


RESEARCH

Open Access



# CpG island turnover events predict evolutionary changes in enhancer activity

Acadia A. Kocher<sup>1,7</sup>, Emily V. Dutrow<sup>1,8</sup>, Severin Uebbing<sup>1,9</sup>, Kristina M. Yim<sup>1</sup>, María F. Rosales Larios<sup>1</sup>, Marybeth Baumgartner<sup>1</sup>, Timothy Nottoli<sup>2,3</sup> and James P. Noonan<sup>1,4,5,6\*</sup> 

\*Correspondence:

james.noonan@yale.edu

<sup>1</sup> Department of Genetics, Yale School of Medicine, New Haven, CT 06510, USA

<sup>2</sup> Department of Comparative Medicine, Yale School of Medicine, New Haven, CT 06510, USA

<sup>3</sup> Yale Genome Editing Center, Yale School of Medicine, New Haven, CT 06510, USA

<sup>4</sup> Department of Ecology and Evolutionary Biology, Yale University, New Haven, CT 06520, USA

<sup>5</sup> Department of Neuroscience, Yale School of Medicine, New Haven, CT 06510, USA

<sup>6</sup> Wu Tsai Institute, Yale University, New Haven, CT 06510, USA

<sup>7</sup> Division of Molecular Genetics and Oncode Institute, Netherlands Cancer Institute, Amsterdam, The Netherlands

<sup>8</sup> Zoetis, Inc, 333 Portage St, Kalamazoo, MI 49007, USA

<sup>9</sup> Genome Biology and Epigenetics, Institute of Biodynamics and Biocomplexity, Department of Biology, Utrecht University, Utrecht, The Netherlands

## Abstract

**Background:** Genetic changes that modify the function of transcriptional enhancers have been linked to the evolution of biological diversity across species. Multiple studies have focused on the role of nucleotide substitutions, transposition, and insertions and deletions in altering enhancer function. CpG islands (CGIs) have recently been shown to influence enhancer activity, and here we test how their turnover across species contributes to enhancer evolution.

**Results:** We integrate maps of CGIs and enhancer activity-associated histone modifications obtained from multiple tissues in nine mammalian species and find that CGI content in enhancers is strongly associated with increased histone modification levels. CGIs show widespread turnover across species and species-specific CGIs are strongly enriched for enhancers exhibiting species-specific activity across all tissues and species. Genes associated with enhancers with species-specific CGIs show concordant biases in their expression, supporting that CGI turnover contributes to gene regulatory innovation. Our results also implicate CGI turnover in the evolution of Human Gain Enhancers (HGEs), which show increased activity in human embryonic development and may have contributed to the evolution of uniquely human traits. Using a humanized mouse model, we show that a highly conserved HGE with a large CGI absent from the mouse ortholog shows increased activity at the human CGI in the humanized mouse diencephalon.

**Conclusions:** Collectively, our results point to CGI turnover as a mechanism driving gene regulatory changes potentially underlying trait evolution in mammals.

**Keywords:** Transcriptional enhancer evolution, Gene regulation, Orphan CpG islands, Comparative genomics

## Background

Genetic variation in transcriptional enhancers has been associated with trait variation across species [1–8]. Sequence changes in enhancers are hypothesized to modify the regulatory information enhancers encode by changing transcription factor binding site (TFBS) composition, thereby altering recruitment of transcription factors and



© The Author(s) 2024. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

co-activators [9–12]. Such molecular changes may alter the spatiotemporal pattern and degree of enhancer activation and lead to corresponding changes in the expression of their target genes [13, 14]. Numerous studies have characterized the contribution of nucleotide substitutions [3, 4, 6, 15–22], transposable elements [23–26], and insertions and deletions [27, 28] to evolutionary changes in enhancer function. However, despite these advances, understanding the specific mechanisms by which genetic variation modifies enhancer activity during evolution remains challenging.

One approach to identifying changes in enhancer function across species relies on comparing levels of histone modifications, such as histone H3 lysine 27 acetylation (H3K27ac), that are strongly associated with enhancer activity [29–32]. These histone modifications can be mapped across the genome in order to identify regions, termed peaks, that show enriched levels of each modification. Several recent studies have compared levels of H3K27ac and other enhancer-associated histone modifications across multiple tissues in multiple mammalian species. These studies identified changes in histone modification levels at thousands of enhancers across species, suggesting abundant evolutionary turnover in enhancer activity across mammals [33–35].

This comparative approach has also been used to identify changes in enhancer activity relevant to human evolution. Two studies in developing human, rhesus macaque, and mouse cortex [36] and limb [31] identified Human Gain Enhancers (HGEs), defined as putative enhancers with higher levels of enhancer-associated histone modifications in human compared to the other two species. Subsequent massively parallel reporter assays (MPRAs) have identified sequence changes within HGEs that drive differential activity between human and chimpanzee orthologs [18]. Chromatin interaction maps in mid-fetal human brain indicate that many HGEs target neurodevelopmental and neuronal genes [37], suggesting that altered HGE activity may have contributed to the evolution of uniquely human brain features.

Although previous studies have hypothesized that such changes in enhancer activity are due to genetic variation that altered transcription factor binding, CpG islands (CGIs) have also recently been found to contribute to enhancer activation. CGIs are genomic intervals with high GC-content and CpG dinucleotide frequency [38]. They are frequently unmethylated, unlike the majority of CpG dinucleotides in the genome, and are associated with 70% of annotated promoters [39]. However, CGIs are also located in intronic and intergenic regions. These “orphan CGIs (oCGIs)” are often located within enhancers that show higher levels of histone modifications, transcription factor binding, and three-dimensional interactions with other genomic regions than non-oCGI containing enhancers [40, 41].

Several mechanisms link CGIs with transcription-factor independent recruitment of chromatin modifiers. First, unmethylated CpG dinucleotides recruit proteins containing ZF-CxxC finger domains, which in turn recruit histone methyltransferase complexes that deposit H3K4me3 [42]. These include CFP1, a subunit of the SET1A/B histone methyltransferase complexes [43], and MLL2, a member of the MLL2 complex [44, 45]. The presence of H3K4me3 promotes open, active chromatin by several mechanisms, including recruitment of histone acetylases, exclusion of factors that deposit repressive histone modifications, recruitment of chromatin remodelers, exclusion of DNA methylation, and direct recruitment of the transcriptional machinery [42]. The recruitment

of H3K4me3-depositing enzymes and consequently all of these downstream effects are mediated by the CpG dinucleotides within CGIs and therefore could be independent of transcription factor binding events.

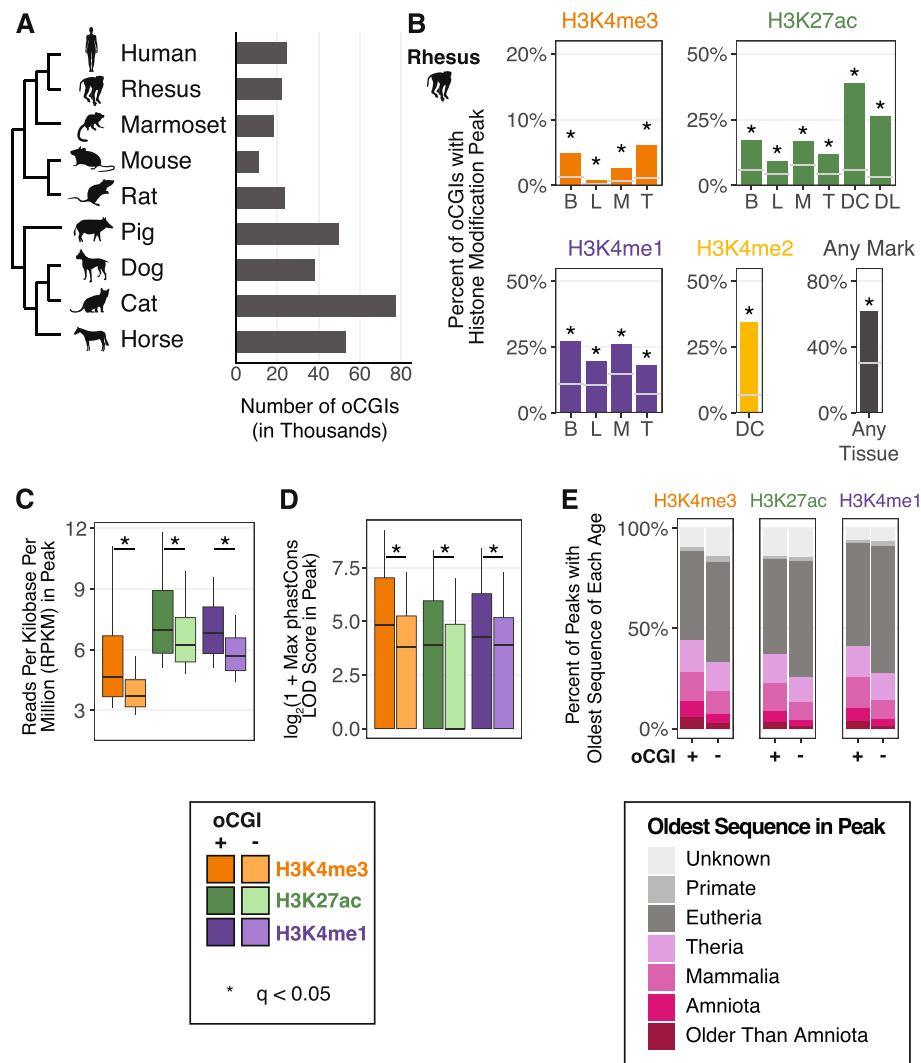
Given the abundance of evidence linking oCGIs to enhancer activity, we chose to examine oCGI turnover as a potential mechanism driving differences in enhancer activity across mammals. We identified oCGIs in nine mammalian genomes and integrated these maps with changes in the deposition of several activity-associated histone modifications. We first found that oCGIs are significantly enriched for H3K27ac, H3K4me3, H3K4me2, and H3K4me1 peaks in multiple tissues and species [31, 34, 36]. We also found extensive turnover of oCGIs across species and that species-specific oCGIs in putative enhancers were associated with species-specific increases in enhancer-associated histone modification levels. Orphan CGI turnover was associated with changes in enhancer activity even within ancient, highly conserved enhancers.

Our findings also support that oCGIs contribute to the increased activity of many HGEs. In light of these results, we selected one HGE to study *in vivo*, which has both an oCGI and prior evidence of increased regulatory activity in human compared to mouse. We generated a humanized mouse line for this locus and found that the oCGI-containing human enhancer exhibits increased H3K27ac and H3K4me3 in the developing mouse brain, although we did not observe associated changes in gene expression. However, using published data, we found that species-specific oCGIs with species-specific histone modifications were generally associated with increased expression of nearby genes. Finally, we found that turnover in transcription factor (TF) binding co-occurred at enhancers with oCGI turnover. We propose that oCGI turnover contributes to enhancer evolution both directly via altering recruitment of chromatin modifiers, and also indirectly by generating permissive chromatin that reveals TFBSs which may then be functionally integrated into nascent or existing enhancers. Our findings identify oCGI turnover as a novel class of sequence change, in addition to copy number changes and studies of evolutionary acceleration focused on nucleotide substitutions, that may be leveraged in comparative studies to identify gene regulatory innovations in mammalian genomes.

## Results

### **oCGIs are significantly enriched for enhancer-associated histone modifications**

We first identified genome-wide CGIs in nine mammalian genomes for which histone modification maps were previously generated (Fig. 1A). We used the canonical definition of CGIs: length  $\geq 200$  bp, GC content  $\geq 50\%$ , and observed/expected CpG dinucleotides  $\geq 0.6$  [38]. We restricted our analysis to oCGIs by excluding CGIs overlapping exons or located within regions 2 kb upstream of a transcription start site, as annotated by RefSeq in every genome [46] (Additional file 1: Fig. S1). For human and mouse, whose genomes have been more extensively annotated, we also removed CGIs overlapping ENCODE blacklist regions and regions 2 kb upstream of transcription start sites identified by the FANTOM Consortium [47, 48]. To further ensure that the identified oCGIs in each mammal were not part of unannotated exons or promoters, we only considered CGIs in each genome that had orthologous sequence in human that did not overlap human RefSeq, ENCODE blacklist, and FANTOM annotations



**Fig. 1** oCGIs are enriched for enhancer-associated histone modifications. **A** The number of oCGIs identified in nine mammalian genomes considered in this study. **B** Percent of oCGIs overlapping a histone modification peak for each indicated histone modification and tissue in rhesus macaque [31, 34, 36]. B = adult brain, L = adult liver, M = adult muscle, T = adult testis, DC = developing cortex, DL = developing limb. Gray horizontal lines indicate the expected overlap and stars indicate significant enrichment ( $q < 0.05$ , BH-corrected, determined by permutation test; see "Methods"). **C** The level of each indicated histone modification in peaks with and without oCGIs, measured in reads per kilobase per million (RPKM). Box plots show the interquartile range and median, and whiskers indicate the 90% confidence interval. Stars indicate a significant difference between peaks with and without an oCGI ( $q < 0.05$ , Wilcoxon rank-sum test, BH-corrected). **D** Maximum phastCons LOD (log-odds) scores in peaks with and without oCGIs. Box plots show the interquartile range and median, and whiskers indicate the 90% confidence interval. Stars indicate a significant difference between peaks with and without an oCGI ( $q < 0.05$ , Wilcoxon rank-sum test, BH-corrected). **E** Evolutionary origin of peaks with and without oCGIs. Bar plots show the percentage of peaks with and without oCGIs whose oldest sequence belongs to each age category. The results shown in panels (C) through (E) were generated using peaks from adult brain in rhesus macaque; see Additional File 1: Figs. S8, S10, and S12 for results from additional species and tissues

(Additional file 1: Fig. S1). After filtering, we found thousands of oCGIs in each mammalian genome, ranging from 11,067 in mouse to 77,199 in cat (Fig. 1A).

We then compared the oCGI sets we used in this study to oCGIs identified using other approaches. Our oCGIs contain thousands more sites than the standard CpG island annotations present in the UCSC Genome Browser, which imposes additional requirements beyond the canonical definition to retain only the strongest, most CpG-dense sites [49] (Additional file 1: Fig. S2A, *left*). Other approaches have been developed to identify CpG islands using probabilistic models rather than strict cut-offs, which can be more sensitive at identifying CpG islands, especially in non-vertebrate genomes [50, 51]. These model-based approaches find many more oCGIs than the UCSC annotations (Additional file 1: Fig. S2A, *middle*) and include several thousand sites we did not consider in this study (Additional file 1: Fig. S2A, *right*). However, most of the sites present in the model-based oCGI sets overlap repeat-masked regions of the genome, which we intentionally excluded from our study based on the difficulty of identifying orthologs for such sequences across species (Additional file 1: Fig. S2B). When these masked sites are removed from the comparison, our oCGI sets identify several thousand oCGIs in the non-repeat genome that are not identified using model-based methods (Additional file 1: Fig. S2C).

Next we sought to interrogate the relationship between these oCGI sets and enhancer activity as measured by local enrichment for specific histone modifications. The association between H3K27ac enrichment and enhancer activity has been previously established [29–32]. We sought to evaluate whether the other histone modifications included in this study (H3K4me3, H3K4me2, and H3K4me1) were also associated with enhancer activity. Furthermore, given that CGIs are associated with CpG-mediated recruitment of histone modifiers, we tested whether all the marks in our study were predictive of enhancer activity for putative enhancers containing oCGIs. We took two orthogonal approaches to assess enhancer activity: (1) activity in a LacZ reporter assay in developing mouse embryos as assessed by the VISTA Enhancer Browser [52], and (2) interactions with target genes as measured by promoter capture HiC in developing mouse liver [53]. We collected putative enhancers (for the first approach, all tested VISTA elements (Additional file 1: Fig. S3A); for the second approach, all ENCODE cCREs categorized as “distal enhancer-like sequences” (Additional file 1: Fig S4A). We then sorted these based on their overlap with histone modification peaks in the matching tissue from ENCODE [30] and their overlap with mouse oCGIs. All histone modifications in the analysis were predictive of transgenic enhancer activity for VISTA elements not containing an oCGI (Fisher’s exact test, Benjamini Hochberg (BH)-corrected; Additional file 1: Fig. S3B (*left*), Additional file 2: Table S1). For VISTA elements that do contain an oCGI, those overlapping a histone modification peak were also more likely to show transgenic reporter activity for H3K27ac, H3K4me2, and H3K4me1, although the sample sizes in this analysis were smaller and did not always reach significance for H3K4me2 and H3K4me1 (Fisher’s exact test, BH-corrected; Additional file 1: Fig. S3B (*right*), Additional file 2: Table S1). We also found that all histone modifications in the analysis were significantly associated with enhancer-promoter interactions for putative enhancers with and without oCGIs (Wilcoxon rank-sum test, BH-corrected; Additional file 1: Fig.

S4B). These results support that the histone modifications included in this study are predictive of enhancer function, both for enhancers with and without oCGIs.

Having established that these histone modifications are predictive of enhancer activity when found at oCGIs, we next quantified the overlap of oCGIs in each species with published maps of each histone modification in multiple tissues [31, 34, 36]. For all datasets, oCGIs were enriched for peaks from all histone modifications (permutation test, see “[Methods](#)”; Fig. 1B, Additional file 1: Fig. S5). In all species, a large percentage (ranging from 36.7% in cat to 78.5% in mouse) of oCGIs overlap histone modification peaks in at least one dataset (“Any Mark” and “Any Tissue” in Fig. 1B, Additional file 1: Fig. S5). This result supports that oCGIs are strongly associated with histone modifications indicative of enhancer activity. To further characterize our oCGIs, we intersected them with candidate cis-regulatory elements (cCREs) annotated by the ENCODE consortium in mouse and human [29]. The majority of cCREs with oCGIs belong to the “distal enhancer-like sequence” category, which are distal to transcription start sites and characterized by DNase accessibility and H3K27ac, and which may also be marked by H3K4me3 (Additional file 1: Fig. S6). The enhancer definitions used by ENCODE reflect growing evidence that H3K4me3 is often found at enhancers, particularly strongly active ones [54–56].

To contextualize the contribution of oCGI-mediated regulatory mechanisms to enhancers, we sought to determine the proportion of putative enhancers that included an oCGI. We calculated the percentage of all histone modification peaks in each species and tissue that contained an oCGI. H3K4me3 peaks showed the highest overlap with oCGIs (in rhesus macaque, ranging from 21.9% of regions in testis to 49.5% in liver), consistent with the role of oCGIs in directly recruiting factors involved in H3K4me3 deposition (Additional file 1: Fig. S7) [42]. Although we found that H3K27ac, H3K4me1, and H3K4me2 peaks were less frequently associated with oCGIs, we still identified thousands of regions for each modification that contained oCGIs, revealing putative enhancers with possible oCGI-dependent functions (Additional file 1: Fig. S7).

We next examined the association between oCGIs and the activity and evolutionary constraint of putative enhancers. We found that histone modification peaks containing an oCGI exhibited higher levels of each modification than peaks without an oCGI (Fig. 1C, Additional file 1: Fig. S8), suggesting that oCGI-containing enhancers may show stronger activity than enhancers lacking an oCGI. Histone modification peaks with an oCGI were also longer than peaks without an oCGI, consistent with an oCGI-associated increased recruitment of histone modifiers across a broader region (Additional file 1: Fig. S9). Peaks with an oCGI were more constrained, as measured using phastCons [57], compared to peaks lacking oCGIs. For example, in adult rhesus macaque brain, H3K4me3 peaks with an oCGI had higher maximum phastCons LOD (log-odds) scores compared to H3K4me3 peaks without an oCGI (log<sub>2</sub>-transformed median of 4.81 versus 3.81; log<sub>2</sub>-transformed upper quartile of 7.03 versus 5.25; significance determined by Wilcoxon rank-sum test, BH-corrected, see “[Methods](#)”) (Fig. 1D, Additional file 1: Fig. S10). Peaks with an oCGI also had higher aggregate phastCons LOD scores and more bases included within constrained regions defined by phastCons (Additional file 1: Fig. S11). We obtained further support for this finding using an age segmentation map of the human genome, in which the evolutionary origin of a human genomic region is inferred

based on the most distantly related species in which an orthologous sequence can be identified [58]. Peaks containing oCGIs were more likely to include ancient sequences compared to peaks lacking oCGIs (Fig. 1E). These trends generalized across all marks tested and across all species in our analysis (Additional file 1: Fig. S12), supporting that oCGIs are components of ancient, constrained, and active enhancers.

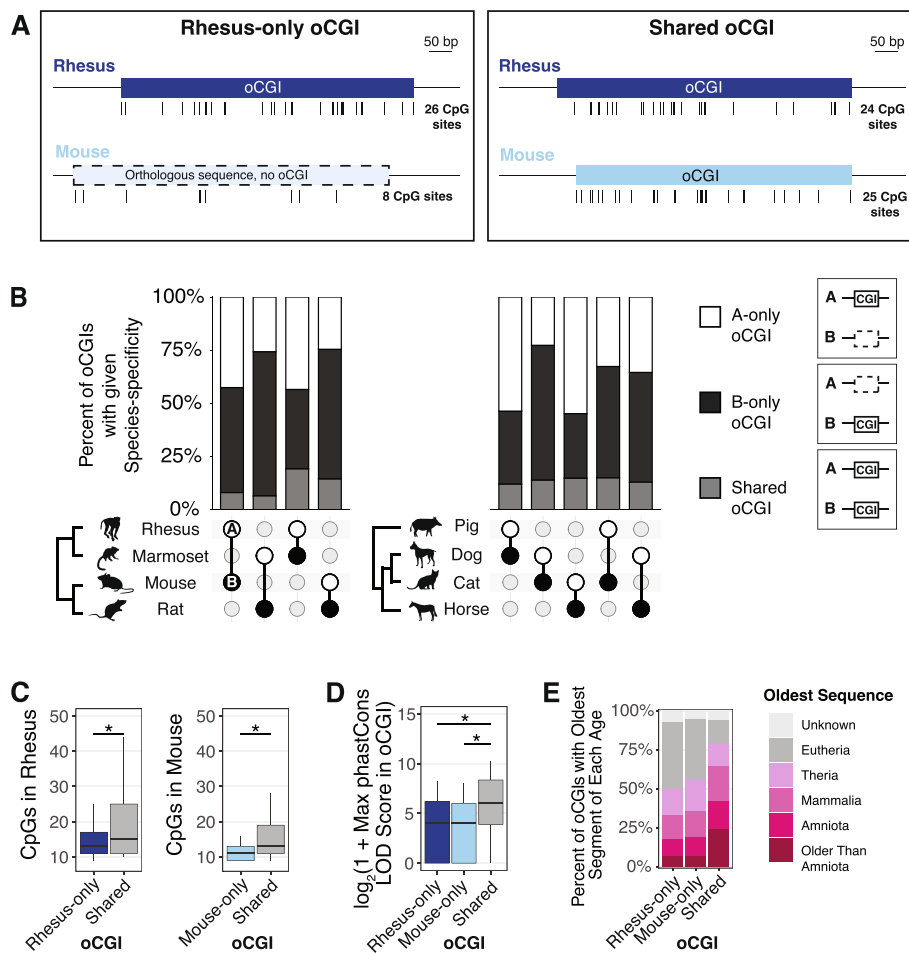
### Extensive turnover of oCGIs across mammalian species

To examine whether oCGI turnover is a widespread mechanism contributing to evolutionary changes in enhancer activity, we first asked how conserved oCGIs are by comparing species pairs within our dataset. For each species pair, we identified all oCGIs located within orthologous sequences. We then classified each oCGI as called in only one species (defined as species-specific within the scope of each pairwise comparison and labeled as “A-only” or “B-only” throughout the subsequent figures) or called in both species in the pairwise comparison (labeled as “shared”) (Fig. 2A, Additional file 1: Fig. S13). We found that, in every species pair we considered, the majority of oCGIs were present in only one species in the pair (Fig. 2B, see Additional file 1: Fig. S14 and Additional file 2: Table S2 for all species pairs). This was true both in closely related species (for example see rhesus macaque (species A) versus marmoset (species B) in Fig. 2B) and more distantly related species (for example see rat versus horse).

We next asked whether species-specific and shared oCGIs differ in their composition and their constraint. As a representative example, we show a comparison of rhesus macaque (species A) versus mouse (species B) in Fig. 2C–E. We found that shared oCGIs had more CpG dinucleotides than rhesus-specific (A-only) or mouse-specific (B-only) oCGIs (Fig. 2C, see Additional file 1: Fig. S15 for all species pairs), were longer (Additional file 1: Fig. S16), and generally, but not in all cases, had a higher ratio of observed/expected CpG dinucleotides (Additional file 1: Fig. S17). Shared oCGIs were more constrained, in that they were more likely to contain a higher scoring phastCons element (Fig. 2D, see Additional file 1: Fig. S18 for all species pairs). Shared oCGIs were more likely to have higher aggregate phastCons LOD scores and a greater percentage of bases covered by a phastCons element (Additional file 1: Fig. S19). Shared oCGIs were also more likely to be located within more ancient sequences compared to species-specific oCGIs (Fig. 2E, see Additional file 1: Fig. S20 for all species pairs).

Although shared oCGIs generally showed evidence of higher constraint, a substantial proportion of species-specific oCGIs did overlap phastCons elements, evidence that they are also located in sequences under constraint. In the rhesus macaque versus mouse comparison, 67.0% of rhesus-specific (A-only) oCGIs and 77.6% of mouse-specific (B-only) oCGIs overlapped a phastCons element, compared to 99.2% of shared oCGIs. Additionally, many species-specific oCGIs were located in ancient enhancers: in the rhesus versus mouse comparison, 17.4% of rhesus-specific (A-only) oCGIs & 18.7% of mouse-specific (B-only) oCGIs were located in enhancers conserved among Amniotes or older clades, compared to 41.9% of shared oCGIs (Fig. 2E). These results support that oCGI turnover has occurred even within ancient, highly constrained enhancers.

In order to obtain insight into the origin of species-specific oCGIs, we performed an analysis using species pairs and an outgroup to polarize species-specific oCGIs as gain or loss events (Additional file 1: Fig. S21). For example, we identified all oCGIs gained in



**Fig. 2** oCGIs show extensive turnover across species. **A** Schematic illustrating how we defined species-specific oCGIs in pairwise comparisons, using rhesus macaque and mouse as an example. *Left*: a rhesus-only oCGI (the sequence is present in both rhesus and mouse, but the oCGI is only present in rhesus). *Right*: a shared oCGI (both the sequence and oCGI are present in both rhesus and mouse). Ticks under each oCGI represent the locations of CpG dinucleotides. **B** Percent of oCGIs across the indicated species pairs (species A versus species B) that are “A-only,” “B-only,” or “shared” as described in the main text. The species pair is shown under each bar, with species A denoted by a white circle and species B denoted by a black circle. Percentages of oCGIs that are species A-only (white), species B-only (black), or shared (gray) are shown. **C** Number of CpG dinucleotides in rhesus-only (dark blue) or mouse-only (light blue) oCGIs compared to shared (gray) oCGIs. Box plots show the interquartile range and median, and whiskers indicate the 90% confidence interval. Stars indicate significant differences ( $q < 0.05$  Wilcoxon rank-sum test, BH-corrected). **D** Maximum phastCons LOD scores in rhesus-only, mouse-only, and shared oCGIs. Box plots show the interquartile range and median, and whiskers indicate the 90% confidence interval. Stars indicate significant differences ( $q < 0.05$ , Wilcoxon rank-sum test, BH-corrected). **E** Evolutionary origins of rhesus-only, mouse-only, and shared oCGI sequences

humans (present in human, absent in the sister group rhesus macaque and the outgroup marmoset) or lost in humans (absent in human, present in rhesus macaque and marmoset) (Additional file 1: Fig. S21A). In general, we identified many more gain events than loss events, although we identified more loss events in sites that overlapped a phastCons element compared to sites that did not (Additional file 1: Fig. S21B). The bias towards gains is likely because they are easier to detect, since the oCGI must be present in only one species, whereas there must be an oCGI present in two species to detect a loss event.

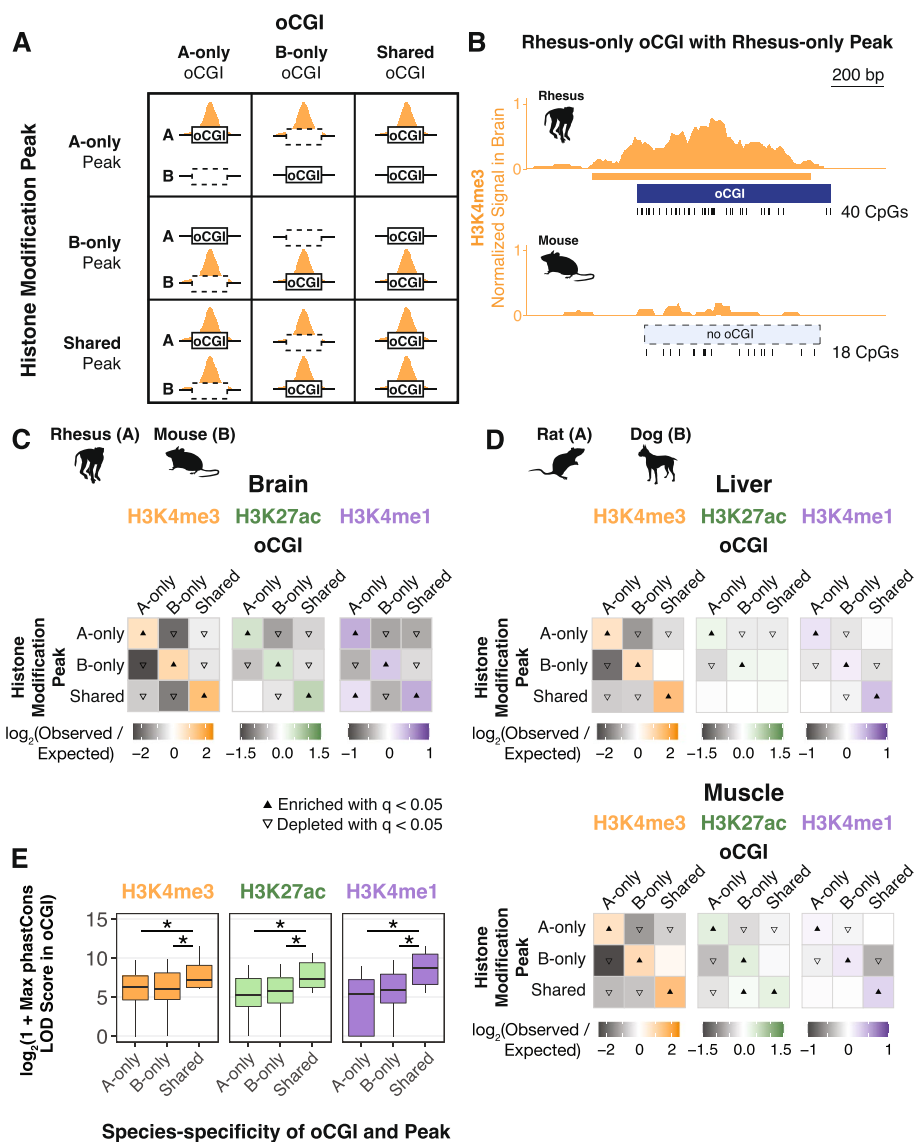


### Species-specific oCGIs are significantly enriched for species-specific histone modification peaks

In light of the relationship we identified between oCGIs and enhancer activity-associated histone modifications, coupled with the extensive turnover of oCGIs we observed, we next examined whether species differences in oCGIs were associated with species differences in histone modification levels. Using pairwise species comparisons as described above, we sorted oCGIs based on their species specificity and the species specificity of co-localized histone modification peaks, performing a separate analysis for each species pair, tissue, and modification. Using a permutation test (Fig. 3A,B, Additional file 1: Fig. S22, “Methods”), we found that species-specific oCGIs were significantly enriched for species-specific histone modification peaks; three representative pairwise comparisons are shown in Fig. 3C, D. We also found that oCGIs specific to one species in the pair were depleted for histone modification peaks specific to the other species, and for peaks that were shared between both species. Shared oCGIs present in both species were enriched for shared histone modification peaks and were depleted for species-specific peaks. We observed the greatest enrichment of species-specific oCGIs for species-specific H3K4me3 peaks, but we observed significant enrichment for species-specific H3K27ac and H3K4me1 peaks as well (Fig. 3C, D). This trend was consistent across all species pairs and tissues, as demonstrated by the representative examples shown in Figs. 3C and 3D and by all 28 species pairwise comparisons in adult brain, liver, muscle, and testis for H3K4me3 (Additional file 1: Fig. S23-24), H3K27ac (Additional file 1: Fig. S25-26), and H3K4me1 (Additional file 1: Fig. S27-28) (Additional file 2: Table S3). We also performed a peak-centric analysis (as opposed to the oCGI-centric analysis described above) in order to evaluate whether the patterns we observed were consistent between the two approaches (Methods). In this reciprocal analysis, we found that

(See figure on next page.)

**Fig. 3** Species-specific oCGIs are significantly enriched for species-specific histone modification peaks. **A** Schematic illustrating how we defined species-specific and shared oCGIs and peaks. In each pairwise species comparison for each histone modification and tissue, we sorted oCGIs based on their species specificity (designated as A-only, B-only, or shared as in Fig. 2) and the species specificity of their histone modification peaks (shown in orange in the schematic). **B** An example of a rhesus macaque-specific oCGI overlapping a rhesus-specific H3K4me3 peak in a pairwise comparison with mouse. Ticks show the location of CpG dinucleotides. Normalized H3K4me3 signal at this locus (orange) is shown as read counts per million in adjacent 10-bp bins. **C** Enrichment and depletion in each indicated comparison of species-specific and shared oCGIs (*top*: A-only, B-only, Shared) and species-specific and shared peaks (*left*: A-only, B-only, Shared), compared to a null expectation of no association between oCGI turnover and peak turnover. Each 3 × 3 grid shows the results for a specific test examining oCGIs and their overlap with three histone modifications in adult rhesus macaque brain. Each box in each grid is colored according to the level of enrichment over expectation (orange for H3K4me3, green for H3K27ac, or purple for H3K4me1) or depletion (gray for all marks) of genome-wide sites that meet the criteria for that box. The color bar below each plot illustrates the level of enrichment or depletion over expectation. Filled upward-pointing triangles denote significant enrichment and open downward-pointing triangles denote significant depletion ( $q < 0.05$ , permutation test, BH-corrected, see Additional file 1: Fig. S22 and “Methods”). **D** Enrichment and depletion in an additional species comparison, rat versus dog, and in additional tissues (liver, *top*, and muscle, *bottom*), shown as described in (C). **E** Maximum LOD score in species-specific oCGIs in species-specific peaks and shared oCGIs in shared peaks, using data from adult rhesus macaque brain. Box plots show the interquartile range and median, and whiskers indicate the 90% confidence interval. Stars indicate significance ( $q < 0.05$ , Wilcoxon rank-sum test, BH-corrected)



**Fig. 3** (See legend on previous page.)

species-specific peaks were enriched for species-specific oCGIs (Additional file 1: Fig. S29, Additional file 2: Table S4), supporting the patterns observed in the oCGI-centric analysis.

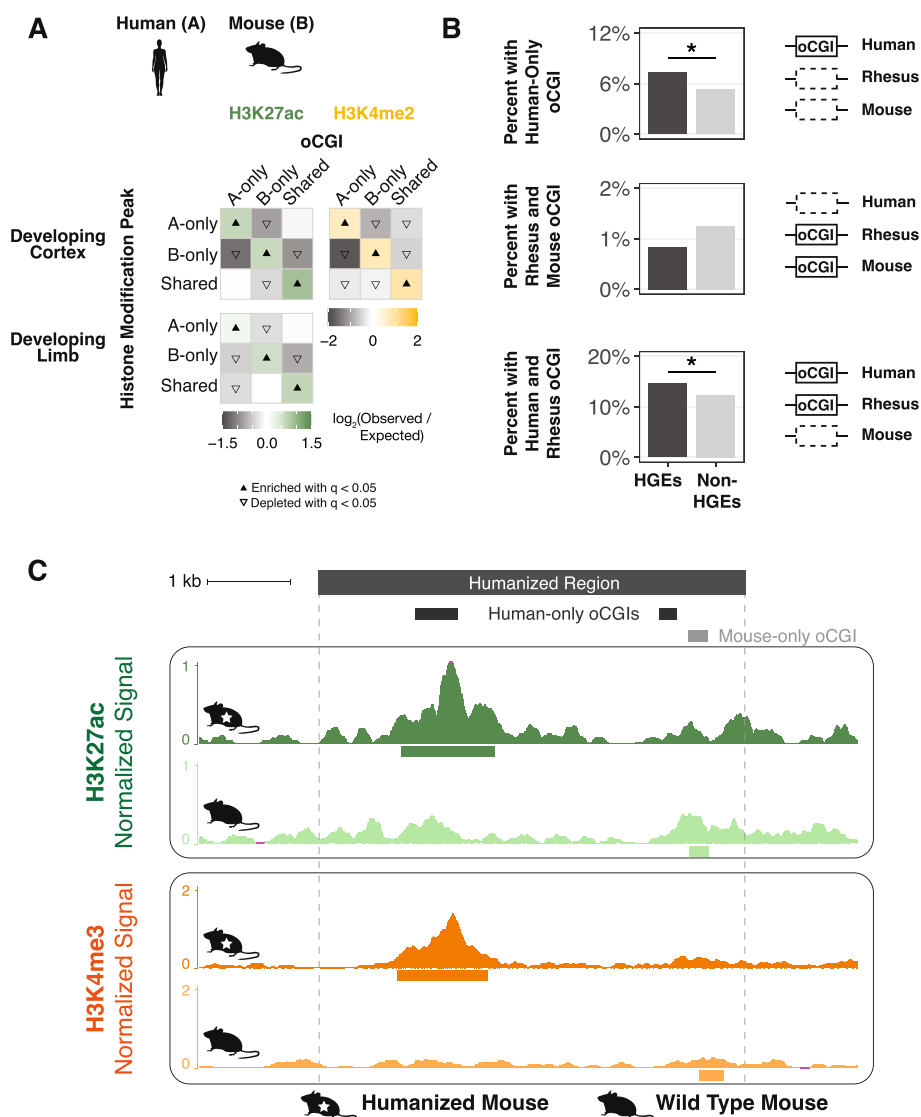
Our approach for identifying species-specific oCGIs involves hard thresholds built into the CGI definition. For example, if an oCGI in species A has a ratio of observed/expected CpG dinucleotides of 0.6, then the site is called as species-specific if the orthologous site in species B has a ratio of less than 0.6. However, this ratio in species B could be anywhere from 0.59 to much lower. Therefore, we examined the effect of imposing successively larger minimal difference requirements for the ratio of observed/expected CpG dinucleotides between each species (Additional file 1: Fig. S30). We found that imposing larger minimal difference thresholds reduces the number of species-specific oCGIs that we identify within each species pair (Additional file 1: Fig. S30A), but also

results in stronger enrichment and depletion results than what we observed using our original criteria (Additional file 1: Fig. S30B).

We next examined whether species-specific oCGIs associated with species-specific changes in histone modification marking were under evolutionary constraint. Overall, shared oCGIs in shared peaks exhibited higher levels of constraint as measured by maximum phastCons LOD scores (Fig. 3E, see Additional file 1: Fig. S31 for more species pairs). Shared oCGIs in shared peaks also more frequently overlapped sequences of greater evolutionary age (Additional file 1: Fig. S32) than species-specific oCGIs within species-specific peaks. However, we also found that a substantial fraction of species-specific oCGIs located within species-specific peaks were under constraint (Fig. 3E). For example, in a comparison of species-specific rhesus macaque and mouse oCGIs and H3K27ac peaks in adult brain, 75.5% of rhesus-specific oCGIs in rhesus-specific peaks and 83.7% of mouse-specific oCGIs in mouse-specific peaks overlapped a constrained region annotated by phastCons, compared to 94.7% of shared oCGIs located in shared peaks. In this same comparison, we also found that 20.9% of rhesus-specific oCGIs in rhesus-specific peaks and 18.2% of mouse-specific oCGIs in mouse-specific peaks overlapped with a sequence conserved within Amniota or more ancient clades, compared to 29.5% of shared active oCGIs in shared peaks, suggesting that many species-specific oCGIs associated with species-specific peaks are components of ancient, constrained enhancers.

In order to assess whether the enrichment of species-specific oCGIs for species-specific peaks we observed in adult tissues was consistent in other contexts, we performed the same analysis on two datasets generated by independent studies of the human, rhesus macaque, and mouse developing cortex [36] and developing limb [31]. For the histone modifications H3K27ac and H3K4me2, and at four developmental time points (Additional file 2: Table S5), we found that species-specific oCGIs were enriched for species-specific histone modification peaks and shared oCGIs were enriched for shared peaks, consistent with our previous findings in adult tissues (Fig. 4A, see Additional file 1: Fig. S33 for all species pairs and time points).

Both of these developmental studies also identified Human Gain Enhancers (HGEs) based on increased histone modification levels in human cortex or limb compared to rhesus macaque and mouse [31, 36]. HGEs have been shown to exhibit human-specific changes in their activity and have been implicated in the regulation of neurodevelopmental genes [37], suggesting they may have contributed to human cortical evolution. Therefore, we examined whether oCGI turnover may have contributed to HGE evolution. We classified HGEs based on whether they had an oCGI in human, rhesus, and mouse, and compared them to histone modification level-matched non-HGE enhancers using a resampling test (Additional file 1: Fig. S34, “[Methods](#)”). We found that cortex HGEs were significantly more likely to contain a human-only oCGI or an oCGI shared between human and rhesus, but absent in mouse (resampling test, BH-corrected; Fig. 4B, Additional file 1: Fig. S35, Additional file 2: Table S6). This finding suggests that changes in oCGI content may have contributed to the increased activity of HGEs in the human cortex. Consistent with this hypothesis, cortex HGEs were depleted for mouse-only oCGIs and oCGIs shared between rhesus and mouse, but absent in human. Limb HGEs were depleted across most oCGI categories and were enriched for having no oCGI



**Fig. 4** Association of species-specific oCGIs with species-specific histone modification peaks and HGEs in the developing human cortex and limb. **A** Enrichment and depletion in each indicated comparison of species-specific and shared oCGIs (top: A-only, B-only, Shared) and species-specific and shared peaks (left: A-only, B-only, Shared), compared to a null expectation of no association between oCGI turnover and peak turnover. Results are shown as in Fig. 3C,D, with enrichment in green for H3K27ac and yellow for H3K4me2, and depletion in gray. One representative comparison is shown for developing cortex (8.5 post-conception weeks (p.c.w.) in human versus embryonic day 14.5 in mouse) and developing limb (embryonic day 41 in human versus embryonic day 12.5 in mouse). **B** Enrichment of specific oCGI species patterns in HGEs compared to non-HGE enhancers in human cortex at 8.5 p.c.w. Bar plots show the percentage of HGEs (left bar) or non-HGE enhancers (right bar) that overlap an oCGI with the species pattern shown on the left. Significance was determined using a resampling test comparing HGEs to non-HGE human enhancers matched for overall histone modification levels (resampling test, BH-corrected; see Additional file 1: Fig. S34 and “Methods”). **C** H3K27ac levels in developing diencephalon at the humanized *hs754* (top tracks) or wild type (bottom tracks) mouse locus at E11.5. Locations of oCGIs within *hs754* and its mouse ortholog are shown at the top (dark gray boxes for two human oCGIs not present in the mouse sequence, and a light gray box for a mouse oCGI not present in the human sequence). Dark green (humanized) and light green (wild type) signal tracks show normalized H3K27ac levels as counts per million reads calculated in adjacent 10-bp bins. Dark orange (humanized) and light orange (wild type) signal tracks show normalized H3K4me3 levels. Peak calls are shown as boxes below the signal tracks. Nominal  $p$ -values were obtained by DESeq2 using a Wald test, then BH-corrected for multiple testing across all peaks genome-wide to generate  $q$ -values (see values in main text and in Additional file 1: Fig. S38)

at all, suggesting that oCGIs may be less relevant to enhancer activity in the developing limb than the developing cortex (Additional file 1: Fig. S35).

#### **Changes in enhancer oCGI content are associated with changes in histone modification levels in a humanized mouse model**

We next sought to study the effect of species differences in oCGI content on enhancer activity using an experimental approach in an in vivo system. We selected one candidate HGE, named hs754 [36]. This enhancer is highly constrained and includes sequences that are inferred to have originated in the stem lineage of jawed vertebrates (Additional file 1: Fig. S36, *bottom*). The human sequence is marked by an H3K27ac peak in the developing human cortex that was not present in the rhesus macaque or mouse cortex. The sequence underlying the human peak contains a 608-bp oCGI that has orthologous sequences in both rhesus and mouse. However, the rhesus orthologous sequence contains a smaller, 306-bp oCGI and the mouse orthologous sequence contains no oCGI (Additional file 1: Fig. S36). In order to determine whether this enhancer is associated with changes in enhancer activity during development, we generated a mouse model in which a 5.5-kb human sequence containing hs754 replaced the mouse orthologous locus via CRISPR/Cas9-mediated homology-directed repair in C57BL/6 mouse embryos (“Methods,” Additional file 1: Fig. S37, Additional file 2: Table S7-8).

We first examined whether the humanized sequence containing an oCGI showed increased levels of H3K27ac and H3K4me3 compared to the mouse sequence. We focused on developing diencephalon, based on a reported regulatory interaction between hs754 and the gene *Irx2* [59], which is expressed during diencephalon development [60]. We included both an early (embryonic day 11.5, E11.5) and late (embryonic day 17.5, E17.5) time point.

The humanized locus showed high levels of H3K27ac overlapping the human oCGI at both E11.5 and E17.5 (Fig. 4C, Additional file 1: Fig. S38A-B). At both time points, the difference in H3K27ac level between the wild type and humanized orthologs was nominally statistically significant as measured using DESeq2, but did not reach significance after genome-wide multiple-testing correction (E11.5:  $p < 1.10 \times 10^{-5}$ ,  $q =$  not significant; E17.5:  $p < 9.83 \times 10^{-4}$ ,  $q =$  not significant; Wald test, BH-corrected). We also found a strong H3K4me3 peak overlapping the human oCGI at both time points, in contrast to the near absence of H3K4me3 in the wild type mouse (Fig. 4C, Additional file 1: Fig. S38C-D). These H3K4me3 peaks at both time points did reach significance after genome-wide correction (E11.5:  $p < 6.14 \times 10^{-20}$ ,  $q < 4.23 \times 10^{-16}$ ; E17.5:  $p < 1.66 \times 10^{-48}$ ,  $q < 1.48 \times 10^{-44}$ ; Wald test, BH-corrected). The full differential analysis results for all genome-wide peaks are shown in Additional file 1: Figure S39 and Additional file 2: Table S9. The overlap of the human oCGI with increased H3K27ac and H3K4me3 levels is consistent with the role of oCGIs in mediating changes in chromatin activity, and consistent with our genome-wide results across species pairs.

In order to determine whether these changes in enhancer activity were associated with downstream effects on gene expression, we carried out RNA-seq on E11.5 and E17.5 diencephalon (4 replicates per genotype at each time point). Using DESeq2, we found only two differentially expressed genes reaching genome-wide significance, both at embryonic day E11.5: the gene *Ppp3cc* which is on a different chromosome than hs754

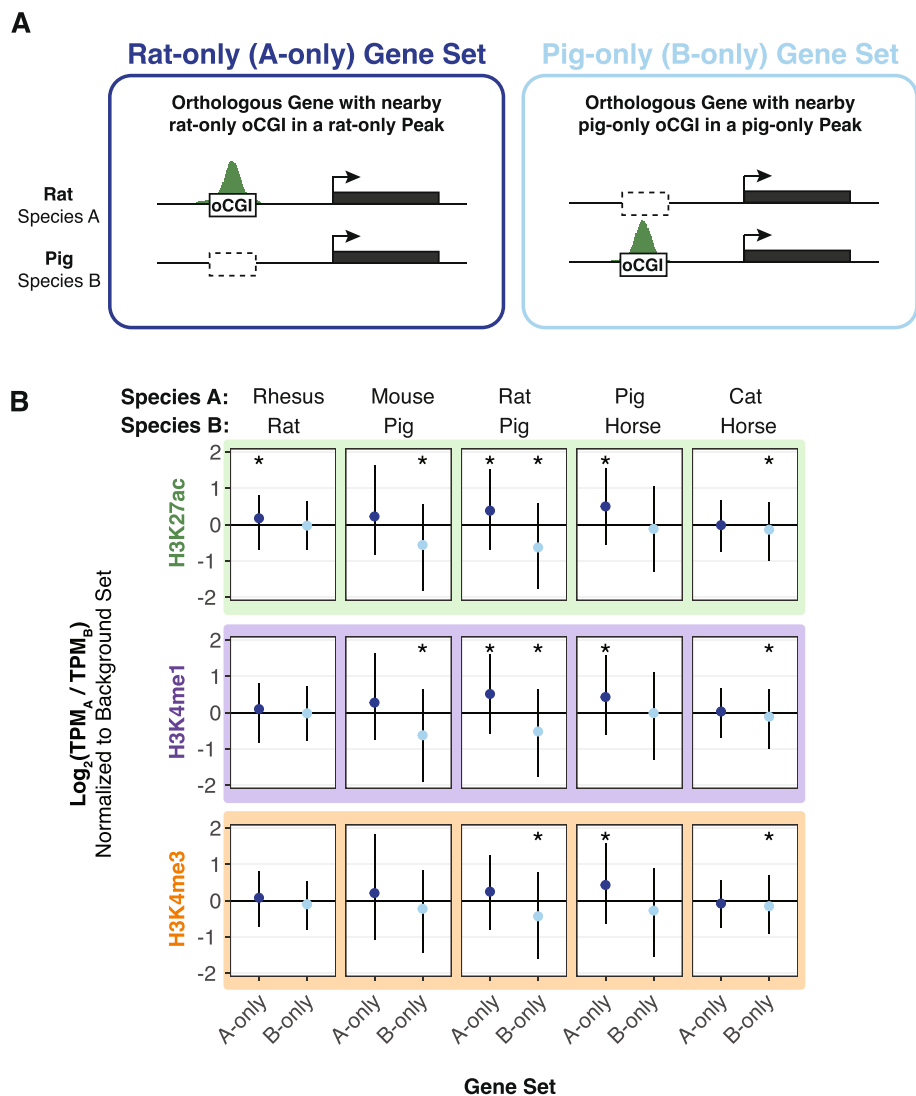
( $q < 1.98 \times 10^{-7}$ ; Wald test, BH-corrected), and the gene *Serinc5* which is 20 Mb away from *hs754* on the same chromosome ( $q < 4.37 \times 10^{-2}$ ; Wald test, BH-corrected) (Additional file 1: Fig. S40, Additional file 2: Table S10). As reported *cis*-regulatory interactions over distances greater than several Mb are rare, as are *trans*-chromosome interactions [61], we consider it unlikely that either gene is a direct target of *hs754*.

Finally, we sought to understand whether the changes in histone modification levels at *hs754* could stem not only from changes in oCGI content, but also from differences in the recruitment of transcription factors. There may be changes to TFBSs, either involving the oCGI or independent of it, in the human ortholog compared to the rhesus and mouse orthologs. To assess this possibility, we performed ChIP-seq for the factor CTCF, which has several predicted motifs in *hs754* and is involved in chromatin looping between enhancers and promoters [62–64]. We identified two additional CTCF binding events in the *hs754* humanized locus (Additional file 1: Fig. S41). These changes in CTCF binding suggest that the gain of putative enhancer activity in the humanized mouse model may be due to a combination of oCGI-mediated mechanisms and changes in TF binding. We will return to the implications of this finding for enhancer evolution in the “Discussion.”

#### **Enhancers exhibiting oCGI and histone modification peak turnover are associated with gene expression changes across species**

Although we did not observe changes in the expression of any potential target gene due to the oCGI-associated increase in enhancer activity in our *hs754* humanized mouse model, such increases may nevertheless be generally associated with changes in enhancer activity. To evaluate this hypothesis, we used transcriptome datasets generated in [34] to determine whether species-specific oCGIs were correlated with increased expression of potential target genes across species pairs. We first identified all species-specific oCGIs in species-specific peaks, performing a separate analysis for each species pair, histone modification, and tissue. We associated these sites with potential target genes based on proximity (“Methods”; Additional file 1: Fig. S42A) [65]. We then focused on the subset of genes that are annotated as 1:1 orthologs between each species pair by Ensembl.

Taking as an example the comparison between rat and pig using adult brain H3K27ac peaks shown in Fig. 5A, we sorted genes into a “rat-only set” which were associated with rat-only oCGIs in rat-only H3K27ac peaks, a “pig-only set” which are associated with pig-only oCGIs in pig-only H3K27ac peaks, and a “background set” not included in either category (Additional file 1: Fig. S42A). For each gene in each set, we compared expression as the ratio of TPM (transcripts per million) in the two species. We then compared the  $\log_2$ -transformed TPM ratios in the rat-only set and the pig-only set to the background set (resampling test matching gene expression level, BH-corrected; see “Methods” and Additional file 1: Fig. S42B-E). We found that genes associated with a rat-only active oCGI in a rat-only H3K27ac peak were more highly expressed in rat (median  $\log_2$ -transformed TPM ratio of 0.39), and genes associated with a pig-only active oCGI in a pig-only H3K27ac peak were more highly expressed in pig (median  $\log_2$ -transformed TPM ratio of -0.64). These results generalized across species pairs and tissues (Fig. 5, Additional file 1: Fig. S43-S45), and reached statistical significance most often for species-specific oCGIs in species-specific H3K27ac peaks (Additional



**Fig. 5** Species-specific oCGIs in species-specific peaks are associated with gene expression changes. **A** Schematic illustrating our method for assigning oCGIs and peaks to genes as described in the text and Additional file 1: Figure S42, using a pairwise comparison of rat and pig as an example. *Left*: A gene associated with a rat-only oCGI in a rat-only H3K27ac peak, which means the gene is assigned to the “rat-only set” (A-only set) of genes. *Right*: A gene associated with a pig-only oCGI in a pig-only H3K27ac peak, which means the gene is assigned to the “pig-only set” (B-only set) of genes. **B** The log<sub>2</sub>-transformed TPM ratio for genes in the A-only set and the B-only set for each indicated species pair and histone modification using data from adult brain. Points indicate median values for the A-only set (dark blue) and the B-only set (light blue) and lines indicate the interquartile range. All values in the A-only set and B-only set were normalized to the median TPM ratio across resampling rounds from the background set. Stars indicate a significant difference between the observed median and the expected median ( $q < 0.05$ , resampling test to compare to the background set, BH-corrected; see Additional file 1: Fig. S42 and “Methods”)

file 1: Fig. S43,  $q < 0.05$  in 62 of 120 tests), followed by those in species-specific H3K4me1 peaks (Additional file 1: Fig. S44,  $q < 0.05$  in 52 of 120 tests) and those in species-specific H3K4me3 peaks (Additional file 1: Fig. S45,  $q < 0.05$  in 24 of 120 tests). This finding suggests that oCGI turnover is not only associated with changes in enhancer activity, but also with changes in gene expression.

### Species-specific oCGIs are significantly enriched for species-specific transcription factor binding events

Orphan CGIs may contribute to enhancer activity by two mechanisms. The first is by CpG-mediated recruitment of ZF-CxxC domain proteins such as CFP1 and MLL2, leading to increases in H3K4me3 and subsequent downstream effects [42]. The second is via the TFBSs that they contain, which enable binding of TFs that then recruit coactivators. Many TFs have motifs with high GC content or that contain CpG dinucleotides and could therefore be components of oCGIs [39].

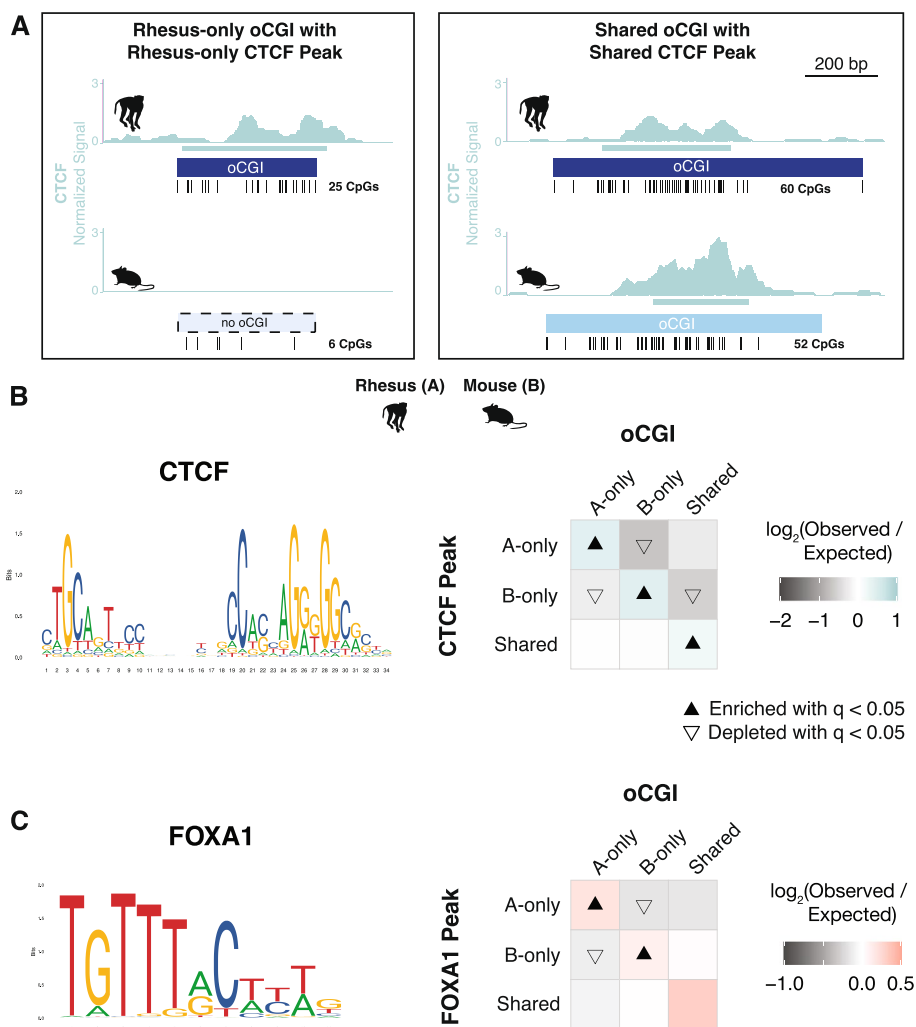
To examine this second mechanism, we assessed whether species-specific oCGIs were enriched for species-specific transcription factor binding events. Using previously generated genome-wide binding data for several transcription factors in adult liver [66–68], we classified oCGIs based on their species specificity and the species specificity of co-localized TF peaks (Fig. 6A). For CTCF, we found that species-specific oCGIs were enriched for species-specific CTCF peaks (Fig. 6B, see Additional file 1: Fig. S46A for all species pairs). Shared oCGIs were also enriched for shared CTCF peaks. This result is consistent with the sequence composition of the CTCF motif, which is GC-rich and contains a CpG dinucleotide, both features of oCGIs. We found similar enrichment patterns for an additional TF with a GC-rich motif (HNF4A, Additional file 1: Fig. S46B) and a TF with a motif containing one CpG dinucleotide (HNF6, Additional file 1: Fig. S46C), also consistent with oCGIs enabling binding of factors with GC-rich motifs and motifs containing CpGs.

We next asked whether oCGIs were associated with the binding of other transcription factors with GC-poor motifs, namely FOXA1, a factor with an AT-rich motif containing no CpG sites. We found that species-specific oCGIs were enriched for species-specific FOXA1 binding (Fig. 6C, see Additional file 1: Fig. S46D for all species pairs). Because the FOXA1 motif is AT-rich and contains no CpG sites, this enrichment is unlikely to be due to the sequence features of the oCGI and suggests that other oCGI-related mechanisms promote FOXA1 binding. Additionally, FOXA1 is a pioneer factor that is able to bind to its motif within closed chromatin [69, 70]. Therefore, FOXA1 would not necessarily require open chromatin, such as the chromatin state generated by oCGIs, in order to bind. However, its binding is still enriched at oCGIs, suggesting that oCGIs do favor FOXA1 binding. Another pioneer factor with an AT-rich motif, CEBPA, was not associated with oCGIs, (Additional file 1: Fig. S46E), suggesting that there is functional heterogeneity in oCGI recruitment of TFs. Nonetheless, our analysis of TF binding turnover suggests that evolutionary changes in the binding of multiple TFs are associated with oCGI turnover, and we will return to this finding in the “Discussion.”

## Discussion

Understanding the genetic and molecular mechanisms that drive evolutionary changes in gene regulation is essential for understanding how such changes contribute to the evolution of novel traits. Previous studies have focused on the role of nucleotide substitutions, transposable elements, and insertions and deletions to identify enhancers that may encode lineage-specific functions [3, 4, 20, 24, 25, 27, 28, 71]. Here, we investigated the contribution of orphan CpG islands (oCGIs) to changes in transcriptional





**Fig. 6** oCGI turnover is associated with changes in transcription factor binding. **A** Schematic illustrating how we compared species-specific oCGIs with species-specific transcription factor binding events in adult liver, using rhesus macaque and mouse as an example case. *Left*: a rhesus-only (species A-only) oCGI with a rhesus-only (species A-only) CTCF peak. *Right*: a shared oCGI with a shared CTCF peak. Ticks show the locations of CpG dinucleotides. **B** *Left*: the consensus motif for CTCF (MA1929.1 from the JASPAR database). *Right*: Enrichment and depletion in each indicated comparison of species-specific and shared oCGIs (*top*: A-only, B-only, Shared) and species-specific and shared CTCF peaks (*left*: A-only, B-only, Shared), compared to a null expectation of no association between oCGI turnover and peak turnover. Each  $3 \times 3$  grid shows the results for a specific test examining oCGIs and their overlap with CTCF peaks. Each box in each grid is colored according to the level of enrichment over expectation (teal) or depletion (gray) of genome-wide sites that meet the criteria for that box. The color bar below each plot illustrates the level of enrichment or depletion over expectation. The filled upward-pointing triangles denote significant enrichment and open downward-pointing triangles denote significant depletion ( $q < 0.05$ , permutation test, BH-corrected; see Additional file 1: Fig. S22 and “Methods”). **C** *Left*: the consensus motif for FOXA1 (MA0148.1 from the JASPAR database). *Right*: Enrichment and depletion in each indicated comparison of species-specific and shared oCGIs (*top*: A-only, B-only, Shared) and species-specific and shared FOXA1 peaks (*left*: A-only, B-only, Shared). Shown as in (B) but with boxes colored according to the level of enrichment over expectation (red) and depletion (gray) of genome-wide sites that meet the criteria for that box

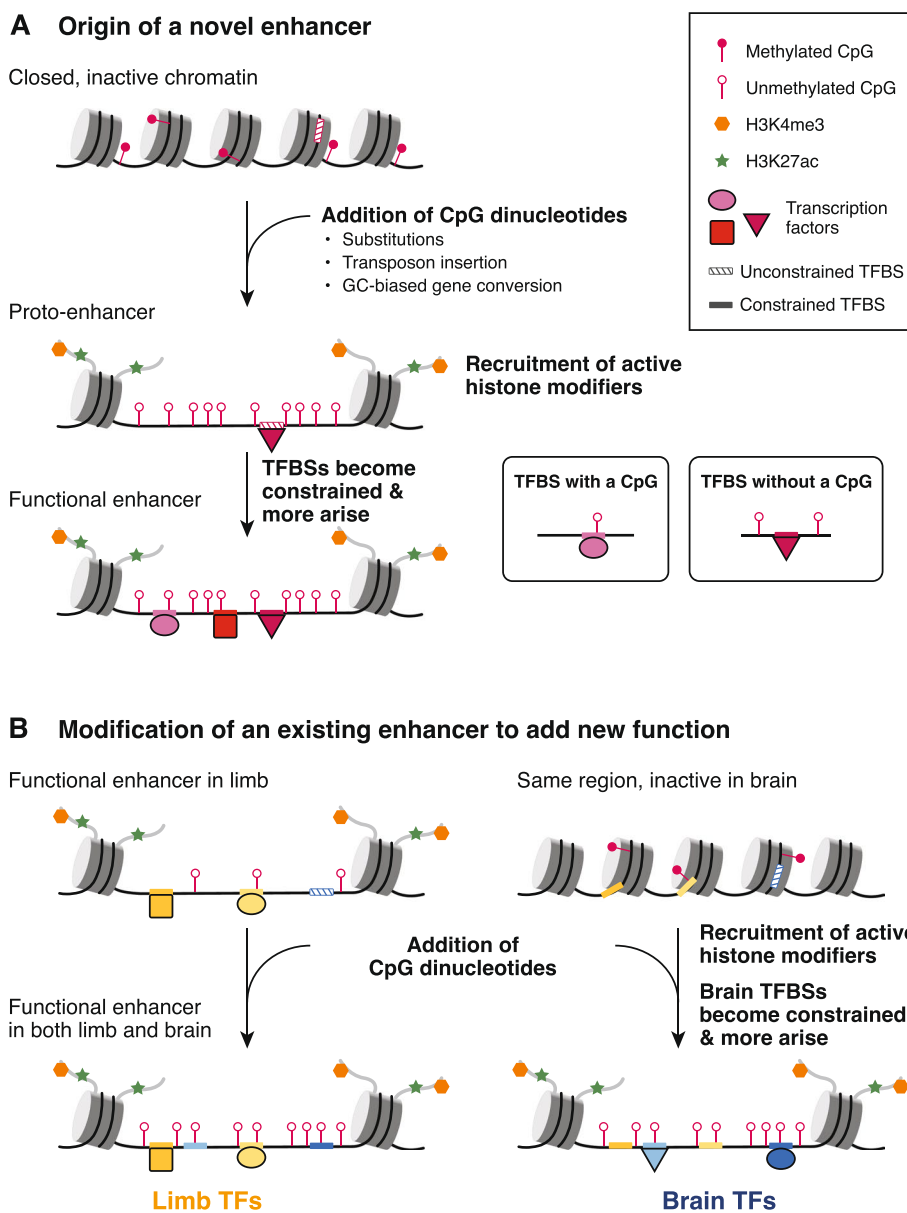
enhancer activity across species. Our findings support that oCGI turnover is associated with changes in the levels of several enhancer-associated histone modifications in mammals. We first found that oCGIs are enriched for histone modification peaks in all mammals we investigated, in line with previous findings in human [40, 41]. Additionally,

we identified extensive turnover of oCGIs across species. We then found that species-specific oCGIs were enriched for species-specific histone modification peaks in multiple developing and adult tissues, and this result was consistent in comparisons of both closely and distantly related species. We also found evidence that oCGI turnover is associated with changes in transcription factor binding events and changes in gene expression. Collectively, our results point to oCGI turnover as a major driver of gene regulatory innovation in mammalian evolution.

Our results also support that oCGIs contribute to the increased activity of Human Gain Enhancers (HGEs) in the developing cortex, which are hypothesized to alter gene expression during human brain development and contribute to uniquely human brain features. We experimentally modeled one such HGE, *hs754*, using humanized mice, and found that changes in oCGI content in the human ortholog were associated with *in vivo* changes in the levels of H3K27ac and H3K4me3, both associated with enhancer activity. These results are consistent with the role of oCGIs in active enhancers, and with the contribution of oCGI turnover to changes in enhancer activity we observed in our comparisons across mammals. However, we also identified correlated changes in CTCF binding at the humanized locus, which will need to be experimentally isolated from the effects of the human oCGI to determine the relative contributions of each to the increased enhancer activity of *hs754*. In general, deconvoluting the relative contributions of CpG dinucleotides and transcription factor binding sites to enhancer activity will require further experimental perturbations, as described in a recent study [72].

Given the oCGI-associated changes in CTCF binding in our humanized mouse model, we sought to assess whether changes in TF binding might be coincident with changes in oCGIs. As might be expected from the sequence characteristics of oCGIs, species-specific oCGIs were enriched for species-specific binding of TFs with GC-rich motifs and motifs containing CpG sites (CTCF, HNF4A, HNF6/ONECUT1). However, species-specific oCGIs were also enriched for FOXA1 binding events. FOXA1 has an AT-rich motif lacking any CpG site. This result suggests that CpG islands are associated with increased TF binding in a manner that is independent of their high GC- and CpG-content. We hypothesize that the gain of a CGI in an enhancer may not only add TFBSs due to its CpG-content, but that it also may promote the incorporation of TFBSs more generally. FOXA1 is a pioneer transcription factor with the ability to bind to its motif even within closed, nucleosomal DNA [69, 70]. However, species-specific FOXA1 binding is still enriched at species-specific oCGIs, suggesting that oCGIs provide a favorable locus for FOXA1 binding to occur even given its function as a pioneer factor. We hypothesize that the presence of oCGIs that promote active, open chromatin is also likely to facilitate the binding of non-pioneer factors with both GC- and AT-rich motifs.

Given these results, we propose a hypothetical model by which gain or loss of oCGIs influences the evolution of both new (Fig. 7A) and existing enhancers (Fig. 7B). Gain of oCGIs may occur via several mechanisms, including individual nucleotide substitution events or transposable element insertions. Transposable element insertion has been proposed as a way that new germ line differentially methylated regions, also CGIs, can arise [73]. Transposable elements could also carry CpG-rich promoters that decay over time while still retaining enhancer function. Another mechanism that may generate oCGIs is GC-biased gene conversion (gBGC), a process by which GC base pairs are preferentially



**Fig. 7** Model of enhancer evolution via oCGI turnover. **A** Evolution of a new enhancer from a locus in a closed chromatin state. This locus may include unconstrained, inaccessible TFBSs (striped boxes on DNA). DNA is depicted as a black line wrapped around cylindrical nucleosomes. After oCGI gain by several potential mechanisms (indicated in the figure), the site now acts as a proto-enhancer located within open, active chromatin recruited by the oCGI [42], which allows TFs to bind previously inaccessible TFBSs. A subset of histone tails (curved gray lines) with H3K4me3 (orange hexagons) and H3K27ac (green stars) modifications are shown. Filled lollipops indicate methylated CpGs, and unfilled lollipops indicate unmethylated CpGs. Over time, TFBSs become constrained (filled boxes on DNA) and additional TFBSs may arise and become fixed, resulting in the evolution of an enhancer with a constrained biological function. **B** Co-option of an existing enhancer in a novel biological context via oCGI gain. In an ancestral species, the enhancer is active in the developing limb and inactive in the developing brain, where the chromatin at the locus is closed. After oCGI gain, CpG-related mechanisms generate open chromatin in the developing brain, which allows existing unconstrained brain TFBSs to be bound. Over time, these and additional TFBSs may gain biological functions and be maintained by selection. The locus becomes a functional enhancer in the developing brain

fixed during meiotic recombination [74]. Increased GC content would increase the number of CpG sites at a locus. In fact, a substantial proportion of oCGIs in each species overlaps a gBGC tract identified using phastBias [75], especially in pig, dog, cat, and horse (Additional file 1: Fig. S47). This analysis is consistent with the possibility that gBGC is a mechanism that generates oCGIs.

We propose that oCGIs may contribute to the evolution of novel enhancers via several mechanisms (Fig. 7A). Our model builds on a previously described model of de novo enhancer birth in the genome, which focused on the role of TFBSs in this process [58] and proposed that enhancers arise from “proto-enhancers,” which are regions of the genome containing TFBSs that are able to recruit TFs and subsequently histone modifiers, leading to deposition of enhancer-associated histone modifications. Although biochemically active, these proto-enhancers do not have biological functions, are not under evolutionary constraint and are rapidly gained or lost over time. However, in some cases these proto-enhancers may serve as nucleation points for novel enhancers to evolve via genetic changes that generate additional TFBSs, producing more complex enhancers with functional effects that may be favored and maintained by selection.

Gain of oCGIs may have similar effects, due to newly introduced CpG dinucleotides recruiting ZF-CxxC domain-containing proteins and associated histone H3K4 methylation machinery, thereby generating an open chromatin environment [42]. Additionally, oCGIs can exclude DNA methylation [42, 76], which may alter the binding of methylation-sensitive TFs [77, 78]. This mechanism may be relevant to our result that species-specific CTCF binding is associated with species-specific oCGIs, because CTCF has been reported to be methylation-sensitive [79, 80]. Orphan CGI-related chromatin changes may themselves contribute to increased enhancer function and recruitment of transcriptional machinery. They may also make existing, previously inaccessible and unconstrained TFBSs at the locus available to TFs, adding new regulatory functions now subject to selection. Additionally, the open chromatin environment generated by novel oCGIs may act as a nucleation point for the evolution of additional TFBSs that contribute new regulatory information. Over time, TFBSs may accumulate and generate an enhancer with biological functions under evolutionary constraint.

Our findings support that highly constrained enhancers have also gained and lost oCGIs (Fig. 3E). Therefore, we propose that oCGI gain may also contribute to the co-option of existing, constrained enhancers in novel biological contexts (Fig. 7B). An existing enhancer may be active in the developing limb, where it is bound by limb TFs that contribute to its function. In developing brain, however, the enhancer remains embedded within closed, inactive chromatin, and the limb TFBSs are not used. In our example, this constrained limb enhancer gains an oCGI and CpG-dependent mechanisms generate active, open chromatin in the brain where the region was previously closed. This new chromatin state may facilitate the evolution of binding sites for brain TFs, leading to a novel function for this enhancer in the brain.

Loss of oCGIs may also contribute to regulatory innovation by reducing enhancer activity. The tendency of CpG dinucleotides to mutate by deamination [38] makes oCGIs susceptible to decay. Decay could occur at weakly constrained enhancers, but also at highly constrained enhancers, since loss of regulatory activity may also lead to novel functions favored by selection. Additionally, we proposed above that oCGIs contribute to

the initial opening of chromatin in a new context that allows existing TFBSs to become constrained and new TFBSs to emerge. It is possible that, once an enhancer has gained enough TFBSs to be strongly active, the CpG dinucleotides of the oCGI are no longer required for function (unless part of a TFBS) and begin to decay. In this way, oCGIs could act transiently to seed enhancers, then decay once no longer needed for function.

Turnover of oCGIs may also impact the evolution of poised enhancers, a class of enhancers whose function is dependent in part on oCGIs for their tethering to target promoters in embryonic stem cells and subsequent activation during development [72, 81]. Poised enhancers are enriched near neurodevelopmental genes and are important for their activation during neural differentiation [82]. Evolutionary turnover of oCGIs could therefore impact the function of poised enhancers and downstream neurodevelopmental processes. Recent work has identified poised enhancers in both mouse and human embryonic stem cells (ESCs) [83]. We note that the enhancer we modeled in humanized mice, hs754, was identified as a poised enhancer only in human but not in mouse, in line with the presence of a large oCGI only in the human sequence. Examination of this locus in published micro-C datasets [84, 85] reveals a contact between hs754 and the gene *Irx2* in human, but not mouse, ESCs (Additional file 1: Fig. S48). Further characterization of the hs754 model may reveal changes in the expression of *Irx2* or other genes at earlier developmental time points when neurodevelopmental genes first become active, which is when poised enhancers are thought to function. Evolutionary changes in oCGI content that impact poised enhancer function during neurodevelopment may thus have implications for brain evolution, including in humans.

## Conclusions

Our study identifies oCGI turnover as a novel mechanism affecting enhancer evolution in a variety of tissue contexts. Given this finding, oCGI content should be considered when assessing the mechanisms driving regulatory evolution across species and its impact on trait evolution. In the context of human and non-human primate evolution, several previous studies have focused on substitution events [15, 17, 21, 86] or deletions [27, 28], many of which are thought to alter TFBS content at enhancers. Our work highlights that an additional class of sequence change, oCGI turnover, is likely to reveal additional enhancers with lineage-specific functions, broadening the set of candidate regulatory innovations that may contribute to the evolution of novel traits.

## Methods

### Generating genome-wide CGI maps

We defined CGIs computationally (Additional file 1: Fig. S1) using a program developed by Andy Law [87, 88]. This program first scans the genome to identify all CpG dinucleotides. It then performs a windowing procedure to identify CGIs. Briefly, it selects the first CpG in the genome and adds downstream CpGs until the interval between two CpGs is at least 200 bp, when it performs a test for the following criteria: GC content of at least 50% and a ratio of observed over expected CpG dinucleotides of at least 0.6 [38]. If the interval meets the criteria, it attempts to build a bigger CGI by continuing the process of adding a CpG and testing whether the criteria are met, then once the criteria are no longer met the interval is output as a CGI. Otherwise if the criteria are not met, the

program drops the initial CpG and continues adding downstream CpGs until the length of 200 bp is met again, at which point it performs the test as above. This program generates more permissive CGI tracks than the default CGI tracks displayed by the UCSC Genome Browser, which imposes additional tests on CGI intervals to output only the largest, most CpG-dense CGIs [49].

We used the following genome versions in this study: rheMac10, calJac4, mm39, rn7, susScr11, canFam6, felCat9, and equCab3 for analyzing adult tissue data [34], and hg19, rheMac2, and mm9 for analyzing developing limb and cortex data [31, 36]. Repeat-masked genomes were used for CGI identification.

### Filtering strategy to identify oCGIs

We restricted our analysis to oCGIs by excluding CGIs overlapping several categories of annotated gene-associated features including promoters (2 kb upstream of transcription start sites, TSSs) and exons, as described in Additional file 1: Fig. S1 and Additional file 2: Table S11, using NCBI and UCSC versions of RefSeq for each genome [46, 89]. We also excluded CGIs falling in two additional annotated feature sets in the human and mouse genomes: promoters (2 kb upstream of TSSs) annotated by the FANTOM Consortium based on CAGE (Cap Analysis of Gene Expression) data from several hundred human and mouse cell lines and tissues [47] and blacklist regions annotated by the ENCODE project in human and mouse [48]. We further restricted our analysis to oCGIs with orthologous sequence in the human genome (hg38 for rheMac10, calJac4, mm39, rn7, susScr11, canFam6, felCat9, and equCab3, or hg19 for rheMac2 and mm9), which we identified using liftOver [90], that did not overlap features annotated in human (Additional file 1: Fig. S1).

### Comparison to other CGI annotations

We obtained CGI maps from two other sources. First, we downloaded maps from the UCSC Genome Browser for a subset of species: human (hg19), rhesus (rheMac10), mouse (mm39), rat (rn7), dog (canFam6), and horse (equCab3). Second, we downloaded CGIs identified using probabilistic models (<https://www.haowulab.org/software/makeCGI/index.html>) [51]. For the second set, we used liftOver (minMatch=0.8) to convert to the genome builds we are using (rheMac2 to rheMac10, mm10 to mm39, rn4 to rn7, canFam2 to canFam6, and equCab2 to equCab3). As done for the primary set of CGIs we used in this study, we removed CGIs that overlapped annotated promoters and gene bodies, including for pseudogenes and ncRNAs, to generate oCGI sets. We then assessed the overlap of oCGIs from each annotation strategy using BEDTools merge to combine across the sets (default settings, preserving information on the source files for each merged oCGI using “-c 4 -o collapse”) [91]. We also assessed the number of N bases in each oCGI using faCount [89].

### Analysis of VISTA enhancer data

We downloaded bed files from the ENCODE Portal (<https://www.encodeproject.org>) containing annotated peaks for the four histone modifications used in this study (H3K27ac, H3K4me3, H3K4me2, H3K4me1) in five E11.5 mouse tissues (forebrain, mid-brain, hindbrain, heart, and limb). We overlapped tested VISTA elements (<https://enhan>

cer.lbl.gov; see “Results”) with each peak set in each tissue using BEDTools intersect [91] and performed a separate analysis for VISTA elements that overlapped and did not overlap an oCGI in mouse (Additional file 1: Fig. S3A). For each test, we then counted the number of VISTA elements falling into each of four categories based on their overlap with a ChIP-seq peak (overlap versus no overlap) and whether they showed transgenic reporter activity in the tissue predicted by the ChIP-seq data (reporter activity versus no reporter activity). We performed Fisher’s exact test on the four categories in this contingency table to evaluate whether there was a significant association between the presence or absence of a ChIP-seq peak and reporter activity (Additional file 1: Fig. S3B, Additional file 2: Table S1). We corrected *p*-values across all tests using the Benjamini–Hochberg procedure [92]. A result was considered significant if the *q*-value (BH-corrected *p*-value) was less than 0.05.

#### **Analysis of cCREs and published enhancer-promoter interaction data**

We downloaded ENCODE cCREs defined in human and mouse from the UCSC Genome Browser [29]. We also downloaded peak files in BED format from the ENCODE Portal as above for H3K27ac, H3K4me3, H3K4me2, and H3K4me1 in E14.5 mouse liver. We sorted mouse intronic and intergenic ENCODE cCREs based on their overlap with each histone modification and mouse oCGIs. Then we measured the number of chromatin interactions involving each cCRE in a promoter capture HiC dataset [53], as listed in the file “FLC\_promoter\_other\_significant\_interactions.txt” associated with that study that we downloaded from the ArrayExpress repository (accession number E-MTAB-2414). We identified interactions by allowing an overlap with the cCRE or a 1-kb window on either side of the cCRE. We also measured the proportion of total cCREs, intronic and intergenic cCREs, and cCREs overlapping oCGIs that belonged to each regulatory element category (annotated in the cCRE files) using BEDTools intersect with default settings.

#### **Analysis of published ChIP-seq data**

We downloaded FASTQ files for histone modification profiles obtained by ChIP-seq in adult tissues [34] from the ArrayExpress repository (accession number E-MTAB-7127). We downloaded FASTQ files for transcription factor profiles obtained by ChIP-seq in adult liver from the ArrayExpress repository (accession numbers E-MTAB-437 for CTCF, E-MTAB-1509 for FOXA1, HNF4A, HNF6/ONECUT1, and CEBPA). We excluded two species for which the available genomes were generated with less than 10X sequencing coverage: opossum (monDom5) and rabbit (oryCun2).

We mapped reads using Bowtie2 [93] (with settings specified by the flag “–very-sensitive”) to unmasked genomes. We removed duplicate and multi-mapping reads using the Sambamba [94] commands “view” and “markdup.” We called peaks using MACS2 [95] with default settings, set to narrow peaks for H3K4me3, H3K27ac, and all transcription factors, and set to broad peaks using “–broad” for H3K4me1.

For histone modifications we discarded outlier replicates, defined as those with  $\geq 20\%$  fewer or  $\geq 50\%$  more peaks than the average of the other replicates for that species, tissue, and histone modification. For transcription factors, we discarded replicates with low peak numbers, as was done previously with this data [68] for one replicate of rhesus

CEBPA, rhesus HNF4A, rhesus FOXA1, and dog HNF6/ONECUT1. We also excluded one replicate with low peak numbers for the following species and TFs: dog FOXA1, rhesus HNF6/ONECUT1, rat HNF6/ONECUT1, and mouse FOXA1. We then identified reproducible peaks by taking the intersection of peaks in each replicate using BEDTools intersect.

For visualization, we used bamCoverage from deepTools [96] to generate bigWig files with the following settings: `--normalizeUsing CPM --binSize 10 --extendReads 300 --centerReads`. These bigWig files summarize ChIP-seq signal data as extended read counts per million (CPM) in bins with a width of 10 bp. We generated bigBed files with peak intervals for each replicate [97].

For the histone modification ChIP-seq in developing cortex, we downloaded processed peak files and bigWig files from [http://noonan.ycga.yale.edu/noonan\\_public/reilly2015/](http://noonan.ycga.yale.edu/noonan_public/reilly2015/) (also available in the Gene Expression Omnibus under accession number GSE63649) [36]. For species pair analysis, we matched time points between species as in Additional file 2: Table S5. These bigWig files summarize ChIP-seq signal data as the number of sequenced fragments (extended to 300 bp) that overlap each base pair, normalized to 1 million aligned reads.

For the histone modification ChIP-seq in developing limb, we downloaded processed peak files and bigWig files from [http://noonan.ycga.yale.edu/noonan\\_public/Limb\\_hub/](http://noonan.ycga.yale.edu/noonan_public/Limb_hub/) (also available from GEO under accession number GSE42413) [31]. For species pair analysis, we matched time points between species as in Additional file 2: Table S5. These bigWig files summarize ChIP-seq signal data as the number of sequenced fragments (extended to 300 bp) that overlap each base pair, normalized to 1 million aligned reads.

### Associating oCGIs and histone modification peaks

We measured the proportion of oCGIs in each species that overlap reproducible histone modification peaks using BEDTools intersect (with flags “-wa -u”). Additionally, we performed a genomic reshuffling test to determine whether the overlap is greater than what would be expected if oCGIs and peaks were distributed independently throughout the genome (Additional file 1: Fig. S5). For each species, we used BEDTools shuffle (with flags “-chrom -noOverlapping”) to randomly redistribute oCGI intervals on each reference genome, excluding all regions with an annotated RefSeq feature (all promoters and all exons for protein-coding genes, lncRNAs, and other ncRNAs, plus introns for pseudogenes and features of an unknown type) and all regions falling in RepeatMasker regions because oCGIs were called on repeat-masked genomes. Since our oCGI sets are additionally filtered by lifting to the human genome and excluding human features, we also lifted each shuffled set to the human genome and filtered out oCGIs overlapping human features, then restricted the shuffled set in the original species based on its status in human. We then counted the percentage of this shuffled and filtered set overlapping a peak for each histone modification we examined in each tissue. We repeated the shuffling procedure 20,000 times and generated an expected value based on the mean percentage of shuffled oCGIs overlapping peaks for each tissue and histone modification across all 20,000 shuffling rounds. We calculated *p*-values by measuring the proportion of shuffling rounds in which the percentage of shuffled oCGIs overlapping peaks was more extreme than the observed percentage of oCGIs overlapping peaks. We corrected



$p$ -values across all species, tissues, and marks and a result was considered significant if the  $q$ -value (BH-corrected  $p$ -value) was less than 0.05.

We also measured the percentage of peaks that overlap oCGIs by first identifying intronic and intergenic peaks for each species, tissue, and histone modification, including lifting to human and filtering based on human feature annotations as for oCGIs. We then measured the percentage of peaks passing this filter that overlapped an oCGI in the original genome (Additional file 1: Fig. S7).

We measured the histone modification levels in peaks with oCGIs and peaks without oCGIs (Fig. 1C, Additional file 1: Fig. S8) by quantifying reads for adult tissue datasets using featureCounts from Subread [98] and normalizing reads per kilobase per million mapped reads (RPKM). We quantified levels from bigWig files from developing tissue datasets using bigWigAverageOverBed, a measure analogous to RPKM.

### **Analysis of sequence conservation and age**

We downloaded 100-way Vertebrate phastCons elements in the human genome from the UCSC Genome Browser, generated using an alignment of 99 vertebrate species to human (hg38). Each phastCons element covers a specified interval in the human genome and is associated with a normalized LOD (log-odds) score reflecting the probability of the element being generated under the constrained phylo-HMM (phylogenetic hidden Markov model) compared to its probability under the non-constrained model [57]. Higher LOD scores indicate a greater degree of inferred constraint. Throughout this study, we overlapped oCGIs and phastCons elements and reported several measures, including maximum LOD score of overlapping phastCons elements, aggregate LOD scores of overlapping phastCons elements (the sum across all elements), and the proportion of bases covered by a phastCons element.

We also used an age segmentation map of the human genome to date oCGIs and histone modification peaks [58]. This map dates intervals within the human genome based on the most distantly related species in the 46-Way MultiZ alignment (hg19) that has alignable sequence in that interval. The ages are as follows: Human, Ape, Primate, Eutheria, Theria, Mammalia, Amniota, Tetrapoda, Gnathostomata, and Vertebrata. Some intervals have an unknown age. We combined the most ancient three categories (Tetrapoda, Gnathostomata, and Vertebrata) into a category called “Older than Amniota.” The map was generated based on the human genome, so we infer sequence ages in the other species based on the age assigned to their orthologous human sequence. As a consequence, we converted the age of some oCGIs and peaks to “Unknown” if their assigned age was for a clade that did not include the species being analyzed; for example, if an oCGI in the pig genome was dated to age “Human,” we converted its age to “Unknown.”

### **Identification of orthologous oCGIs across species pairs**

For each species pair in the dataset (“species A” and “species B”), we identified oCGIs that are mappable between both species (using human as an intermediate), regardless of their oCGI status in both species. In other words, the underlying sequence was required to be present in both species, even if the oCGI was only called in one species. The procedure is described in Additional file 1: Figure S13.

We collected several additional pieces of information for each orthologous oCGI. We counted the number of CpG dinucleotides in each oCGI using `faCount` from the UCSC Genome Browser (Fig. 2C, Additional file 1: Fig. S15). We also intersected the human coordinates of orthologous oCGIs between each species pair with `phastCons` elements (Fig. 2D, Additional file 1: Fig. S18-19) and the age segmentation map of the human genome (Fig. 2E, Additional file 1: Fig. S20), both of which are described above.

#### **Analysis of oCGI gain and loss**

We lifted oCGIs in all species to the human genome, then merged them using `BEDTools merge` using default settings, but with the flags “-c 4 -o collapse” to preserve the species of origin for each site, which was annotated in the fourth column of the original files. We then counted the number of oCGIs with each of the presence-absence patterns depicted in Additional file 1: Fig. S21A in order to determine the number gained or lost on a subset of terminal branches of the tree (human, rhesus, mouse, rat, dog, and cat).

#### **Enrichment analysis of species-specific oCGIs and species-specific peaks**

For each species pair, tissue, and histone modification or transcription factor, we intersected the set of species-specific and shared oCGIs with reproducible ChIP-seq peaks in both species in the pair, requiring 1 bp of overlap. We sorted species-specific (A-only, B-only) and shared oCGIs based on their overlap with species-specific (A-only, B-only) and shared peaks. It was also possible for an orthologous oCGI to overlap with no peak in either species; however, these sites were excluded from the analysis.

We determined enrichment and depletion in each category using a procedure described in Additional file 1: Figure S22. We calculated  $p$ -values for enrichment and depletion using a permutation test (Additional file 1: Fig. S22). A result was considered significant if the  $q$ -value (BH-corrected  $p$ -value) was less than 0.05.

We also performed a peak-centric version of this test (rather than oCGI-centric) (Additional file 1: Fig. S29). Instead of first identifying orthologous oCGIs between a species pair, we first identified orthologous histone modification or transcription factor peaks. We then counted their overlap with oCGIs and performed enrichment analysis and a permutation test as described above for the oCGI-centric analysis, although in this case the peak species specificity labels were randomly permuted and the number falling into each oCGI category was counted. Statistical significance was determined as above, and a result was considered significant if the  $q$ -value (BH-corrected  $p$ -value) was less than 0.05.

#### **Association between oCGIs and HGEs**

For the human developing cortex and limb datasets [31, 36], at each human time point (Additional file 2: Table S5), we sorted all intronic and intergenic histone modification peaks based on their status as HGEs or as non-HGE human enhancers. We counted the proportion of HGEs with an oCGI species pattern in each of the following categories: human-only oCGI (absent in rhesus and mouse), rhesus-only oCGI (absent in human and mouse), mouse-only oCGI (absent in human and rhesus), oCGI shared between human and rhesus only (absent in mouse), oCGI shared between human and mouse (absent in rhesus), oCGI shared between rhesus and mouse (absent in human), oCGI

in all three species, and oCGI in none of the three species. We compared these proportions (the observed values) to proportions in the non-HGE human enhancer set using a resampling test (Fig. 4B, Additional file 1: Fig. S35, Additional file 2: Table S6), described in Additional file 1: Figure S34. A result was considered significant if the  $q$ -value (BH-corrected  $p$ -value) was less than 0.05.

### Mouse line generation and validation

The hs754 humanized mouse line was generated at the Yale Genome Editing Center by injecting the editing plasmid, purified Cas9 RNA, and sgRNA into C57BL/6 J embryos, as previously described [99]. The sequence of the CRISPR guide RNA was 5' GAACCA AATATGGTGGGGAC. Coordinates of human and mouse sequences are in Additional file 2: Table S7. F0 edited mice were backcrossed with C57BL/6 J mice from Jackson Laboratory for several generations.

Genotyping primers for the humanized and wild type locus are provided in Additional file 2: Table S8, along with cloning primers for amplification of human and mouse genomic DNA for the generation of the editing construct and Sanger sequencing to verify the integrity of the edited locus. We amplified across both homology arms and the humanized region in both humanized and wild type mice, then cloned the products into pUC19 and used Sanger sequencing to verify each product. We also amplified 6.1 kb upstream and 3.8 kb downstream of the humanized locus, followed by Sanger sequencing, to establish that the extended locus was intact (Additional file 1: Fig. S37B-C). We also established that the homozygous line carried two copies of the humanized haplotype using quantitative PCR (qPCR) (Additional file 1: Fig. S37D). We performed copy number qPCR using genomic DNA from three humanized and three wild type individuals using Light Cycler 480 SYBR Green I Master Mix (Roche #04707516001). The biological replicates for each individual were run in triplicate and error bars show the standard deviation from these technical replicates. All Ct values were normalized to a control region on chromosome 5, and normalized again to the values for the first wild type individual. Copy number primers are provided in Additional file 2: Table S8.

### Chromatin immunoprecipitation and ChIP-seq

We collected tissue from developing diencephalon at E11.5 and E17.5 from homozygous humanized crosses and homozygous wild type crosses. Each biological replicate came from 6 pooled embryos (at E11.5) or 3 pooled embryos (at E17.5) from a single litter, and each experiment involved two (at E11.5) or three (at E17.5) biological replicates. We crosslinked and sonicated tissue, then immunoprecipitated chromatin with antibodies for H3K27ac (Active Motif #91193, 15ug chromatin at E11.5 or 20  $\mu$ g chromatin at E17.5 and 5  $\mu$ g of antibody), H3K4me3 (Cell Signaling #9751S, 5ug chromatin and 5ug of antibody), or CTCF (Diagenode #C15410210-50, 15  $\mu$ g chromatin and 5  $\mu$ g of antibody) as previously described [100]. Using the MODified peptide array from Active Motif (#13005 and #13006), we have found that the lot of H3K27ac antibody that we used is cross-reactive for histone H2B lysine 5 acetylation, in addition to H3K27 acetylation (Additional file 1: Fig. S49). This modification has also been shown to predict active enhancers [101]. Samples and their matched inputs were sequenced on an Illumina NovaSeq 6000 to generate paired 150-bp reads.

### RNA extraction and RNA-seq

We collected diencephalon tissue from 6 pooled embryos (at E11.5) or 3 pooled embryos (at E17.5) from a single litter per biological replicate, from 4 biological replicates per genotype. We purified RNA using the Qiagen miRNeasy Kit (#74106). The Yale Center for Genome Analysis prepared libraries using polyA-selection (Roche Kapa mRNA Hyper Prep Cat #KR1352) and sequenced them on an Illumina NovaSeq 6000 to generate paired 150-bp reads.

### Analysis of ChIP-seq data from the hs754 humanized mouse

We generated Bowtie2 indexes for the mouse genome (mm39) and for the humanized mouse genome, made by replacing the sequence of the mouse locus with the human sequence. Coordinates are shown in Additional file 2: Table S7. We mapped ChIP-seq reads to the appropriate genome using Bowtie2 (v2.4.2) using the settings “-sensitive -no-unal.” We removed multi-mapping and duplicate reads using the Sambamba commands “view” and “markdup.” We called peaks using MACS2 with default settings set to narrow peaks for H3K27ac, H3K4me3, and CTCF.

For each time point and ChIP target, we generated a set of merged peaks across the genotypes in order to perform differential peak calling. The humanized allele of hs754 is 5531 bp, compared to 5141 bp for the replaced region in mouse, and this size discrepancy meant that an adjustment procedure was required in order to assign orthology and merge peaks between genotypes, which is described in Additional file 1: Fig. S50. We counted reads in all merged peaks using HTSeq [102], then called differential peaks in R using DESeq2 [103]. Results for all peaks are shown in Additional file 2: Table S9. A peak was considered significantly different if the *q*-value (BH-corrected *p*-value) was less than 0.05. We found several significant differential peaks on chromosome 19, which are all located within two known copy-number variants in the C57BL/6 mouse line (Additional file 1: Fig. S51) [104]. We hypothesize that in some comparisons, the samples used for wild type and humanized tissue contained different copy numbers of these loci. Because our testing for differential ChIP-seq peaks using DESeq2 uses raw read counts from the histone modification IP tracks, which are not normalized to input counts, the copy number variant led to calling these regions as differentially marked between the genotypes.

We generated bigWig files for visualization using bamCoverage from deepTools with the following settings: -normalizeUsing CPM -binSize 10 -extendReads 300 -centerReads. These bigWig files show counts per million (CPM) in bins with a width of 10 bp. We also generated bigBed files with peak intervals for each replicate.

### Analysis of RNA-seq data from the hs754 humanized mouse

We generated a STAR index for mm39 using the unmasked genome and the basic gene annotation GTF from GENCODE release 31 [105]. We mapped reads using STAR (v2.7.9a) and used “-quantMode GeneCounts” to directly output counts. We then used DESeq2 in R to perform differential expression analysis using an adjusted *p*-value cutoff of 0.05. Full results are shown in Additional file 2: Table S10.

### Analysis of published RNA-seq data

In order to analyze differences in gene expression between species pairs, we analyzed genes annotated as 1:1 orthologs by Ensembl [106]. This required mapping reads to Ensembl genomes for counting based on Ensembl GTF files. However, in order to integrate gene expression data with the orthologous oCGI sets which were called in UCSC genome versions, we restricted this analysis to species for which Ensembl release 106 offered GTF files built on the same NCBI genome assemblies as the UCSC genome versions (rhesus macaque, mouse, rat, pig, cat, and horse; see Additional file 2: Table S12). We converted the chromosome names in these Ensembl GTF files to be compatible with UCSC chromosome names by adding “chr” in front of the numbers used by Ensembl.

We downloaded FASTQ files for RNA-seq from the ArrayExpress repository (E-MTAB-8122). We mapped reads using STAR to unmasked genomes (rheMac10, mm39, rn7, susScr11, felCat9, and equCab3) with the setting `-sjdbOverhang 149` [107]. We specified the Ensembl GTF files described above at this mapping step. STAR output the counts for each gene in each species, which we used to calculate TPM (transcripts per million) [108], using the R package GenomicRanges to calculate exon lengths for each gene from GTF files [109].

### Association between species-specific oCGIs in species-specific histone modification peaks and gene expression

Within each species, we defined the regulatory domain of each gene annotated in the Ensembl GTF using GREAT rules [65] in order to associate the gene with its putative enhancers. First, we assigned a basal regulatory domain to each protein-coding gene, which was 5 kb upstream and 1 kb downstream of the TSS. We then extended each basal regulatory domain to the next nearest basal regulatory domain upstream and downstream, up to a maximum distance of 1 Mb. We then sorted genes into four categories based on species-specific oCGIs in species-specific peaks that fell into their regulatory domains (performing a separate analysis for each species pair, tissue, and histone modification): the “A-only set” of genes associated with an A-only oCGI in an A-only peak, the “B-only set” of genes associated with a B-only oCGI in a B-only peak, the “mixed set” of genes associated with both types of oCGIs which was excluded from the analysis, and the “background set” containing all other genes (Additional file 1: Fig. S42A). Note that this assignment procedure considers only species-specific oCGIs in species-specific peaks, ignoring all other oCGIs and peaks that may be associated with the gene.

We then compared the TPM ratio (TPM in Species A / TPM in Species B) for every gene in the A-only set and the B-only set to the background set using a resampling test (Additional file 1: Fig. S42B-E). A result was considered significant if the *q*-value (BH-corrected *p*-value) was less than 0.05.

### Analysis of GC-biased gene conversion tracts

GC-biased gene conversion tracts were defined using phastBias [75], which finds regions where weak-to-strong substitutions are enriched compared to strong-to-weak substitutions. The program was run separately for each of the nine species in the analysis, taking that species as the foreground branch, using the setting “`--output-tracts`.”

A whole-genome alignment for all nine species was extracted from a 120-mammal alignment [110] using `maf_parse`. The neutral rate was determined based on fourfold degenerate sites. These sites were extracted from the 120-mammal alignment using `msa_view` with settings “--4d” and “--features” from the human Gencode annotation (v41). Branch lengths were determined using `phyloFit` with the setting “--subst-mod REV” [111].

### Visualization of Micro-C data

We visualized the Micro-C data shown in Additional file 1: Fig. S48 [84, 85] using the 4DN Visualization Workspace hosted by 4D Nucleome. Specifically, we visualized files 4DNFI9GMP2J8.mcool (human H1 ESCs) and 4DNFINNZDDXV.mcool (mouse JM8. N4 ESCs).

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-024-03300-z>.

Additional file 1: Supplementary figures S1-S51.

Additional file 2: Supplementary tables S1-S12.

Additional file 3: Review history.

### Acknowledgements

We thank members of the Noonan lab, B. Lesch and lab members, A. Louvi, G. Wagner, and I-H. Park for their feedback on the manuscript, and L. Pennacchio for assistance with interrogating the VISTA Enhancer Browser database.

### Review history

The review history is available as Additional File 3.

### Peer review information

Veronique van den Berghe was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

### Authors' contributions

AAK, EVD, and JPN conceived of and designed the study; AAK performed all computational analyses related to oCGIs, with assistance from SU and KMY; EVD designed the humanized mouse line and performed validation experiments; EVD and TN generated the humanized mouse line; AAK performed humanized mouse line validation experiments, ChIP-seq experiments, and RNA-seq experiments with assistance from MFRL and MB; AAK analyzed data from the humanized mouse line; AAK and JPN wrote the manuscript with input from all authors.

### Funding

This work was supported by a grant from the National Institute of General Medical Sciences (NIGMS) (R01 GM094780, to JPN), a grant from the Eunice Kennedy Shriver National Institute of Child Health and Human Development (R01 HD102030 to JPN) and funds from the Yale School of Medicine (to JPN). AAK was supported by an NSF Graduate Research Fellowship (DGE-1122492). EVD was supported in part by NIGMS training grant T32 GM007499. KMY was supported by an NSF Graduate Research Fellowship (DGE-1752134). SU was supported in part by a Research Fellowship (352711928) from the Deutsche Forschungsgemeinschaft (DFG). MB was supported by an NIH F32 Postdoctoral Fellowship (NICHD) (F32 HD108935). This research program and related results were also made possible by the support of the NOMIS foundation (to JPN).

### Availability of data and materials

ChIP-seq and RNA-seq data for the humanized mouse model have been deposited under GEO accession GSE231307 [112]. We used the following previously generated datasets: ChIP-seq for histone modifications in adult tissue [113], RNA-seq in adult tissue [114], ChIP-seq for transcription factors in adult liver [115, 116], ChIP-seq for histone modifications in developing brain [117] and developing limb [118], and Capture Hi-C in developing liver [119]. ChIP-seq data for histone modifications in developing mouse tissues were downloaded from ENCODE. <https://www.encodeproject.org> (2020): E11.5 tissues: ENCF958GPD, ENCF651IYR, ENCF445DQR, ENCF908XCE, ENCF928TGM, ENCF897EEM, ENCF566DFK, ENCF203QTV, ENCF236UMU, ENCF016BEF, ENCF147OKD, ENCF202HIO, ENCF098IGX, ENCF218FKJ, ENCF255WOG, ENCF047OVD, ENCF132QFU, ENCF835VQG, ENCF316YSJ, ENCF703AUU; E14.5 mouse liver data: ENCF384FJW, ENCF510NON, ENCF627CYS, ENCF880DXF.

We used Micro-C data in embryonic stem cells from human [120] and mouse [121].

All code for the study is available under an MIT license at Github [122] and <https://doi.org/10.5281/zenodo.11236169> [123].

## Declarations

### Ethics approval and consent to participate

All animal work was performed according to approved Yale IACUC protocols (#2020–07271, #2019–11167, #2022–11167).

### Consent for publication

Not applicable.

### Competing interests

The authors declare that they have no competing interests.

Received: 1 June 2023 Accepted: 4 June 2024

Published online: 13 June 2024

## References

1. Reilly SK, Noonan JP. Evolution of Gene Regulation in Humans. *Annu Rev Genom Hum G.* 2016;17(1):45–67.
2. Whalen S, Pollard KS. Enhancer Function and Evolutionary Roles of Human Accelerated Regions. *Annu Rev Genet.* 2022;56(1):423–39.
3. Dutrow EV, Emera D, Yim K, Uebbing S, Kocher AA, Krenzer M, et al. Modeling uniquely human gene regulatory function via targeted humanization of the mouse genome. *Nat Commun.* 2022;13(1):304.
4. Aldea D, Atsuta Y, Kokalari B, Schaffner SF, Prasasya RD, Aharoni A, et al. Repeated mutation of a developmental enhancer contributed to human thermoregulatory evolution. *Proc National Acad Sci.* 2021;118(16):e2021722118.
5. Boyd JL, Skove SL, Rouanet JP, Pilaz LJ, Bepler T, Gordân R, et al. Human-Chimpanzee Differences in a FZD8 Enhancer Alter Cell-Cycle Dynamics in the Developing Neocortex. *Curr Biol.* 2015;25(6):772–9.
6. Mangan RJ, Alsina FC, Mosti F, Sotelo-Fonseca JE, Snellings DA, Au EH, et al. Adaptive sequence divergence forged new neurodevelopmental enhancers in humans. *Cell.* 2022;185(24):4587–4603.e23.
7. Shapiro MD, Marks ME, Peichel CL, Blackman BK, Nereng KS, Jónsson B, et al. Genetic and developmental basis of evolutionary pelvic reduction in threespine sticklebacks. *Nature.* 2004;428(6984):717–23.
8. Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, et al. The genomic basis of adaptive evolution in threespine sticklebacks. *Nature.* 2012;484(7392):55–61.
9. Long HK, Prescott SL, Wysocka J. Ever-Changing Landscapes: Transcriptional Enhancers in Development and Evolution. *Cell.* 2016;167(5):1170.
10. Rebeiz M, Tsiantis M. Enhancer evolution and the origins of morphological novelty. *Curr Opin Genetics Dev.* 2017;45:115.
11. Peter IS, Davidson EH. Evolution of Gene Regulatory Networks Controlling Body Plan Development. *Cell.* 2011;144(6):970–85.
12. Levine M. Transcriptional enhancers in animal development and evolution. *Curr Biol.* 2010;20(17):R754.
13. Wray GA. The evolutionary significance of cis-regulatory mutations. *Nat Rev Genet.* 2007;8(3):206.
14. Carroll SB. Evolution at Two Levels: On Genes and Form. *Plos Biol.* 2005;3(7): e245.
15. Prabhakar S, Noonan JP, Pääbo S, Rubin EM. Accelerated Evolution of Conserved Noncoding Sequences in Humans. *Science.* 2006;314(5800):786–786.
16. Prabhakar S, Visel A, Akiyama JA, Shoukry M, Lewis KD, Holt A, et al. Human-Specific Gain of Function in a Developmental Enhancer. *Science.* 2008;321(5894):1346–50.
17. Pollard KS, Salama SR, King B, Kern AD, Dreszer T, Katzman S, et al. Forces Shaping the Fastest Evolving Regions in the Human Genome. *Plos Genet.* 2006;2(10):e168.
18. Uebbing S, Gockley J, Reilly SK, Kocher AA, Geller E, Gandotra N, et al. Massively parallel discovery of human-specific substitutions that alter enhancer activity. *Proc National Acad Sci.* 2021;118(2):e2007049118.
19. Li S, Hannenhalli S, Ovcharenko I. De novo human brain enhancers created by single-nucleotide mutations. *Sci Adv.* 2023;9(7):eadd2911.
20. Whalen S, Inoue F, Ryu H, Fair T, Markenscoff-Papadimitriou E, Keough K, et al. Machine learning dissection of human accelerated regions in primate neurodevelopment. *Neuron.* 2023;111(6):857–873.e8.
21. Keough KC, Whalen S, Inoue F, Przytycki PF, Fair T, Deng C, et al. Three-dimensional genome rewiring in loci with human accelerated regions. *Science.* 2023;380(6643):eabm1696.
22. Frankel N, Erezylmaz DF, McGregor AP, Wang S, Payre F, Stern DL. Morphological evolution caused by many subtle-effect substitutions in regulatory DNA. *Nature.* 2011;474(7353):598–603.
23. Chuong EB, Elde NC, Feschotte C. Regulatory evolution of innate immunity through co-option of endogenous retroviruses. *Sci New York N Y.* 2016;351(6277):1083–7.
24. Chuong EB, Rumi MAK, Soares MJ, Baker JC. Endogenous retroviruses function as species-specific enhancer elements in the placenta. *Nat Genet.* 2013;45(3):325–9.
25. Lynch VJ, Nnamani MC, Kapusta A, Brayer K, Plaza SL, Mazur EC, et al. Ancient transposable elements transformed the uterine regulatory landscape and transcriptome during the evolution of mammalian pregnancy. *Cell Rep.* 2015;10(4):551–61.
26. Trizzino M, Park Y, Holsbach-Beltrame M, Aracena K, Mika K, Caliskan M, et al. Transposable elements are the primary source of novelty in primate gene regulation. *Genome Res.* 2017;27(10):1623–33.
27. Mclean CY, Reno PL, Pollen AA, Bassan AI, Capellini TD, Guenther C, et al. Human-specific loss of regulatory DNA and the evolution of human-specific traits. *Nature.* 2011;471(7337):216–9.
28. Xue JR, Mackay-Smith A, Mouri K, Garcia MF, Dong MX, Akers JF, et al. The functional and evolutionary impacts of human-specific deletions in conserved elements. *Science.* 2023;380(6643):eabn2253.

29. Abascal F, Acosta R, Adleman NJ, Adrian J, Afzal V, Ai R, et al. Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature*. 2020;583(7818):699–710.
30. Gorkin DU, Barozzi I, Zhao Y, Zhang Y, Huang H, Lee AY, et al. An atlas of dynamic chromatin landscapes in mouse fetal development. *Nature*. 2020;583(7818):744–51.
31. Cotney J, Leng J, Yin J, Reilly SK, DeMare LE, Emera D, et al. The Evolution of Lineage-Specific Regulatory Activities in the Human Embryonic Limb. *Cell*. 2013;154(1):185–96.
32. Cotney J, Leng J, Oh S, DeMare LE, Reilly SK, Gerstein MB, et al. Chromatin state signatures associated with tissue-specific gene expression and enhancer activity in the embryonic limb. *Genome Res*. 2012;22(6):1069–80.
33. Villar D, Berthelot C, Aldridge S, Rayner TF, Lukk M, Pignatelli M, et al. Enhancer evolution across 20 mammalian species. *Cell*. 2015;160(3):554–66.
34. Roller M, Stamper E, Villar D, Izuogu O, Martin F, Redmond AM, et al. LINE retrotransposons characterize mammalian tissue-specific and evolutionarily dynamic regulatory regions. *Genome Biol*. 2021;22(1):62.
35. Berthelot C, Villar D, Horvath JE, Odom DT, Flicek P. Complexity and conservation of regulatory landscapes underlie evolutionary resilience of mammalian gene expression. *Nat Ecol Evol*. 2017;2(1):152–63.
36. Reilly SK, Yin J, Ayoub AE, Emera D, Leng J, Cotney J, et al. Evolutionary changes in promoter and enhancer activity during human corticogenesis. *Science*. 2015;347(6226):1155–9.
37. Won H, Huang J, Opland CK, Hartl CL, Geschwind DH. Human evolved regulatory elements modulate genes involved in cortical expansion and neurodevelopmental disease susceptibility. *Nat Commun*. 2019;10(1):2396.
38. Gardiner-Garden M, Frommer M. CpG Islands in vertebrate genomes. *J Mol Biol*. 1987;196(2):261–82.
39. Deaton AM, Bird A. CpG islands and the regulation of transcription. *Gene Dev*. 2011;25(10):1010–22.
40. Bell JSK, Vertino PM. Orphan CpG islands define a novel class of highly active enhancers. *Epigenetics*. 2017;12(6):449–64.
41. Steinhaus R, Gonzalez T, Seelou D, Robinson PN. Pervasive and CpG-dependent promoter-like characteristics of transcribed enhancers. *Nucleic Acids Res*. 2020;48(10):5306–17.
42. Hughes AL, Kelley JR, Klose RJ. Understanding the interplay between CpG island-associated gene promoters and H3K4 methylation. *Biochimica Et Biophysica Acta Bba - Gene Regul Mech*. 2020;1863(8): 194567.
43. Thomson JP, Skene PJ, Selfridge J, Clouaire T, Guy J, Webb S, et al. CpG islands influence chromatin structure via the CpG-binding protein Cfp1. *Nature*. 2010;464(7291):1082–6.
44. Denissov S, Hofemeister H, Marks H, Kranz A, Ciotta G, Singh S, et al. Mll2 is required for H3K4 trimethylation on bivalent promoters in embryonic stem cells, whereas Mll1 is redundant. *Development*. 2014;141(3):526–37.
45. Hu D, Gao X, Cao K, Morgan MA, Mas G, Smith ER, et al. Not All H3K4 Methylations Are Created Equal: Mll2/COMPASS Dependency in Primordial Germ Cell Specification. *Mol Cell*. 2017;65(3):460–475.e6.
46. O'Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res*. 2016;44(D1):D733–45.
47. Forrest ARR, Kawaji H, Rehli M, Baillie JK, de Hoon MJL, FANTOM Consortium and the RIKEN PMI and CLST (DGT), et al. A promoter-level mammalian expression atlas. *Nature*. 2014;507(7493):462–70.
48. Amemiya HM, Kundaje A, Boyle AP. The ENCODE Blacklist: Identification of Problematic Regions of the Genome. *Sci Rep-uk*. 2019;9(1):9354.
49. UCSC Genome Browser. CpG Islands [Internet]. 2006. Available from: [http://genomewiki.ucsc.edu/index.php/CpG\\_Islands](http://genomewiki.ucsc.edu/index.php/CpG_Islands). cited 2023 Feb 25.
50. Irizarry RA, Wu H, Feinberg AP. A species-generalized probabilistic model-based definition of CpG islands. *Mamm Genome*. 2009;20(9–10):674.
51. Wu H, Caffo B, Jaffee HA, Irizarry RA, Feinberg AP. Redefining CpG islands using hidden Markov models. *Biostatistics*. 2010;11(3):499–514.
52. Visel A, Minovitsky S, Dubchak I, Pennacchio LA. VISTA Enhancer Browser—a database of tissue-specific human enhancers. *Nucleic Acids Res*. 2007;35(Database issue):D88.
53. Schoenfelder S, Furlan-Magaril M, Mifsud B, Tavares-Cadete F, Sugar R, Javierre BM, et al. The pluripotent regulatory circuitry connecting promoters to their long-range interacting elements. *Genome Res*. 2015;25(4):582–97.
54. Core LJ, Martins AL, Danko CG, Waters CT, Siepel A, Lis JT. Analysis of nascent RNA identifies a unified architecture of initiation regions at mammalian promoters and enhancers. *Nat Genet*. 2014;46(12):1311–20.
55. Henriques T, Scruggs BS, Inouye MO, Muse GW, Williams LH, Burkholder AB, et al. Widespread transcriptional pausing and elongation control at enhancers. *Gene Dev*. 2018;32(1):26–41.
56. Pekowska A, Benoukraf T, Zacarias-Cabeza J, Belhocine M, Koch F, Holota H, et al. H3K4 tri-methylation provides an epigenetic signature of active enhancers. *EMBO J*. 2011;30(20):4198–210.
57. Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res*. 2005;15(8):1034–50.
58. Emera D, Yin J, Reilly SK, Gockley J, Noonan JP. Origin and evolution of developmental enhancers in the mammalian neocortex. *Proc National Acad Sci*. 2016;113(19):E2617–26.
59. Rajarajan P, Borrmann T, Liao W, Schrodde N, Flaherty E, Casiño C, et al. Neuron-specific signatures in the chromosomal connectome associated with schizophrenia risk. *Science*. 2018;362(6420):eaat4311.
60. Mallika C, Guo Q, Li JYH. Gbx2 is essential for maintaining thalamic neuron identity and repressing habenular characters in the developing thalamus. *Dev Biol*. 2015;407(1):26.
61. de Laat W, Duboule D. Topology of mammalian developmental enhancers and their regulatory landscapes. *Nature*. 2013;502(7472):499–506.
62. Merkschlager M, Nora EP. CTCF and Cohesin in Genome Folding and Transcriptional Gene Regulation. *Annu Rev Genom Hum G*. 2015;17(1):1–27.
63. Ren G, Jin W, Cui K, Rodriguez J, Hu G, Zhang Z, et al. CTCF-Mediated Enhancer-Promoter Interaction Is a Critical Regulator of Cell-to-Cell Variation of Gene Expression. *Mol Cell*. 2017;67(6):1049–1058.e6.
64. Kubo N, Ishii H, Xiong X, Bianco S, Meitinger F, Hu R, et al. Promoter-proximal CTCF-binding promotes long-range-enhancer dependent gene activation. *Nat Struct Mol Biol*. 2021;28(2):152–61.



65. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, et al. GREAT improves functional interpretation of cis-regulatory regions. *Nat Biotechnol.* 2010;28(5):495–501.
66. Schmidt D, Wilson MD, Ballester B, Schwalie PC, Brown GD, Marshall A, et al. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Sci New York N.Y.* 2010;328(5981):1036–40.
67. Schmidt D, Schwalie PC, Wilson MD, Ballester B, Gonçalves Â, 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.
68. Ballester B, Medina-Rivera A, Schmidt D, González-Porta M, Carlucci M, Chen X, et al. Multi-species, multi-transcription factor binding highlights conserved control of tissue-specific biological pathways. *Elife.* 2014;3:e02626.
69. Cirillo LA, Lin FR, Cuesta I, Friedman D, Jarnik M, Zaret KS. Opening of Compacted Chromatin by Early Developmental Transcription Factors HNF3 (FoxA) and GATA-4. *Mol Cell.* 2002;9(2):279–89.
70. Iwafuchi-Doi M, Zaret KS. Cell fate control by pioneer transcription factors. *Development.* 2016;143(11):1833–7.
71. Kunarso G, Chia NY, Jeyakani J, Hwang C, Lu X, Chan YS, et al. Transposable elements have rewired the core regulatory network of human embryonic stem cells. *Nat Genet.* 2010;42(7):631–4.
72. Pachano T, Sánchez-Gaya V, Ealo T, Mariner-Fauli M, Bleckwehl T, Asenjo HG, et al. Orphan CpG islands amplify poised enhancer regulatory activity and determine target gene responsiveness. *Nat Genet.* 2021;53(7):1036–49.
73. Suzuki S, Shaw G, Kaneko-Ishino T, Ishino F, Renfree MB. The Evolution of Mammalian Genomic Imprinting Was Accompanied by the Acquisition of Novel CpG Islands. *Genome Biol Evol.* 2011;3:1276–83.
74. Duret L, Galtier N. Biased Gene Conversion and the Evolution of Mammalian Genomic Landscapes. *Annu Rev Genom Hum G.* 2009;10(1):285–311.
75. Capra JA, Hubisz MJ, Kostka D, Pollard KS, Siepel A. A Model-Based Analysis of GC-Biased Gene Conversion in the Human and Chimpanzee Genomes. *Plos Genet.* 2013;9(8):e1003684.
76. Krebs AR, Dessus-Babus S, Burger L, Schübeler D. High-throughput engineering of a mammalian genome reveals building principles of methylation states at CG rich regions. *Elife.* 2014;3:e04094.
77. Kaluscha S, Domcke S, Wirbelauer C, Stadler MB, Durdu S, Burger L, et al. Evidence that direct inhibition of transcription factor binding is the prevailing mode of gene and repeat repression by DNA methylation. *Nat Genet.* 2022;54(12):1895–906.
78. Yin Y, Morgunova E, Jolma A, Kaasinen E, Sahu B, Khund-Sayeed S, et al. Impact of cytosine methylation on DNA binding specificities of human transcription factors. *Science.* 2017;356(6337):eaaj2239.
79. Hashimoto H, Wang D, Horton JR, Zhang X, Corces VG, Cheng X. Structural Basis for the Versatile and Methylation-Dependent Binding of CTCF to DNA. *Mol Cell.* 2017;66(5):711–720.e3.
80. Wang H, Maurano MT, Qu H, Varley KE, Gertz J, Pauli F, et al. Widespread plasticity in CTCF occupancy linked to DNA methylation. *Genome Res.* 2012;22(9):1680–8.
81. Rada-Iglesias A, Bajpai R, Swigut T, Brugmann SA, Flynn RA, Wysocka J. A unique chromatin signature uncovers early developmental enhancers in humans. *Nature.* 2011;470(7333):279.
82. Cruz-Molina S, Respuela P, Tebartz C, Kolovos P, Nikolic M, Fueyo R, et al. PRC2 Facilitates the Regulatory Topology Required for Poised Enhancer Function during Pluripotent Stem Cell Differentiation. *Cell Stem Cell.* 2017;20(5):689–705.e9.
83. Crispatzu G, Rehimi R, Pachano T, Bleckwehl T, Cruz-Molina S, Xiao C, et al. The chromatin, topological and regulatory properties of pluripotency-associated poised enhancers are conserved in vivo. *Nat Commun.* 2021;12(1):4344.
84. Hsieh THS, Cattoglio C, Slobodyanyuk E, Hansen AS, Rando OJ, Tjian R, et al. Resolving the 3D Landscape of Transcription-Linked Mammalian Chromatin Folding. *Mol Cell.* 2020;78(3):539–553.e8.
85. Krietenstein N, Abraham S, Venev SV, Abdennur N, Gibcus J, Hsieh THS, et al. Ultrastructural Details of Mammalian Chromosome Architecture. *Mol Cell.* 2020;78(3):554–565.e7.
86. Kostka D, Holloway AK, Pollard KS. Developmental Loci Harbor Clusters of Accelerated Regions That Evolved Independently in Ape Lineages. *Mol Biol Evol.* 2018;35(8):2034–45.
87. Law A. CpG islands - Andy Law track for galGal3 [Internet]. 2006. Available from: <https://github.com/ucscGenomeBrowser/kent/blob/master/src/hg/makeDb/doc/galGal3.txt#L649>.
88. UCSC Genome Browser. Pre-process chromosome files for Andy Law CpG island program [Internet]. 2022. Available from: [https://hgwdev-gperez2.gi.ucsc.edu/~gperez2/mlq/mlq\\_29758/](https://hgwdev-gperez2.gi.ucsc.edu/~gperez2/mlq/mlq_29758/).
89. Nassar LR, Barber GP, Benet-Pagès A, Casper J, Clawson H, Diekhans M, et al. The UCSC Genome Browser database: 2023 update. *Nucleic Acids Res.* 2022;51(D1):D1188–95.
90. Hinrichs AS, Karolchik D, Baertsch R, Barber GP, Bejerano G, Clawson H, et al. The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res.* 2006;34(suppl\_1):D590–8.
91. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.* 2010;26(6):841–2.
92. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J Royal Statistical Soc Ser B Methodol.* 1995;57(1):289–300.
93. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.* 2012;9(4):357–9.
94. Tarasov A, Vilella AJ, Cuppen E, Nijman IJ, Prins P. Sambamba: fast processing of NGS alignment formats. *Bioinformatics.* 2015;31(12):2032–4.
95. Zhang Y, Liu T, Meyer CA, Eeckhoutte J, Johnson DS, Bernstein BE, et al. Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* 2008;9(9):R137.
96. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* 2016;44(W1):W160–5.
97. Kent WJ, Zweig AS, Barber G, Hinrichs AS, Karolchik D. BigWig and BigBed: enabling browsing of large distributed datasets. *Bioinformatics.* 2010;26(17):2204–7.
98. Liao Y, Smyth GK, Shi W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics.* 2014;30(7):923–30.
99. Price NL, Rotllan N, Zhang X, Canfrán-Duque A, Nottoli T, Suarez Y, et al. Specific Disruption of Abca1 Targeting Largely Mimics the Effects of miR-33 Knockout on Macrophage Cholesterol Efflux and Atherosclerotic Plaque Development. *Circ Res.* 2019;124(6):874–80.

100. Cotney JL, Noonan JP. Chromatin Immunoprecipitation with Fixed Animal Tissues and Preparation for High-Throughput Sequencing. *Cold Spring Harb Protoc.* 2015;2015(2):pdb.prot084848.
101. Narita T, Higashijima Y, Kilic S, Liebner T, Walter J, Choudhary C. Acetylation of histone H2B marks active enhancers and predicts CBP/p300 target genes. *Nat Genet.* 2023;55(4):679–92.
102. Putri GH, Anders S, Pyl PT, Pimanda JE, Zanini F. Analysing high-throughput sequencing data in Python with HTSeq 2.0. *Bioinformatics.* 2022;38(10):2943–5.
103. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* 2014;15(12):550.
104. Watkins-Chow DE, Pavan WJ. Genomic copy number and expression variation within the C57BL/6J inbred mouse strain. *Genome Res.* 2008;18(1):60–6.
105. Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, et al. GENCODE 2021. *Nucleic Acids Res.* 2020;49(D1):gkaa1087.
106. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, et al. Ensembl 2022. *Nucleic Acids Res.* 2021;50(D1):D988–95.
107. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics.* 2013;29(1):15–21.
108. Wagner GP, Kin K, Lynch VJ. Measurement of mRNA abundance using RNA-seq data: RPKM measure is inconsistent among samples. *Theor Biosci.* 2012;131(4):281–5.
109. Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, et al. Software for Computing and Annotating Genomic Ranges. *Plos Comput Biol.* 2013;9(8): e1003118.
110. Hecker N, Hiller M. A genome alignment of 120 mammals highlights ultraconserved element variability and placenta-associated enhancers. *Gigascience.* 2020;9(1):giz159.
111. Siepel A, Haussler D. Phylogenetic Estimation of Context-Dependent Substitution Rates by Maximum Likelihood. *Mol Biol Evol.* 2004;21(3):468–88.
112. Kocher AA, Noonan JP. CpG island turnover is associated with evolutionary changes in enhancer activity. *Datasets GSE231307.* Gene Expression Omnibus. 2023. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE231307>.
113. Roller M, Stamper E. H3K4me3, H3K27ac and H3K4me1 ChIP-seq in 4 tissues of 10 mammals. Dataset E-MTAB-7127. *ArrayExpress.* 2020. <https://www.ebi.ac.uk/arrayexpress/experiments/E-MTAB-7127>.
114. Roller M. RNA-seq in 4 tissues of 10 mammals. Dataset E-MTAB-8122. *ArrayExpress.* 2019. <https://www.ebi.ac.uk/arrayexpress/E-MTAB-8122>.
115. Schwalie PC. CTCF binding evolution in mammals. Dataset E-MTAB-437. *ArrayExpress.* 2011. <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-437>.
116. Wilson M, Ballester B, Schmidt D, Stefflova K, Watt S, Brown G, Lukk M, Flicek P, Odom D. Combinatorial transcription factor binding evolution in five placental mammals. Dataset E-MTAB-1509. *ArrayExpress.* 2014. <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-1509>.
117. Reilly SK, Yin J, Ayoub AE, Emera D, Leng J, Cotney J, Sarro R, Rakic P, Noonan JP. Evolutionary changes in promoter and enhancer activity during human corticogenesis. Dataset GSE63649. 2015. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63649>.
118. Cotney J, Noonan JP. The evolution of lineage-specific regulatory activities in the human embryonic limb. Dataset GSE42413. 2013. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE42413>.
119. Mifsud B, Schoenfelder S, Fraser P. Promoter capture mESC and fetal liver. Dataset E-MTAB-2414. *ArrayExpress.* 2015. <https://www.ebi.ac.uk/biostudies/arrayexpress/studies/E-MTAB-2414>.
120. Dekker J. MicroC Libraries from H1ESC cells from 2 biological replicates, FA and DSG cross-linked, MNase digestion. Dataset 4DNES21D8SP8. 4DNucleome. 2019. <https://data.4dnucleome.org/experiment-set-replicates/4DNES21D8SP8>.
121. Darzacq X. Micro-C on JM8.N4 WT. Dataset 4DNES14CNC1I. 4DNucleome. 2019. <https://data.4dnucleome.org/experiment-set-replicates/4DNES14CNC1I>.
122. Kocher, AA. CpG island turnover events predict evolutionary changes in enhancer activity. Github. 2023. [https://github.com/NoonanLab/Kocher\\_et\\_al\\_oCGI\\_turnover](https://github.com/NoonanLab/Kocher_et_al_oCGI_turnover).
123. Kocher AA. CpG island turnover events predict evolutionary changes in enhancer activity. 2024. Zenodo. <https://doi.org/10.5281/zenodo.11236169>.

## Publisher's Note

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