding in metazoans-mechanisms and functional implications. *Nat Rev Genet.* 2014; 15(4):221–33.
28. Matzke MA, Mette MF, Matzke AJ. Transgene silencing by the host genome defense: implications for the evolution of epigenetic control mechanisms in plants and vertebrates. *Plant Mol Biol.* 2000;43(2–3):401–15.
29. Yoder JA, Walsh CP, Bestor TH. Cytosine methylation and the ecology of intragenomic parasites. *Trends Genet.* 1997;13(8):335–40.
30. Slotkin RK, Martienssen R. Transposable elements and the epigenetic regulation of the genome. *Nat Rev Genet.* 2007;8(4):272–85.
31. Huda A, Mariño-Ramírez L, Jordan IK. Epigenetic histone modifications of human transposable elements: genome defense versus exaptation. *Mob DNA.* 2010;1:2.
32. Lowe CB, Haussler D. 29 mammalian genomes reveal novel exaptations of mobile elements for likely regulatory functions in the human genome. *PLoS One.* 2012;7(8):e43128.
33. Bejerano G, Lowe CB, Ahituv N, King B, Siepel A, Salama SR, et al. A distal enhancer and an ultraconserved exon are derived from a novel retroposon. *Nature.* 2006;441(7089):87–90.
34. Kidwell MG, Lisch DR. Transposable elements and host genome evolution. *Trends Ecol Evol.* 2000;15(3):95–9.
35. Kanduri C, Pant V, Loukinov D, Pugacheva E, Qi CF, Wolffe A, et al. Functional association of CTCF with the insulator upstream of the H19 gene is parent of origin-specific and methylation-sensitive. *Curr Biol.* 2000;10(14): 853–6.
36. Kurukuti S, Tiwari VK, Tavosoidana G, Pugacheva E, Murrell A, Zhao Z, et al. CTCF binding at the H19 imprinting control region mediates maternally inherited higher-order chromatin conformation to restrict enhancer access to Igf2. *Proc Natl Acad Sci.* 2006;103(28):10684–9.
37. Shen JC, Rideout WM, Jones PA. The rate of hydrolytic deamination of 5-methylcytosine in double-stranded DNA. *Nucleic Acids Res.* 1994;22(6):972–6.
38. Bird AP. DNA methylation and the frequency of CpG in animal DNA. *Nucleic Acids Res.* 1980;8(7):1499–504.
39. Ebersole T, Kim JH, Samoshkin A, Kouprina N, Pavlicek A, White RJ, et al. tRNA genes protect a reporter gene from epigenetic silencing in mouse cells. *Cell Cycle.* 2011;10(16):2779–91.
40. Tempera I, Klichinsky M, Lieberman PM. EBV Latency Types Adopt Alternative Chromatin Conformations. *PLoS Pathog.* 2011;7(7):e1002180.
41. Pentland I, Parish JL. Targeting CTCF to control virus gene expression: a common theme amongst diverse DNA viruses. *Viruses.* 2015;7(7):3574–85.
42. Lee JS, Raja P, Pan D, Pesola JM, Coen DM, Knipe DM. CCCTC-binding factor acts as a heterochromatin barrier on herpes simplex viral latent chromatin and contributes to poised latent infection. *MBio.* 2018;9:e02372–17.
43. Satou Y, Miyazato P, Ishihara K, Yaguchi H, Melamed A, Miura M, et al. The retrovirus HTLV-1 inserts an ectopic CTCF-binding site into the human genome. *Proc Natl Acad Sci U S A.* 2016;113(11):3054–9.
44. Lang F, Li X, Vladimirova O, Hu B, Chen G, Xiao Y, et al. CTCF interacts with the lytic HSV-1 genome to promote viral transcription. *Sci Rep.* 2017;7: 39861.
45. Kenteropoulou E, Aitken SJ, Feig C, Stefflova K, Ibarra-Soria X, Odom DT, et al. Clustered CTCF binding is an evolutionary mechanism to maintain topologically associating domains. *bioRxiv.* 2019. <https://doi.org/10.1101/668855>.
46. Zhang Y, Li T, Preissl S, Grinstein J, Farah E, Destici E, et al. Transcriptionally active HERV-H retrotransposons demarcate topologically associating domains in human pluripotent stem cells. *Nat Genet.* 2019;51(9):1380–8.
47. Meredith RW, Janečka JE, Gates J, Ryder OA, Fisher CA, Teeling EC, et al. Impacts of the cretaceous terrestrial revolution and KPg extinction on mammal diversification. *Science.* 2011;334(6055):521–4.
48. Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, et al. Initial sequencing and comparative analysis of the mouse genome. *Nature.* 2002;420(6915):520–62.
49. Haeussler M, Schönig K, Eckert H, Eschstruth A, Mianné J, Renaud JB, et al. Evaluation of off-target and on-target scoring algorithms and integration into the guide RNA selection tool CRISPOR. *Genome Biol.* 2016;17(11):148.
50. Moreno-Mateos MA, Vejnar CE, Beaudoin JD, Fernandez JP, Mis EK, Khokha MK, et al. CRISPRscan: designing highly efficient sgRNAs for CRISPR/Cas9 targeting in vivo. *Nat Methods.* 2015;12(10):982–8.
51. Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic genomes. *Mob DNA.* 2015;6(1):11.
52. Smit A, Hubley R, Green P. RepeatMasker Open-4.0.6 2013–2015. 2017; <http://www.repeatmasker.org>.
53. Rice P, Longden L, Bleasby A. EMBOSS: The European Molecular Biology Open Software Suite. *Trends Genet.* 2000;16(6):276–7.
54. Choudhary MNK, Friedman RZ, Wang JT, Jang HS, Zhuo X, Wang T. Co-opted transposons help perpetuate conserved higher-order chromosomal structures. *Datasets. Gene Expression Omnibus.* 2019; <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE141550>.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

