

1 Tandem repeat variation within and between species reveals signatures  
2 of selection in humans and chimpanzees

3

4 Carolina L. Adam<sup>1\*</sup>, Joana L. Rocha<sup>2,3</sup>, Peter H. Sudmant<sup>2+</sup>, Rori V. Rohlf<sup>1,4\*+</sup>

5

6 <sup>1</sup>Institute of Ecology and Evolution, University of Oregon, Eugene, OR 97403, USA

7 <sup>2</sup>Department of Integrative Biology, University of California - Berkeley, Berkeley, CA 94720, USA

8 <sup>3</sup>Department of Biology, New York University, New York, NY 10003, USA

9 <sup>4</sup>Department of Data Science, University of Oregon, Eugene, OR 97403, USA

10

11 \* Corresponding authors e-mails: [carolinaladam@gmail.com](mailto:carolinaladam@gmail.com), [rori@uoregon.edu](mailto:rori@uoregon.edu)

12 + These authors contributed equally

### 13 Abstract

14 Tandem repeats (TRs) are highly mutable DNA elements that influence gene regulation<sup>1</sup>, protein  
15 structure<sup>2</sup>, and disease<sup>3</sup>. Until recently, their repetitive nature has hindered accurate TR  
16 sequencing and genotyping, resulting in sparse comparative data across species. In addition, we  
17 lack population-aware approaches to analyze TR conservation, divergence, and mutational  
18 dynamics. Here, leveraging telomere-to-telomere primate genomes and long-read data from 46  
19 humans and 23 chimpanzees, we constructed a catalog of homologous TR loci, and developed an  
20 analytical framework to jointly analyze TR variation within- and between-species. Across  
21 primates, TR diversity and conservation vary strongly with genomic context, with coding and 5'  
22 UTR TRs exhibiting reduced polymorphism and constraint across species, consistent with  
23 stabilizing selection. Yet, while TRs are depleted in coding sequence, they are enriched in 5'  
24 UTRs, suggesting functional roles that outweighs mutational risks. TR heterozygosity varies  
25 across motif lengths and is concordant with both evolutionary and trio-based mutation rate  
26 estimates<sup>4</sup>. Introducing an HKA-like approach to control for locus-specific mutation rates, we  
27 identified TRs with signatures of directional and balancing selection. These candidates are  
28 significantly enriched in genes involved in nervous system development and synaptic function,  
29 highlighting TRs as potential contributors to neural evolution. Further, TR divergence correlates  
30 with gene expression divergence, particularly for promoter-related TRs and expression in  
31 organoids related to neurodevelopment, implicating a subset of regulatory TRs as candidates for  
32 adaptive expression evolution. Finally, trait-associated TRs display longer alleles and higher  
33 diversity in humans compared to chimpanzees, consistent with lineage-specific runaway  
34 mutations and/or directional selection<sup>5</sup>. Together, our results establish a comparative framework  
35 for TR evolutionary analyses, revealing how mutational processes and selection jointly shape  
36 repeat variation, and supporting their role as both conserved functional elements and as drivers  
37 of evolutionary innovation.

## 38 Main

39 Tandem repeats (TRs) are ubiquitous across metazoans and account for nearly 8% of the  
40 human genome<sup>6</sup>. Their repetitive structure makes them prone to replication slippage, leading to  
41 mutation rates that are orders of magnitude higher than single-nucleotide variants (SNVs)<sup>7</sup>. On a  
42 per-locus basis, TR variants are more likely than SNVs to impact functional traits<sup>8</sup>, and human  
43 genetics has repeatedly demonstrated their impact on development and disease, from  
44 neurological disorders to cancer<sup>9-11</sup>. These properties suggest that TRs may act as fuel for rapid  
45 adaptation<sup>12,13</sup>.

46 Despite multiple lines of evidence supporting the role of TRs in adaptation<sup>13-16</sup>, technical  
47 limitations in sequencing and genotyping long repetitive DNA have hindered systematic TR  
48 analyses across species. Now, advances in long-read sequencing and telomere-to-telomere (T2T)  
49 assemblies have resolved previously inaccessible repeat regions<sup>6</sup>, fundamentally reshaping our  
50 understanding of genome architecture and structural variation<sup>17</sup>. As a result, recent comparative  
51 studies have uncovered lineage-specific TRs associated with evolutionary divergence among  
52 primates, particularly in neuron-specific regulatory mechanisms<sup>18-20</sup>. Characterizing TR variation  
53 between species can provide insight into how TRs impact molecular processes by revealing  
54 conserved sites under stabilizing selection, as well as sites under directional selection that  
55 underlie adaptive lineage-specific traits. Still, we lack a framework that enables comparisons of  
56 long-read TR genotypes both between and within species, an approach that has proven  
57 transformative in SNV-based evolutionary studies<sup>21,22</sup>.

58 Here, we leveraged T2T reference genomes to identify millions of homologous TR loci  
59 between humans and non-human primates (NHPs), revealing broad patterns of TR constraint  
60 over evolutionary time. We integrated long-read, assembly-level data from 46 humans, and 23  
61 newly sequenced chimpanzees, described in a companion paper<sup>23</sup>, providing a unique resource  
62 for joint analysis of within- and between-species variation to investigate the mutational and  
63 selective processes shaping TR evolution. Further, we propose an analytical HKA-like  
64 framework for comparative TR analyses, enabling the identification of locus-specific deviations  
65 in divergence-diversity ratios and candidate loci under distinct selective regimes. Finally, we  
66 examine the contribution of TR variation to expression divergence and trait variation,  
67 highlighting how extreme patterns of divergence and diversity may help prioritize putatively  
68 functional loci.

69

## 70 Results

### 71 Species-specific TR catalogs, homology to humans, and genomic distribution

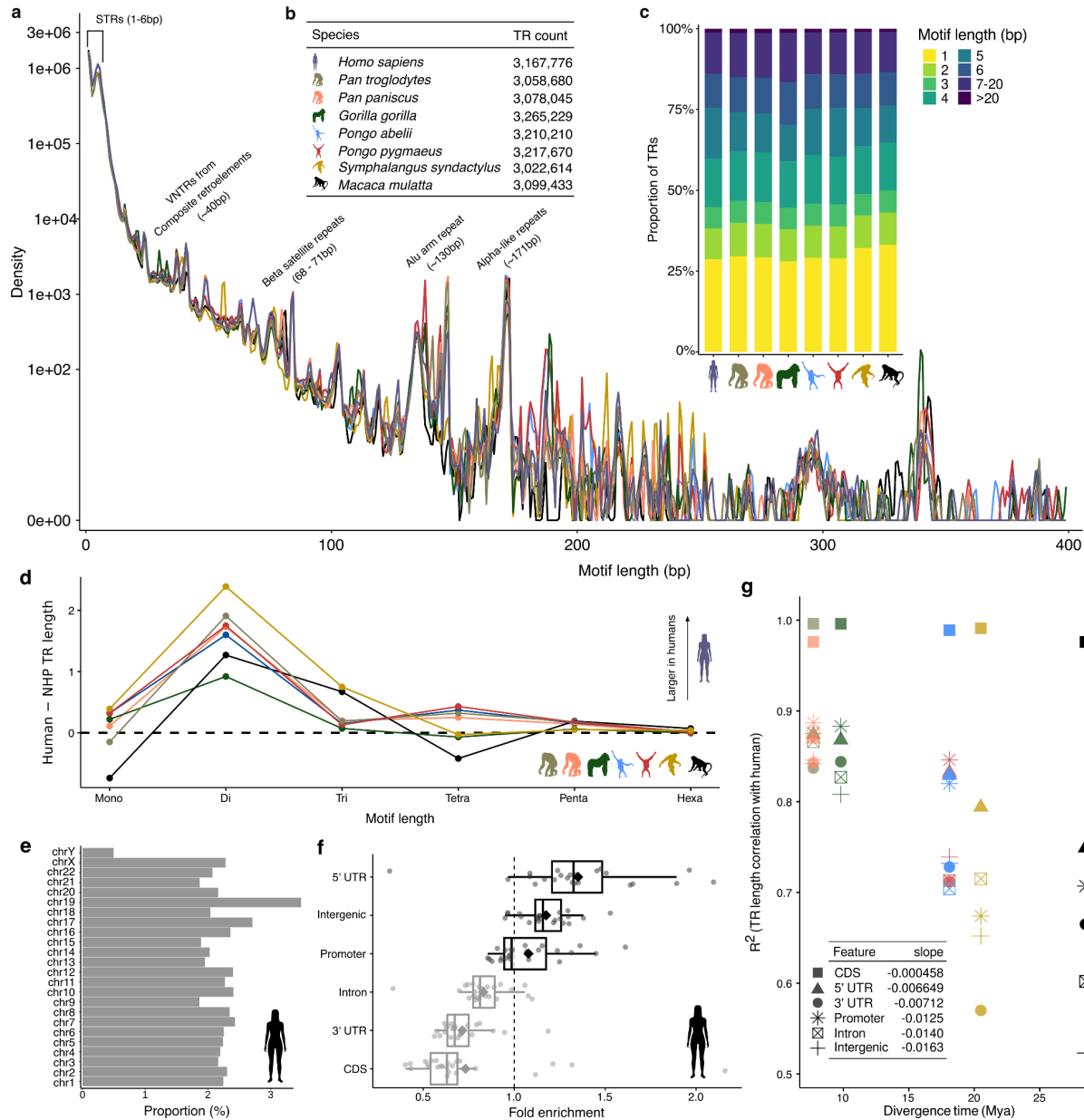
72 We developed the TRACK pipeline to create species-specific TR catalogs and identify  
73 homology between human and each NHP reference genome. TRs are subdivided based on their  
74 motif length, where short tandem repeats (STRs) exhibit motifs  $\leq 6$  bp, and variable number  
75 tandem repeats (VNTRs) display motifs  $\geq 7$ bp. Their distributions across motif lengths and  
76 catalog sizes were comparable across species ([Fig. 1a-c](#)). Among STRs, abundance across motif  
77 lengths was highly concordant with other studies ([Fig. 1c](#))<sup>24-26</sup>. Discordances, such as hexamer

78 abundance<sup>24,25</sup> where we observed a mild depletion, likely stem from differences in TR  
79 identification and filtering. For instance, while some studies consider only TRs with high or  
80 complete sequence constancy<sup>24-26</sup>, our analysis allowed a lower threshold (>60%), capturing  
81 more variable repeats. Motifs larger than 20 bp were considerably less common and have not  
82 been systematically explored in prior studies. When queried in the Dfam database, several  
83 common long motifs matched known repetitive elements, including Alu elements, and VNTRs  
84 embedded within composite retrotransposons ([Fig. 1a](#)). Young Alu elements have been reported  
85 to insert adjacent to older Alus, potentially through reuse of nearby LINE-1 cleavage sites<sup>27</sup>. The  
86 GC content of TR motifs also varied with motif lengths, but showed broadly consistent patterns  
87 across species (Supplementary Fig. 2).

88 Total STR length varies significantly across motif sizes and between species ([Fig. 1d](#);  
89 Supplementary Table S1). The largest differences were observed for dimers, which were  
90 consistently longer in humans than in all NHP species, consistent with earlier  
91 human-chimpanzee results<sup>28</sup>. The persistence of this pattern across a broader primate phylogeny  
92 is consistent with a shift on the human lineage for an increased rate of dinucleotide expansions.

93 STRs were generally uniformly distributed along chromosome lengths (Extended Data  
94 Fig. 1). In humans, chromosomes 19 and 17 showed the highest relative TR densities ([Fig. 1e](#)),  
95 consistent with earlier studies<sup>29,30</sup>. Chromosome Y displayed the lowest density of non-CenSat  
96 TRs after filtering (see Methods), reflecting the exclusion of its extensive satellite-rich regions<sup>31</sup>.  
97 In contrast to STRs, human VNTRs are concentrated in subtelomeric regions (Extended Data  
98 Fig. 1a), in line with previous observations that subtelomeric structural variation in humans is  
99 largely composed of VNTRs with motifs >15 bp<sup>32</sup>. In several NHP, however, telomeric and  
100 subtelomeric regions are depleted of TRs after filtering (Extended Data Fig. 1b-h). This reflects  
101 the presence of lineage-specific satellite-like repeat structures at chromosome ends, including the  
102 StSats arrays described in non-human great apes<sup>33</sup> and the large SiaRep repeats found in siamang  
103 gibbons<sup>34</sup>, which were masked during filtering. The density of TRs across homologous genomic  
104 regions shows decreasing correlation with increasing time since the most recent common  
105 ancestor (TMRCA), indicating progressive divergence of repeat landscapes over evolutionary  
106 time (Extended Data Fig. 2).

107 In humans, TRs are depleted in coding (CDS) regions, 3' untranslated regions (UTRs),  
108 and introns, while they are enriched in 5' UTRs, intergenic regions, and promoters ([Fig. 1f](#)).  
109 When stratified by motif length, distinct distributions emerged: (Supplementary Fig. 3,  
110 Supplementary Table S2), in agreement with a large STR population panel<sup>12</sup>. Monomers and  
111 dimers were strongly depleted in coding regions (0.01 and 0.05-fold, respectively), while repeats  
112 with motif lengths in multiples of three were overrepresented (Supplementary Fig. 3,  
113 Supplementary Table S2), reflecting the coding-frame constraint<sup>25,30</sup>. A similar trend was  
114 observed in the 5' UTRs. The GC content of TR motifs also varied across genomic features,  
115 being enriched in coding and 5' UTRs (Supplementary Fig. 4).



116

117 Fig. 1. **Tandem repeats across T2T primate genomes.** **a**, Distributions of TR motif lengths across primate  
118 genomes. Labeled peaks indicate TRs with motif lengths >20 bp that show strong matches to entries in the Dfam  
119 database. **b**, TR count per species. **c**, Proportional distributions of STR motif lengths across primates. **d**, Difference  
120 in total STR lengths between human and each NHP species across motif lengths. **e**, Proportion of each human  
121 chromosome contained in TRs. **f**, Fold enrichment of total TR length (bp) overlapping annotated genomic features  
122 per chromosome. Diamonds indicate mean enrichment across the genome. **g**, Divergence time (in millions of years)  
123 and the coefficient of determination ( $R^2$ ) for reference TR length comparisons between humans and each non-human  
124 primate species. The inset table shows the slope of the linear regression line for each genomic feature.  
125

126 The varying abundance of homologous TRs between human and NHPs reflects  
127 phylogenetic distances, with chimpanzees retaining the largest overlap (Supplementary Fig. 5).  
128 TR length comparison between species reveals conservation following the phylogeny, with  
129 correlations weakening as TMRCA increases, with higher conservation for TRs within CDS

130 (Fig. 1g, Extended Data Fig. 3). In coding regions, the maintenance of TR length has been shown  
131 to be crucial in several cases to preserve functional protein structure<sup>35,36</sup>, leading to deep  
132 conservation of coding TRs across mammals<sup>37</sup>. We also observe notable length conservation for  
133 TRs in 5' UTRs (Fig. 1g, Extended Data Fig. 3). In 5' UTRs, extreme conservation has been  
134 linked to RNA-mediated translational control in essential developmental genes<sup>38</sup>. In accordance,  
135 genes with TRs in their 5'UTRs are gene ontology (GO) enriched with biological processes  
136 spanning broad terms such as developmental growth, to more specialized terms such as  
137 axonogenesis, with overall enrichment for core neurodevelopmental processes (Supplementary  
138 Fig. 6). Additionally, variants in UTRs are known to disrupt transcription initiation and, in some  
139 cases, contribute to disease<sup>39</sup>.

140

#### 141 **TR genotypes, comparative length variation, and mutation process**

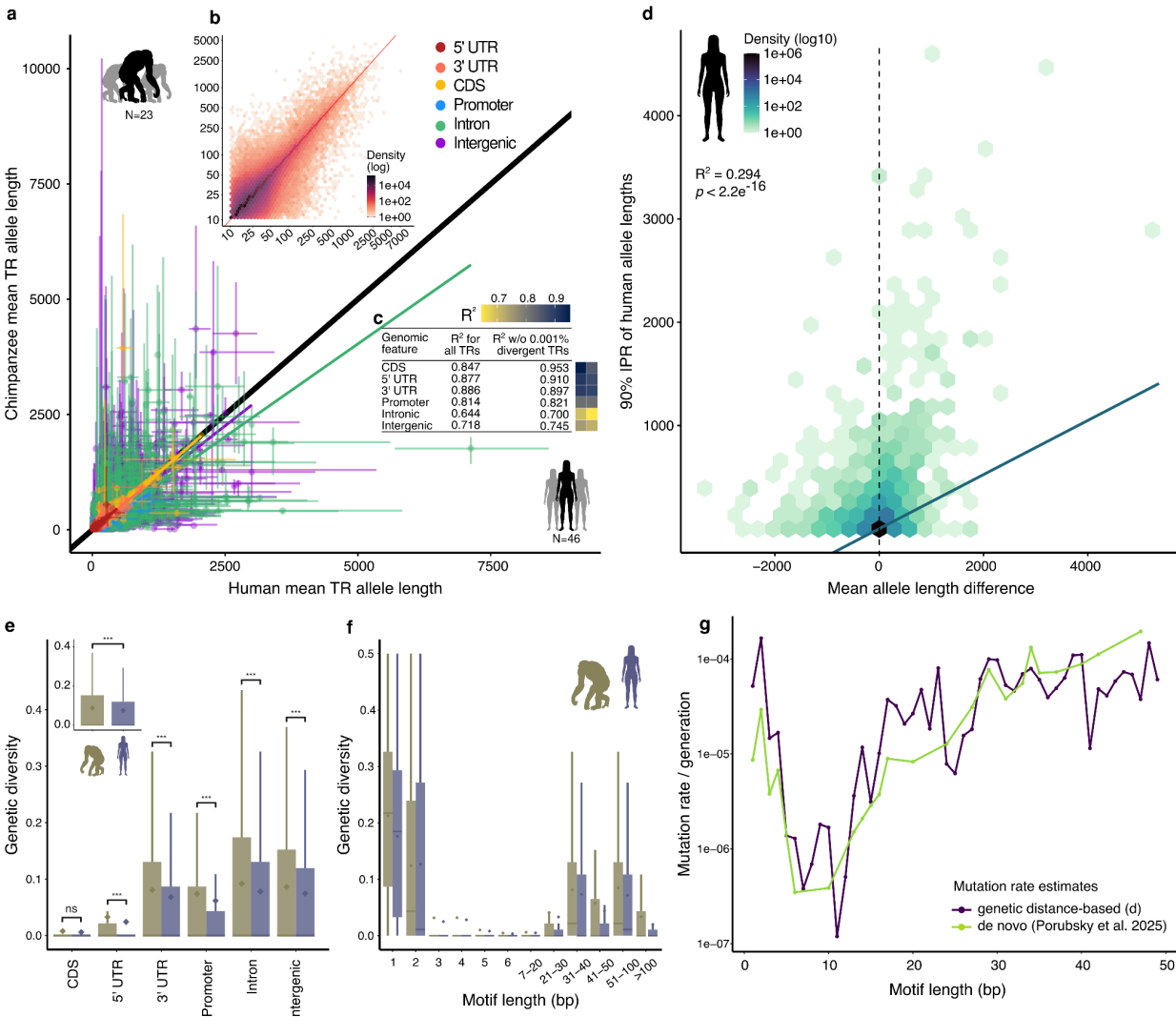
142 After quality control and filtering (see Methods), we identified 1,905,929 homologous  
143 TRs between humans and chimpanzees. Of these, 1,033,666 (54.2%) are invariant, i.e.,  
144 monoallelic across samples, and 872,237 (45.7%) exhibit at least two different alleles in either  
145 species, here denoted TR variants, or TRVs. Invariant TRs display longer motifs and have an  
146 overrepresentation of coding and 5' UTR repeats (Supplementary Fig. 7).

147 Homologous TRs are generally short (mean allele length  $\leq 22$  bp for 62.1% of TRs) and  
148 highly correlated between species (Fig. 2a-c; Supplementary Fig. 8). In coding and UTR regions,  
149 concordance is exceptionally high and mirrors previous assembly-level comparisons of  
150 orthologous human-chimpanzee STRs<sup>40</sup>, suggesting stronger stabilizing selection on functional  
151 TR variation. Highly divergent TRs also tend to exhibit greater allelic diversity within the  
152 species harboring the longer alleles (Fig. 2d; Supplementary Fig. 9). This pattern is consistent  
153 with longer alleles being subject to higher mutation rates and the potential for lineage-specific  
154 runaway repeat expansions<sup>5,41</sup>.

155 Overall, chimpanzees show higher TR heterozygosity (Fig. 2e), reflecting lower human  
156 genetic diversity found across diverse markers<sup>22,42,43</sup>. Consistent with other studies<sup>12,44</sup>, TR  
157 genetic diversity is reduced in 5' UTRs and coding regions in both species, suggesting that  
158 strong purifying selection limits functional TR variation across lineages. Heterozygosity also  
159 varied across motif lengths, highest for monomers and dimers, decreasing with motif lengths up  
160 to 20 bp, and increasing for longer motifs (Fig. 2f). This pattern holds for TRs across genomic  
161 features, except for coding and, to some extent, 5' UTR regions, which have lower  
162 heterozygosity even for monomers and dimers (Supplementary Fig. 10). Thus, even when  
163 harboring short motifs known to have elevated mutation rates due to frequent slippage<sup>7</sup>, variation  
164 in coding, and to a lesser degree 5' UTR, TRs are likely constrained by strong purifying  
165 selection. In 5' UTRs, TRs exhibit remarkably low heterozygosity, reflecting the abundance of  
166 3-6bp motifs, which display lower genetic diversity compared to monomers or dimers  
167 (Supplementary Figs. 3, 10). Nonetheless, TRs are overrepresented in 5' UTRs compared to the  
168 genome-wide distribution (Fig. 2c; Odds ratio = 1.37,  $p < 2.2 \times 10^{-16}$ ), and display even stronger  
169 enrichment in hyperconserved 5' UTRs<sup>38</sup> (Odds ratio = 3.54,  $p < 2.2 \times 10^{-16}$ ), suggesting that they

170 serve functional roles that offset the risk of mutations.

171



172

173 Fig. 2. **Comparative analyses of homologous tandem repeats between humans and chimpanzees.** **a**, Mean TR  
 174 allele lengths between humans (x-axis) and chimpanzees (y-axis) for TRs shared between species, including iTRs  
 175 and TRVs. Each point represents a single TR locus, color-coded by its genomic annotation in the CHM13 genome.  
 176 Whiskers represent the 5th-95th quantiles of allele lengths in humans (horizontal bars) and chimpanzees (vertical  
 177 bars). **b**, Heatmap of the same TR loci showing the joint distribution of human and chimpanzee mean allele lengths.  
 178 Color intensity reflects the density of TRs within each bin (50 bins). **c**, Correlation between human and chimpanzee  
 179 TR length grouped by genomic feature. **d**, Heatmap of mean allele length divergence between chimpanzees and  
 180 humans (x-axis) and the 90% interpercentile range of human allele lengths per locus. **e,f**, Expected heterozygosity  
 181 across different (**e**) genomic features and (**f**) across motif lengths. \* $p \leq 0.05$ , \*\* $p \leq 0.01$  and \*\*\* $p \leq 0.001$  denote  
 182 statistically significant differences given by Wilcoxon rank-sum tests. **g**, Estimates of *de novo* (light green) and  
 183 genetic distance-based (dark purple) mutation rates averaged across motif lengths.

184

185 Comparative TR variation data also facilitates model-based estimates of locus-specific  
 186 mutation rates. We estimated per locus mutation rates under the stepwise mutation model (SMM)  
 187 and selective neutrality<sup>45</sup> based on observed genetic distances between species. When stratified

188 by TR motif length, our estimates are highly concordant with *de novo* mutation rates observed in  
189 human pedigree data<sup>4</sup> ([Fig. 2g](#)). This concordance is particularly remarkable since SMM-based  
190 estimates overestimate mutation rates for loci which undergo multi-step mutations, and for the *de*  
191 *novo* observations we consider only mutation occurrence, not magnitude of change<sup>46,47</sup>. For short  
192 motifs where we have the most data, our SMM-based mutation rate estimates are higher than *de*  
193 *novo* observations. This may reflect frequent multistep mutations, as have been observed  
194 particularly for monomers and dimers<sup>46,48</sup>, or a TR mutation rate slow-down on the human  
195 lineage. Both mutation rate estimates generally decrease with motif length until about 20bp,  
196 when they increase. This pattern suggests a shift in the dominant mutational process, potentially  
197 from replication slippage in short motifs to recombination or gene conversion in longer motifs<sup>49</sup>.  
198 However, when considering only coding and regulatory TRs, model-based estimates drop below  
199 *de novo* mutation estimates, with concordance decreasing by an order of magnitude (Extended  
200 Data Fig. 4; MSE =  $2.5 \times 10^{-8}$  for intergenic TRs versus  $1.4 \times 10^{-7}$  for CDS TRs). This likely  
201 reflects violation of the assumption of selective neutrality at TRs subject to stronger stabilizing  
202 selection over evolutionary time. Further, we observe concordance between heterozygosity and  
203 mutation rate estimates over motif lengths, underscoring the role of mutational dynamics in  
204 shaping TR variation.

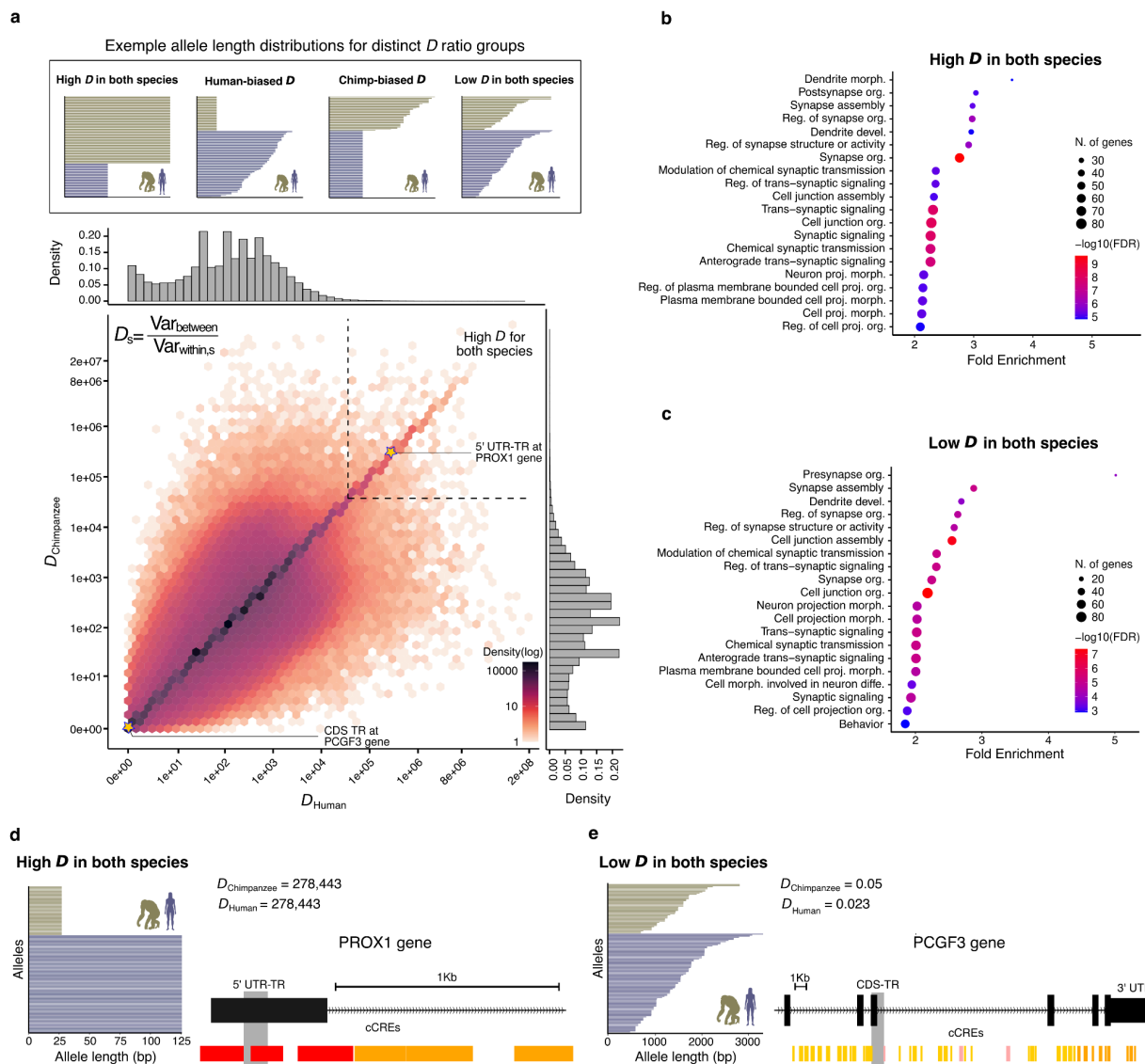
205

#### 206 **TRs with exceptional divergence or diversity**

207 Since TR mutation rate varies so widely across loci, allele length divergence between  
208 species of any given TR can arise from either species-specific selection or high mutation rates.  
209 We sought to account for this using an HKA-like approach by computing a species-specific  
210 divergence-diversity ratio ( $D$ ) for each TR (See Methods), defined as the variance in TR length  
211 between species divided by the variance within the focal species. This approach controls for  
212 per-locus mutation rate variation by identifying loci with extreme divergence between species  
213 relative to within-species diversity.  $D$  values were broadly concordant between species ([Fig. 3a](#);  
214 Pearson's  $r = 0.15$ ,  $p < 2.2 \times 10^{-16}$ ; Spearman's  $\rho = 0.84$ ;  $p < 2.2 \times 10^{-16}$ ), consistent with a  
215 general pattern of similar TRV mutational dynamics and selective pressures between humans and  
216 chimpanzees. To explore outlier TRVs as candidates for selection, we considered the 1,000 most  
217 extreme genic TRVs in each of three groups: 1) those with high  $D$  in both species, reflecting  
218 elevated divergence between species relative to within-species variation, consistent with  
219 species-specific directional selection; 2) those with low  $D$  in both species, characterized by  
220 relatively reduced divergence compared to within-species diversity, which may reflect balancing  
221 selection or strong stabilizing selection with a high mutation rate; and 3) those with asymmetric  
222  $D$  patterns, where between-species divergence is coupled with higher diversity in either  
223 chimpanzees (chimp-biased  $D$ ) or humans (human-biased  $D$ ), which may be a product of  
224 directional selection<sup>50,51</sup> or species-specific runaway mutations<sup>36,52</sup> (Supplementary Table S3).  
225 Across groups, highly divergent genic TRVs were longer and with larger motifs (Wilcoxon  
226 rank-sum test, FDR-adjusted  $p < 0.0001$ ; Supplementary Fig. 11a-d). They also displayed greater  
227 GC content when compared to a set of background TRVs with the same TR length and motif

length distribution (Wilcoxon rank-sum test, FDR-adjusted  $p < 0.0001$ ; Supplementary Fig. 11c) and were highly enriched in coding regions compared to the distribution of all genic TRVs (Odds ratio = 6.6,  $p < 2.2 \times 10^{-16}$ ; Supplementary Fig. 11e).

To investigate the biological processes associated with TRVs with extreme divergence or diversity, we performed GO enrichment analysis on the genes intersecting TRVs from each ratio group. While all TR-containing genes are modestly enriched for broad developmental processes (~1.2-fold; Supplementary Fig. 12), genes containing both high and low  $D$  TRs showed >2-fold enrichment for terms underlying nervous system development, synaptic organization and cell signaling (Fig. 3b and c), despite their contrasting evolutionary patterns. For comparison, random sets of TR-containing genes showed no GO enrichment (see Supplementary Methods).



238

Fig. 3. Divergence-diversity landscapes of TRs in humans and chimpanzees. a, Heatmap showing the joint distribution of TR Divergence-Diversity Ratios ( $D$ ) in humans and chimpanzees. Marginal histograms display the distribution of  $D$  for each species. Black dashed line shows the boundaries of high  $D$  in both species. Stars indicate

242 example TRs belonging to extreme ratio categories in both species. **b,c**, Gene Ontology enrichment analysis for  
243 biological processes terms associated with genes intersecting the top 1,000 genic TRs with **(b)** high  $D$  in both  
244 species and low  $D$  in both species. The set of all TR-containing genes was used as background. **d,e**, Genomic  
245 location and allele length distribution for two example TRs classified as **(d)** high  $D$  in both species and **(e)** low  $D$  in  
246 both species. Grey vertical bars represent TR location. Candidate cis-regulatory elements (cCREs) indicate  
247 promoter-like signature (red), proximal (orange), and distal (yellow) enhancer-like signature, and DNase-H3K4me3  
248 elements (pink).

249

250 Of the 32 genes containing coding TRVs in the high  $D$  group, 13 are zinc-fingers (ZNF)  
251 genes, a functionally diverse family involved in processes such as transcriptional regulation and  
252 DNA repair<sup>53</sup>. Several ZNFs are known to undergo rapid lineage-specific divergence and positive  
253 selection on DNA-binding domains between humans and chimpanzees<sup>54,55</sup>. The strongest  
254 divergence signals from genic TRVs occur in introns of genes associated with neural and sensory  
255 systems, epithelial integrity, and signal transduction (Supplementary Table S4). One compelling  
256 candidate overlaps the 5' UTR of PROX1 (Fig. 3d), a homeobox transcription factor essential for  
257 embryonic and central nervous system development<sup>56</sup>.

258 At the opposite extreme, six of the ten genic TRVs with the lowest  $D$  ratios occur in  
259 immune-related genes (Supplementary Table S4), including PRKCE, XRCC4, and CALCR,  
260 which have been implicated in innate immunity, antibody diversification, and immune-associated  
261 signaling, respectively<sup>57-60</sup>. One striking example exhibits nearly all unique alleles in the coding  
262 regions of PCGF3 (Fig. 3e), which acts primarily as a transcriptional activator required for  
263 mesodermal differentiation<sup>61</sup>, but also contributes to antiviral immunity by promoting  
264 interferon-responsive gene transcription<sup>62</sup>. Although extreme diversity compared to divergence  
265 alone could be caused by stabilizing selection coupled with high mutation rate<sup>63</sup>, the  
266 overrepresentation of immune genes is consistent with long-term balancing selection<sup>64,65</sup>.

267 We also consider genic TRVs with asymmetric  $D$  values, which reflect diversity in one  
268 species relative to the level of between-species divergence. Human-biased TRVs were enriched  
269 for cell morphogenesis and nervous system development terms, particularly neurogenesis  
270 (Extended Data Fig. 5a), and were primarily located in introns of cell signaling and intracellular  
271 trafficking genes, several with brain-specific activity (Supplementary Table S5). In contrast,  
272 chimp-biased TRVs were enriched only for cell-cell adhesion (Extended Data Fig. 5b), and occur  
273 largely in introns of genes involved in RNA processing, cell signaling and cytoskeletal  
274 organization (Supplementary Table S5). While such asymmetric  $D$  patterns can be explained by  
275 species-specific elevated mutation rates rapidly introducing new alleles, the genomic context of  
276 several of these TRVs suggests a contribution from directional selection, or a combination of  
277 both processes.

278

### 279 **Differential gene expression and TR variation**

280 Given the predicted causal role of TRVs in gene expression variation in humans<sup>1,8</sup>, we  
281 analyzed 7,050 orthologous genes expressed in humans and chimpanzees across six tissues<sup>66</sup>.  
282 Genes harboring promoter TRVs did not exhibit greater expression divergence (Wilcoxon  
283 rank-sum test, FDR-adjusted  $p > 0.05$  across all tissues; Supplementary Fig. 13). However, we

284 identified a weak but significant positive correlation between absolute gene expression  
285 divergence and the average absolute TRV allele length divergence (Pearson's  $r = 0.02$ ,  
286 FDR-adjusted  $p = 0.006$ ; Supplementary Fig. 14). This is consistent with studies suggesting that  
287 evolutionary changes in TR lengths can modulate expression levels of nearby genes<sup>1,16</sup>.  
288 Furthermore, this result agrees with studies showing that regions of elevated expression tend to  
289 exhibit an increase in the rate of mutations<sup>67,68</sup>, thus presumably harboring higher variation and  
290 divergence. When stratified by tissue type and genomic feature, significant positive correlations  
291 were observed in a subset of comparisons, including brain x intron TRs, kidney x 5' UTR TRs,  
292 liver x intron TRs, and testis x promoter TRs (Supplementary Fig. 15). To further investigate this  
293 relationship, we focused on 4,125 differentially expressed (DE) genes between humans and  
294 chimpanzees in at least one of the six tissues, identified using *limma*<sup>66</sup> (see Methods). While  
295 TRVs are generally not overrepresented in DE genes (Odds ratio = 0.84,  $p = 0.987$ ), DE genes  
296 are significantly enriched for TRs with strong evidence to impact gene expression (eTRs<sup>1,69</sup>)  
297 (Odds ratio = 1.45,  $p = 0.012$ ). Thus, while human-chimpanzee differential expression does not  
298 appear to be primarily driven by genomic TR species length variation, limited context-dependent  
299 associations may be due to a subset of regulatory TR.

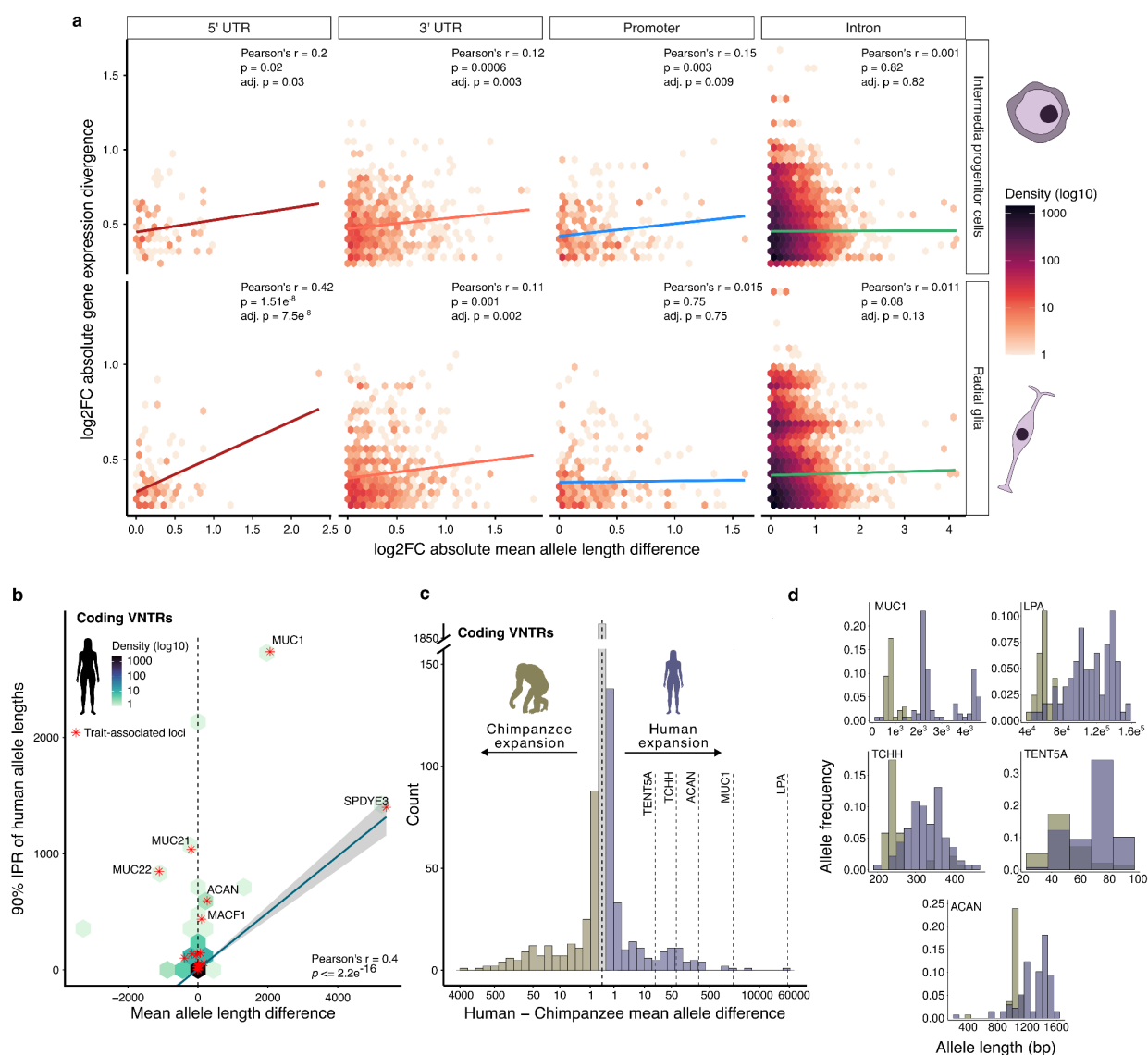
300 Given longstanding evidence of the role of TRs in neuronal development and brain  
301 function<sup>70,71</sup>, we focused on a set of 738 DE genes between human and chimpanzee in organoids  
302 with telencephalon identity<sup>20</sup>. Again, we found a weak but significant positive correlation  
303 between mean absolute TR allele length divergence averaged over genes and absolute gene  
304 expression divergence (Supplementary Fig. 16a; Pearson's  $r = 0.074$ , FDR-adjusted  $p = 0.0009$ ).  
305 When stratified by cell type and genomic feature, positive correlations were significant for  
306 intermediate progenitor cells and radial glia cells in UTR and promoter TRVs (Fig. 4a; Extended  
307 Data Fig. 6). We also observe a modest but significant enrichment of TRs with larger mean  
308 lengths in humans among genes upregulated in humans relative to those upregulated in  
309 chimpanzees (Supplementary Fig. 16b; Odds ratio = 1.11,  $p = 1.05 \times 10^{-12}$ ). These results suggest  
310 a connection between TR expansions and regulatory divergence in primates, in line with  
311 previous findings that TR divergence causes expression changes<sup>19</sup>.

312 By incorporating human and macaque cells from primary telencephalon samples into the  
313 DE analysis, Pollen et al., (2019) identified 261 candidate genes with human-specific regulatory  
314 changes during cortical development, which were significantly enriched for TRVs (Odds ratio =  
315 2.15,  $p = 4.25 \times 10^{-6}$ ). The majority of these candidate genes are up-regulated in humans ( $n=207$ )  
316 in one or more cell types, and several contain divergent TRVs (Supplementary Fig. 17). These  
317 patterns of gene expression divergence and allele length variation suggest that a subset of TRs  
318 might affect the regulation of genes involved in cortical neurogenesis, potentially contributing to  
319 human-specific aspects of neural development. Examples include VPS53 and PTPRS, which  
320 harbor highly divergent intronic VNTRs overlapping or near regulatory regions with high  
321 chromatin accessibility (Supplementary Fig. 18). Point mutations in VPS53 are linked to  
322 progressive cerebello-cerebral atrophy<sup>72</sup> and increased levels of PTPRS caused by a point  
323 mutation are associated with decreased risk of Alzheimer's disease<sup>73</sup>.

324

### 325 Functional TR variation

326 Beyond their regulatory potential, TRs can directly affect protein structure and function.  
 327 For example, a large-scale analysis of UK Biobank data identified a set of 25 coding VNTRs  
 328 with significant trait association<sup>74</sup>. We found that of these 25 trait-associated TRs, 19 exhibit  
 329 significant differences in length distribution between humans and chimpanzees, with 11 showing  
 330 longer alleles in humans (Supplementary Fig. 19; Wilcoxon rank-sum test, FDR-adjusted  $p <$   
 331 0.0001). Notably, nine of these fall within the top 2% of the  $\sim 219,593$  genic TRs showing higher  
 332 mean lengths in humans, drawn from 837,612 homologous genic TRs.



333

334 Fig. 4 - **Gene expression, divergence, and trait-associated TRs.** **a**, Heatmaps showing the correlation between  
 335 absolute log fold change in mean TR allele length across TRs overlapping 5'UTR, 3'UTR and promoter, and  
 336 intronic regions (x-axis) and absolute log fold change in gene expression divergence between humans and  
 337 chimpanzee (y-axis) across intermediate progenitor cells and radial glia cells. **b**, TR mean allele length divergence  
 338 between chimpanzees and humans (x-axis) versus the 90% interpercentile range of human allele lengths per locus.

339 Red asterisks denote loci with evidence of genotype-phenotype association. The LPA locus is omitted. **c**,  
340 Distribution of mean allele length difference between humans and chimpanzees for coding VNTRs, highlighting five  
341 trait-causal VNTRs<sup>74</sup>. **d**, Chimpanzee and human allele frequency distributions for the same five loci.

342

343 These trait-associated VNTRs display not only high divergence but also increased human  
344 diversity (Fig. 4b). In a companion paper<sup>23</sup>, we show a similar trend in pathogenic TRs, which  
345 are consistently expanded and highly diverse in humans compared to chimpanzees. Notably, the  
346 25 trait-associated VNTRs from Mukamel et al<sup>74</sup> and 61 expansion-disorder TRs from the  
347 STRipy database<sup>75</sup> are overrepresented among TRVs with exceptional relative length divergence  
348 and diversity in humans (top 10% in each)(Odds ratio = 10.66, CI: 5.13-20.51,  $p = 1.12 \times 10^{-8}$ ).  
349 Focusing on five fine-mapped TR loci with strong evidence that length polymorphism, rather  
350 than linked SNVs or indels, directly drives trait association (posterior probability > 0.95<sup>74</sup>), we  
351 show that mean allele length difference between humans and chimpanzees was significantly  
352 greater than expected under the null distribution for all coding VNTRs (permutation test,  
353 1,000,000 runs,  $p \leq 1 \times 10^{-6}$ , and  $P = 0.003$  when excluding the long and highly divergent VNTR  
354 in LPA) (Fig. 4c, d; Supplementary Fig. 20). Repeat length variation in these loci directly affects  
355 protein structure by altering the length of protein domains (see Fig. S1 from Mukamel et al.<sup>74</sup>).  
356 This pattern is expected for TRs with species-specific elevated mutability, particularly when  
357 longer expanded alleles have higher mutation rates<sup>5,41</sup>. Selection may also play a role if  
358 intermediate-length alleles are too weakly deleterious to overcome high mutation rates<sup>76,77,78</sup>, or if  
359 they are beneficial in an ecological context until further expansion causes fitness reduction, in  
360 the case of expansion disorders.

361 These observations support a hypothesis that highly expanded and diverse TRs are strong  
362 candidates for functional variation. This may be because increased TR length and variability  
363 provide a greater opportunity to alter genome structure and therefore function, potentially  
364 leading to a larger phenotypic impact<sup>79,80</sup>. In addition, by having disproportionately large trait  
365 effect sizes, the impact of these TRs may be more detectable, leading to ascertainment bias.  
366 Regardless of the underlying mechanism, our results suggest that these TRs are important targets  
367 for future large-scale comparative studies aimed at understanding the contribution of repeat  
368 variation to genome function, phenotypic diversity, and adaptation across evolutionary  
369 timescales.

## 370 Discussion

371 TRs have long been hypothesized to play a significant role in phenotypic variation and  
372 evolutionary change. However, their evolutionary dynamics have remained underexplored due to  
373 limitations in resolving repeat variation at scale across species. Here, using assembly-level  
374 long-read data in a comparative framework, we present a comprehensive survey of the  
375 heterogeneous landscape of TR variation in primates, with focus on humans and chimpanzees.  
376 We recovered loci spanning a continuum of selective regimes, from pervasive length  
377 conservation in functional regions, consistent with stabilizing selection, to high divergence  
378 consistent with directional selection. We also observe loci with exceptional diversity, resulting  
379 either from balancing selection or elevated mutation rates and strong stabilizing selection,

380 particularly outside functional or regulatory contexts.

381 TRs in CDS, and to a lesser extent 5' UTRs, exhibit reduced polymorphism and extreme  
382 length conservation across evolutionary time (~28.5 My; [Fig. 1g](#), [Fig. 2a-c](#)), signatures of  
383 negative selection. However, while TRs are depleted in CDS, they are enriched in 5' UTRs ([Fig.](#)  
384 [1f](#)). This suggests that in these regions the high cost of their mutational risk is outweighed by a  
385 selected functional role, possibly in RNA folding<sup>81</sup>, translational regulation<sup>82</sup>, or modulation of  
386 transcription factor binding affinity<sup>83</sup>. Hence, TRs in 5' UTR repeats are strong candidates for  
387 functional molecular impact.

388 We estimated TR-specific mutation rates based on TR divergence between humans and  
389 chimpanzees under a simple single-step mutation model (SMM) and selective neutrality, yielding  
390 results similar to those from trio-based average *de novo* mutation rate estimates<sup>4</sup>. This similarity  
391 suggests that, in broad strokes, on average the neutral SMM reasonably describes TR evolution,  
392 while groups of loci where the estimates depart suggest deviation from the SMM or neutrality. At  
393 the same time, we observe that species-specific expansions are associated with elevated  
394 within-species diversity, which could be caused by locus and species-specific increases in  
395 expansion-biased mutational processes.

396 Leveraging TR divergence-diversity ratios, we identified putative candidates for  
397 directional and balancing selection. TRs with both exceptional relative divergence and  
398 exceptional relative diversity were enriched in genes involved in nervous system development  
399 and synaptic processes, agreeing with prior associations between TR variation and  
400 neurodevelopment<sup>8,71,84</sup>. Notably, most known pathogenic repeat expansions are associated with  
401 neurological or neurodegenerative disorders, further supporting the particular sensitivity of  
402 nervous system processes to TR variation. We also identified TRVs with high diversity compared  
403 to divergence, including several extreme examples in immune-related genes. Although this  
404 pattern of variation may be due to high mutation rates coupled with stabilizing selection, it is  
405 also expected under long-term balancing selection, particularly at immune loci<sup>21,85</sup>.

406 We found associations between brain organoid expression divergence and divergence of  
407 TRs located in promoters and UTRs, where repeat variation can directly alter chromatin  
408 organization, nucleosome occupancy, and affect alternative splicing<sup>16,86</sup>. Further, genes with  
409 human-specific regulatory changes during cortical development are enriched in TRs, some of  
410 which are highly divergent. More broadly, trait and disease-causal TRs are associated  
411 bidirectionally with longer human alleles and higher diversity. These patterns may reflect  
412 elevated, lineage-specific mutational processes leading to runaway expansions with functional  
413 consequences, or possibly the action of directional selection shaping repeat length distributions.  
414 While the underlying mechanism remains unresolved, we hypothesize that exceptional length  
415 divergence and species-specific diversity may serve as useful criteria for prioritizing TRs for  
416 functional studies.

417 While we provide a comprehensive overview of TR variation in primates, several  
418 limitations should be noted. First, estimates of within-species diversity are restricted to humans  
419 and chimpanzees. Second, we focused on loci shared between species. Future studies including

420 species-specific TRs may provide additional insight into forces impacting TR birth and death.  
421 Third, current TR-phenotype association evidence is largely derived from human short-reads,  
422 which have limited resolution for long TRs and do not capture trait associations in other species.  
423 Comparable long-read studies across NHP are essential to assess whether similar patterns in  
424 trait-causal loci extend beyond humans. Finally, our analyses are limited to TR allele length  
425 polymorphism and do not capture other forms of variation capable of affecting phenotypes, such  
426 as sequence interruptions and changes in motif composition<sup>87</sup>.

427 In summary, our results show that TR variation reflects the interplay of mutational  
428 processes and selective pressures, acting both as conserved functional elements and as fuel for  
429 evolutionary innovation. Ongoing advances in TR genotyping and mutation inference, along  
430 with increasingly accessible population-scale cross-species variant catalogs, will enable  
431 systematic characterization of TR mutational dynamics and their contribution to adaptation  
432 across evolutionary timescales.

433

## 434 **Methods**

### 435 **TR reference catalogs**

436 Using the TRACK pipeline<sup>88</sup>, we generated TR catalogs for eight ape T2T reference  
437 genomes<sup>6,89</sup> ([Fig. 1b](#)). Briefly, TRs were identified with Tandem Repeat Finder v.4.09<sup>90</sup> using the  
438 following parameters: *matchscore* 2, *mismatchscore* 5, *indelscore* 7, *pm* 80, *pi* 10, *minscore* 24,  
439 *maxperiod* 2000, *-l* 6. Resulting annotations were filtered by total length (>11 bp and <10 Kbp),  
440 copy number (>2.5), and constancy score, i.e., sequence similarity between adjacent repeat units  
441 (>60%). When TRs overlapped by 5 bp or less, we retained the repeat with the shortest motif  
442 length. We then normalized motif sequences to their smallest periodic units. For example, the  
443 motif “ATATAT” was reduced to “AT”, and the copy number was recalculated accordingly. This  
444 normalization step was implemented because TRF reports consensus motifs that maximize  
445 alignment scores, which do not always correspond to the minimal motifs. As a result, TRF output  
446 may contain nested motifs composed of repeated instances of a shorter underlying unit. Finally,  
447 we queried our catalog against the Dfam database to identify instances where TRs intersect  
448 known repetitive element families<sup>91</sup>. After this initial characterization of the catalogs, we applied  
449 additional filtering to exclude centromeric regions and regions containing alpha satellite DNA  
450 (cenSat) and subterminal satellites (StSat). Annotations were obtained from the CHM13 and  
451 T2T-Primate Consortium to identify and remove complex, high-order repeat (HOR)-rich regions  
452 before homology assessment.

453

### 454 **Annotation of TR genomic features**

455 To classify TRs by genomic feature, we used the GENCODE GFF2 gene annotation for  
456 CHM13. We retained only transcripts annotated as the APPRIS principal isoform to ensure  
457 high-confidence gene models<sup>92</sup>. Exons, coding sequences (CDS), and introns were directly  
458 extracted from the annotation. Transcription start sites (TSS) were defined as the 5' end of each  
459 transcript, and promoter regions were defined as the 1 Kb upstream of the TSS, accounting for

460 strand orientation. Untranslated regions (UTRs) were inferred based on exon coordinates: 5'  
461 UTRs were inferred as the exon segments upstream of the start codon, and 3' UTRs as those  
462 downstream of the stop codon, considering strand orientation. To avoid overestimating UTRs,  
463 we required inferred regions to not overlap CDS regions. The TR catalog was intersected to these  
464 annotations and a hierarchical classification was applied to resolve edge cases where TRs  
465 overlap more than one feature: CDS > 5' UTR > 3' UTR > promoter > intron > intergenic.

466

#### 467 **Homology assessment**

468 To identify homologous TRs between reference genomes, we used a multi-step  
469 alignment-based pipeline implemented in TRACK<sup>88</sup>. TR coordinates from a target genome were  
470 lifted to a query genome using the UCSC LiftOver tool<sup>93</sup> with *-minMatch* 0.1 and *-bedPlus=3*  
471 *-tab* to preserve metadata tracking of TRs across genome builds. Lifted TRs were intersected  
472 with the native TR catalog of the query genome using *bedtools intersect*, requiring a reciprocal  
473 overlap of at least 10%.

474 For each overlapping TR pair, we extracted and compared their motifs. To ensure  
475 strand-invariant and phase-independent comparison, TRACK computes the lexicographically  
476 smallest cyclic permutation of each motif and its reverse complement. This step allows the direct  
477 comparison of motif sequences regardless of strand orientation or rotational phase. To assess  
478 motif similarity, we performed global pairwise alignments between each candidate homologous  
479 TR pair with EMBOSS Needle<sup>94</sup>, using the following parameters: *-gapopen 10* and *-gapextend*  
480 *0.5*. Each motif pair was aligned in both forward and reverse-complement orientation, and we  
481 retained the alignment with the highest similarity score. TR pairs with the best alignment  
482 similarity score  $\geq 95\%$  were retained as confidently homologous. Homology detection was  
483 performed bidirectionally, with both genomes in the pair used as target and query. The final set  
484 of homologous TRs was defined as the intersection of high-similarity TR pairs identified in both  
485 directions. This reciprocal filtering strategy ensures that homologs are robust to mapping  
486 artifacts or asymmetric genome annotations.

487 To quantify the correlation of normalized TR density between species we computed  
488 density in non-overlapping 1 Mb windows in the human genome, and windows were lifted over  
489 to the corresponding NHP genomes. Lifted windows were filtered based on length, retaining  
490 regions between 0.8-2 Mb for most species, and 0.2-2 Mb for the more divergent taxa, i.e.  
491 siamang gibbon and macaque, where liftover performance is reduced. For each species, TR  
492 density was calculated as the proportion of bps covered by TRs within each lifted genomic  
493 window and normalized to sum to 1 prior to comparison.

494

#### 495 **Determining TR genotypes**

496 TRs were genotyped using PacBio HiFi data from 46 humans in the Human Pangenome  
497 Research Consortium (HPRC)<sup>95</sup> and 23 near-T2T chimpanzee genomes, described in a  
498 companion paper<sup>23</sup>, with Tandem Repeat Genotyping Tool v.3.0 (TRGT)<sup>96</sup>. Variants were filtered  
499 for missing data (*--max-missing 1*), minimum allele spanning depth ( $>3$ ) and allele constancy

500 score ( $>0.6$ ).

501

### 502 **Estimating TR mutation rates**

503 Per-locus mutation rates were estimated under the stepwise mutation model (SMM) and  
504 neutrality based on the observed genetic distance between humans and chimpanzees.  
505 Specifically, we computed the TR genetic distance, originally termed  $\delta\mu^2$ , but we refer to as  $d^{45}$ :

$$506 \quad d = \left( \sum_i a_i p_i - \sum_j a_j p_j \right)^2$$

507 where  $a_i$  is allele copy number and  $p_i$  is the allele frequency for allele  $i$  in humans, and  $j$  indexed  
508 quantities in chimpanzees. Under neutrality and the SMM, the expected value of  $d$  increases  
509 linearly with the mutation rate ( $\mu$ ) and the divergence time ( $t$ )<sup>45</sup>. We therefore estimated  
510 per-generation mutation rates as

$$511 \quad \hat{\mu} = \frac{d}{t_{total}}$$

512 where  $t_{total}$  is the sum of the human and chimpanzee branch lengths. We assumed branch lengths  
513 of 7.7 million years for each lineage and generation times of 29 years for humans and 25 years  
514 for chimpanzees, yielding  $t_{total} = 5.8 \times 10^5$  generations<sup>97</sup>.

515

### 516 **Divergence-diversity ratio ( $D$ )**

517 TRV per locus divergence-diversity ratio ( $D$ ) was computed as the ratio of  
518 between-species to within-species variance in allele length, computed separately for each species  
519 after removing 3 outlier alleles per species:

$$520 \quad D_s = \frac{Var_{between}}{Var_{within,s}}$$

521 where between-species variance for each TR was defined as the weighted squared difference  
522 between each species-specific mean length and the overall mean length across both species:

$$523 \quad Var_{between} = \sum_{s=1}^2 n_s (\bar{x}_s - \bar{x})^2, \quad \bar{x} = \frac{\sum_{s=1}^2 n_s \bar{x}_s}{\sum_{s=1}^2 n_s}$$

524 where  $n_s$  is the number of alleles in species  $s$ , and  $\bar{x}$  is the overall mean length across both  
525 species. One characteristic of  $D_s$  is that elevated within-species diversity can also inflate between  
526 species variance, reducing its power to detect loci with exceptionally high relative diversity.

527 TRs with high  $D_s$  in both species were defined as the top 1,000 loci in both species, while  
528 TRs with low  $D_s$  in both species were defined as the bottom 1,000 loci after excluding repeats  
529 with within-species variance below 1. We implemented this threshold as a conservative noise

530 floor in order to eliminate small TRs with low variance that produce small ratios likely not  
531 biologically relevant (Supplementary Fig. 1). Species-specific extremes were defined based on  
532 deviation from the diagonal in the human-chimpanzee *D* comparison, representing loci with  
533 pronounced divergence relative to diversity in a single lineage.

534

### 535 **Gene ontology enrichment analysis**

536 We first performed Gene Ontology (GO) enrichment analysis to test whether genes  
537 containing TRs were associated with specific biological process terms, using the full set of  
538 human genes as a background set. Next, we selected genes intersecting the top 1,000 TRVs  
539 within each group of extreme *D* values, using the full set of genes intersecting TRs as the  
540 background set. Analyses were conducted in ShinyGO v.0.82<sup>98</sup> with a minimum pathway size of  
541 15 and a maximum pathway size of 1,000 genes. Significance was assessed using the false  
542 discovery rate (FDR), and only pathways with FDR-corrected  $p < 0.01$  were considered  
543 enriched. To increase confidence in the results, we performed enrichment analysis with 10  
544 random sets of 1,000 TR-containing genes, which showed no GO enrichment.

545

### 546 **TR divergence and gene expression divergence**

547 Gene expression data were obtained from Brawand et al.<sup>66</sup>, which provides normalized  
548 RPKM (reads per kilobase of exon per million mapped reads) values for six tissues (brain,  
549 cerebellum, heart, kidney, liver, and testis) from humans and chimpanzees. To reduce noise from  
550 lowly expressed genes, we applied an RPKM threshold of  $>1$  in at least one individual for each  
551 tissue and species. Differential expression (DE) between species was assessed per tissue using  
552 *limma*<sup>99</sup> on log-transformed RPKM values, with empirical Bayes moderation of gene-wise  
553 standard errors. Genes with FDR-adjusted  $p$ -values  $< 0.05$  were considered DE. TR length  
554 divergence was quantified as the log<sub>2</sub> fold change of mean allele length per locus and compared  
555 to gene expression divergence using Pearson's correlation. To focus on genes containing TRs  
556 known to be associated with expression (eTRs), we retained those associated with fine-mapped  
557 expression STR<sup>1</sup> and expression-associated VNTRs<sup>69</sup>.

558

### 559 **Code Availability**

560 All code used in the paper can be obtained at [https://github.com/caroladam/tr\\_analyses](https://github.com/caroladam/tr_analyses).

561

### 562 **Data availability**

563 The HPRC data can be obtained at <https://humanpangenome.org/data/>. The chimpanzee  
564 assemblies can be obtained at [pending]. Chimpanzee-human homologous TR catalog and VCF  
565 files for humans and chimpanzees can be obtained at 10.5281/zenodo.20616249 [pending  
566 publication].

567

### 568 **Acknowledgements**

569 We thank the Rohlfs lab for helpful discussion.

570

571 **Funding statement**

572 This work was supported by the National Science Foundation (NSF) CAREER award 2144878  
573 to R.V.R., the NIH National Institute of General Medicine award R35GM142916 to P.H.S., and  
574 the NIH National Human Genome Research Institute award R01HG013017 to P.H.S.

575

576 **Author contributions**

577 R.V.R., P.H.S. conceived the design of the study, acquired funding, and supervised the research.  
578 C.L.A, J.L.R., R.V.R., P.H.S. processed and analyzed the data. C.L.A, R.R. drafted the  
579 manuscript. C.L.A, J.L.R., R.V.R., P.H.S. reviewed and edited the manuscript.

580

581 **Competing interests**

582 The authors declare no competing financial interests.

583

584 **Additional information**

585 Supplementary information is available for this paper.

586 Correspondence should be addressed to R.V.R. and C.L.A.

## 587 References

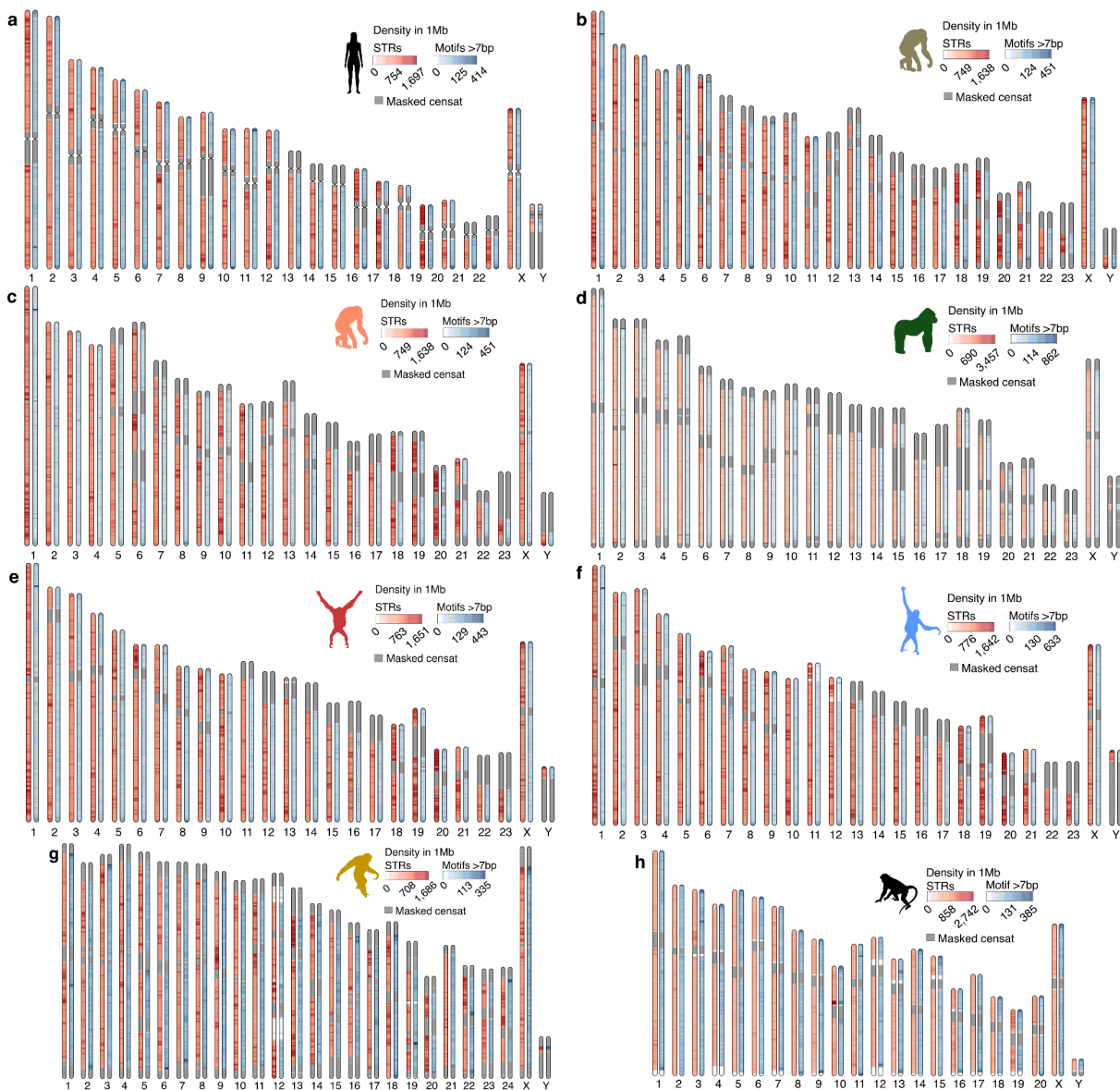
- 588 1. Fotsing, S. F. *et al.* The impact of short tandem repeat variation on gene expression. *Nat. Genet.* **51**,  
589 1652–1659 (2019).
- 590 2. Doyle, L. *et al.* Rational design of  $\alpha$ -helical tandem repeat proteins with closed architectures. *Nature*  
591 **528**, 585–588 (2015).
- 592 3. Hannan, A. J. Tandem repeat polymorphisms: modulators of disease susceptibility and candidates for  
593 ‘missing heritability’. *Trends Genet.* **26**, 59–65 (2010).
- 594 4. Porubsky, D. *et al.* Human de novo mutation rates from a four-generation pedigree reference. *Nature*  
595 **643**, 427–436 (2025).
- 596 5. Brinkmann, B., Klitschar, M., Neuhuber, F., Hühne, J. & Rolf, B. Mutation rate in human  
597 microsatellites: influence of the structure and length of the tandem repeat. *Am. J. Hum. Genet.* **62**,  
598 1408–1415 (1998).
- 599 6. Nurk, S. *et al.* The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
- 600 7. Fan, H. & Chu, J.-Y. A brief review of short tandem repeat mutation. *Genomics Proteomics*  
601 *Bioinformatics* **5**, 7–14 (2007).
- 602 8. Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene expression variation in  
603 humans. *Nat. Genet.* **48**, 22–29 (2016).
- 604 9. Erwin, G. S. *et al.* Recurrent repeat expansions in human cancer genomes. *Nature* **613**, 96–102  
605 (2023).
- 606 10. Schloissnig, S. *et al.* Long-read sequencing and structural variant characterization in 1,019 samples  
607 from the 1000 Genomes Project. *bioRxiv.org* (2024) doi:[10.1101/2024.04.18.590093](https://doi.org/10.1101/2024.04.18.590093).
- 608 11. Song, J. H. T., Lowe, C. B. & Kingsley, D. M. Characterization of a human-specific tandem repeat  
609 associated with bipolar disorder and schizophrenia. *Am. J. Hum. Genet.* **103**, 421–430 (2018).
- 610 12. Huang, Y. *et al.* Short tandem repeats in populations of the Qinghai-Tibet Plateau and adjacent  
611 regions provide insights into high-altitude adaptation. *Sci. Adv.* **11**, eadx1590 (2025).
- 612 13. Zhou, K., Aertsen, A. & Michiels, C. W. The role of variable DNA tandem repeats in bacterial  
613 adaptation. *FEMS Microbiol. Rev.* **38**, 119–141 (2014).
- 614 14. Fondon, J. W., 3rd & Garner, H. R. Molecular origins of rapid and continuous morphological  
615 evolution. *Proc. Natl. Acad. Sci. U. S. A.* **101**, 18058–18063 (2004).
- 616 15. Villanea, F. A. *et al.* The MUC19 gene: An evolutionary history of recurrent introgression and  
617 natural selection. *Science* **389**, eadl0882 (2025).
- 618 16. Vences, M. D., Legendre, M., Caldara, M., Hagihara, M. & Verstrepen, K. J. Unstable tandem repeats  
619 in promoters confer transcriptional evolvability. *Science* **324**, 1213–1216 (2009).
- 620 17. L Rocha, J., Lou, R. N. & Sudmant, P. H. Structural variation in humans and our primate kin in the  
621 era of telomere-to-telomere genomes and pangenomics. *Curr. Opin. Genet. Dev.* **87**, 102233 (2024).
- 622 18. Liu, Q. & Tian, W. Association of human-specific expanded short tandem repeats with  
623 neuron-specific regulatory features. *Sci. Adv.* **11**, eadp9707 (2025).
- 624 19. Sulovari, A. *et al.* Human-specific tandem repeat expansion and differential gene expression during  
625 primate evolution. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 23243–23253 (2019).
- 626 20. Pollen, A. A. *et al.* Establishing cerebral organoids as models of human-specific brain evolution. *Cell*  
627 **176**, 743–756.e17 (2019).
- 628 21. Leffler, E. M. *et al.* Multiple instances of ancient balancing selection shared between humans and  
629 chimpanzees. *Science* **339**, 1578–1582 (2013).

- 630 22. Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature* **499**, 471–475  
631 (2013).
- 632 23. Rocha, J. *et al.* A Pan-pangenome illuminates complex structural variation and selection in humans,  
633 chimpanzees, and bonobos. *bioRxiv* (2026) doi:[10.64898/2026.06.06.730619](https://doi.org/10.64898/2026.06.06.730619).
- 634 24. Sharma, A. & Sowpati, D. T. Analysis of tandem repeats in seven telomere-to-telomere primate  
635 genomes. *J. Genet.* **104**, 14 (2025).
- 636 25. Srivastava, S., Avvaru, A. K., Sowpati, D. T. & Mishra, R. K. Patterns of microsatellite distribution  
637 across eukaryotic genomes. *BMC Genomics* **20**, 153 (2019).
- 638 26. Verbiest, M. *et al.* Mutation and selection processes regulating short tandem repeats give rise to  
639 genetic and phenotypic diversity across species. *J. Evol. Biol.* **36**, 321–336 (2023).
- 640 27. El-Sawy, M. & Deininger, P. Tandem insertions of Alu elements. *Cytogenet. Genome Res.* **108**,  
641 58–62 (2005).
- 642 28. Webster, M. T., Smith, N. G. C. & Ellegren, H. Microsatellite evolution inferred from  
643 human-chimpanzee genomic sequence alignments. *Proc. Natl. Acad. Sci. U. S. A.* **99**, 8748–8753  
644 (2002).
- 645 29. Grimwood, J. *et al.* The DNA sequence and biology of human chromosome 19. *Nature* **428**, 529–535  
646 (2004).
- 647 30. Subramanian, S., Mishra, R. K. & Singh, L. Genome-wide analysis of microsatellite repeats in  
648 humans: their abundance and density in specific genomic regions. *Genome Biol.* **4**, R13 (2003).
- 649 31. Rhie, A. *et al.* The complete sequence of a human Y chromosome. *Nature* **621**, 344–354 (2023).
- 650 32. Linthorst, J. *et al.* Extreme enrichment of VNTR-associated polymorphicity in human subtelomeres:  
651 genes with most VNTRs are predominantly expressed in the brain. *Transl. Psychiatry* **10**, 369  
652 (2020).
- 653 33. Cechova, M. *et al.* High satellite repeat turnover in great apes studied with short- and long-read  
654 technologies. *Mol. Biol. Evol.* **36**, 2415–2431 (2019).
- 655 34. Koga, A., Hirai, Y., Hara, T. & Hirai, H. Repetitive sequences originating from the centromere  
656 constitute large-scale heterochromatin in the telomere region in the siamang, a small ape. *Heredity*  
657 (*Edinb.*) **109**, 180–187 (2012).
- 658 35. Madsen, B. E., Villesen, P. & Wiuf, C. Short tandem repeats in human exons: a target for disease  
659 mutations. *BMC Genomics* **9**, 410 (2008).
- 660 36. Usdin, K., House, N. C. M. & Freudenreich, C. H. Repeat instability during DNA repair: Insights  
661 from model systems. *Crit. Rev. Biochem. Mol. Biol.* **50**, 142–167 (2015).
- 662 37. Schaper, E., Gascuel, O. & Anisimova, M. Deep conservation of human protein tandem repeats  
663 within the eukaryotes. *Mol. Biol. Evol.* **31**, 1132–1148 (2014).
- 664 38. Byeon, G. W. *et al.* Functional and structural basis of extreme conservation in vertebrate 5'  
665 untranslated regions. *Nat. Genet.* **53**, 729–741 (2021).
- 666 39. Wieder, N. *et al.* The role of untranslated region variants in Mendelian disease: a review. *Eur. J.*  
667 *Hum. Genet.* **33**, 1096–1105 (2025).
- 668 40. Kronenberg, Z. N. *et al.* High-resolution comparative analysis of great ape genomes. *Science* **360**,  
669 eaar6343 (2018).
- 670 41. Neff, B. D. & Gross, M. R. Microsatellite evolution in vertebrates: inference from AC dinucleotide  
671 repeats. *Evolution* **55**, 1717–1733 (2001).
- 672 42. Bilgin Sonay, T. *et al.* Tandem repeat variation in human and great ape populations and its impact on  
673 gene expression divergence. *Genome Res.* **25**, 1591–1599 (2015).

- 674 43. Li, W. H. & Sadler, L. A. Low nucleotide diversity in man. *Genetics* **129**, 513–523 (1991).
- 675 44. Willems, T. *et al.* The landscape of human STR variation. *Genome Res.* **24**, 1894–1904 (2014).
- 676 45. Goldstein, D. B., Ruiz Linares, A., Cavalli-Sforza, L. L. & Feldman, M. W. Genetic absolute dating  
677 based on microsatellites and the origin of modern humans. *Proc. Natl. Acad. Sci. U. S. A.* **92**,  
678 6723–6727 (1995).
- 679 46. Gymrek, M., Willems, T., Reich, D. & Erlich, Y. Interpreting short tandem repeat variations in  
680 humans using mutational constraint. *Nat. Genet.* **49**, 1495–1501 (2017).
- 681 47. Sainudiin, R., Durrett, R. T., Aquadro, C. F. & Nielsen, R. Microsatellite mutation models: insights  
682 from a comparison of humans and chimpanzees. *Genetics* **168**, 383–395 (2004).
- 683 48. Mitra, I. *et al.* Patterns of de novo tandem repeat mutations and their role in autism. *Nature* **589**,  
684 246–250 (2021).
- 685 49. Lai, Y. & Sun, F. The relationship between microsatellite slippage mutation rate and the number of  
686 repeat units. *Mol. Biol. Evol.* **20**, 2123–2131 (2003).
- 687 50. Beaumont, M. A. & Balding, D. J. Identifying adaptive genetic divergence among populations from  
688 genome scans. *Mol. Ecol.* **13**, 969–980 (2004).
- 689 51. Excoffier, L., Foll, M. & Petit, R. J. Genetic consequences of range expansions. *Annu. Rev. Ecol.*  
690 *Evol. Syst.* **40**, 481–501 (2009).
- 691 52. Ellegren, H. Microsatellites: simple sequences with complex evolution. *Nat. Rev. Genet.* **5**, 435–445  
692 (2004).
- 693 53. Cassandri, M. *et al.* Zinc-finger proteins in health and disease. *Cell Death Discov.* **3**, 17071 (2017).
- 694 54. Jovanovic, V. M. *et al.* Positive selection in gene regulatory factors suggests adaptive pleiotropic  
695 changes during human evolution. *Front. Genet.* **12**, 662239 (2021).
- 696 55. Nowick, K. *et al.* Gain, loss and divergence in primate zinc-finger genes: a rich resource for  
697 evolution of gene regulatory differences between species. *PLoS One* **6**, e21553 (2011).
- 698 56. Elsir, T., Smits, A., Lindström, M. S. & Nistér, M. Transcription factor PROX1: its role in  
699 development and cancer. *Cancer Metastasis Rev.* **31**, 793–805 (2012).
- 700 57. Altman, A. & Kong, K.-F. Protein kinase C enzymes in the hematopoietic and immune systems.  
701 *Annu. Rev. Immunol.* **34**, 511–538 (2016).
- 702 58. Soulas-Sprauel, P. *et al.* Role for DNA repair factor XRCC4 in immunoglobulin class switch  
703 recombination. *J. Exp. Med.* **204**, 1717–1727 (2007).
- 704 59. Wang, S., Wang, W. & Zeng, J. Role of CALCR expression in liver cancer: Implications for the  
705 immunotherapy response. *Mol. Med. Rep.* **31**, 41 (2025).
- 706 60. Maleitzke, T. *et al.* The calcitonin receptor protects against bone loss and excessive inflammation in  
707 collagen antibody-induced arthritis. *iScience* **25**, 103689 (2022).
- 708 61. Zhao, W. *et al.* Polycomb group RING finger proteins 3/5 activate transcription via an interaction  
709 with the pluripotency factor Tex10 in embryonic stem cells. *J. Biol. Chem.* **292**, 21527–21537  
710 (2017).
- 711 62. Da, G. *et al.* Nuclear PCGF3 inhibits the antiviral immune response by suppressing the  
712 interferon-stimulated gene. *Cell Death Discov.* **10**, 429 (2024).
- 713 63. Zhang, X.-S. & Hill, W. G. Genetic variability under mutation selection balance. *Trends Ecol. Evol.*  
714 **20**, 468–470 (2005).
- 715 64. Bitarello, B. D. *et al.* Signatures of long-term balancing selection in human genomes. *Genome Biol.*  
716 *Evol.* **10**, 939–955 (2018).
- 717 65. Minias, P. & Vinkler, M. Selection balancing at innate immune genes: Adaptive polymorphism

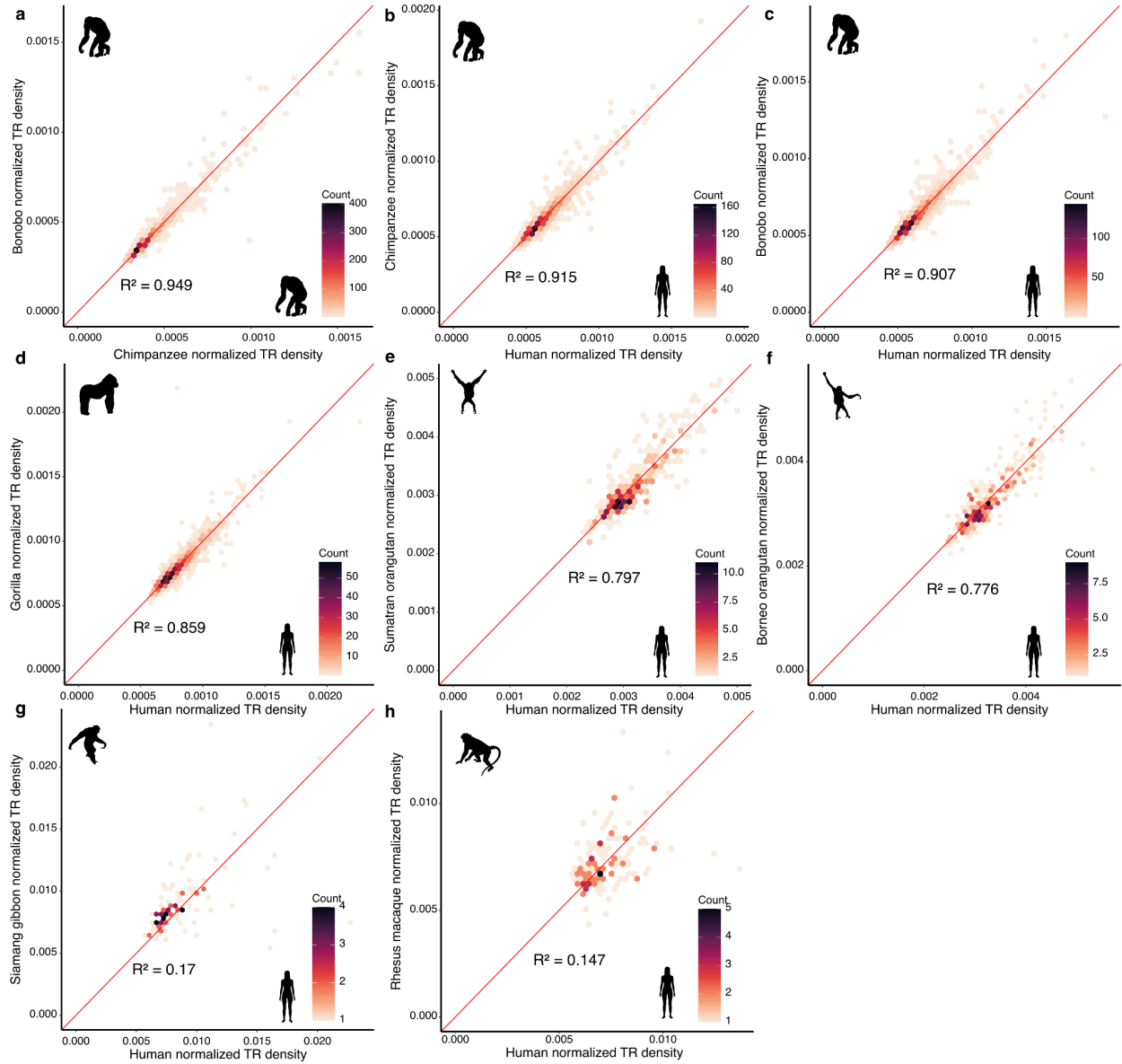
- 718 maintenance in Toll-like receptors. *Mol. Biol. Evol.* **39**, msac102 (2022).
- 719 66. Brawand, D. *et al.* The evolution of gene expression levels in mammalian organs. *Nature* **478**,  
720 343–348 (2011).
- 721 67. Park, C., Qian, W. & Zhang, J. Genomic evidence for elevated mutation rates in highly expressed  
722 genes. *EMBO Rep.* **13**, 1123–1129 (2012).
- 723 68. Bachl, J., Carlson, C., Gray-Schopfer, V., Dessing, M. & Olsson, C. Increased transcription levels  
724 induce higher mutation rates in a hypermutating cell line. *J. Immunol.* **166**, 5051–5057 (2001).
- 725 69. Eslami Rasekh, M., Hernández, Y., Drinan, S. D., Fuxman Bass, J. I. & Benson, G. Genome-wide  
726 characterization of human minisatellite VNTRs: population-specific alleles and gene expression  
727 differences. *Nucleic Acids Res.* **49**, 4308–4324 (2021).
- 728 70. Hammock, E. A. D. & Young, L. J. Microsatellite instability generates diversity in brain and  
729 sociobehavioral traits. *Science* **308**, 1630–1634 (2005).
- 730 71. Xiao, X. *et al.* Revisiting tandem repeats in psychiatric disorders from perspectives of genetics,  
731 physiology, and brain evolution. *Mol. Psychiatry* **27**, 466–475 (2022).
- 732 72. Feinstein, M. *et al.* VPS53 mutations cause progressive cerebello-cerebral atrophy type 2 (PCCA2).  
733 *J. Med. Genet.* **51**, 303–308 (2014).
- 734 73. Poirier, A. *et al.* PTPRS is a novel marker for early Tau pathology and synaptic integrity in  
735 Alzheimer’s disease. *Sci. Rep.* **14**, 14718 (2024).
- 736 74. Mukamel, R. E. *et al.* Protein-coding repeat polymorphisms strongly shape diverse human  
737 phenotypes. *Science* **373**, 1499–1505 (2021).
- 738 75. Halman, A., Dolzhenko, E. & Oshlack, A. STRipy: A graphical application for enhanced genotyping  
739 of pathogenic short tandem repeats in sequencing data. *Hum. Mutat.* **43**, 859–868 (2022).
- 740 76. Pajic, P. & Gokcumen, O. Evolutionary balancing of genetic consequence and innovation in  
741 mammals through variable number tandem repeats. *Genome Biol. Evol.* **18**, evaf250 (2026).
- 742 77. Press, M. O., Hall, A. N., Morton, E. A. & Queitsch, C. Substitutions are boring: Some arguments  
743 about parallel mutations and high mutation rates. *Trends Genet.* **35**, 253–264 (2019).
- 744 78. Ibañez, K. *et al.* Increased frequency of repeat expansion mutations across different populations. *Nat.*  
745 *Med.* **30**, 3357–3368 (2024).
- 746 79. Gemayel, R., Vinces, M. D., Legendre, M. & Verstrepen, K. J. Variable tandem repeats accelerate  
747 evolution of coding and regulatory sequences. *Annu. Rev. Genet.* **44**, 445–477 (2010).
- 748 80. Manigbas, C. A. *et al.* A phenome-wide association study of tandem repeat variation in 168,554  
749 individuals from the UK Biobank. *Nat. Commun.* **15**, 10521 (2024).
- 750 81. Kinney, N., Pathak, D., Evans, E. & Arias, P. Short tandem repeat variants are possibly associated  
751 with RNA secondary structure and gene expression. *PLoS One* **20**, e0326355 (2025).
- 752 82. Leppek, K., Das, R. & Barna, M. Functional 5’ UTR mRNA structures in eukaryotic translation  
753 regulation and how to find them. *Nat. Rev. Mol. Cell Biol.* **19**, 158–174 (2018).
- 754 83. Horton, C. A. *et al.* Short tandem repeats bind transcription factors to tune eukaryotic gene  
755 expression. *Science* **381**, eadd1250 (2023).
- 756 84. Cui, Y. *et al.* Multi-omic quantitative trait loci link tandem repeat size variation to gene regulation in  
757 human brain. *Nat. Genet.* **57**, 369–378 (2025).
- 758 85. Ferrer-Admetlla, A. *et al.* Balancing selection is the main force shaping the evolution of innate  
759 immunity genes. *J. Immunol.* **181**, 1315–1322 (2008).
- 760 86. Hamanaka, K. *et al.* Genome-wide identification of tandem repeats associated with splicing variation  
761 across 49 tissues in humans. *Genome Res.* **33**, 435–447 (2023).

- 762 87. Rajan-Babu, I.-S., Dolzhenko, E., Eberle, M. A. & Friedman, J. M. Sequence composition changes  
763 in short tandem repeats: heterogeneity, detection, mechanisms and clinical implications. *Nat. Rev.*  
764 *Genet.* **25**, 476–499 (2024).
- 765 88. Adam, C. L., Rocha, J., Sudmant, P. & Rohlf, R. TRACKing tandem repeats: a customizable  
766 pipeline for identification and cross-species comparison. *Bioinform. Adv.* **5**, vbaf066 (2025).
- 767 89. Yoo, D. *et al.* Complete sequencing of ape genomes. *bioRxiv.org* (2024)  
768 doi:[10.1101/2024.07.31.605654](https://doi.org/10.1101/2024.07.31.605654).
- 769 90. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* **27**,  
770 573–580 (1999).
- 771 91. Hubley, R. *et al.* The Dfam database of repetitive DNA families. *Nucleic Acids Res.* **44**, D81–9  
772 (2016).
- 773 92. Rodriguez, J. M. *et al.* APPRIS: annotation of principal and alternative splice isoforms. *Nucleic*  
774 *Acids Res.* **41**, D110–7 (2013).
- 775 93. Hinrichs, A. S. *et al.* The UCSC Genome Browser Database: update 2006. *Nucleic Acids Res.* **34**,  
776 D590–8 (2006).
- 777 94. Madeira, F. *et al.* The EMBL-EBI Job Dispatcher sequence analysis tools framework in 2024.  
778 *Nucleic Acids Res.* **52**, W521–W525 (2024).
- 779 95. Liao, W.-W. *et al.* A draft human pangenome reference. *Nature* **617**, 312–324 (2023).
- 780 96. Dolzhenko, E. *et al.* Characterization and visualization of tandem repeats at genome scale. *Nat.*  
781 *Biotechnol.* **42**, 1606–1614 (2024).
- 782 97. Shao, Y. *et al.* Phylogenomic analyses provide insights into primate evolution. *Science* **380**, 913–924  
783 (2023).
- 784 98. Ge, S. X., Jung, D. & Yao, R. ShinyGO: a graphical gene-set enrichment tool for animals and plants.  
785 *Bioinformatics* **36**, 2628–2629 (2020).
- 786 99. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and  
787 microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).



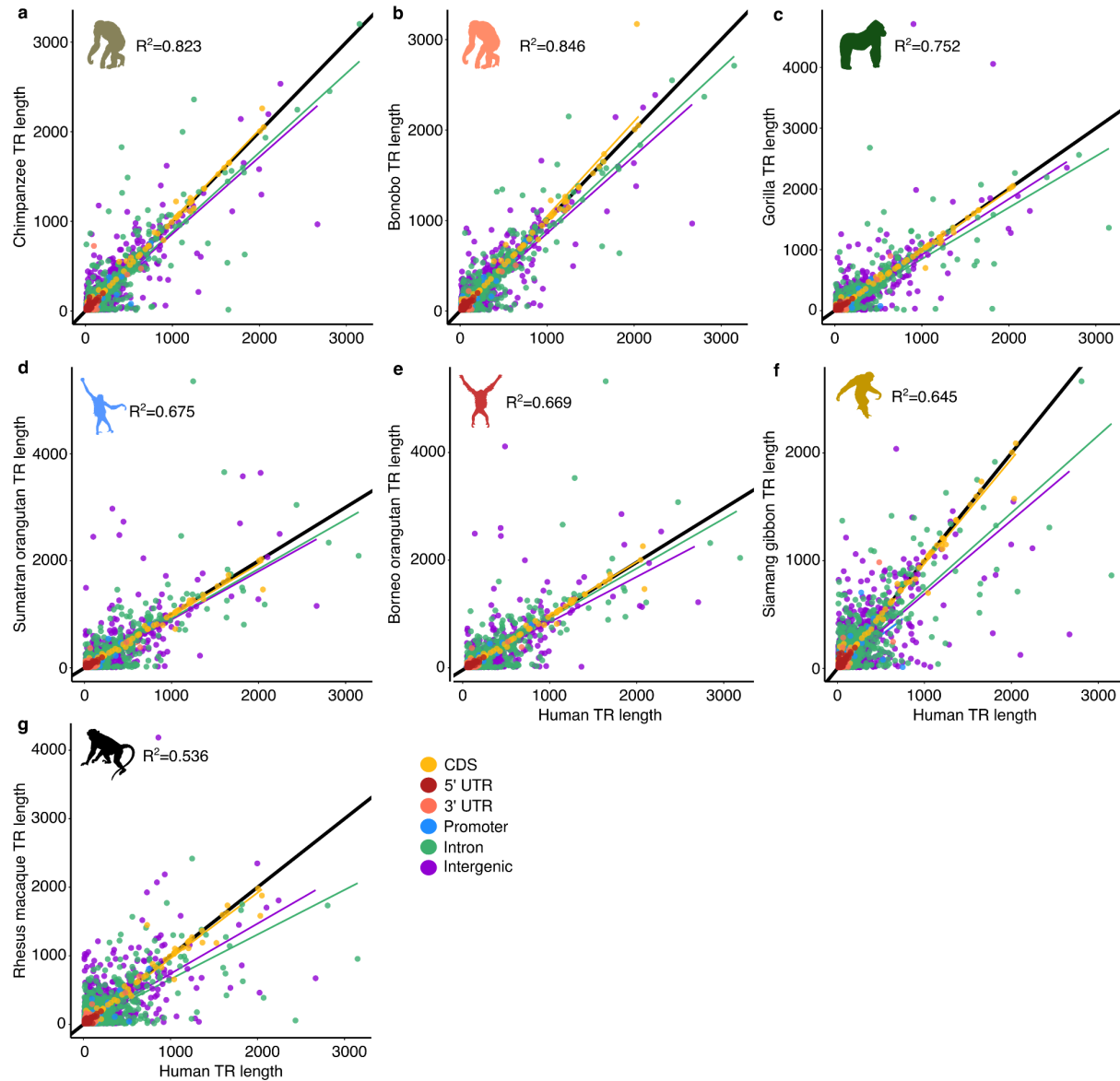
788

789 Extended Data Fig. 1 - **Genomic TR distributions in non-human primate genomes.** Ideograms  
 790 showing the density of STRs in red and VNTRs in blue across non-overlapping 1 Mb windows. Masked  
 791 CenSat regions are marked in gray. TR densities are shown for **a**, *Homo sapiens*; **b**, *Pan troglodytes*; **c**,  
 792 *Pan paniscus*; **d**, *Gorilla gorilla*; **e**, *Pongo pygmaeus*; **f**, *Pongo abelii*; **g**, *Symphalangus syndactylus*; and  
 793 **h**, *Macaca mulatta*.



794

795 Extended Data Fig. 2 - **Correlation of TR density in homologous regions.** Correlation between  
796 normalized density of TRs between **a**, chimpanzee and bonobo; and between human and **b**, chimpanzee;  
797 **c**, bonobo; **d**, gorilla; **e**, Sumatran orangutan; **f**, Borneo orangutan; **g**, siamang gibbon; and **h**, macaque.



798

799 **Extended Data Fig. 3 - Landscape of TR length variation across human and non-human primate**  
 800 **T2T reference genomes.** Scatterplot of reference TR lengths between human (x-axis) and seven  
 801 non-human T2T reference primates (y-axis): **a**, *Pan troglodytes*; **b**, *Pan paniscus*; **c**, *Gorilla gorilla*; **d**,  
 802 *Pongo pygmaeus*; **e**, *Pongo abelii*; **f**, *Symphalangus syndactylus*; and **g**, *Macaca mulatta*. Each point  
 803 represents a single TR locus, color-coded according to CHM13 genomic annotation.

804

805

806

807

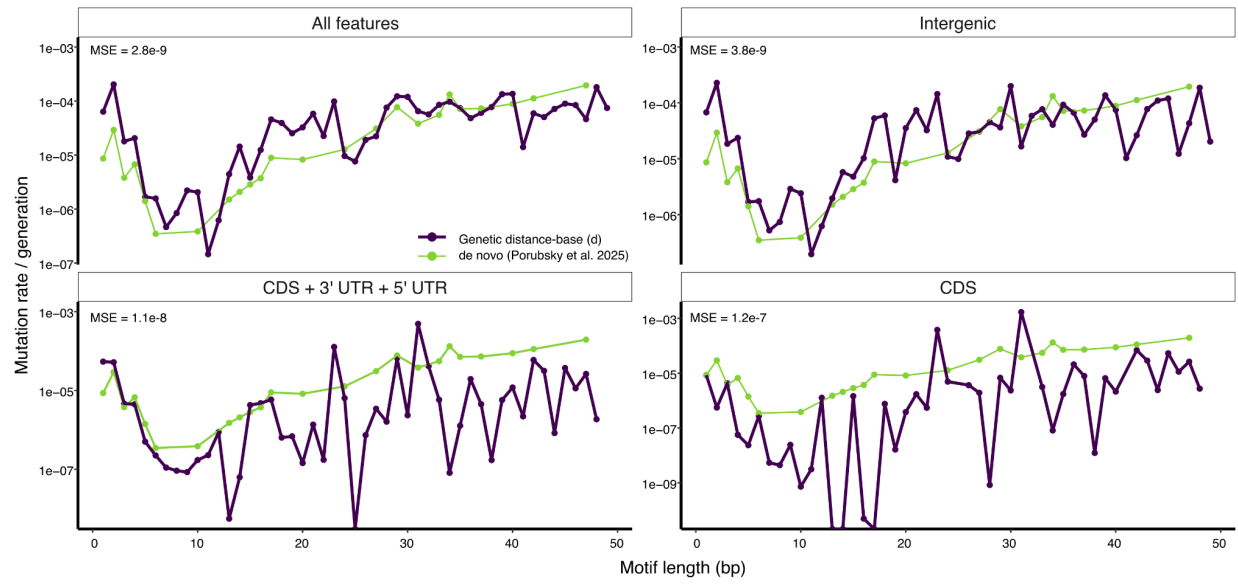
808

809

810

811

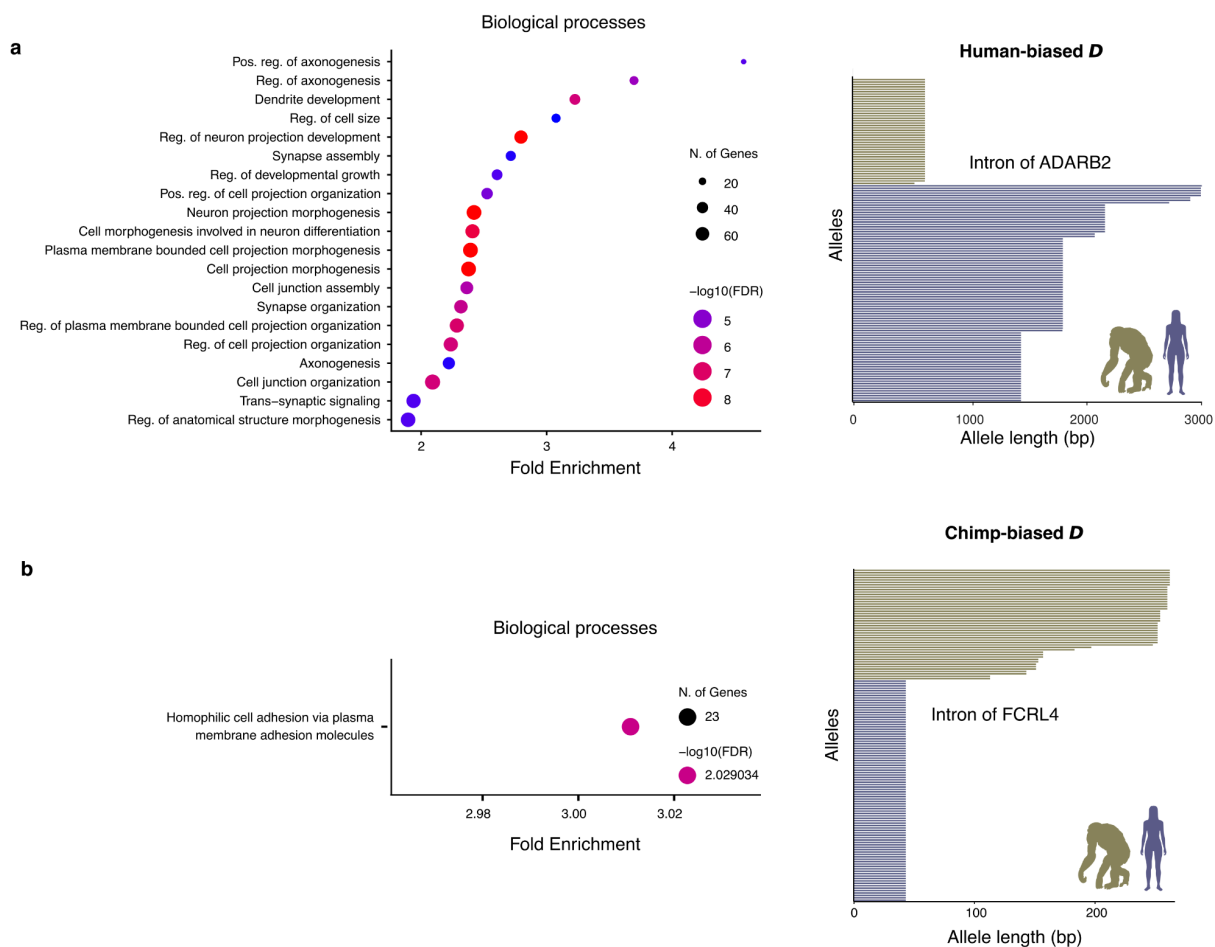
812



813

814 Extended Data Fig. 4 - Estimates of *de novo* and genetic distance-based (*d*) TR mutation rates averaged

815 across motif lengths for different subsets of loci.



816

817 Extended Data Fig. 5 - **Biological processes associated with asymmetric divergence-diversity (*D*)**

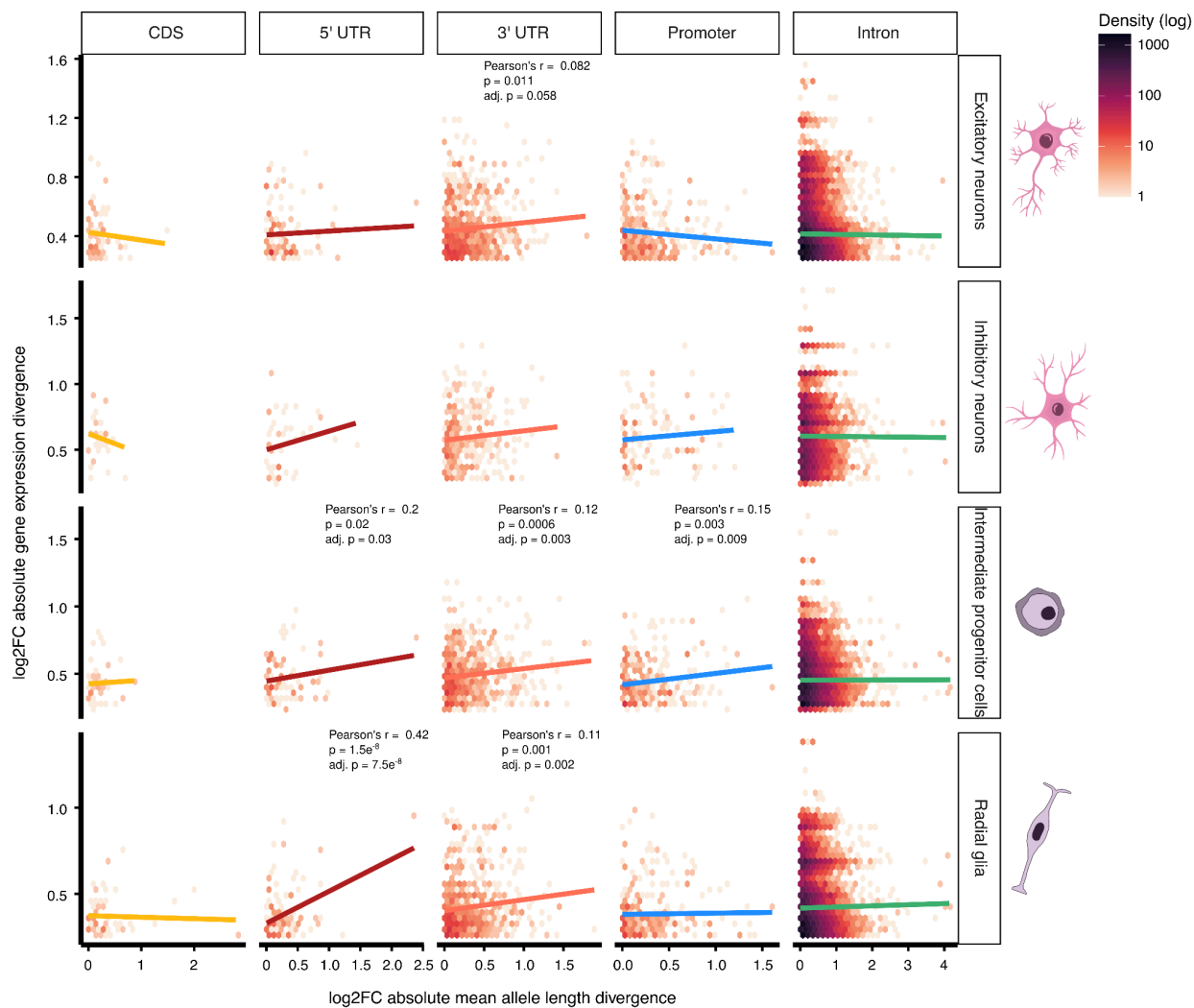
818 **ratios.** Gene Ontology enrichment analysis for biological processes terms associated with genes

819 intersecting the top 1,000 genic TRs with asymmetric *D*, and allele length distribution of one

820 representative genic TR with **a**, human-biased *D*, with an intronic TR at ADARB2; and **b**, chimp-biased

821 *D*, with an intronic TR at FCRL4.

822



823

824 Extended Data Fig. 6 - Heatmaps showing the correlation between absolute log fold change in mean TR  
 825 allele length across genes in each genomic feature (x-axis) and absolute log fold change in gene  
 826 expression divergence between humans and chimpanzees (y-axis) across organoids with telencephalon  
 827 identity from Pollen et al. (2019).