



Epigenomic differences in the human and chimpanzee genomes are associated with structural variation

Xiaoyu Zhuo, Alan Y Du, Erica C Pehrsson, et al.

Genome Res. published online December 10, 2020

Access the most recent version at doi:[10.1101/gr.263491.120](https://doi.org/10.1101/gr.263491.120)

P<P	Published online December 10, 2020 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Epigenomic differences in the human and chimpanzee genomes are**
2 **associated with structural variation**

3 Xiaoyu Zhuo^{1,2}, Alan Y. Du^{1,2}, Erica C. Pehrsson^{1,2}, Daofeng Li^{1,2} and Ting Wang^{1,2,3}

4 **1** Washington University School of Medicine in St. Louis, Department of Genetics, Saint Louis,
5 Missouri, USA

6 **2** The Edison Family Center for Genome Sciences and Systems Biology, Washington University
7 School of Medicine, St Louis, MO, USA

8 **3** McDonnell Genome Institute, Washington University School of Medicine, St Louis, MO, USA

9

10

11

12

13

14

15

16

17 **ABSTRACT**

18 Structural variation (SV), including insertions and deletions (indels), is a primary mechanism of
19 genome evolution. However, the mechanism by which SV contributes to epigenome evolution
20 is poorly understood. In this study, we characterized the association between lineage-specific
21 indels and epigenome differences between human and chimpanzee to investigate how SVs
22 might have shaped the epigenetic landscape. By intersecting medium-to-large human-
23 chimpanzee indels (20bp-50kb) with putative promoters and enhancers in cranial neural crest
24 cells (CNCC) and repressed regions in induced pluripotent cells (iPSC), we found that 12% indels
25 overlap putative regulatory and repressed regions (RRRs), and 15% of these indels are
26 associated with lineage-biased RRRs. Indel-associated putative enhancer and repressive regions
27 are ~1.3 and ~3 times as likely to be lineage-biased, respectively, as those not associated with
28 indels. We found a 2-fold enrichment of medium-sized indels (20bp to 50bp) in CpG island
29 (CGI)-containing promoters than expected by chance. Lastly, from human-specific transposable
30 element insertions, we identified putative regulatory elements, including NR2F1-bound
31 putative CNCC enhancers derived from SVAs and putative iPSC promoters derived from LTR5s.
32 Our results demonstrate that different types of indels are associated with specific epigenomic
33 diversity between human and chimpanzee.

34

35

36 INTRODUCTION

37 The question of what makes us uniquely human has long been of interest (Darwin 1871).
38 Comparative genomics has sought the genetic basis of human-specific traits (King and Wilson
39 1975; The Chimpanzee Sequencing and Analysis Consortium 2005; Kronenberg et al. 2018;
40 Rogers and Gibbs 2014; Wall 2013), including human-specific gene gain/loss or regions under
41 accelerated evolution (Pollard et al. 2006; Zhu et al. 2007; Enard et al. 2002; Atkinson et al.
42 2018; Florio et al. 2018; Fiddes et al. 2018; Suzuki et al. 2018; Franchini and Pollard 2017). In
43 addition, epigenetic and transcriptomic differences also contribute to human-specific
44 phenotypes (Prescott et al. 2015; Gallego Romero et al. 2015; Trizzino et al. 2017; Ward et al.
45 2018; Danko et al. 2018; Pai et al. 2011; Hernando-Herraez et al. 2013; Eres et al. 2019).
46 However, how structural variations (SVs) affect human-specific functions is just beginning to be
47 explored (Gordon et al. 2016; Fudenberg and Pollard 2019).

48 SVs, which include deletions, duplications, inversions, insertions, and translocations, are
49 responsible for the majority of genetic differences within populations and between species. The
50 1000 Genomes Project estimated that an individual carries a median of 8.9 Mb of SVs versus 3.6
51 Mb of single nucleotide variants (SNVs) (Sudmant et al. 2015). Long-read sequencing of two
52 haploid human genomes revealed that the majority of SVs were novel, suggesting that the
53 impact of SVs is underestimated (Huddleston et al. 2017). SVs also contribute to inter-species
54 divergence. In 2005, the Chimpanzee Sequencing and Analysis Consortium reported ~90 Mb of
55 insertions or deletions (indels) between human and chimpanzee; in contrast, SNVs constituted
56 only ~35 Mb (The Chimpanzee Sequencing and Analysis Consortium 2005).

57 Non-coding cis-regulatory elements (CREs) play a critical role in gene regulation (The ENCODE
58 Project Consortium 2004). One powerful method to identify putative functional elements is
59 epigenomic profiling (Ernst and Kellis 2012; Roadmap Epigenomics Consortium et al. 2015). For
60 example, H3K4me3 is usually associated with promoters, and H3K27ac is associated with both
61 active promoters and active enhancers. In contrast, H3K9me3 is associated with
62 heterochromatin, a repressed state characterized by densely packed DNA and low gene
63 expression (Becker et al. 2015). By applying epigenomic profiling to related organisms
64 (“comparative epigenomics”), we can compare epigenetic signature across species between
65 syntenic regions and investigate the birth and death of regulatory elements during evolution
66 (Xiao et al. 2012; Lowdon et al. 2016). Although studies have investigated enhancer evolution
67 between human and chimpanzee (Trizzino et al. 2017; Gallego Romero et al. 2015; Ward et al.
68 2018; Prescott et al. 2015), the impact of SV on these elements has not been studied.

69 Here, we developed a novel computational strategy to define syntenic regions that contain
70 indels. Using publicly available epigenomic datasets from human and chimpanzee, we defined
71 the association between indels and epigenetic differences between the two species. We
72 explored both epigenomic conservation and innovation in association with medium to large
73 indels (20bp to 50kb), and how lineage-specific transposable element (TE) insertions contribute
74 to new putative functional elements. Our findings indicate that SVs and epigenomic changes
75 between human and chimpanzee are significantly interrelated.

76

77 **RESULTS**

78 **Development of a novel method to find orthologous regions overlapping large** 79 **indels**

80 Conventional comparative genomic/epigenomic methods rely on tools such as UCSC liftOver to
81 retrieve syntenic regions between species based on alignments between genome assemblies
82 (Kuhn et al. 2013). However, these tools usually fail to return syntenic regions when the
83 synteny is disrupted by medium to large SVs. To overcome this obstacle, we developed a new
84 pipeline that combines CrossMap (Zhao et al. 2014), an alternative to liftOver, with our newly
85 developed tool called OrthoINDEL (Methods). Instead of filtering syntenic regions using the
86 minimum percentage of bases that can be converted to the new assembly, CrossMap outputs
87 the syntenic region as multiple blocks split by alignment gaps. OrthoINDEL then concatenates
88 the fragmented orthologs output by CrossMap if they are continuous or separated only by
89 indels (Fig. 1). This way, our pipeline stringently converts regions from the source genome to
90 their syntenic coordinates in the target genome even if they overlap large indels. In contrast,
91 without sacrificing specificity, the UCSC liftOver dismisses regions with a large fraction absent in
92 the target genome. To illustrate this improvement, we performed the same genomic
93 coordinates conversion from human to chimpanzee using either OrthoINDEL or liftOver. We
94 found OrthoINDEL successfully converted ~1000 more regions that were dismissed as "partially
95 deleted" or "split in new" by liftOver (Supplemental_Fig_S1, Supplemental_Table_S1). Thus,
96 OrthoINDEL is better suited to retrieve syntenic regions overlapping indels.

97 **Indel-associated enhancers and H3K9me3 regions are more likely to be lineage-**
98 **biased**

99 By analyzing two human-chimpanzee comparison ChIP-seq datasets (Prescott et al. 2015, Ward
100 et al. 2018), we identified ~15,000 putative promoters (H3K4me3 ChIP-seq peaks), ~27,000
101 putative enhancers (H3K27ac peaks outside of H3K4me3 peaks) and ~31,000 H3K9me3 regions
102 (H3K9me3 broad peaks) in each species that constitute ~40Mb, 67Mb and 300Mb, respectively
103 (Fig. 2A). Together, we define them as putative regulatory and repressed regions (RRRs). We
104 found support for ~88% of putative CNCC promoters by overlapping them with annotated
105 FANTOM5 or GENCODE promoters (~84% and ~85% in GENCODE and FANTOM5, respectively)
106 (Frankish et al. 2019; The FANTOM Consortium and the RIKEN PMI and CLST (DGT) 2014).

107 Next, we applied the OrthoINDEL pipeline to establish the syntenic relationship between
108 human and chimpanzee RRRs. We further classified regions as human-biased, chimpanzee-
109 biased, or invariant based on both peak-calling and ChIP-seq reads difference defined by
110 DESeq2 (Methods; Supplemental_Table_S2) (Landt et al. 2012; Love et al. 2014). A datahub
111 displaying the processed data over human and chimpanzee syntenic regions is available at the
112 WashU Epigenome Browser (Li et al. 2019)
113 (<https://epigenomegateway.wustl.edu/browser/?sessionFile=https://wangftp.wustl.edu/~xzhu>
114 [o/CNCC/publicationSession.json](https://epigenomegateway.wustl.edu/browser/?sessionFile=https://wangftp.wustl.edu/~xzhu)).

115 We annotated all indels larger than 20bp between the human (hg38) and chimpanzee
116 (panTro5) reference genomes using the DASVC pipeline (Methods) (Gordon et al. 2016). In
117 total, we defined 193,180 medium-to-large indels (20bp to 50kbp) encompassing 95.8Mb (~3%

118 of the human haploid genome). We selected ~127,000 of them (42.2Mb) with a defined
119 ancestral state (using gorilla genome as an outgroup) and located within nonrepetitive regions
120 for epigenomic analysis (Methods) (Supplemental_Fig_S3). We also annotated TE-derived
121 insertions within these indels (Methods) (Fig. 2B, Supplemental_Table_S3). The overall number
122 and length of our indels were in excellent agreement with previously published results (The
123 Chimpanzee Sequencing and Analysis Consortium 2005; Kronenberg et al. 2018)
124 (Supplemental_Fig_S4).

125 Next, we characterized the association between indels and putative RRRs. We identified
126 ~15,000 indel-overlapping putative RRRs by intersecting their coordinates using BEDTools
127 (Quinlan and Hall 2010) (Supplemental_Table_S4). Indels are slightly depleted in putative RRRs
128 instead of being uniformly distributed in the genome (Fisher's exact test enrichment ratio 0.94,
129 P-value 2.4×10^{-9}). We found that ~88% (112,433/127,350) of indels do not overlap any
130 putative RRRs in the study (Supplemental_Table_S5), 10% of indels overlap with invariant
131 elements, while the remaining 1.8% (2267 indels) are associated with lineage-biased elements
132 (Fig. 2C). An association between a human-specific indel and a human-biased putative RRR (135
133 with enhancer and 801 with H3K9me3 heterochromatin) could suggest that the indel may have
134 created the RRR in the human lineage. On the other hand, association between a human-
135 specific indel and a chimpanzee-biased element (149 with enhancer and 173 with H3K9me3
136 heterochromatin) could suggest that the indel may have disrupted an ancestral RRR in the
137 human lineage. The same logic applies to chimpanzee lineage indels (Fig. 3,
138 Supplemental_Table_S5).

139 In accordance with previously reported findings, we found that all except five putative
140 promoters are invariant between the two species (Fig. 2C). In contrast, ~85% of putative
141 enhancer and repressed regions are invariant (Prescott et al. 2015; Ward et al. 2018) (Fig. 2C).
142 Compared with those not associated with indels, putative enhancers associated with indels are
143 ~30% more likely to be lineage-biased (Fisher's exact test P-value 4.7×10^{-6}), and H3K9me3
144 regions associated with indels are about three times as likely to be lineage-biased (Fisher's
145 exact test P-value 10^{-760}) (Fig. 2C). This result suggests that indels have a moderate association
146 with putative enhancers and a strong association with H3K9me3 regions.

147 **The enrichment of different indels with putative RRRs**

148 We sought to understand if different size categories of indels had different association with
149 putative RRRs. We divided non-TE-derived indels to four groups (20-50bp, 50-500bp, 500-5kb,
150 5k-50kb) and separated TE-derived insertions by TE class. We defined indels from 20 to 50bp as
151 medium-sized indels, and indels ≥ 50 bp as large indels following conventions (The 1000
152 Genomes Project Consortium 2015; Kronenberg et al. 2018). We calculated the enrichment of
153 the intersection between these indel categories with putative RRRs over the genomic
154 background using Fisher's exact test (Fig. 4). We plotted the enrichment ratio and p-value of all
155 pairs with at least one intersection (Fig. 4).

156 We noticed three main trends. First, medium-sized indels (20-50bp) are enriched in invariant
157 putative CNCC promoters. In contrast, indels larger than 500bp are depleted in invariant
158 putative promoters. Second, we found indels longer than 5kb are depleted in invariant putative
159 CNCC enhancers. Third, in H3K9me3 repressed regions, lineage-specific sequences longer than

160 500bp (insertions > 500bp in that lineage and deletions > 500bp in the other lineage) are
161 enriched for H3K9me3 regions from the same lineage (Fig. 4). As an example, we found that
162 human lineage insertions and chimpanzee lineage deletions > 500bp are enriched in human-
163 biased H3K9me3 regions. None of the lineage-biased putative CNCC promoters overlap indel
164 (Fig. 4), and the enrichment/depletion of indels with putative invariant promoters is almost
165 identical to their enrichment/depletion with all putative CNCC promoters.

166 We separated TE-derived insertions by TE class and performed the same enrichment analysis
167 (Fig. 4). Lineage-specific SVA insertions are significantly enriched in both putative lineage-biased
168 CNCC enhancers and iPSC H3K9me3 repressed regions for both species, which implies that
169 newly inserted SVA elements may have provided CNCC-specific enhancers and been targeted
170 by the repressive marks in both species (Fig. 4). ERV insertions are only enriched as lineage-
171 biased H3K9me3 repressed regions in the chimpanzee lineage, contributed mainly by the
172 H3K9me3-modified chimpanzee-specific PTERV regions (Supplemental_Table_S4) (Yohn et al.
173 2005). Previous studies in fly and yeast suggest that repressive marks can spread beyond the
174 heterochromatin boundary and affect nearby genes (Elgin and Reuter 2013; Obersriebnig et al.
175 2016; Greenstein et al. 2018). However, we found here that H3K9me3 marks rarely expand
176 beyond 2kb outside newly inserted TE boundaries (Supplemental_Fig_S5). In general, out of the
177 SVs associated with biased chromatin marks, insertions tend to be associated with creation,
178 rather than destruction, of putative RRRs.

179 We further explore the enrichment of indels from 20-50bp with putative CNCC promoters and
180 the enrichment of SVA insertions with putative CNCC enhancers in the following sections.

181 **Enrichment of medium-sized indels within CpG islands**

182 Promoters are considered conserved elements due to their low nucleotide substitution rate
183 (The FANTOM Consortium and the RIKEN PMI and CLST (DGT) 2014). We also found both
184 putative CNCC promoters and GENCODE-annotated promoters are conserved using phastCons
185 score (Supplemental_Fig_S6A, B). In contrast, using Fisher's exact test, we found that medium-
186 sized (20-50bp) indels are highly enriched within putative CNCC promoters, whereas indels
187 >500bp are depleted within putative CNCC promoters (Fig. 4). To characterize the relationship
188 between indel size and their enrichment in promoters at a finer resolution, we separated indels
189 <500bp by size at 50bp intervals and tested their enrichment with putative CNCC promoters
190 using Fisher's exact test. We found that 50-100bp indels are barely enriched, and all indel bins
191 >100bp are not statistically enriched (Supplemental_Fig_S7). To validate our observation, we
192 further calculated indel frequency around annotated genes for 20-50bp and 50-100bp indels.
193 Again, we found that 20-50bp indels, but not 50-100bp indels, have elevated frequency
194 immediately upstream of transcription start sites (Fig. 5A). We also observed slightly lower
195 indel frequency in more conserved promoters (phastCons >0.2) than in less conserved
196 promoters (phastCons <0.2), indicating that conserved promoters have fewer indels
197 (Supplemental_Fig_S6C, D). To test whether the enrichment of 20-50bp indels in promoters
198 also exists in the human population, we extracted variants from the 1000 Genomes Project
199 Phase 3 release (The 1000 Genomes Project Consortium 2015; Sudmant et al. 2015) and
200 repeated the analysis. Consistent with our inter-species observation, only medium-sized indels
201 (20-50 bp) within the human population are enriched in the promoters (Fig. 5B).

202 Seventy percent of human gene promoters contain CpG islands (CGIs) (Saxonov et al. 2006).
203 Many indels in promoters are located within UCSC annotated CGIs, suggesting that indels within
204 CGIs could be driving the observed high indel rate in promoters (Kuhn et al. 2013). We
205 separated promoters into CGI and non-CGI promoters, and found that 20-50bp indels are
206 enriched only in CGI promoters (enrichment ratio 2.60, P-value 10^{-762}) but not in non-CGI
207 promoters (enrichment ratio 1.06, P-value 0.02).

208 Since indels > 150bp are not enriched within putative CNCC promoters (Supplemental_Fig_S7),
209 we analyzed the enrichment of 20-50bp, 50-100bp, and 100-150bp indels with promoters and
210 CGIs. To verify the Fisher's exact test results (Fig. 4 and Supplemental_Fig_S7), we directly
211 compared observed number of intersections with the expected distribution based on random
212 sampling (Fig. 5C) (Methods). CGIs outside promoters are enriched for all three sizes of indels (P
213 < 0.001, permutation test), and the enrichment of indels in CGIs within promoters decreases as
214 a function of indel size from 20bp to 150bp. In contrast, promoters without CGIs are barely
215 enriched for any indel category (Fig. 5C). Thus, consistent with previous findings, our data
216 suggests that CGIs are hot spots for indels in evolution (Tian et al. 2011; Kiktev et al. 2018).

217 Indels from 100bp to 150bp are less enriched in CGIs within promoters than the smaller 20-
218 50bp and 50-100bp indels, potentially because they may drastically influence regulatory activity
219 and have thus been selected against during evolution. Medium-sized indels (20-50bp), on the
220 other hand, seem to have been tolerated, as promoters with CGIs are highly epigenetically
221 conserved between human and chimpanzee (Fig. 4).

222 Lastly, we characterized the enrichment of indels in CGIs within the human population using the
223 variants from 1000 Genomes Project Phase 3 release and found the same pattern
224 (Supplemental_Fig_S8) (The 1000 Genomes Project Consortium 2015). We also characterized
225 the CGI enrichment of small indels (< 20bp). Small indels (1-20bp), especially indels with length
226 <= 6bp, are depleted in CGIs. Since 98% of 1000 Genome Project indels are smaller than 20bp,
227 our conclusion is consistent with the previous observation that indels are depleted within CGIs
228 (Neininger et al. 2019), but highlights the novel finding of enrichment of medium sized indels in
229 CGIs.

230 **Lineage-specific TEs give rise to putative promoters and enhancers and** 231 **gradually become repressed during evolution**

232 Lineage-specific SVA insertions are highly enriched in both putative enhancers and repressed
233 regions (Fig. 4). For the 37 SVA elements that overlap with human-biased putative enhancer
234 regions, the aggregated H3K27ac ChIP-seq profile shows areas of elevated signal similar to that
235 in putative CNCC enhancers (Supplemental_Fig_S9A, B). To better illustrate their enrichment
236 pattern, we plotted aggregated ChIP-seq profiles around all human-specific SVAs and compared
237 them to the profiles of their orthologous pre-insertion sites in chimpanzee (Fig. 6A). In addition
238 to the previously described ChIP-seq data, we also included iPSC H3K27ac from both species
239 and human iPSC H3K4me3 (The ENCODE Project Consortium 2004; Gallego Romero et al. 2015).
240 We also analyzed human-specific LTR5 insertions in a similar fashion (Supplemental_Fig_S9C,
241 D).

242 We used the mappability score as a measurement of repetitiveness and the propensity of a
243 genomic region to produce uniquely mappable reads (Derrien et al. 2012). Lineage-specific SVA
244 and LTR5 have low mappability scores (Fig. 6A). Therefore, their corresponding ChIP-seq signals
245 were likely underestimated. Indeed, with the exception of H3K9me3 ChIP-seq, which was
246 sequenced using 100bp paired-end reads, all other ChIP-seq datasets generated using single-
247 end reads have close to zero signal over the low mappability regions. Nevertheless, we found a
248 strong H3K27ac signature suggesting putative enhancer activity on the 3' flanking region of SVA
249 in CNCC and a strong promoter signature (H3K4me3 and H3K27ac) on both flanking regions of
250 LTR5 in iPSCs in human but not in chimpanzee. Conversely, we found H3K4me3 ChIP-seq signal
251 3' of LTR5 insertions in CNCC, but these were found in both human and chimpanzee, suggesting
252 that in this case the epigenomic mark may be independent of the TE insertion (Fig. 6A).

253 To understand how epigenetic profiles might evolve over time, we extended our analysis to
254 related TE insertions that are shared by both species. SVA proliferated and diverged in the
255 human genome in a similar fashion as the amplification of L1 (Khan et al. 2006, Hancks and
256 Kazazian 2016). SVAs in the human genome were classified into six subfamilies (SVA-A to SVA-
257 F). The expansion of subfamilies from SVA-A to SVA-D predates the human-chimpanzee split,
258 and SVA-E and SVA-F expanded after human-chimpanzee divergence (Wang et al. 2005). We
259 defined all human-specific SVAs as SVA-human and classified human-chimpanzee shared SVAs
260 based on their subfamilies (from SVA-A to SVA-D). We plotted H3K27ac, H3K9me3 and
261 mappability profiles of different SVA subfamilies in Fig. 6B. We found that the 3' boundary
262 H3K27ac signal in CNCC decreases as SVA subfamily ages, with the greatest signal in human-

263 specific subfamilies. In contrast, H3K9me3 signal intensifies with increasing SVA age in iPSC (Fig.
264 6B).

265 Next, we similarly classified LTR5s based on lineage-specificity and subfamilies (from oldest to
266 youngest: LTR5A-shared, LTR5B-shared, and LTR5Hs-shared to LTR5-human) and performed the
267 same analysis using H3K4me3 data in iPSC. Similar to SVAs, the strength of the LTR5 promoter
268 signature in iPSC is negatively correlated with age. However, LTR5s are not marked by H3K9me3
269 in iPSC (Fig. 6C).

270 **NR2F1 binding is correlated with enhancer signature on the 3' end of SVA in** 271 **CNCC**

272 Our data thus far predicted a putative enhancer within SVAs close to their 3' end. The H3K27ac
273 signal is likely “hidden” due to low mappability, but part of it extends into mappable regions, as
274 we observed in Fig. 6. However, we cannot observe the boundary ChIP-seq signal at the 5' end
275 of a full-length (1.6kb) SVA, which might be too long for the enhancer signal to extend beyond.
276 To determine the precise location of the putative enhancer, we extracted all SVAs with
277 complete 3' ends in the human genome, sorted them based on size, anchored them at the 3'
278 end, and annotated them with CNCC H3K27ac signal (Fig. 7A). Again, we observed the elevated
279 H3K27ac signal on the 3' end of SVAs (Fig. 7A). However, once the SVA length was reduced to
280 300-500 bp, a similar boundary H3K27ac signal on the 5' end emerged. This result is consistent
281 with the hypothesis that the enhancer signature originating within SVAs can extend beyond the
282 low mappability region. The boundary H3K27ac signal disappeared on both ends of shorter SVA
283 copies, suggesting that further truncation resulted in a loss of the internal enhancer (Fig. 7A, B).

284 Zooming in onto the 5' end of these shorter SVA fragments revealed a strong NR2F1 binding
285 motif (Fig. 7B, C). Importantly, the disappearance of the boundary H3K27ac signal correlated
286 with the truncation of the NR2F1 motif (Fisher's exact test P-value 7×10^{-5}) (Fig. 7C). NR2F1 is a
287 critical regulator in CNCC (Rada-Iglesias et al. 2012). Indeed, the NR2F1 and H3K27ac ChIP-seq
288 signals co-occur in SVAs in CNCC (Fig. 7B) (Prescott et al. 2015). However, not every putative
289 human-biased CNCC SVA enhancer has a related NR2F1 peak (Fig. 7D, E), suggesting the
290 possible involvement of other transcription factors in these SVA-derived putative enhancers.

291

292 **DISCUSSION**

293 In this study, we systematically characterized how medium to large indels are correlated with
294 differences in the epigenome in the human and chimpanzee lineages. We found that indels are
295 enriched in putative lineage-biased enhancers and H3K9me3 repressed regions. We should
296 note that all genomic enrichments were estimated using a whole genome random distribution
297 as the background. However, genomic features are not randomly distributed and the
298 assumption may not always be appropriate. Yokoyama et al. described a substitution-based
299 framework to model birth-death of lineage-specific functional elements (Yokoyama et al. 2014).
300 We demonstrated here that in addition to substitutions, indels can also contribute to the birth-
301 death of regulatory and repressed elements. Our strategy is readily applicable to other
302 comparative epigenomic datasets, and we have made our processed data available in our
303 comparative browser (Li et al. 2019).

304 Mutation rate varies depending on genomic region. It has been reported that the substitution
305 rate in closed chromatin regions is higher than the rate in open chromatin regions (Makova and
306 Hardison 2015; Fortin and Hansen 2015). Lunter et al. found the highest indel rates in regions
307 with extremely high and extremely low GC content (Lunter et al. 2006). Using macaque as an
308 outgroup, Kvikstad et al. reported a curvilinear relationship between human lineage indel rate
309 and GC content, and they also found weak anti-correlation between insertion rate and number
310 of CGIs at 1Mb genomic windows (Kvikstad et al. 2007). Specific to CGIs, the substitution rate at
311 CpG sites is higher than the mutation rate at other sites because of the high deamination rate
312 of 5-methyl cytosine (Coulondre et al. 1978). Cohen et al. found that the lack of methylation of
313 CGIs can explain their maintenance without implying purifying selection in primates (Cohen et
314 al. 2011). However, CGIs are often associated with genome instability (Deaton and Bird 2011;
315 Du et al. 2014).

316 We found that CGIs are hotspots for medium/large indels in hominids. Kiktev et al.
317 demonstrated that the high-GC region in yeast has a high deletion/duplication rate resulting
318 from DNA polymerase slippage (Kiktev et al. 2018). Thus, analogous to well-characterized SVs in
319 human coding exons (Montgomery et al. 2013; Challis et al. 2015), high-GC regions are prone to
320 forming a single stranded DNA loop due to polymerase slippage during DNA replication (Tian et
321 al. 2011), which likely results in the increased indel rate. However, GC rich regions are prone to
322 sequencing error and we cannot completely rule out the possibility that some indels we called
323 were caused by the difficulty of sequencing and calling variants in these GC rich regions.

324 By comparing sequence conservation with epi-conservation, Xiao et al. reported elevated epi-
325 conservation of H3K27ac, H3K27me3, and methylated CpGs (but not H3K4me3) in rapidly
326 evolving sequences and proposed that epi-conservation could buffer some deleterious
327 mutations (Xiao et al. 2012). Here, we report the conservation of H3K4me3 marks between
328 human and chimpanzee despite an elevated rate of medium-sized indels, further supporting
329 the concept of epi-conservation and its potential role in buffering the impact of mutations.

330 Britten and Davidson proposed the gene battery model in the 1970s to explain the evolution of
331 regulatory networks (Britten and Davidson 1971). They proposed that a single “activator gene”
332 can control a “battery of genes” by interacting with diffused repetitive sequences throughout
333 the genome. Since then, TEs have been repeatedly demonstrated to contribute novel
334 regulatory elements (Feschotte 2008; Wang et al. 2007; Lynch et al. 2011; Chuong et al. 2017).
335 However, in primates, especially in the human lineage, there have been conflicting reports
336 about the regulatory role of TEs. Trizzino et al. found specific TE subfamilies enriched in liver
337 enhancers (Trizzino et al. 2017). On the other hand, Ward et al. could not find significant
338 contribution of TEs to gene regulation in pluripotent stem cells (Ward et al. 2018). We report
339 clear signals of TE-derived, tissue-specific putative enhancers and promoters unique to human
340 or chimpanzee. We identified LTR5 as putative promoters in iPSC, while Fuentes et al. found
341 them to have enhancer activity in human embryonal carcinoma NCCIT cells (Fuentes et al.
342 2018). We could not find a large impact on nearby gene expression associated with these TE-
343 derived enhancers with our limited dataset. It is possible that these new enhancers do not
344 regulate the closest genes. Alternatively, they may provide functional redundancy instead of
345 inventing new regulation (Osterwalder et al. 2018; Choudhary et al. 2020). We also report that

346 TE-associated heterochromatin displays limited spreading (Supplemental_Fig_S5). We found a
347 rapid conversion of TE epigenetic modification from active to repressive states as a function of
348 age in young TEs. This discovery echoes a previous report of the transition of repressive marks
349 from cytosine methylation to H3K9me3 as ERVs age in the human genome (Ohtani et al. 2018).
350 These data suggest that although many TEs carry regulatory elements, the host rapidly and
351 continuously silences such activity during evolution.

352 Most genomic analyses rely on second-generation sequencing, which produces short reads,
353 restricting our ability to detect signals from repetitive, low mappability regions. To overcome
354 this limitation, we investigated not only the epigenetic signal from within TEs but also from
355 flanking regions. By comparing TE insertions with orthologous pre-insertion sites in another
356 species, we can infer that the epigenetic signal originates from these highly repetitive regions.
357 Our approach expands the application of second-generation sequencing and reveals that there
358 are more potentially functional elements hidden in unmapped territories. However, the
359 sensitivity of our method is limited by the distance of the element to the boundary, read length,
360 DNA fragment size and other factors.

361 Only four different TEs have been actively transposing in the human lineage. Of the four, we
362 found that two are associated with putative enhancers and promoters using data from only two
363 cell types. Although most TEs are neutrally evolving in the genome, our discovery suggested
364 that many CREs carried by TEs were active upon insertion. Our finding begs more thorough
365 investigation of CREs derived from recently inserted TEs in more cell types and between
366 different species.

367 **METHODS**

368 **Indel identification**

369 We applied the DASVC tool to annotate indels between human and chimpanzee (Gordon et al.
370 2016). The DASVC tool was downloaded from (<https://github.com/zeeev/DASVC>) and was used
371 with default parameters to identify 20bp to 50kb indels between the hg38 and panTro5
372 genomes. We further processed DASVC output using a Python script (`refine_calledSV.py`) to
373 remove segmental replacements and extract the exact coordinates in both species. To identify
374 indels that occurred in mappable regions, we calculated the average 75bp-mappability score of
375 the flanking 200bp of all indels in both species, and selected those indels with a mappability
376 score > 0.7 in both flanking regions in both species (Derrien et al. 2012).

377 To differentiate deletion in one lineage from insertion in the other, we identified orthologous
378 coordinates for all indels in the gorilla reference genome (`gorGor5`). We considered an indel to
379 be a deletion if the indel region is present in the gorilla genome, and an insertion if the region is
380 absent in the gorilla.

381 A new chimpanzee reference genome `panTro6` was published shortly after we started this
382 project (Kronenberg et al. 2018). `PanTro6` closed 52% of remaining gaps, but it does not refute
383 the high-quality `panTro5` reference genome in assembled regions. Therefore, the conclusions
384 we reach here using `panTro5` should remain valid.

385 **TE-derived insertion annotation**

386 To annotate TE-derived insertions within the identified indels, we intersected the indel list with
387 RepeatMasker annotations for both the hg38 and panTro5 genomes using BEDTools (Smit et al.;
388 Quinlan and Hall 2010). To avoid calling fragmented TEs as separate TE insertion events, we
389 defined an indel as a TE-derived insertion if it was derived from a single TE insertion event.
390 Indels containing only an *Alu* element, a full-length endogenous retrovirus (ERV), or a solo long
391 terminal repeat (LTR) were counted as *Alu*/ERV insertions. Solitary LTRs were included because
392 they are derived from an ERV insertion followed by non-homologous recombination (Mager
393 and Stoye 2015). Due to the prevalence of 5' end truncation during target-primed reverse
394 transcription (TPRT) (Luan et al. 1993), we also tolerated incomplete 5' ends in defining L1 and
395 SVA insertions. One limitation to our rigorous approach is that we did not annotate lineage-
396 specific solo-LTRs, where part of the solo-LTR aligned with the 5' end LTR of the full-length ERV
397 and part aligned with the 3' end, as TE-derived insertions. In addition to the known actively
398 transposing TE subfamilies described above, we also identified a few lineage-specific LTR12C
399 elements (Supplemental_Table_S3).

400 Since TE insertions are homoplasmy-free and unidirectional (Bashir et al. 2005; Ray et al. 2006),
401 the orthologous regions corresponding to most human-/chimpanzee-specific TE insertions
402 should be found as pre-insertion sites in the gorilla genome. As expected, the orthologous
403 locations of 98% (15,803 of 16,068) of lineage-specific TE insertions are pre-insertion sites in
404 the gorilla genome (Supplemental_Fig_S4B), whereas the remaining 2% were present as TE

405 insertions. These cases could be explained by incomplete lineage sorting, as demonstrated
406 before (Kronenberg et al. 2018; Ray et al. 2006).

407 **Peak calling and cross-species comparison**

408 We downloaded both CNCC and iPSC raw read FASTQ files from NCBI GEO repository,
409 accessions GSE70751, GSE61343 and GSE96712 (Prescott et al. 2015; Gallego Romero et al.
410 2015; Ward et al. 2018). Human H3K27ac iPSC ChIP-seq data and the associated input BAM files
411 mapped to hg38 were downloaded from the ENCODE portal experiment ENCSR729ENO. We
412 called ChIP-seq peaks using ENCODE recommendations with MACS2 and IDR thresholding
413 (Supplemental Methods) (Li et al. 2011; Landt et al. 2012; Li 2013).

414 We applied CrossMap to identify orthologous segments for each peak in the other species
415 (Zhao et al. 2014). For each peak region, CrossMap outputs all fragmented orthologous loci
416 separated by any indel >1 bp. We developed a new tool, OrthoINDEL, that processes
417 fragmented syntenic regions from the CrossMap output file. OrthoINDEL uses two parameters
418 to filter fragmented regions. A maximum distance of 50kb was used to filter out fragments with
419 too large a separation. We used a minimum distance of 50bp to define continuous fragments.
420 To be defined as an indel, the split fragments we required to be continuous in one species.
421 Since we focus on indels, our pipeline removes other SVs including inversions. Lastly, we
422 filtered out regions with average 50bp-mappability < 0.7 in either species (compared with indel
423 identification, we used more stringent mappability criteria here to eliminate false positive
424 lineage-biased RRRs) (Derrien et al. 2012). The last filtering step is critical to filter out false
425 positive lineage-biased regions (Supplemental_Fig_S10). With this pipeline, we defined

426 stringent 1:1 syntenic regions between human and chimpanzee tolerating indels up to 50kb
427 that can be converted reciprocally using OrthoINDEL.

428 Peak calling is sensitive to sequencing coverage and background signal. To better differentiate
429 invariant peaks from lineage-biased peaks, we counted the number of reads mapped to each
430 peak in both species and applied DESeq2 to quantify peak intensity difference (Love et al. 2014).
431 We classified regions as "human-biased" if a peak was called only in the human genome by
432 MACS2 and the number of ChIP-seq reads in the human peak is significantly higher than the
433 number in its orthologous region in the chimpanzee genome (DESeq2, $q < 0.0001$); "chimpanzee-
434 biased" regions were defined in a similar fashion. Lastly, we classified regions as "invariant" if
435 ChIP-seq peaks were called in both species or if a peak was called in only one species, but the
436 difference of ChIP-seq reads number between syntenic regions is not significant.

437 **Enrichment analysis**

438 We calculated the number of indel-RRR overlaps using the BEDTools intersect function (Quinlan
439 and Hall 2010). To perform Fisher's exact test with the hg38 genome as background, we used
440 the BEDTools fisher function. The permutation test used BEDTools to shuffle indel coordinates
441 and intersect with CGIs and promoters (Supplemental Methods). Briefly, we counted the
442 number of intersections of human-chimpanzee indels ranging from 20-50bp, 50-100bp and 100-
443 150bp with CGIs within promoters, CGIs outside of promoters and promoters without CGIs,
444 respectively. We then randomly shuffled indel coordinates 1000 times and repeated the
445 intersection.

446 **Identification of transcription factor binding motifs**

447 We used FIMO from the MEME suite to find potential transcription factor binding sites within
448 SVA elements (Bailey et al. 2009; Grant et al. 2011).

449 **Data visualization**

450 We generated bigWig files using methylQA and displayed them on the WashU Epigenome
451 Browser (Li et al. 2015). All data are visualized on the WashU Epigenome Browser (Li et al.
452 2019). All ChIP-seq data were normalized to reads per genomic content (RPGC) using deepTools
453 bamcoverage (Ramirez et al. 2016). Binding profiles and heatmaps were generated using
454 deepTools2 (Supplemental Methods) (Ramirez et al. 2016).

455 **DATA ACCESS**

456 All processed data are accessible on the WashU Comparative Epigenome browser:
457 (<https://epigenomegateway.wustl.edu/browser/?sessionFile=https://wangftp.wustl.edu/~xzhuo/CNCC/publicationSession.json>). Our pipeline and scripts generated in this study are available
458 on GitHub (https://github.com/xzhuo/indel_epi_landscape) and as Supplemental Code.
459

460 **ACKNOWLEDGMENTS**

461 We thank all members of the Wang lab for their helpful suggestions; Dr. Zev Kronenberg from
462 PacBio for his help with the DASVC pipeline and sharing with us their apes indel variant-calling
463 result; and Silas Hsu from UIUC for building the updated WashU Epigenome Browser. Finally,

464 we want to thank three reviewers for their constructive criticisms. This manuscript is dedicated
465 to the memory of Yimin Yang (1960--2018).

466 X.Z. is supported in part by 5R25DA027995.

467 A.Y.D. is supported by a grant from NHGRI (no. T32 HG000045).

468 E.C.P. is supported by a Postdoctoral Fellowship, PF-17-201-01, from the American Cancer
469 Society.

470 T.W. is supported by NIH grants R01HG007175, U24ES026699, U01CA200060, U01HG009391,
471 and U41HG010972, and by the American Cancer Society Research Scholar grant RSG-14-049-01-
472 DMC.

473 **AUTHOR CONTRIBUTIONS**

474 X.Z. and T.W. conceived and implemented the study. X.Z. performed the analysis and wrote the
475 paper; A.Y.D. contributed to data analysis; D.L. contributed to data visualization and
476 interpretation; E.C.P. and A.Y.D. edited the paper; and T.W. supervised the study. All authors
477 read and approved the final manuscript.

478 **DISCLOSURE DECLARATION**

479 The authors declare no competing interests.

480 **FIGURE LEGENDS**

481 **Figure 1: Comparison of UCSC liftOver with OrthoINDEL.** Briefly, CrossMap splits syntenic
482 regions into fragments separated by any gap in the alignment. OrthoINDEL concatenates the
483 fragments split by indels and returns syntenic regions containing these indels. UCSC liftOver
484 does not convert the third example (yellow), where a large portion from the species1 region is
485 absent in species2. OrthoINDEL enables us to retrieve syntenic regions with large indels and
486 filter out other SVs such as inversions. Rectangles represent one-to-one alignments from
487 species1 to species2. Diamonds represents insertions in species1. Triangles represents an
488 inverted region between the two species.

489 **Figure 2: All putative RRRs and indels between human and chimpanzee and the number of**
490 **overlaps between them.** A. Length distribution of all putative RRRs and their orthologs. Violin
491 plots show the length of putative promoter, enhancer and H3K9me3 repressed regions. For
492 regions called in only one species, the length of their syntenic regions are plotted in the other
493 species. Size in the human genome is shown on the left; size in the chimpanzee on the right.
494 The average length of each distribution is marked and labeled. B. Size distribution of indels
495 between human and chimpanzee. The number of insertions and deletions in each lineage is
496 plotted in a back-to-back histogram with indel length on the x-axis and the number of indels of
497 different lengths on the y-axis. Colors distinguish indels based on TE classification ("noTE", not
498 derived from a TE insertion). C. The number of lineage-biased/invariant putative CNCC
499 promoters, CNCC enhancers, and iPSC H3K9me3 heterochromatin regions with or without indel
500 association. Regions were separated into those without an indel or with one of the four types of

501 indels. Colors distinguish putative RRR invariant between the two species or biased in either
502 lineage. The percentage of each category is displayed in a stacked histogram with the number
503 of occurrences labeled.

504 **Figure 3: Examples of indels associated with different CREs on the WashU Epigenome**

505 **Browser.** Human-chimpanzee track represents pairwise alignment between the human (blue)
506 and chimpanzee (pink) genomes. A. A human-specific insertion associated with a chimpanzee-
507 biased enhancer. B. A chimpanzee-specific deletion associated with a human-biased enhancer.
508 C. A chimpanzee-specific insertion associated with a human-biased H3K9me3 heterochromatin
509 region. D. A human-specific deletion associated with a chimpanzee-biased H3K9me3
510 heterochromatin region.

511 **Figure 4: Enrichment of indel categories with putative RRR categories.** Each dot in the matrix
512 represents the enrichment of one type of indel within a specific putative RRR. Indels were first
513 separated into human insertions, human deletions, chimpanzee insertions, or chimpanzee
514 deletions and then further subdivided by size or TE classification. In addition to lineage-specific
515 HERVK(HML2), chimpanzee-specific ERV insertions also include PTERV insertions absent from
516 human genome. Putative CNCC promoter, CNCC enhancer, and iPSC H3K9me3 heterochromatin
517 regions are presented from left to right. Each putative RRR is further classified horizontally into
518 human-chimpanzee invariant, human-biased, and chimpanzee-biased regions. Human-biased
519 and chimpanzee-biased putative CNCC promoters are not shown in the figure because they do
520 not intersect with any indel. The enrichment P-value (BEDTools Fisher's exact test) is displayed

521 as the background greyscale in the matrix, and the enrichment ratio is displayed using both the
522 color and size of each dot. Combinations with no intersections are crossed out.

523 **Figure 5: Enrichment of 20-50bp indels in promoters/CGI.** A. Frequency of human-chimpanzee
524 indels of sizes 20-50bp and 50-100bp in and around annotated human genes. Metaplots display
525 genes with introns removed and 5 kb flanking regions surrounding the transcription start site
526 and end site. B. Frequency of human population indels found in the 1000 Genomes Project
527 around the same promoter regions described in Fig. 5A. C. Comparing the observed number of
528 indels intersecting with CpG island within promoters, CpG island outside of promoters and
529 promoters without CpG island with the same intersections between randomly shuffled indels
530 and the three CpG island or promoter regions. All indels were shuffled 1000 times. The density
531 distributions of the numbers of shuffled indel intersections are illustrated by black lines. The
532 observed numbers are indicated with a red vertical line with the observation/expectation ratio
533 (O/E) and P value at the top of each graph.

534 **Figure 6: ChIP-seq read count, normalized to reads per genomic content, surrounding**
535 **repetitive TE insertions reveals putative hidden CREs.** A. ChIP-seq signal profiles and 50bp
536 mappability over all human-specific SVA and LTR5 insertions in the human genome (top) and
537 their orthologous pre-insertion sites in the chimpanzee genome (bottom), with 1kb flanking
538 regions. B. Profile of CNCC H3K27ac ChIP-seq, iPSC H3K9me3 ChIP-seq, and 50bp mappability
539 around SVA insertions in the human genome, with 1 kb flanking regions. SVA subfamilies are
540 distinguished by color with copy numbers inside parenthesis. C. Profile of iPSC H3K4me3 ChIP-
541 seq, iPSC H3K9me3 ChIP-seq, and 50bp mappability around LTR5 insertions in the human

542 genome, with 1kb flanking regions. LTR5 subfamilies are distinguished by color with copy
543 numbers labeled.

544 **Figure 7: NR2F1 binding profile in CNCCs coincides with the putative SVA enhancer profile.** A.
545 Heatmaps of H3K27ac (50bp single-end) and NR2F1 (202bp paired-end) ChIP-seq signal in CNCC
546 over all SVA elements with complete 3' ends (outlined by dotted line) and 3kb flanking regions
547 in the human genome sorted by length (top to bottom, longest to shortest). B. Zoomed-in view
548 of the heatmaps displaying 53 5' truncated SVA elements from 587 bp to 300 bp boxed in Fig.
549 7a. High H3K27ac and NR2F1 signal are visible on both ends, and flanking H3K27ac ChIP-seq
550 signal correlates with NR2F1 ChIP-seq signal. C. Nucleotide sequence alignment of the 10
551 truncated SVA elements boxed in Fig. 7B. Elements with strong CNCC H3K27ac signals are
552 marked by checkmarks on the right; elements without CNCC H3K27ac signal are marked by a
553 cross. The NR2F1 binding motif is provided at the bottom of the alignment. D. A human SVA
554 insertion associated with a human-biased enhancer and NR2F1 binding. E. A human SVA
555 insertion associated with a human-biased enhancer but no NR2F1 binding. Note the low
556 mappability, indicating high repetitiveness, over the SVA insertions and the difference in NR2F1
557 ChIP-seq peaks between E and F.

558

559 REFERENCES

560

561 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic variation.
562 *Nature* **526**: 68–74.

- 563 Atkinson EG, Audesse AJ, Palacios JA, Bobo DM, Webb AE, Ramachandran S, Henn BM. 2018.
564 No Evidence for Recent Selection at FOXP2 among Diverse Human Populations. *CELL* **174**:
565 1424–1435.e15.
- 566 Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. 2009.
567 MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Research* **37**: W202–8.
- 568 Bashir A, Ye C, Price AL, Bafna V. 2005. Orthologous repeats and mammalian phylogenetic
569 inference. *Genome Research* **15**: 998–1006.
- 570 Becker JS, Nicetto D, Zaret KS. 2015. H3K9me3-Dependent Heterochromatin: Barrier to Cell
571 Fate Changes. *Trends in genetics*: *TIG* **32**: 29–41.
- 572 Britten RJ, Davidson EH. 1971. Repetitive and non-repetitive DNA sequences and a speculation
573 on the origins of evolutionary novelty. *Q Rev Biol* **46**: 111–138.
- 574 Challis D, Antunes L, Garrison E, Banks E, Evani US, Muzny D, Poplin R, Gibbs RA, Marth G, Yu F.
575 2015. The distribution and mutagenesis of short coding INDELS from 1,128 whole exomes.
576 *BMC Genomics* **16**: 143.
- 577 The Chimpanzee Sequencing and Analysis Consortium. 2005. Initial sequence of the
578 chimpanzee genome and comparison with the human genome. *Nature* **437**: 69–87.
- 579 Choudhary MN, Friedman RZ, Wang JT, Jang HS, Zhuo X, Wang T. 2020. Co-opted transposons
580 help perpetuate conserved higher-order chromosomal structures. *Genome Biol* **21**: 16.
- 581 Chuong EB, Elde NC, Feschotte C. 2017. Regulatory activities of transposable elements: from
582 conflicts to benefits. *Nat Rev Genet* **18**: 71–86.
- 583 Cohen NM, Kenigsberg E, Tanay A. 2011. Primate CpG islands are maintained by heterogeneous
584 evolutionary regimes involving minimal selection. *CELL* **145**: 773–786.
- 585 Coulondre C, Miller JH, Farabaugh PJ, Gilbert W. 1978. Molecular basis of base substitution
586 hotspots in *Escherichia coli*. **274**: 775–780.
- 587 Danko CG, Choate LA, Marks BA, Rice EJ, Wang Z, Chu T, Martins AL, Dukler N, Coonrod SA, Tait
588 Wojno ED, et al. 2018. Dynamic evolution of regulatory element ensembles in primate
589 CD4+ T cells. *Nature Ecology & Evolution* **2**: 537–548.
- 590 Darwin C. 1871. *The Descent of Man, and Selection in Relation to Sex*. Createspace Independent
591 Publishing Platform.
- 592 Deaton AM, Bird A. 2011. CpG islands and the regulation of transcription. *Genes Dev* **25**: 1010–
593 1022.

- 594 Derrien T, Estellé J, Marco Sola S, Knowles DG, Raineri E, Guigó R, Ribeca P. 2012. Fast
595 computation and applications of genome mappability. ed. C.A. Ouzounis. *PLoS ONE* **7**:
596 e30377.
- 597 Du X, Gertz EM, Wojtowicz D, Zhabinskaya D, Levens D, Benham CJ, Schäffer AA, Przytycka TM.
598 2014. Potential non-B DNA regions in the human genome are associated with higher rates
599 of nucleotide mutation and expression variation. *Nucleic Acids Research* **42**: 12367–12379.
- 600 Elgin SCR, Reuter G. 2013. Position-effect variegation, heterochromatin formation, and gene
601 silencing in *Drosophila*. *Cold Spring Harb Perspect Biol* **5**: a017780–a017780.
- 602 Enard W, Przeworski M, Fisher SE, Lai CSL, Wiebe V, Kitano T, Monaco AP, Pääbo S. 2002.
603 Molecular evolution of FOXP2, a gene involved in speech and language. *Nature* **418**: 869–
604 872.
- 605 The ENCODE Project Consortium. 2004. The ENCODE (ENCyclopedia Of DNA Elements) Project.
606 *Science* **306**: 636–640.
- 607 Eres IE, Luo K, Hsiao CJ, Blake LE, Gilad Y. 2019. Reorganization of 3D genome structure may
608 contribute to gene regulatory evolution in primates. ed. H.S. Malik. *PLoS Genet* **15**:
609 e1008278.
- 610 Ernst J, Kellis M. 2012. ChromHMM: automating chromatin-state discovery and
611 characterization. *Nat Methods* **9**: 215–216.
- 612 The FANTOM Consortium and the RIKEN PMI and CLST (DGT). 2014. A promoter-level
613 mammalian expression atlas. *Nature* **507**: 462–470.
- 614 Feschotte C. 2008. Transposable elements and the evolution of regulatory networks. *Nat Rev*
615 *Genet* **9**: 397–405.
- 616 Fiddes IT, Lodewijk GA, Mooring M, Bosworth CM, Ewing AD, Mantalas GL, Novak AM, van den
617 Bout A, Bishara A, Rosenkrantz JL, et al. 2018. Human-Specific NOTCH2NL Genes Affect
618 Notch Signaling and Cortical Neurogenesis. *CELL* **173**: 1356–1369.e22.
- 619 Florio M, Heide M, Pinson A, Brandl H, Albert M, Winkler S, Wimberger P, Huttner WB, Hiller M.
620 2018. Evolution and cell-type specificity of human-specific genes preferentially expressed in
621 progenitors of fetal neocortex. *eLife Sciences* **7**.
- 622 Fortin J-P, Hansen KD. 2015. Reconstructing A/B compartments as revealed by Hi-C using long-
623 range correlations in epigenetic data. *Genome Biology* **16**: 180.
- 624 Franchini LF, Pollard KS. 2017. Human evolution: the non-coding revolution. *BMC Biol* **15**: 89.

- 625 Frankish A, Diekhans M, Ferreira A-M, Johnson R, Jungreis I, Loveland J, Mudge JM, Sisu C,
626 Wright J, Armstrong J, et al. 2019. GENCODE reference annotation for the human and
627 mouse genomes. *Nucleic Acids Res* **47**: D766–D773.
- 628 Fudenberg G, Pollard KS. 2019. Chromatin features constrain structural variation across
629 evolutionary timescales. *Proceedings of the National Academy of Sciences* **116**: 2175–2180.
- 630 Fuentes DR, Swigut T, Wysocka J. 2018. Systematic perturbation of retroviral LTRs reveals
631 widespread long-range effects on human gene regulation. *eLife Sciences* **7**.
- 632 Gallego Romero I, Pavlovic BJ, Hernando-Herraez I, Zhou X, Ward MC, Banovich NE, Kagan CL,
633 Burnett JE, Huang CH, Mitrano A, et al. 2015. A panel of induced pluripotent stem cells from
634 chimpanzees: a resource for comparative functional genomics. *eLife Sciences* **4**: e07103.
- 635 Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM, Malig M, Raja A,
636 Fiddes I, Hillier LW, et al. 2016. Long-read sequence assembly of the gorilla genome. *Science*
637 **352**: aae0344–aae0344.
- 638 Grant CE, Bailey TL, Noble WS. 2011. FIMO: scanning for occurrences of a given motif.
639 *Bioinformatics* **27**: 1017–1018.
- 640 Greenstein RA, Jones SK, Spivey EC, Rybarski JR, Finkelstein IJ, Al-Sady B. 2018. Noncoding RNA-
641 nucleated heterochromatin spreading is intrinsically labile and requires accessory elements
642 for epigenetic stability. *eLife Sciences* **7**: e0159292.
- 643 Hancks DC, Kazazian HH. 2016. Roles for retrotransposon insertions in human disease. *Mobile*
644 *DNA* **7**: 9.
- 645 Hernando-Herraez I, Prado-Martinez J, Garg P, Fernandez-Callejo M, Heyn H, Hvilsom C,
646 Navarro A, Esteller M, Sharp AJ, Marques-Bonet T. 2013. Dynamics of DNA methylation in
647 recent human and great ape evolution. ed. Y. Gilad. *PLoS Genet* **9**: e1003763.
- 648 Huddleston J, Chaisson MJP, Steinberg KM, Warren W, Hoekzema K, Gordon D, Graves-Lindsay
649 TA, Munson KM, Kronenberg ZN, Vives L, et al. 2017. Discovery and genotyping of structural
650 variation from long-read haploid genome sequence data. *Genome Research* **27**: 677–685.
- 651 Khan H, Smit A, Boissinot S. 2006. Molecular evolution and tempo of amplification of human
652 LINE-1 retrotransposons since the origin of primates. **16**: 78–87.
- 653 Kiktev DA, Sheng Z, Lobachev KS, Petes TD. 2018. GC content elevates mutation and
654 recombination rates in the yeast *Saccharomyces cerevisiae*. *Proceedings of the National*
655 *Academy of Sciences* **115**: E7109–E7118.
- 656 King MC, Wilson AC. 1975. Evolution at two levels in humans and chimpanzees. *Science* **188**:
657 107–116.

- 658 Kronenberg ZN, Fiddes IT, Gordon D, Murali S, Cantsilieris S, Meyerson OS, Underwood JG,
659 Nelson BJ, Chaisson MJP, Dougherty ML, et al. 2018. High-resolution comparative analysis
660 of great ape genomes. *Science* **360**: eaar6343.
- 661 Kuhn RM, Haussler D, Kent WJ. 2013. The UCSC genome browser and associated tools. *Briefings*
662 *in Bioinformatics* **14**: 144–161.
- 663 Kvikstad EM, Tyekucheva S, Chiaromonte F, Makova KD. 2007. A Macaque’s-Eye View of Human
664 Insertions and Deletions: Differences in Mechanisms. *PLOS Computational Biology* **3**:
665 e176.
- 666 Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglou S, Bernstein BE, Bickel P,
667 Brown JB, Cayting P, et al. 2012. ChIP-seq guidelines and practices of the ENCODE and
668 modENCODE consortia. *Genome Research* **22**: 1813–1831.
- 669 Li D, Hsu S, Purushotham D, Sears RL, Wang T. 2019. WashU Epigenome Browser update 2019.
670 *Nucleic Acids Research* **47**: W158–W165.
- 671 Li D, Zhang B, Xing X, Wang T. 2015. Combining MeDIP-seq and MRE-seq to investigate genome-
672 wide CpG methylation. *Methods* **72**: 29–40.
- 673 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
674 *arxiv.org*.
- 675 Li Q, Brown JB, Huang H, Bickel PJ. 2011. Measuring reproducibility of high-throughput
676 experiments. *Ann Appl Stat* **5**: 1752–1779.
- 677 Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for
678 RNA-seq data with DESeq2. *Genome Biol* **15**: 550.
- 679 Lowdon RF, Jang HS, Wang T. 2016. Evolution of Epigenetic Regulation in Vertebrate Genomes.
680 *Trends Genet.*
- 681 Luan DD, Korman MH, Jakubczak JL, Eickbush TH. 1993. Reverse transcription of R2Bm RNA is
682 primed by a nick at the chromosomal target site: a mechanism for non-LTR
683 retrotransposition. *CELL* **72**: 595–605.
- 684 Lunter G, Ponting CP, Hein J. 2006. Genome-wide identification of human functional DNA using
685 a neutral indel model. *PLoS Computational Biology* **2**: e5.
- 686 Lynch VJ, Leclerc RD, May G, Wagner GP. 2011. Transposon-mediated rewiring of gene
687 regulatory networks contributed to the evolution of pregnancy in mammals. *Nat Genet* **43**:
688 1154–1159.
- 689 Mager DL, Stoye JP. 2015. Mammalian Endogenous Retroviruses. *Microbiol Spectr* **3**: MDNA3–
690 0009–2014.

- 691 Makova KD, Hardison RC. 2015. The effects of chromatin organization on variation in mutation
692 rates in the genome. *Nature Reviews Genetics* **16**: 213–223.
- 693 Montgomery SB, Goode DL, Kvikstad E, Albers CA, Zhang ZD, Mu XJ, Ananda G, Howie B,
694 Karczewski KJ, Smith KS, et al. 2013. The origin, evolution, and functional impact of short
695 insertion-deletion variants identified in 179 human genomes. *Genome Research* **23**: 749–
696 761.
- 697 Neininger K, Marschall T, Helms V. 2019. SNP and indel frequencies at transcription start sites
698 and at canonical and alternative translation initiation sites in the human genome. *PLoS*
699 *ONE* **14**: e0214816.
- 700 Obersriebnig MJ, Pallesen EMH, Sneppen K, Trusina A, Thon G. 2016. Nucleation and spreading
701 of a heterochromatic domain in fission yeast. *Nat Commun* **7**: 11518.
- 702 Ohtani H, Liu M, Zhou W, Liang G, Jones PA. 2018. Switching roles for DNA and histone
703 methylation depend on evolutionary ages of human endogenous retroviruses. *Genome*
704 *Research* **28**: 1147–1157.
- 705 Osterwalder M, Barozzi I, Tissières V, Fukuda-Yuzawa Y, Mannion BJ, Afzal SY, Lee EA, Zhu Y,
706 Plajzer-Frick I, Pickle CS, et al. 2018. Enhancer redundancy provides phenotypic
707 robustness in mammalian development. *Nature* **554**: 239–243.
- 708 Pai AA, Bell JT, Marioni JC, Pritchard JK, Gilad Y. 2011. A genome-wide study of DNA
709 methylation patterns and gene expression levels in multiple human and chimpanzee
710 tissues. ed. G. Gibson. *PLoS Genet* **7**: e1001316.
- 711 Pollard KS, Salama SR, King B, Kern AD, Dreszer T, Katzman S, Siepel A, Pedersen JS, Bejerano G,
712 Baertsch R, et al. 2006. Forces shaping the fastest evolving regions in the human genome.
713 *PLoS Genet* **2**: e168.
- 714 Prescott SL, Srinivasan R, Marchetto MC, Grishina I, Narvaiza I, Selleri L, Gage FH, Swigut T,
715 Wysocka J. 2015. Enhancer Divergence and cis-Regulatory Evolution in the Human and
716 Chimp Neural Crest. *CELL* **163**: 68–83.
- 717 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.
718 *Bioinformatics* **26**: 841–842.
- 719 Rada-Iglesias A, Bajpai R, Prescott S, Brugmann SA, Swigut T, Wysocka J. 2012. Epigenomic
720 annotation of enhancers predicts transcriptional regulators of human neural crest. *Cell*
721 *Stem Cell* **11**: 633–648.
- 722 Ramirez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S, Dündar F, Manke T.
723 2016. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic*
724 *Acids Research* **44**: W160–5.

- 725 Ray DA, Xing J, Salem A-H, Batzer MA. 2006. SINEs of a nearly perfect character. *Syst Biol* **55**:
726 928–935.
- 727 Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-
728 Moussavi A, Kheradpour P, Zhang Z, Wang J, et al. 2015. Integrative analysis of 111
729 reference human epigenomes. *Nature* **518**: 317–330.
- 730 Rogers J, Gibbs RA. 2014. Comparative primate genomics: emerging patterns of genome
731 content and dynamics. *Nat Rev Genet*.
- 732 Saxonov S, Berg P, Brutlag DL. 2006. A genome-wide analysis of CpG dinucleotides in the human
733 genome distinguishes two distinct classes of promoters. *Proc Natl Acad Sci USA* **103**: 1412–
734 1417.
- 735 Smit A, Hubley R, Green P. 2013–2015. RepeatMasker Open-4.0.
736 <<http://www.repeatmasker.org>>.
- 737 Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye K, Jun G,
738 Fritz MH-Y, et al. 2015. An integrated map of structural variation in 2,504 human genomes.
739 *Nature* **526**: 75–81.
- 740 Suzuki IK, Gacquer D, Van Heurck R, Kumar D, Wojno M, Bilheu A, Herpoel A, Lambert N,
741 Cheron J, Polleux F, et al. 2018. Human-Specific NOTCH2NL Genes Expand Cortical
742 Neurogenesis through Delta/Notch Regulation. *CELL* **173**: 1370–1384.e16.
- 743 Tian X, Strassmann JE, Queller DC. 2011. Genome nucleotide composition shapes variation in
744 simple sequence repeats. *Molecular Biology and Evolution* **28**: 899–909.
- 745 Trizzino M, Park Y, Holsbach-Beltrame M, Aracena K, Mika K, Caliskan M, Perry GH, Lynch VJ,
746 Brown CD. 2017. Transposable elements are the primary source of novelty in primate gene
747 regulation. *Genome Research* **27**: 1623–1633.
- 748 Wall JD. 2013. Great ape genomics. *ILAR J* **54**: 82–90.
- 749 Wang H, Xing J, Grover D, Hedges DJ, Han K, Walker JA, Batzer MA. 2005. SVA elements: a
750 hominid-specific retroposon family. *Journal of Molecular Biology* **354**: 994–1007.
- 751 Wang T, Zeng J, Lowe CB, Sellers RG, Salama SR, Yang M, Burgess SM, Brachmann RK, Haussler
752 D. 2007. Species-specific endogenous retroviruses shape the transcriptional network of the
753 human tumor suppressor protein p53. *Proc Natl Acad Sci USA* **104**: 18613–18618.
- 754 Ward MC, Zhao S, Luo K, Pavlovic BJ, Karimi MM, Stephens M, Gilad Y. 2018. Silencing of
755 transposable elements may not be a major driver of regulatory evolution in primate iPSCs.
756 *eLife Sciences* **7**: 166.

- 757 Xiao S, Xie D, Cao X, Yu P, Xing X, Chen C-C, Musselman M, Xie M, West FD, Lewin HA, et al.
758 2012. Comparative epigenomic annotation of regulatory DNA. *CELL* **149**: 1381–1392.
- 759 Yohn CT, Jiang Z, McGrath SD, Hayden KE, Khaitovich P, Johnson ME, Eichler MY, McPherson JD,
760 Zhao S, Pääbo S, et al. 2005. Lineage-specific expansions of retroviral insertions within the
761 genomes of African great apes but not humans and orangutans. *Plos Biol* **3**: e110.
- 762 Yokoyama KD, Zhang Y, Ma J. 2014. Tracing the Evolution of Lineage-Specific Transcription
763 Factor Binding Sites in a Birth-Death Framework. *PLoS Computational Biology* **10**:
764 e1003771.
- 765 Zhao H, Sun Z, Wang J, Huang H, Kocher J-P, Wang L. 2014. CrossMap: a versatile tool for
766 coordinate conversion between genome assemblies. *Bioinformatics* **30**: 1006–1007.
- 767 Zhu J, Sanborn JZ, Diekhans M, Lowe CB, Pringle TH, Haussler D. 2007. Comparative genomics
768 search for losses of long-established genes on the human lineage. *PLoS Comput Biol* **3**:
769 e247.
- 770
- 771

Fig. 1

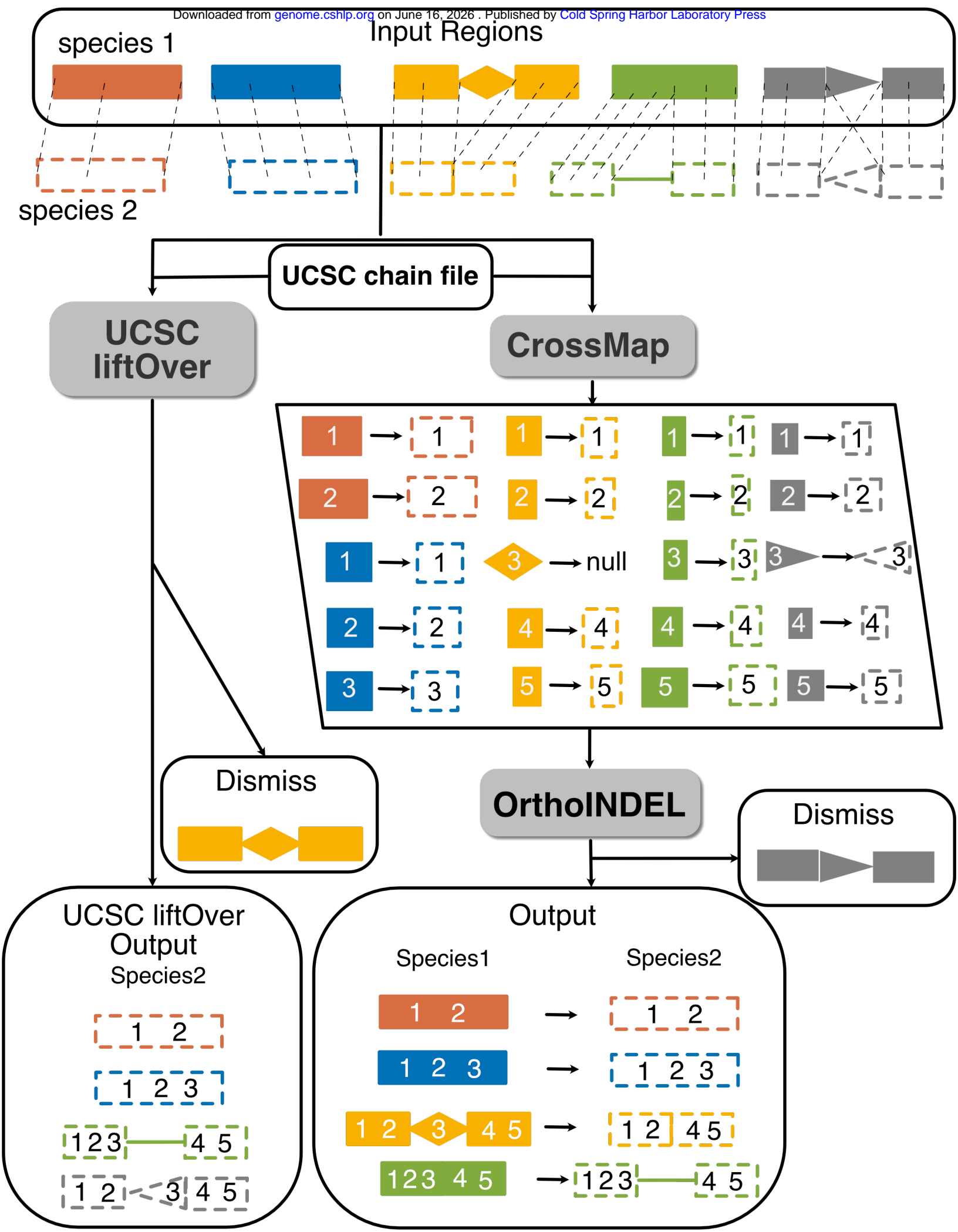
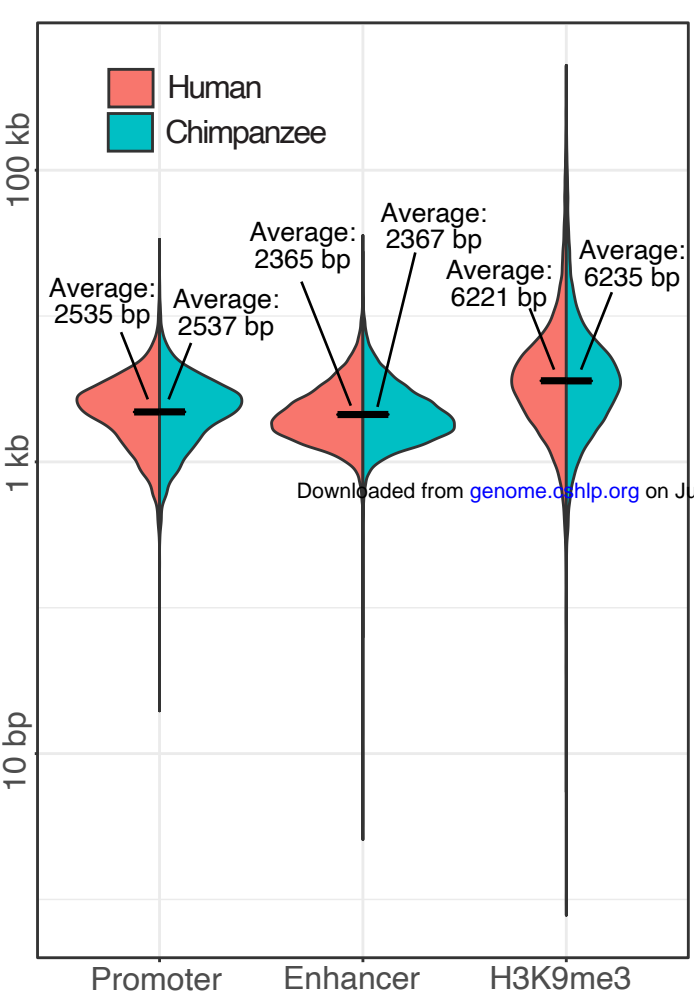
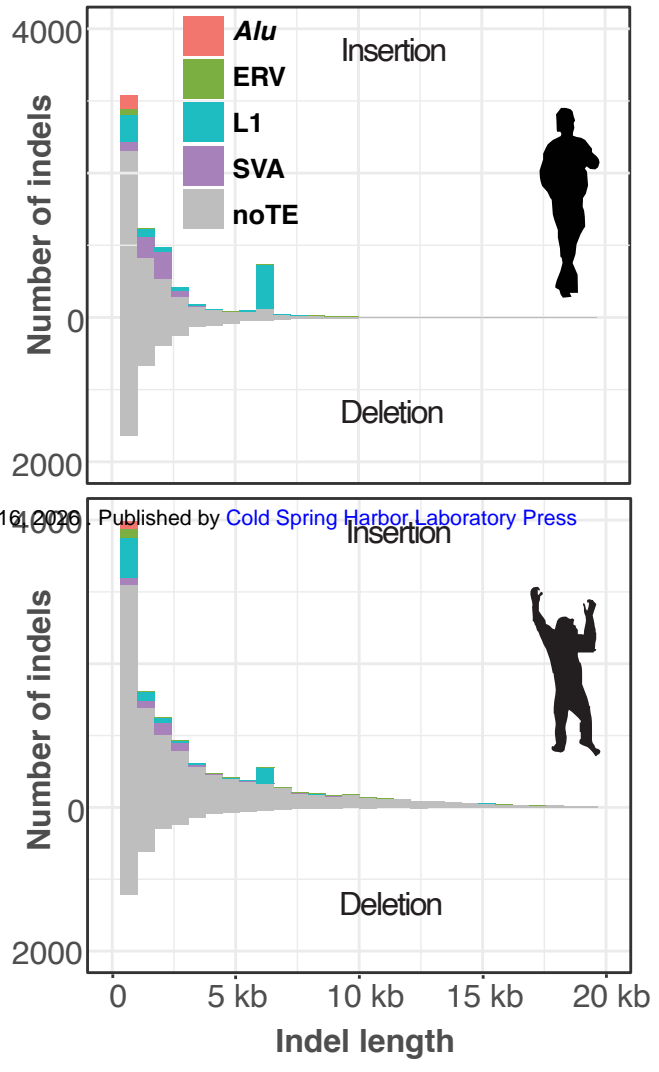


Fig.2

A Length distribution of regulatory/repressive regions and their orthologs



B Indel length distribution



C

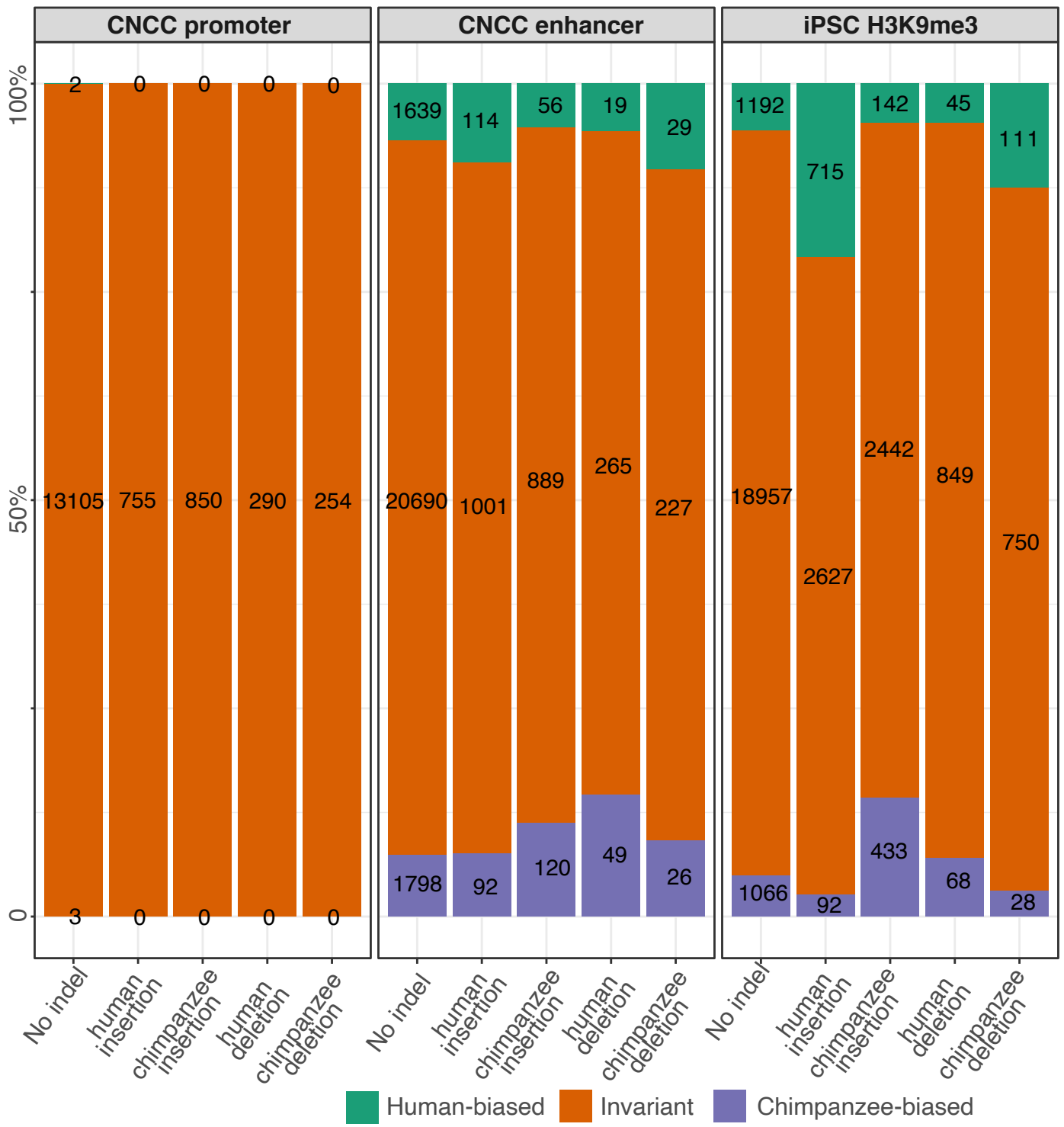


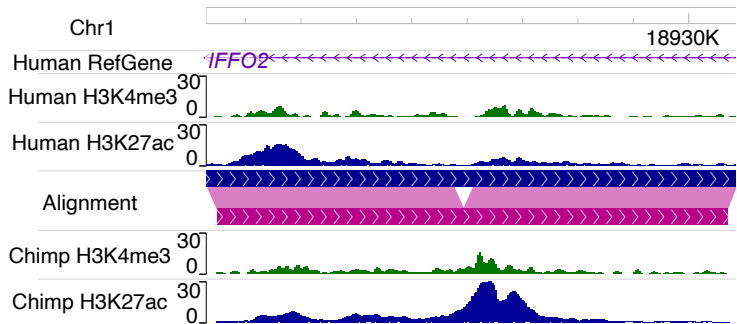
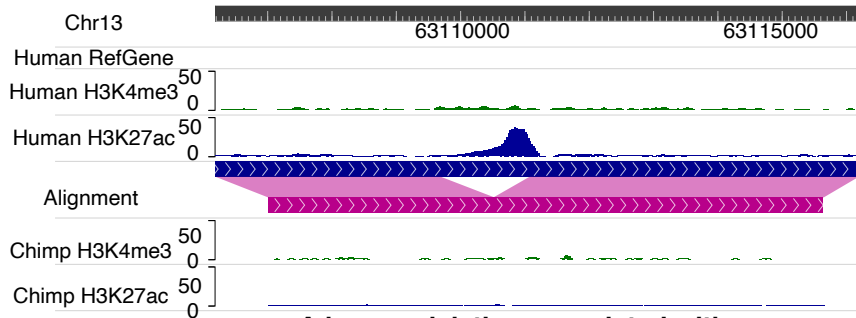
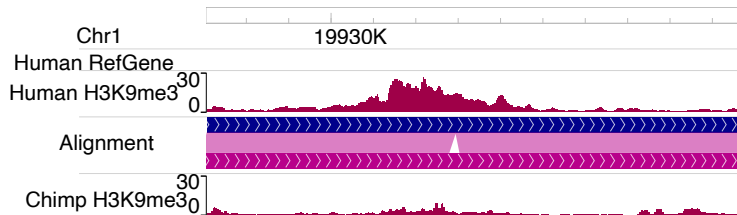
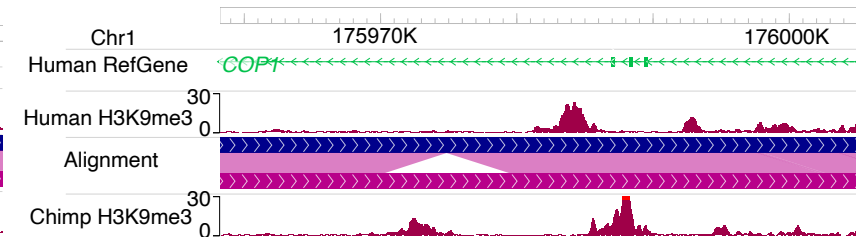
Fig. 3**A****A human-specific insertion associated with chimpanzee-biased enhancer****B****A chimpanzee-specific deletion associated with human-biased enhancer****C****A chimpanzee insertion associated with human-biased H3K9me3 region****D****A human deletion associated with chimpanzee-biased H3K9me3 region**

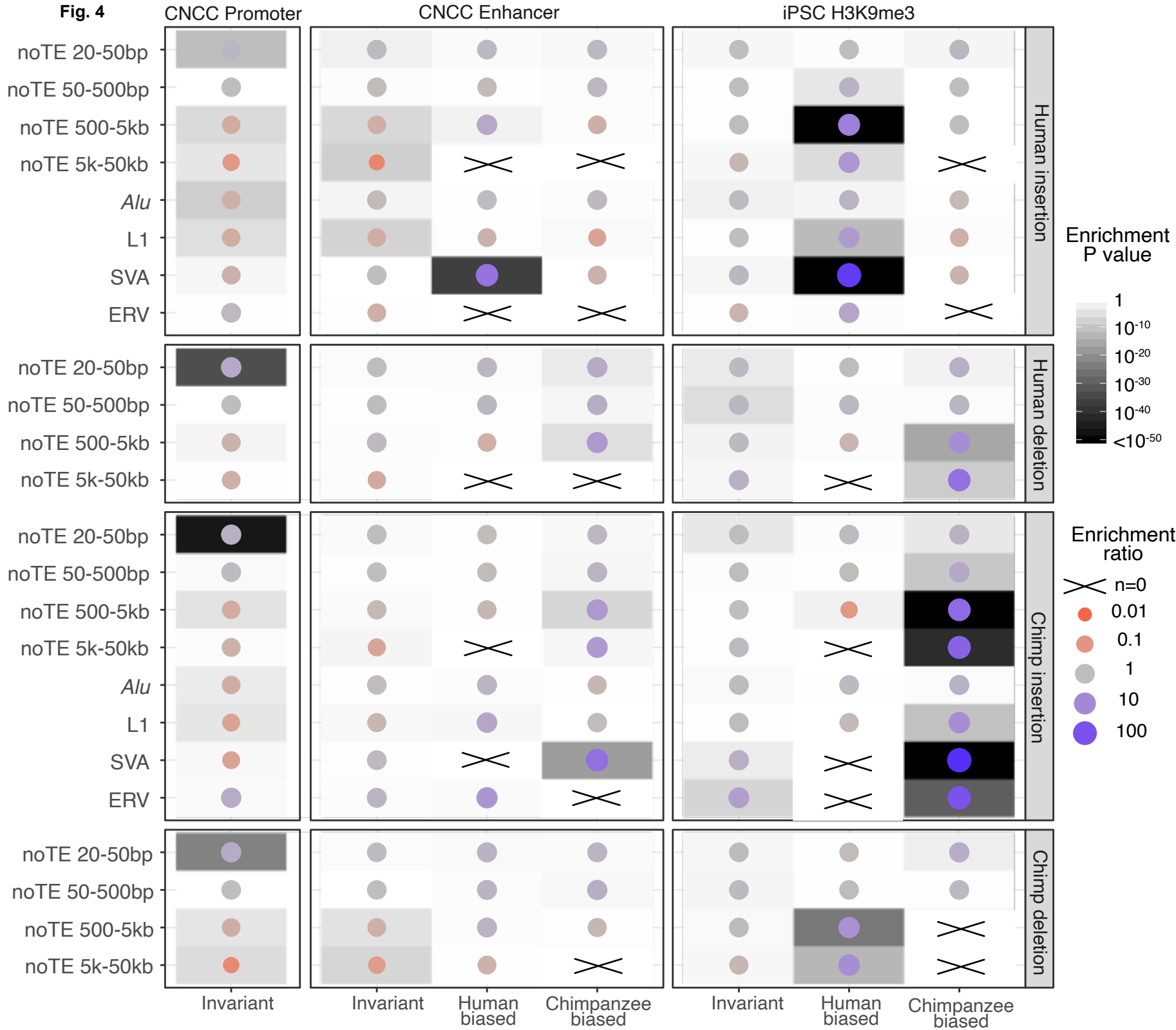
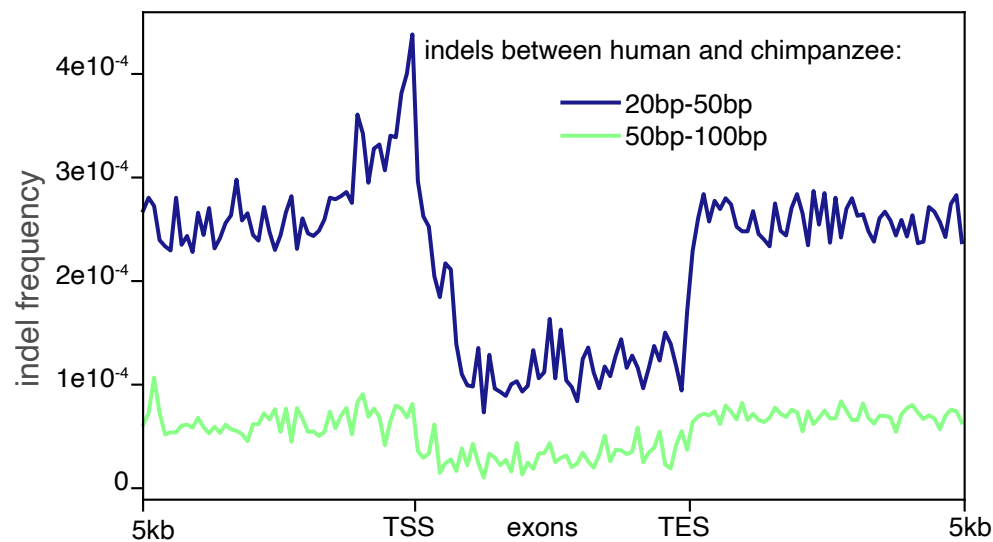
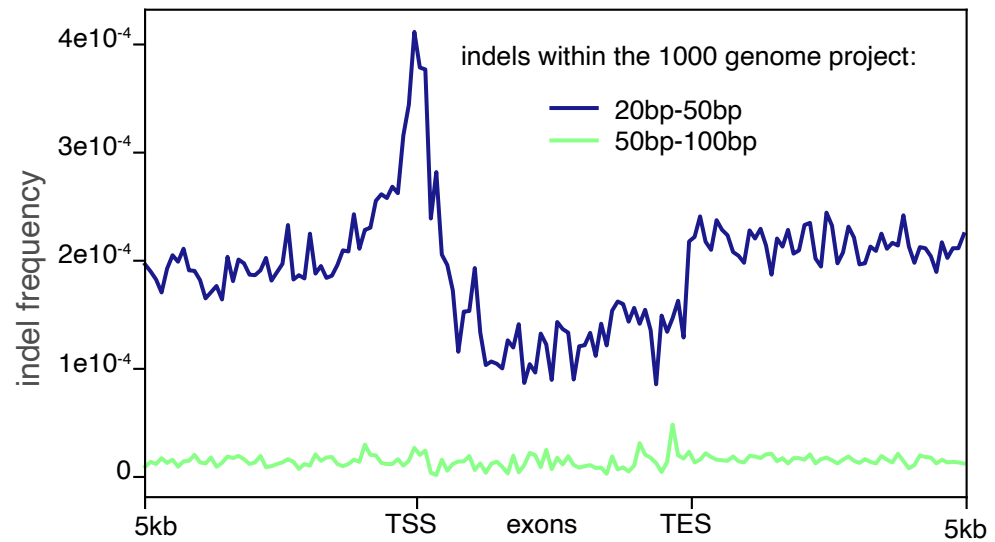
Fig. 4

Fig. 5

A



B



C

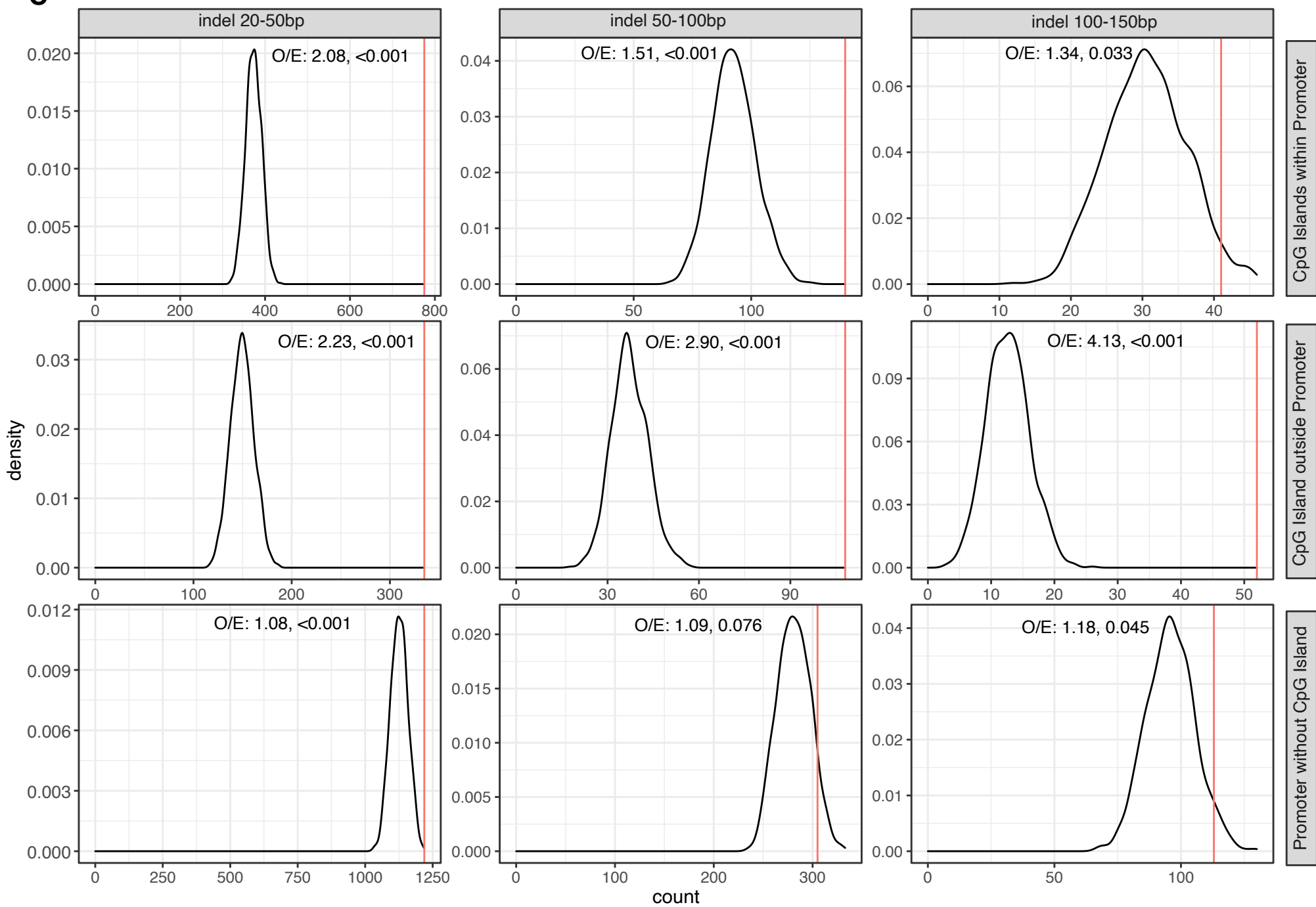


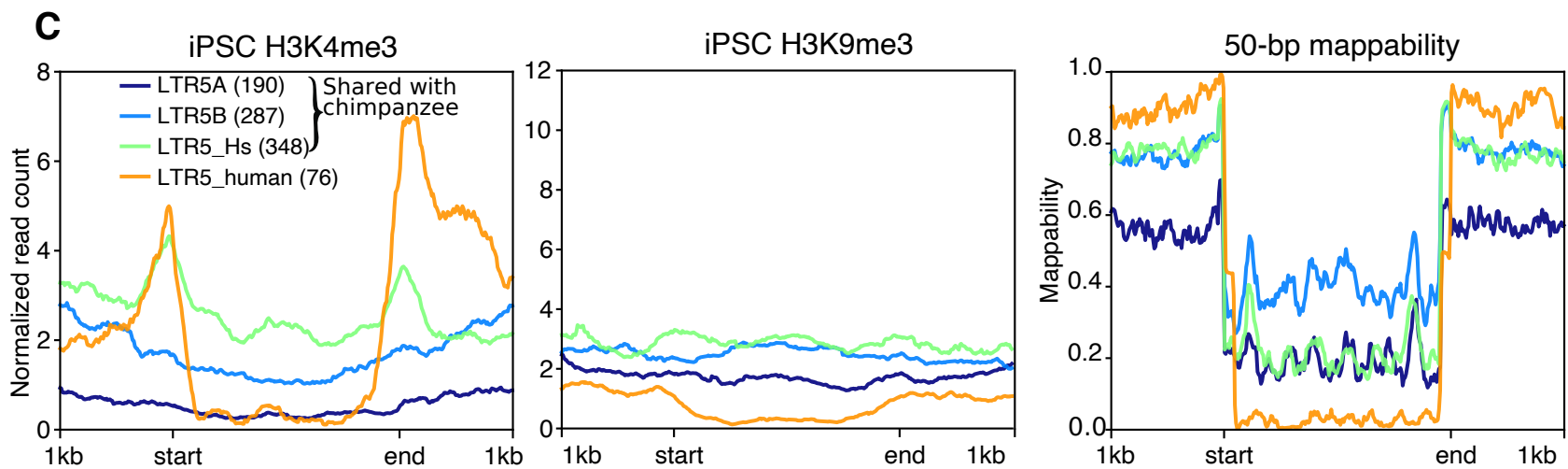
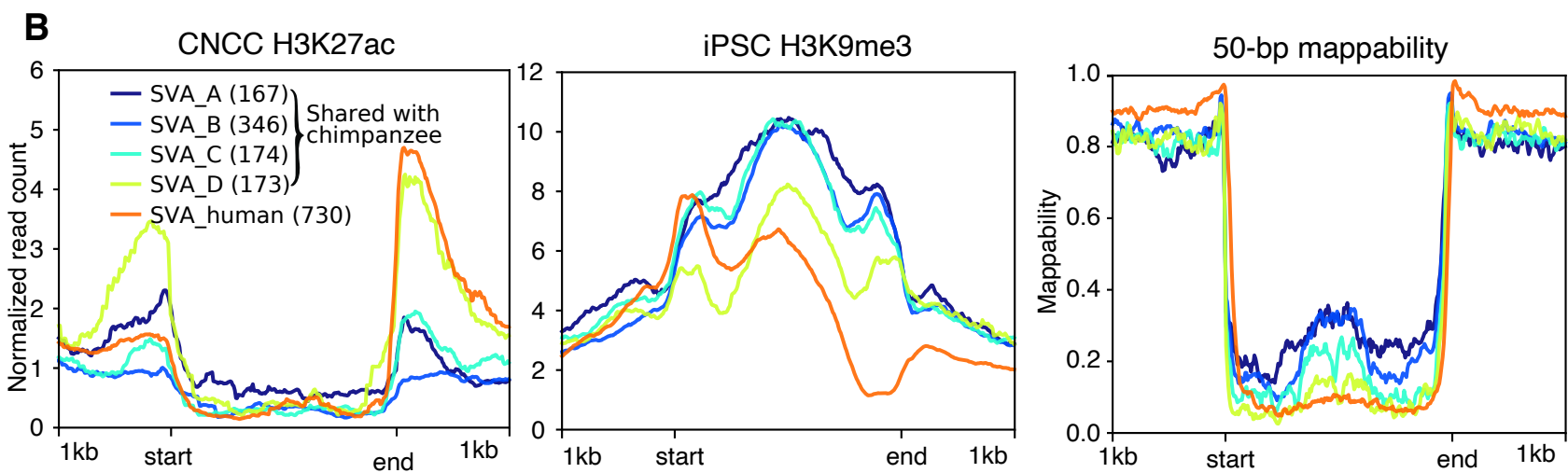
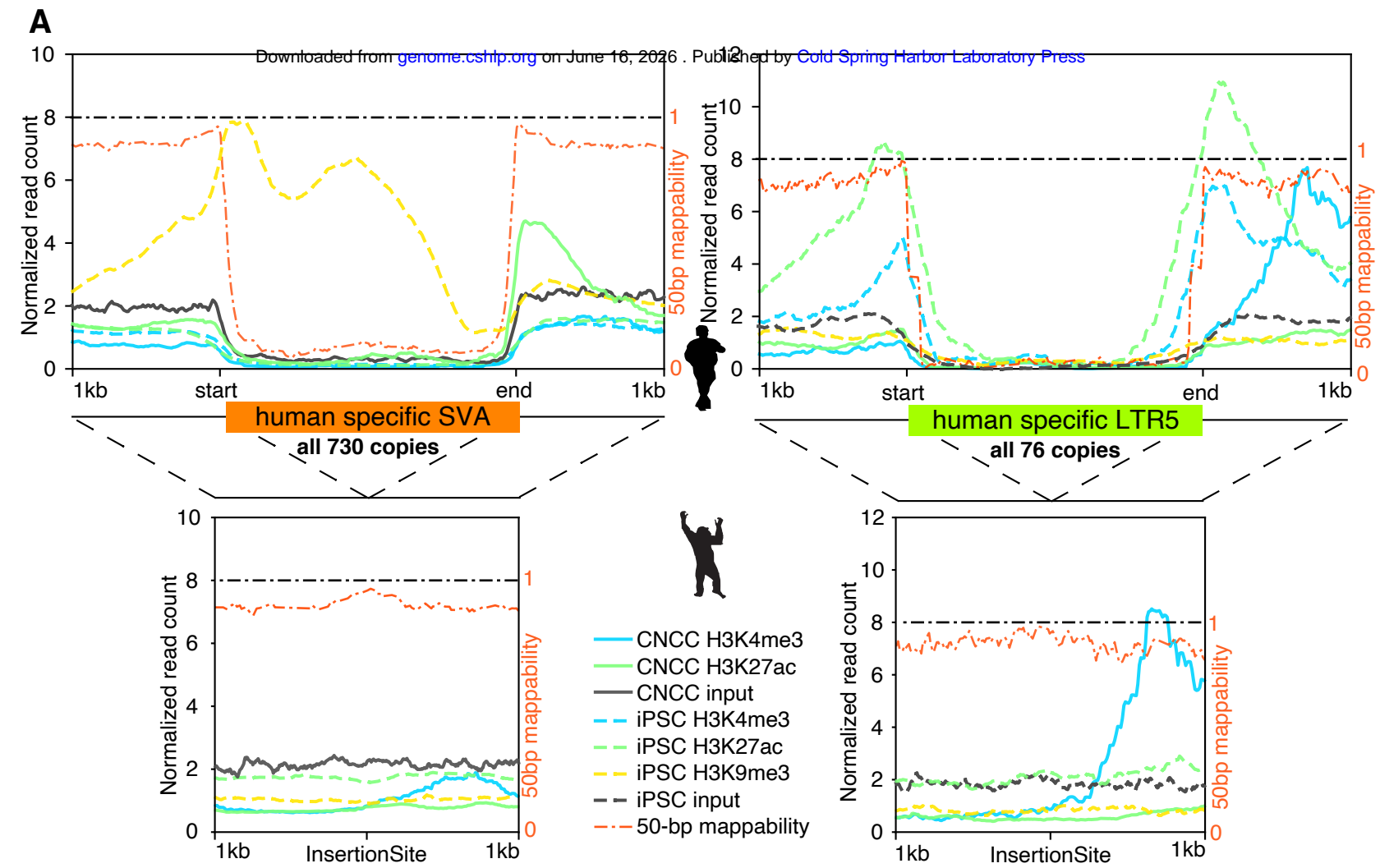
Fig. 6

Fig. 7