**Title**

**Integrated analysis sheds light on evolutionary trajectories of young transcription start sites in the human genome**

Cai Li[1,#], Boris Lenhard[2,3,4], Nicholas M. Luscombe[1,5,6,#]

[1] The Francis Crick Institute, 1 Midland Road, London, NW1 1AT, UK

[2] Computational Regulatory Genomics, MRC London Institute of Medical Sciences, Du Cane Road, London, W12 0NN, UK

[3] Institute of Clinical Sciences, Faculty of Medicine, Imperial College London, Du Cane Road, London, W12 0NN, UK

[4] Sars International Centre for Marine Molecular Biology, University of Bergen, Thormøhlensgate 55, N-5008 Bergen, Norway

[5] Okinawa Institute of Science & Technology Graduate University, Okinawa, 904-0495, Japan

[6] UCL Genetics Institute, University College London, Gower Street, London, WC1E 6BT, UK

# Corresponding authors: cai.li@crick.ac.uk and nicholas.luscombe@crick.ac.uk

22 **Abstract**

23 Previous studies revealed widespread transcription initiation and fast turnover of

24 transcription start sites (TSSs) in mammalian genomes. Yet how new TSSs originate

25 and how they evolve over time remain poorly understood. To address these questions,

26 we analyzed ~200,000 human TSSs by integrating evolutionary and functional

27 genomic data, particularly focusing on TSSs that emerged in the primate lineages. We

28 found that intrinsic factors of repetitive sequences and their proximity to established

29 regulatory modules (extrinsic factors) contribute significantly to origin of new TSSs.

30 In early periods, young TSSs experience rapid sequence evolution driven by

31 endogenous mutational mechanisms that reduce the instability of associated repetitive

32 sequences. In later periods, the regulatory functions of young TSSs are gradually

33 modified, and with evolutionary changes subject to temporal (fewer regulatory

34 changes in younger TSSs) and spatial constraints (fewer regulatory changes in more

35 isolated TSSs). These findings advance our understanding of how regulatory

36 innovations arise in the genome throughout evolution and highlight the roles of

37 repetitive sequences in these processes.

38

## 1. Introduction

Many studies revealed that transcription is pervasive in prokaryotic and eukaryotic genomes[1,2]. One recent study found that three-quarters of the human genome can be transcribed[3], indicating a much more complex transcriptional landscape than previously thought. Transcription Start Sites (TSSs) are the genomic loci where transcription initiation occurs and thus are a critical class of regulatory element for transcriptional control. By harnessing diverse high-throughput sequencing technologies, studies in the past few years have greatly improved TSS annotation in model organism genomes, especially human, and uncovered new characteristics of transcriptional initiation[4-6]. One intriguing phenomenon about TSSs is that they occur widely throughout the genome, not only in typical promoters of annotated genes, but also in other regions such as intergenic or intronic loci. For example, some enhancers also contain TSSs, producing so-called enhancer RNAs[7-9].

Many previous studies about TSS evolution focused on cross-species comparisons and revealed interesting macro-evolutionary patterns[10-14]. For example, by comparing human and mouse TSSs, a recent study found that >56% of protein-coding genes have experienced TSS turnover events since humans and mice diverged[13]. Genes with TSS turnover were also found to experience adaptive evolution in their coding regions and expression levels[13]. Unlike macro-evolution, however, micro-evolutionary processes (i.e. intra-species evolution) of TSSs are relatively poorly understood. Given the high turnover rate of TSSs[13], population genomic data can provide a more detailed view of TSS evolution. Although some previous studies made use of population genomic data, they pooled all TSSs together to compare with non-TSS elements[15] or focused on purifying selection[13,16]. Since different TSSs could have distinct evolutionary histories, pooling all TSSs together could bury the interesting characteristics of a specific TSS categories. A recent comprehensive study in *Drosophila melanogaster* populations investigating the relationship between genetic variations and TSS usage identified thousands of genetic variants affecting transcript levels and promoter shapes, providing important new insights into TSS evolution at the population level[17].

Despite extensive investigation, many questions about TSSs are yet to be addressed. Importantly, the evolutionary origin of new TSSs and evolutionary trajectories of newly emerged TSSs remain unresolved. Previous studies have suggested that

3

71    repetitive sequences are a rich source of new TSSs[13,18], but the underlying

72    mechanisms of how these sequences contribute to novel transcription initiation remain

73    unclear. For instance, why do some repetitive elements initiate transcription and

74    others not? How does the host genome handle the potential conflicts arising from the

75    inherent instability of repetitive elements associated with new TSSs? Furthermore, the

76    subsequent changes of newly emerged TSSs and their evolutionary fates have not

77    been systematically investigated. Only by addressing these questions can we begin to

78    understand how regulatory innovations arise in the genome throughout evolution and

79    how they contribute to biological diversity and adaptation.

80    To gain detailed insights into evolution of young TSSs and the underlying regulatory

81    mechanisms, we analyzed ~200,000 published human TSSs by integrating both

82    evolutionary (inter-species and intra-species) and functional genomic approaches,

83    with an emphasis on evolutionarily young TSSs that emerged in the primate lineages.

84    We show that 1) intrinsic factors of repetitive sequences and extrinsic chromatin

85    environments contribute significantly to the origin of novel transcription initiation; 2)

86    after emerging in the genome, young TSSs undergo rapid sequence evolution which is

87    likely due to several endogenous mutational mechanisms; and 3) regulatory outcomes

88    of young TSSs are gradually modified in subsequent periods and tend to be subject to

89    temporal and spatial constraints.

90    ## 2. Results

91    ### 2.1 Identification of evolutionarily young TSSs in the human genome

92    Using the cap analysis of gene expression (CAGE) sequencing technologies, the

93    FANTOM 5 project[4] generated the most comprehensive TSS annotation to date,

94    covering major primary cell types and tissues in human. To identify evolutionarily

95    young TSSs, we took advantage of the 'robust' human TSS dataset from FANTOM

96    project, which consists of 201,873 high-confidence TSSs. After filtering TSSs that

97    could confound downstream analysis (see Methods for details), we grouped the

98    remaining 151,902 TSSs into categories of different evolutionary ages. Since there is

99    no large-scale CAGE TSS annotation in the other primate genomes, it is impossible to

100   define the evolutionary ages of TSSs by comparing TSS annotations. However,

101   previous studies revealed that sequence-intrinsic properties of many promoters can

102    drive transcription initiation autonomously[19,20], indicating that the sequence itself is

103    an important determinant of promoter capacity. Moreover, Young et al. (2015) found

104    that, of those human TSSs that could be aligned to an orthologous sequence in the

105    mouse, more than 80% have detectable transcriptional initiation in mouse[13]. This

106    implies that if the orthologous sequence of a human TSS can be found in another

107    genome, it probably exhibits initiation in that species.

108    Therefore, to estimate the evolutionary ages of human TSS loci, we investigated the

109    sequence presence/absence patterns based on sequence alignments between human

110    and other 16 genomes (10 primate species representing major primate lineages and 6

111    non-primate mammalian species as outgroups). A human TSS locus is considered

112    present in another genome if the corresponding pairwise alignment satisfies: 1) a

113    mapping ratio of the human TSS peak (i.e. a CAGE tag cluster region predicted by

114    decomposition-based peak identification method in FANTOM) in another genome of

115    ≥90% and 2) a mapping ratio of the TSS peak±100 bp (considered as core promoter

116    region in this study) of ≥50% (see Methods and **Supplementary Tables 1-3** for more

117    details). Based upon the presence/absence patterns in alignments, we categorized the

118    human TSSs into four groups of different sequence ages (**Fig. 1a**): 1) TSSs whose

119    sequence loci can be found in at least one non-primate mammalian genome,

120    consisting of 141,117 TSSs (92.9% of all surveyed TSSs, named 'mammalian' group;

121    **Fig. 1a**); 2) TSSs whose sequences occurred during early primate evolution but before

122    the last common ancestor of Old World anthropoids, consisting of 6,668 TSSs (4.4%,

123    named 'primate' group; **Fig. 1a**); 3) TSSs whose sequences occurred during the

124    evolution of Old World anthropoids but before the last common ancestor of hominids,

125    consisting of 3,318 TSSs (2.2%, named 'OWA' group; **Fig. 1a**); 4) TSSs whose

126    sequences occurred since emergence of hominids, consisting of 799 TSSs (0.5%,

127    named 'hominid' group; **Fig. 1a**). The relatively large numbers of TSSs in three recent

128    periods corroborate the "frequent birth" phenomenon reported previously[13], and

129    enable us to perform detailed comparative analysis between these periods. Hereafter

130    we considered TSSs in the 'mammalian' group as evolutionarily old TSSs and those

131    in other three groups as evolutionarily young TSSs. For instance, in the gene *BAAT*

132    locus shown in **Fig. 1b**, there are two old TSSs present in both primate and non-

133    primate mammalian genomes, and one young TSS established during the evolution of

134    OWAs. The young TSS is located in a region overlapping one long terminal repeat

135  (LTR) element (**Fig. 1b**), suggesting that it originated from an LTR insertion. This
136  young TSS is expressed in many cell types where the old TSSs are expressed,
137  suggesting it may undertake part of the transcription task of old TSSs or up-regulate
138  the expression level of *BAAT* in some conditions.

139  We first examined some general features among TSS groups. We found that old TSSs
140  are mainly associated with mRNAs (59%), while many young TSSs are associated
141  with lncRNAs (54%~60%), indicating a compositional bias in the TSS groups (**Fig.**
142  **1c**). As TSSs become older, the proportion of mRNA TSSs becomes larger, and the
143  opposite happens to the intergenic lncRNA TSSs (**Fig. 1c**). Relative to older TSSs,
144  younger TSSs generally have narrower TSS peaks (**Fig. 1d**) and comprise more
145  TATA-box containing TSSs (**Fig. 1e**) and fewer CpG island (CGI)-associated TSSs
146  (**Fig. 1f**). This is consistent with previous observations about broad and sharp
147  promoters in mammalian genomes[4,21], which found that CGI promoters are usually
148  broad and associated with housekeeping genes, while TATA-box promoters are sharp
149  and associated with less conserved tissue-specific genes. Both old and young TSSs
150  exhibit elevated GC content and CpG content in TSS-proximal positions
151  (**Supplementary Fig. 1**), although relative to young TSSs, old TSSs tend to be more
152  GC-rich. We also noticed that the 'hominid' TSS group has higher average GC and
153  CpG content relative to 'OWA' and 'primate' groups (**Supplementary Fig. 1**), which
154  could be partly due to fewer historical deamination events of methylated cytosines in
155  very young TSS loci (see also later sections about DNA methylation).

156  **2.2 Sources of young TSSs**

157  **2.2.1 Intrinsic factors of repetitive sequences contribute to novel transcription**
158  **initiation**

159  Based upon the defined TSSs groups of different ages, next we systematically
160  investigated how new TSSs originate and how they evolve over time. Previous
161  analyses from earlier FANTOM projects showed that many mammalian transcripts
162  initiate within repetitive elements, especially retrotransposons[13,18]. Given the
163  extensive retrotransposition during mammalian evolution, retrotransposon-derived
164  TSSs could be an important source of novel TSSs. In addition, tandem repeats, which
165  are highly mutable loci, were found to be abundant in promoter regions and have
166  significant impact on gene expression[22,23]. With these observations in mind, we

167    examined the repetitive sequences (or 'repeats' for short hereafter) in all TSS loci,

168    including transposable elements (TEs, i.e. retrotransposons and DNA transposons)

169    and tandem repeats, based on annotations of RepeatMasker[24], TRF[25] and STRcat[26].

170    We found that ~70% of young TSSs have at least one repeat element within core

171    promoter regions (±100 bp of TSSs), but only 24% among old TSSs (**Fig. 2a**).

172    Whereas a large fraction (43%) of repetitive sequences associated with old TSSs are

173    tandem repeats, many young TSS loci are associated with retrotransposons, including

174    LTRs, long intersperse nuclear elements (LINEs) and short interspersed nuclear

175    elements (SINEs) (**Fig. 2a**). Because some tandem repeats could derive from

176    retrotransposons, we performed an alternative analysis considering only the nearest

177    retrotransposon element (**Supplementary Table 4 & Supplementary Fig. 2**). LTRs

178    are the most abundant retrotransposon class associated with young TSSs, with ~30%

179    of young TSSs are associated with LTRs. 14% and 8% of young TSSs are associated

180    with LINEs and SINEs, respectively. The large number of retrotransposons associated

181    with young TSSs suggests a major role of retrotransposition in forming new TSS loci.

182    Faulkner et al. (2009) revealed that many TE-derived TSSs are unevenly distributed

183    along TE element consensus sequences, and many TE-derived TSSs are not present in

184    the canonical 5' promoters of TE elements[18]. However, how these TE-derived

185    sequences contribute to transcription initiation was not discussed in detail and thus

186    remain poorly understood. To gain more detailed insight into this question, we first

187    mapped TSSs to the TE consensus sequences like Faulkner et al. (2009), and analyzed

188    the distributions of TSSs along repeat elements. The distributions obtained from our

189    analysis are similar to those in Faulkner et al. (2009), but also exhibit some

190    differences. The differences are likely due to the upgraded CAGE protocols[27] and

191    improvements in the TSS calling method[28], which largely overcame some previous

192    issues such as 'multimapping' and 'exon painting' in early CAGE datasets used in

193    Faulkner et al. (2009).

194    We found that the TSSs associated with LTR elements are mainly in the sense strand

195    of LTRs and clustered within narrow regions (**Fig. 2b** for the THE1B subfamily and

196    **Supplementary Fig. 3** for more subfamilies). Since LTR elements contain the

197    promoters for endogenous retroviral elements (ERVs), the sense-biased distributions

198    of TSSs suggest that transcription initiation events in these regions are mainly

199    contributed by the original ERV promoter activities within LTRs. These patterns were

200   not observed in Faulkner et al. (2009), as they only investigated the distributions of
201   TSSs along LTR superfamilies but not the subfamilies. We also found that a large
202   fraction (~50%) of young TSSs associated with LTRs contain a TATA-box motif
203   starting at 25~35 bp upstream of the dominant TSSs (**Supplementary Fig. 4**),
204   whereas the ratio drops to ~30% for the old TSSs associated with LTRs, suggesting a
205   substantial fraction of TATA-box promoters derived from LTRs might have turned
206   into TATA-less promoters during evolution.

207   LINE-1(L1) is the most abundant LINE family in the human genome (covering ~20%
208   of human genome). The overall distribution of TSSs along L1 elements (**Fig. 2c**) is
209   similar to that in Faulkner et al. (2009). However, we further observed many
210   differences in the TSS distributions between different L1 subfamilies
211   (**Supplementary Fig. 5**). For some subfamilies, transcription initiation occurs mainly
212   at the region of 5'end antisense promoters (e.g. L1PB1, L1PBa1) which were
213   discussed in Faulkner et al. (2009), whereas for other subfamilies the initiation occurs
214   mainly at the 3'end (e.g. L1MB7) or rather randomly (e.g. L1M4). Although the
215   background distribution of L1 subfamilies in the human genome can explain such
216   difference to some degree, it is apparently not the only reason (**Supplementary Fig.
217   5**). This suggests that sequences from different L1 subfamilies have very variable
218   propensity to drive transcription initiation.

219   Alu elements comprise the most abundant SINE family in the human genome
220   (covering ~10% of human genome). Although Alus are frequently inserted in
221   promoter-proximal and intronic regions, previous research found that they generally
222   lack capacity for driving autonomous transcription[20]. In the FANTOM5 dataset,
223   initially we observed many new TSSs located around the 3' poly(A) region and the A-
224   rich linker region, but later we found that these TSSs probably resulted from the
225   technical artifacts in the CAGE sequencing in FANTOM5 and thus filtered out the
226   related TSSs (**Supplementary Fig. 6**, see Methods for more details). The remaining
227   Alu-associated TSSs tend to be enriched at the 5'end of Alu in the antisense strand
228   (**Supplementary Fig. 6**), but how these sequences help drive transcription initiation is
229   unclear.

230   We found that ~9% of young TSSs contain tandem repeats which are not associated
231   with TEs. Unlike the tandem repeats derived from new TE insertions, the flank

232    regions of these tandem repeats tend to be conserved among mammals and have

233    higher GC content (**Supplementary Fig. 7**), suggesting that some new TSSs in these

234    regions are likely due to autonomous expansions of tandem repeats located in

235    proximal regions of pre-existing promoters (some examples provided in

236    **Supplementary Fig. 7**). This is consistent with previously reported enrichment of

237    tandem repeats in primate promoters[13,22,29,30].

238    Taken together, these findings suggest that repetitive sequences significantly

239    contribute to novel TSSs in multiple ways. Among the repetitive sequences,

240    retrotransposons (especially LTRs) are the biggest contributor for generating new

241    TSSs.

### 2.2.2 Extrinsic factors contribute to novel transcription initiation

243    Although previous studies and our analyses indicate that some sequence-intrinsic

244    features of repeats can promote transcription initiation, the majority of repeats

245    harboring such proto-TSS sequences do not exhibit initiation signals. For instance,

246    fewer than 1% of LTR elements in the human genome are associated with CAGE-

247    defined TSSs, implying that there are extrinsic factors that could affect the

248    transcription initiation in these regions. One reason for this is that most repeat

249    elements tend to be highly suppressed by the host defense mechanisms, such as DNA

250    methylation and methylation of H3 lysine 9[31]. In addition, we reasoned that proximity

251    of some proto-TSSs to established transcription units might be an extrinsic factor for

252    promoting novel transcription initiation, because such proximity could allow them to

253    access the transcription machinery of other TSSs for initiation. To test this hypothesis,

254    we first examined the *cis*-proximity of the LTR proto-TSSs to old TSSs. Indeed, we

255    found that young TSS-associated LTRs are closer to old TSSs compared to other

256    LTRs that are not associated with TSSs and random genomic intervals (**Fig. 2d**). We

257    further took advantage of published ChIA-PET data which identifies spatially

258    proximal regulatory regions in the genome. We focused on the ChIA-PET data for

259    CTCF and RAD21 (a subunit of cohesin), which are important for chromatin

260    architecture and linking regulatory modules for transcriptional regulation[32]. CTCF

261    binding sites were also found to be highly conserved during evolution[32]. We examined

262    the distances of LTRs to the mammalian-conserved ChIA-PET interaction loci (see

263    Methods) and found that TSS-associated LTRs are closer to CTCF or RAD21

264 interaction loci compared to non-TSS-associated LTRs (**Fig. 2e**). We suggest that

265 proximity to CTCF/cohesin anchoring loci may enable some proto-TSSs to be

266 spatially proximal to other transcription units and utilize their transcription machinery

267 for initiation.

268 The spatial proximity of young TSSs to old TSSs may also help to explain the

269 evolution of the number of TSSs per gene. We noticed that the number of TSSs per

270 gene in the human genome approximates to an exponential distribution – the number

271 of genes with a specific number of TSSs decreases exponentially with increase of the

272 number of TSSs per gene (**Fig. 2f**). The exponential relationship appears to be

273 independent of gene lengths, because the it still exists when looking at genes within a

274 specific length range (**Supplementary Fig. 8**). The exponential distribution indicates

275 that most genes have few TSSs, whereas a small fraction of genes have large number

276 of TSSs. A similar relationship is also seen for newly emerged TSSs (**Fig. 2g**), which

277 implies that a small fraction of genes gain many new TSSs during a specific period.

278 We also observed a positive correlation between number of pre-existing TSSs per

279 gene and number of newly gained TSSs per gene (Pearson's r=0.24, p < 2.2e-16,

280 **Supplementary Fig. 9**) - genes that have more existing TSSs are more likely to gain

281 new TSSs in a later period. Based upon the above observations, we suggest that most

282 of new TSSs derived from repeats arise opportunistically, partly due to their sequence-

283 intrinsic properties and proximity to other transcription units. As time goes by, some

284 newly emerged TSSs could be exapted by proximal genes to form alternative

285 promoters. On the other hand, these observations also suggest that the existing

286 transcriptional landscape to some extent constrains the emergence and evolution of

287 new TSSs.

288 **2.3 Rapid sequence evolution of young TSSs**

289 **2.3.1 Young TSSs undergo rapid sequence evolution**

290 Next we investigated the subsequent changes of young TSSs after they appear in the

291 genome. One important aspect is the evolutionary rate, which reflects the general

292 trend of sequence evolution. A previous study based on TSSs of early FANTOM

293 projects[14] showed that evolutionary rates in promoter regions vary between lineages

294 and that the primate lineages appear to have increased rates in promoter regions;

295 however evolutionarily young and old promoters were not separately analyzed. Here

296   we focused on the evolutionary rates for TSS groups of different ages in comparison
297   with the genomic background. To do this, we utilized genomic alignments to infer
298   evolutionary sequence changes around TSS loci for two recent periods (from the last
299   common ancestor of OWAs to the last common ancestor of hominids and from the last
300   common ancestor of hominids to present, as indicated by the phylogeny in **Fig. 3a**),
301   using a maximum likelihood method (see Methods). Based on inferred sequence
302   changes, we calculated the relative rates of substitutions and small insertions/deletions,
303   which were normalized by genomic average. We found that proximal positions of old
304   TSSs have lower substitution rates compared with surrounding regions and genomic
305   average (**Fig. 3a**), suggesting that they were subject to purifying selection in these
306   periods. In contrast, proximal positions of young TSSs exhibit elevated evolutionary
307   rates compared to the surrounding regions as well as genomic average (**Fig. 3a**),
308   suggesting that young TSS loci underwent rapid sequence evolution. Interestingly, for
309   the 'primate' TSS group the substitution rates during the early period are higher than
310   in the later period (**Fig. 3a**), suggesting that newly emerged TSSs evolve rapidly at
311   first and then slow down later. Although this pattern is not observed in the
312   insertion/deletion rates (**Supplementary Fig. 10**), it might be due to saturated
313   insertion/deletion mutations and some ancestral insertion/deletion events not being
314   accurately inferred using alignments of extant species. Additionally, by examining the
315   population polymorphism data from the 1000 genomes project, we found that the
316   young TSSs also have elevated variant densities relative to surrounding regions
317   (**Supplementary Fig. 11**), further supporting that young TSSs undergo rapid
318   sequence evolution.

319   **2.3.2 Endogenous mutational processes contribute to rapid evolution of young**
320   **TSSs**

321   We then asked how the young TSSs evolve rapidly after appearing in the genome.
322   Since many young TSSs are associated with repetitive sequences, we reasoned that
323   some mutational processes associated with repeats could contribute to the rapid
324   evolution.

325   One contributing factor could be DNA methylation, which is one of main mechanisms
326   for repressing TE activities[31]. We found that the younger TSSs have significantly
327   higher levels of CpG methylation in the germline compared to older TSSs (**Fig. 3b**

11

328 and **Supplementary Fig. 12)**. In addition, TE-associated TSSs tend to have higher

329 levels of CpG methylation compared to non-TE TSSs within each TSS group (**Fig.**

330 **3b**). Because methylated cytosine (mC) can frequently mutate to thymine (T) via

331 deamination, the DNA hypermethylation around young TSSs in the germline

332 represents an important contributor for the elevated evolutionary rates. This is further

333 supported by the substitution patterns in the human population genomic data, in which

334 the C > T is the most common substitution type (~40% of all substitutions) in all TSS

335 groups and ~17% of C to T mutations occur in the CpG context (**Fig. 3c**).

336 Another contributing factor is recombination, which has been found to be associated

337 with mutations and GC-biased gene conversion[33]. We found that LTR-associated TSSs

338 have significantly higher recombination rates relative to genomic average (**Fig. 3d**).

339 Higher recombination rates are also observed in non-TE-associated young TSSs (**Fig.**

340 **3d**). Consistently, older LTR-associated TSSs have more solitary LTRs (**Fig. 3e**),

341 which are known to result from allelic or non-allelic homologous recombination[34]. As

342 recombination hotspots evolve rapidly[35] and ancient recombination events are

343 difficult to detect, it is possible that recombination had also contributed to the rapid

344 evolution of SINE/LINE-associated TSSs.

345 A third contributing factor is the instability of tandem repeats. Previous research

346 revealed that the mutability of microsatellites (also known as short tandem repeats)

347 increases with their length and long microsatellites tend to be shortened or interrupted

348 by mutations over time[36,37]. Indeed, we found that tandem repeats associated with

349 younger TSSs tend to be shorter than those in older TSSs (**Fig. 3f**), implying that they

350 are more likely to mutate.

351 **2.3.3 Consequences of rapid evolution in young TSSs**

352 A direct consequence of the rapid evolution around young TSSs is that they

353 accumulated many changes, which could reduce or eliminate the transposition

354 capacity of TEs or the mutability of tandem repeats around TSSs, resulting in a more

355 stable genomic environment. Therefore these mutational processes probably help to

356 resolve the genomic conflicts caused by the inherent instability of associated repeats

357 around young TSSs. In addition, we suspect that rapid evolution may lead to deaths of

358 some young TSSs, because some sequence changes could disrupt critical promoter

359 components required for transcription initiation. In the example shown in **Fig. 3g**, a

360   LTR locus with transcription initiation signal in human has been deleted from rhesus

361   and baboon. However, because we lack large-scale CAGE-defined TSSs in other

362   primate species and there could be polymorphisms in TSS loci, we are currently

363   unable to perform detailed analysis regarding the evolutionary deaths of young TSSs.

364   **2.4 Functional impact of young TSSs**

365   **2.4.1 TSSs of different evolutionary ages exhibit distinct functional signatures**

366   Previous comparison between human and mouse CAGE-defined TSSs revealed that

367   lineage-specific TSSs tend to have tissue-restricted expression profiles, often in

368   samples associated with testis, immunity or brain[13]. Yet how the regulatory functions

369   of these lineage-specific TSSs are gradually established in organisms remain unclear.

370   We sought to investigate the resulting regulatory impact of newly emerged TSSs and

371   how their impact changes over time. We first took advantage of published functional

372   genomic data from ENCODE and other projects to compare related functional

373   signatures between TSS groups, including DNase I hypersensitivity (DHS), histone

374   modifications, DNA methylation, transcription factor (TF) binding and chromatin

375   interactions. Intriguingly, we found that TSSs of different ages exhibit segregating

376   functional signatures (**Fig. 4** for GM12878 cell line) and such patterns are observed in

377   different cell lines (**Supplementary Fig. 13** for K562 and H1-hESC cell lines).

378   Relative to older TSSs, younger TSSs tend to have lower chromatin accessibility

379   (DHS, **Fig. 4a**), lower levels of activating histone modifications (e.g. H3K4me3,

380   H3K27ac, H3K4me1, H3k9ac, **Fig. 4b** and **Supplementary Fig. 14**) and higher CpG

381   methylation (**Fig. 4c**), suggesting younger TSSs are under a more repressed chromatin

382   environment. By examining ChIP-seq data for TFs in ENCODE cell lines, we found

383   that older TSS loci tend to have more binding regions (i.e. more surrounding

384   sequences overlapping ChIP-seq peaks) relative to younger TSSs (**Fig. 4d,** and

385   **Supplementary Fig. 15** for meta-profiles of individual TF ChIP-seq datasets in

386   GM12878). We also observed a similar trend for computationally predicted TFBS

387   (**Supplementary Fig. 16**). We further analyzed the published ChIA-PET interaction

388   data for RNA polymerase II (RNAP II), which are usually formed within

389   CTCF/cohesin looped structures and considered to reflect promoter-enhancer

390   interactions[38]. We found that younger TSSs have fewer RNAP II chromatin

391   interactions compared with older TSS (**Fig. 4e**), suggesting that younger TSSs tend to

392     lack connections to other regulatory modules. This is consistent with the observations

393     in TF binding (**Fig. 4d**), as TF binding is important for forming promoter-enhancer

394     interactions. As for expression output, younger TSSs tend to display lower expression

395     than older TSSs (**Fig. 4f**), which is consistent with a previous observation that

396     evolutionarily volatile promoters tend to have lower expression levels[13]. Taken

397     together, these observations indicate that the evolution of TSSs leave footprints in the

398     functional signatures of TSSs; namely that younger TSSs tend to have smaller

399     regulatory impact on a genome and that the impact increases with time.

400     By comparing the TSS subgroups defined by the transcript types, we also observed

401     heterogeneity of functional signatures within TSS groups. Within a similarly-aged

402     group, TSSs associated with mRNAs tend to have higher DHS, more activating

403     histone modifications, more TF binding and more chromatin interactions than other

404     TSSs (**Fig. 4h-m** and **Supplementary Fig. 17**), indicating they are more

405     transcriptionally active. Consistently, mRNA TSSs tend to have higher expression

406     levels than other TSSs within the same group (**Fig. 4n**). Furthermore, TSSs of

407     proximal lncRNAs appear to be more transcriptionally active compared to that of

408     intergenic lncRNAs, likely because they are more proximal to other transcription units.

409     Overall, these findings suggest that locations of young TSSs in gene annotation

410     context could influence the regulatory outcomes.

411     **2.4.2 Evolution of regulatory functions of young TSSs appears to be subject to**

412     **temporal and spatial constraints**

413     The segregating functional signatures of TSSs of different ages strongly imply that the

414     regulatory outcomes of young TSSs are gradually changed over time. Yet it remains

415     unclear how regulatory changes of young TSSs take place in organisms during

416     evolution, e.g. in what tempo and mode. The regulatory impacts of historical and

417     fixed sequence changes around TSSs are difficult to assess, however, there are many

418     ongoing changes around TSSs within human populations, whose regulatory effects

419     have been widely studied by combining functional and population genomic

420     approaches[39]. Two common strategies are to identify regulatory quantitative trait loci

421     (rQTLs, e.g. TF binding QTLs, histone modification QTLs and DHS QTLs) and

422     variants associated with regulatory allelic specificities (AS, e.g. allele-specific TF

423     binding, allele-specific methylation). Although no QTL or AS study has been

14

424    specifically performed for human CAGE-defined TSSs, we can apply data from

425    genome-wide rQTL and AS studies of other molecular traits. A previous study[40]

426    revealed that expression levels of CAGE-defined TSSs are highly correlated with

427    other functional signatures such as TF binding, histone modifications and DHS in

428    surrounding regions, and can be largely predicted by those functional signatures ($R^2 >$

429    0.7). Therefore we reasoned that changes in the regulatory outcomes of TSSs can be

430    approximated by changes in related functional signatures in surrounding regions. By

431    examining rQTLs and AS variants (together called regulatory variants) in TSS loci of

432    different ages, we can gain insights into the tempo and mode of regulatory evolution

433    of TSSs at different life stages.

434    In our analysis we focused only on the *cis*-regulatory variants around TSS loci, as

435    published *trans*-regulatory variants are rare and of relatively low-quality. Previous

436    expression QTL studies found that the density of *cis*-regulatory variants drops rapidly

437    with increased distances to target TSSs[41], we restricted our analysis to only regulatory

438    variants within ±1 kb of TSSs. By re-analyzing data from multiple independent

439    studies, including DHS, methylation, histone marks and TF binding, we found that

440    younger TSSs tend to have fewer regulatory variants compared with older TSSs (see

441    **Fig. 5a-d** for four representative datasets and **Supplementary Fig. 18** for more

442    datasets). The trend is especially clear for variants associated with DHS, methylation

443    and TF binding. This is interesting because it suggests that although young TSS loci

444    evolve rapidly, many of the sequence changes appear to have none or limited impact

445    on transcriptional regulation. Since some TSSs are closely spaced, regulatory variants

446    could be counted multiple times in the above analysis (though it may be possible for a

447    variant to affect multiple adjacent TSSs). We still observed similar patterns even after

448    excluding all the TSSs separated by less than 2 kb (**Supplementary Fig. 19**).

449    Moreover, similar trends are observed when only including regulatory variants with

450    high derived allele frequencies (**Supplementary Fig. 20**), changes in which are more

451    likely to be fixed in populations in the future. Overall, these observations imply that

452    regulatory evolution of young TSSs is subject to a temporal constraint - younger TSSs

453    have a slower tempo in regulatory evolution (**Fig. 5e**), which might be due to the

454    strong repression in early periods.

455    Separating similarly aged TSSs according to transcript type, mRNA and proximal

456    lncRNA TSSs tend to have more regulatory variants compared with intergenic

457 lncRNA TSSs (**Fig. 5a-d**). Since mRNA and proximal lncRNA TSSs also have more

458 ChIA-PET interactions than other TSSs (**Fig. 4l**), we propose that there is a spatial

459 constraint on the regulatory evolution of young TSSs. Generally, younger TSSs have

460 less connectivity to other regulatory modules (i.e. spatially isolated) than older TSSs

461 (**Fig. 4e**), which likely limits their functional impact. In the subsequent evolution,

462 sequence changes in the young TSSs which are proximal to other regulatory modules

463 tend to have more regulatory effects and these TSSs may be incorporated in the

464 existing regulatory network more quickly (i.e. a higher tempo of regulatory evolution

465 in these TSSs). In contrast, relatively isolated TSSs tend to have a slower tempo of

466 regulatory evolution and are more difficult to be co-opted by the host.

467 Examples of evolving *cis*-proximal and *trans*-proximal young TSSs are shown in **Fig.**

468 **5f-g**. In the gene *RNFT2 locus* shown in **Fig. 5f**, an 'OWA' TSS, which lies on the

469 antisense strand of a newly inserted L1 element, is *cis*-proximal to an upstream old

470 TSS. In the surrounding regions of the 'OWA' TSS, there are multiple polymorphic

471 sites in current populations, two of which are regulatory variants affecting PU.1

472 binding and H3K4me3 respectively (**Fig. 5f**). In the example shown in **Fig. 5g**, a

473 'primate' TSS within an LTR element is ~70 kb away from *TAGAP* locus. However,

474 this young TSS is *trans*-proximal to the TSSs of *TAGAP*, as supported by several

475 CTCF and RNAPII ChIA-PET interaction pairs (**Fig. 5g**). This LTR is a solitary LTR

476 and thus lack capacity for retrotransposition. Six regulatory variants are within ±1 kb

477 of the young TSS (**Fig. 5g**). More examples are given in **Supplementary Fig. 21**.

478 **3 Discussion**

479 Given the large number of identified TSSs in the mammalian genomes and the high

480 TSS turnover rate, it is important to understand where the new TSSs come from, how

481 they evolve over time, and their functional impact on transcripts. By performing

482 evolutionary and functional analyses, we gain several important insights into the

483 evolution of newly emerged TSSs. We summarize our main findings in an integrative

484 model as shown in **Fig. 6**.

485 First, our analyses revealed several sequence-intrinsic and extrinsic factors that

486 promote the emergence of new TSSs (**Fig. 6**). Intrinsic factors are mainly associated

487 with the expansion of repetitive sequences, among which retrotransposons represent a

488 major source of new TSSs. In addition to sequence-intrinsic properties, chromatin

16

489   organization and spatial chromosomal interactions are likely important extrinsic
490   factors. New TSSs are usually proximal in *cis* or *trans* to other established
491   transcriptional units providing easier access to the transcriptional machinery, whereas
492   unexpressed proto-TSSs are more isolated. This dependence on extrinsic chromatin
493   environment partly explains why only a small fraction of proto-TSSs have detectable
494   initiation signals.

495   Secondly, resolving genomic conflicts is likely the main theme in the early period of
496   young TSSs (**Fig. 6**). Our evolutionary rate analysis revealed that young TSSs
497   experienced rapid sequence evolution in early periods, which appear to be associated
498   with several endogenous mutational processes, including DNA methylation,
499   recombination and tandem repeat mutagenesis. We suggest that such rapid evolution
500   can reduce the genomic conflicts caused by the instability of repetitive sequences
501   associated with young TSSs, as the TSS loci became more stable after they mutated.
502   We suspect that a considerable fraction of new TSSs may die during the rapid
503   evolution in early periods, as sequence changes could disrupt critical promoter
504   components required for transcription initiation.

505   Thirdly, by analyzing functional genomic data, we found that in early periods young
506   TSSs tend to have limited transcriptional competency, likely due to the highly
507   repressive environment and lack of connectivity to other functional modules.
508   However, their regulatory potential appear to be gradually enhanced over time (**Fig.
509   6**). Interestingly, by examining regulatory variants around TSS loci, we revealed that
510   the evolution of regulatory functions of young TSSs appears to be subject to temporal
511   and spatial constraints. The temporal constraint - that younger TSSs have fewer
512   regulatory variants within a period (slower tempo) despite faster sequence evolution -
513   is probably due to the genomic conflicts caused by the novel transcription and
514   associated unstable repetitive sequences. Young TSSs tend to be strongly repressed at
515   first and require time to resolve the genomic conflicts caused by associated repeats.
516   The spatial constraint – that TSSs with fewer chromosomal contact display a slower
517   tempo of regulatory evolution - likely limits the regulatory impact of young TSSs in
518   early stages and affects the evolutionary trajectories of young TSSs depending on
519   their genomic context. Based upon these observations and proposed constraints, we
520   speculate that younger and (or) more isolated TSSs are more likely to die out during
521   evolution.

522    Many studies have reported the contribution of repetitive sequences to regulatory
523    innovation[34]. Our detailed analysis on evolutionary trajectories of young human TSSs
524    provide new strong evidence. We have shown that the repeat-derived TSSs are tightly
525    constrained in the beginning and have limited functional impact, but after resolving
526    genomic conflicts some are successfully incorporated into the existing regulatory
527    network, turning "conflicts" into "benefits"[34]. In the long run, the repeat-derived TSSs
528    contribute significantly to regulatory innovation. Interestingly, a similar evolutionary
529    pattern was also observed in Alu exonization in primate genomes[42], implying a
530    commonly used strategy in genome evolution. Given the pervasiveness of repetitive
531    sequences and the similarity of chromatin structures in eukaryotic genomes, the
532    observed evolutionary processes involved in newly emerged TSSs in primate
533    genomes could also exist in other eukaryotic groups. These evolutionary patterns also
534    suggest the importance of balancing evolvability and robustness in genome
535    evolution[43].

536

537 **Methods**

538 **Human TSS annotation dataset**

539 We used the FANTOM 5 TSS dataset because it is the most comprehensive TSS

540 annotation to date, cataloguing/encompassing the genome-wide TSS profiling of most

541 major primary cell types and tissues in human. The high-confidence, "robust" TSSs

542 from the latest FANTOM CAT annotation (http://fantom.gsc.riken.jp/cat/, part of

543 FANTOM 5)[28] were used for our analyses, particularly as each TSS has been assigned

544 a RNA-seq-defined transcript. Coding status and transcript classification of transcripts

545 were defined as in the FANTOM CAT. To facilitate analysis and interpretation, we

546 merged three lncRNA classes ("lncRNA_antisense", "lncRNA_divergent" and

547 "lncRNA_sense_intronic") in the FANTOM CAT annotation into a class called

548 "proximal lncRNA", because these lncRNAs are proximal to other transcript units.

549 We also merged several minor classes ("sense_overlap_RNA", "short_ncRNA",

550 "small_RNA", "structural_RNA" and "uncertain_coding") into a class called "other

551 RNA". For TSSs which are associated with multiple types of transcripts, we assigned

552 them hierarchically to the five categories: mRNA > proximal_lncRNA >

553 intergenic_lncRNA > pseudogene > other_RNA. As CAGE TSS peaks (i.e. tag

554 clusters) usually span more than 1 bp, unless specified otherwise, we used the

555 dominant TSS position (i.e. the most frequently used initiation site) of each TSS peak

556 provided in the FANTOM annotation for most analyses.

557 **Categorization of human TSSs by sequence age**

558 To categorize human TSSs by the evolutionary age of the sequence, we made use of

559 whole genome alignments between human (hg19) and 16 other mammalian genomes

560 (**Supplementary Table 1**) from UCSC genome browser[44]. To estimate the sequence

561 ages of human TSS loci, the UCSC liftOver tool was used to determine presence or

562 absence of each human TSS sequence in other non-human genomes based on

563 available pairwise chain alignment files from UCSC. We required a minimum

564 mapping ratio of 90% for CAGE TSS peaks (~23bp in length on average), which

565 usually covers Initiator (Inr) elements of promoters. The sequence proximal to Inr

566 element has previously been found to be conserved in mammalian promoters[14]. In

567 addition, we required a minimum mapping ratio of 50% for TSS peaks±100 bp, which

568 we considered as "core promoter" regions in our study and are usually under high

19

569 selective constraint[14], although there is no standard definition for "core promoter"
570 currently. To reduce potential false positives resulting from alignments of paralogous
571 loci in two genomes, we further required a minimum alignment chain size of 10 kb for
572 both target and query genomes. A human TSS locus satisfying the above criteria for
573 the pairwise alignment was considered as having the orthologous sequence in the
574 surveyed genome, and its sequence age should be equal to or larger than the age of
575 last common ancestor of two species. The presence/absence patterns of TSSs were
576 then used for defining the four TSS groups as described in the main text. We also tried
577 multiple sets of thresholds for liftOver which did not result in notable variation in the
578 grouping results (**Supplementary Table 2**), mainly because many newly emerged
579 TSS loci were associated with TE insertions, which usually span more than 200 bp.

580 As some genomic regions are highly repetitive and could lead to poor assemblies and
581 erroneous alignments, we filtered out any TSS whose ±1 kb regions overlapping the
582 blacklisted genomic regions (see **Supplementary Table 3**) defined in the ENCODE
583 project and two other studies[45,46]. Because CAGE reads are usually short (20~70bp)[27]
584 and can be mapped to the genome multiple times, we made use of the Duke 20-bp
585 uniqueness track from UCSC browser to filter out the TSS peaks that have an average
586 uniqueness score of <0.5 (a 20-bp uniqueness score of <0.5 means that a 20-mer can
587 be mapped to the human genome more than twice). After excluding these blacklist
588 regions, we still observed that some TSS loci, which are usually associated with low-
589 complexity tandem repeats, exhibited suspiciously high read depths in some
590 functional datasets, suggesting they might be artifacts due to poor mappability for
591 short reads in those regions. Therefore we further filtered out any TSS harboring more
592 than 10% (200 bp) of tandem repeats in the 2 kb region centered on the TSS. In
593 addition, TSSs of chrM and chrY were excluded from all analyses because some
594 genome assemblies or functional datasets lack data for these genomic sequences.

595 When analyzing the remaining TSSs, we further found two significant sources of
596 putative false positives. One is the pseudogene-associated TSSs. Pseudogenes
597 (especially processed pseudogenes) were reported as a notable source of false
598 positives for CAGE-defined TSSs because of their high sequence similarity to
599 original gene loci and the short lengths of CAGE reads[47]. For the GM12878 cell line,
600 only 3.7% of the pseudogene TSSs in primate lineages from FANTOM 5 can be found
601 in the previously published GRO-cap-defined TSSs (**Supplementary Fig. 6**)[5].

20

602 Therefore we excluded all pseudogene TSSs from downstream analyses. Another
603 source of false positives is the TSSs associated with poly(A) or poly(T) tracts. We
604 initially found many young TSSs in FANTOM 5 located around the 3' poly(A) region
605 and the A-rich linker region of Alu elements. However, in the GM12878 cell line,
606 only 5.2% of the poly(dA:dT)-associated TSSs in primate lineages from FANTOM 5
607 can be found in the GRO-cap-defined TSSs (**Supplementary Fig. 6**). On the other
608 hand, a much larger fraction (43%) of the TSSs that are not associated with
609 pseudogenes and poly(dA:dT) tracts can be found in the GRO-cap-defined TSS
610 dataset. Such a large difference in the overlapping ratio suggests that the TSSs
611 associated with poly(dA:dT) tracts have a high fraction of false positives. A recent
612 study also suggested that Alu sequences generally lack the capacity to drive
613 autonomous transcription[20]. Therefore we filtered out the TSSs flanked by a tandem
614 repeat with A content of >50 % or T content of >50 % within ±100 bp.

615 **Analysis of TATA-box and CpG islands (CGI)**

616 The data of CGI annotation in the human genome was from Cohen et al. (2011)[48]. A
617 TSS was considered as CGI-associated if its core promoter region (TSS±100 bp)
618 overlaps a CGI. TATA-box hits were predicted by R package "seqPattern" using the
619 TBP position-weighted matrix with a minimum score of 80%. A TSS was considered
620 as TATA-box-associated if the start of a TATA-box motif is located at 25~35 bp
621 upstream of the TSS.

622 **Analysis of repeats associated with TSSs**

623 The annotation of transposable elements in our analysis was based on RepeatMasker
624 annotation of the hg19 assembly, downloaded from http://www.repeatmasker.org
625 (Repeat Library 20140131)[24]. In addition, as young TSS loci are frequently associated
626 with tandem repeats, tandem repeats annotated by TRF (downloaded from UCSC) and
627 STRcat[26] were also used. The "Simple repeat", "Low complexity" and "Satellite"
628 families in RepeatMasker were considered as tandem repeats in our analysis. The
629 tandem repeats from RepeatMasker, TRF and STRcat were merged into a union
630 dataset. For overlapping tandem repeats in these three datasets, the priority order for
631 being included in the union dataset was STRcat > TRF > RepeatMasker.

632   To investigate the repeat content around TSS loci, we first identified the nearest repeat

633   element to each TSS and counted how many TSSs harbored repeat elements within

634   TSS±100 bp regions (i.e. core promoter regions in this study). Since retrotransposons

635   and tandem repeats were the main types of TSS-associated repeats and many tandem

636   repeats were derived from retrotransposons, for each TSS group defined by sequence

637   age, we further defined four TSS subgroups ('SINE-associated', 'LINE-associated',

638   'LTR-associated' and 'Others') based on the nearest retrotransposon within 100 bp of

639   the TSS. The statistics of subgroups defined by transcript types and associated

640   retrotransposons are given in **Supplementary Table 4.**

641   To analyze the distributions of TSSs along repeat elements, we calculated the relative

642   distances of TSSs to the 5' (corresponding to 0% of the full-length) of corresponding

643   repeat subfamily consensus sequences based on the alignment information provided in

644   RepeatMasker annotation. When investigating distances of young TSSs to ChIA-PET

645   interaction loci of CTCF or RAD21, we only considered the interaction pairs whose

646   sequences could be found in at least one of the six non-primate mammalian genomes

647   listed in **Supplementary Table 1**, based on the liftOver mapping with parameters "-

648   minMatch=0.5 -minChainT=10000 -minChainQ=10000". The chromatin interactions

649   in these mammalian-conserved loci are likely established before emergence of

650   primates and conserved among mammals.

651   **Evolutionary rate analysis**

652   To investigate the evolutionary rates around TSS loci, we extracted alignments of

653   human and 14 other mammalian genomes for all TSS and their surrounding 2 kb

654   regions from the 100-way MULTIZ genome alignments from UCSC (all species used

655   for analysis are listed in **Supplementary Table 1**; tarSyr1 and micMur1 were not in

656   the 100-way alignments and thus not included in this analysis). To improve the

657   alignment quality, the extracted MULTIZ alignments were re-aligned using PRANK

658   with parameter "+F", which was found to generate more accurate gapped alignments

659   for evolutionary analysis [49]. The re-alignment results were then used to infer ancestral

660   sequences for each TSS locus using FASTML[50] with parameters "--SubMatrix HKY -

661   jointReconstruction no --indelReconstruction ML". FASTML produced posterior

662   probabilities for each position of inferred ancestral sequences. Positions with low-

663   confidence inferred sequences (maximum marginal probability of <0.8) were

664 excluded for subsequent analyses. Evolutionary sequence changes (substitutions,
665 insertions and deletions) in TSS loci in different periods were identified by comparing
666 inferred ancestral sequences and derived sequences, and these changes were used to
667 calculate substitution, insertion, and deletion rates for each period respectively. To
668 estimate the genomic average evolutionary rates, we generated 10,000 random 2-kb
669 intervals from the human genome, and ran the same analysis pipeline as described
670 above for the TSS loci. The relative rates of substitutions, insertions and deletions in
671 TSS loci were then obtained by dividing the original rates by genomic average rates
672 estimated from random intervals.

**Analysis of mutational mechanisms**

674 Because spontaneous deamination of methylated cytosines (causing cytosine to
675 thymine substitutions) was found to be a major source of mutations during evolution,
676 we analyzed the germline DNA methylation levels to investigate the impact of
677 methylation on the evolutionary rates of different TSS groups. We used the published
678 germline DNA methylation data from Guo et al. (2015)[51] and focused on the CpG
679 methylation events. The methylome of male primordial germ cells of 7-weeks old
680 embryos was used in our analysis, because this sample exhibited a high degree of
681 methylation across the genome, as shown in that study. Data of recombination rates in
682 human populations was from the HapMap project[52]. The completeness status (solitary
683 or non-solitary) of LTRs was predicted by REannotate[53] with parameters "-n -c",
684 using the RepeatMasker annotation as input.

**Analysis of functional signatures of TSSs**

686 Processed data (files of normalized signals and called peaks) of Dnase I-seq, ChIP-seq
687 and DNA methylation (WGBS) of ENCODE cell lines (GM12878, K562 and H1-
688 hESC) were downloaded from ENCODE website and ENSEMBL database. Analysis
689 and visualization of functional genomics data on the TSS groups and subgroups were
690 performed with BEDtools, R, seqplots[54] and deeptools[55].

691 ChIA-PET data for CTCF and RNAPII in GM12878 were from Tang et al. (2015)[38].
692 ChIA-PET data for RAD21 in GM12878 were from Grubert et al. (2015)[56]. ChIA-
693 PET data for RNAPII in K562 cell line were downloaded from the ENCODE website.

**Regulatory variant analysis**

23

695    The ongoing genomic changes (polymorphic sites) affecting the regulatory outcomes

696    of TSSs in human populations can be considered as a snapshot of regulatory evolution

697    of TSSs. Investigation of these regulatory variants would help to understand how the

698    regulatory impact of TSSs changes over time. Therefore, we analyzed published

699    regulatory variants which affect transcription-related molecular traits, such as TF

700    binding, histone marks, DNA methylation and DNase I hypersensitivity (DHS) from

701    several genome-wide studies.

702    Regulatory variants for allele-specific DHS in multiple cell types were from Maurano

703    et al. (2015)[57]. Regulatory variants for allele-specific CpG methylation in multiple

704    cell types were from Schultz et al. (2015)[58]. Three types of histone mark QTLs

705    (H3K4me3, H3K4me1 and H3K27ac) of lymphoblastoid cell lines (LCLs) were from

706    Grubert et al. (2015)[56]. For the data from Grubert et al. (2015), we only used the

707    regulatory variants that are located within the corresponding regulated histone peak

708    regions for analysis. Binding QTLs of 5 TFs (JunD, NF- kb, Pou2f1, PU.1 and Stat1)

709    and H3K4me3 QTLs in LCLs were from Tehranchi et al. (2016)[59]. The derived allele

710    frequencies (DAFs) of variants were based on the data of 1000 genomes project phase

711    3 release and only variants with known ancestral alleles were used for analysis. For

712    each type of regulatory variant, we calculated the proportion of TSSs harboring at

713    least one regulatory variant within 1 kb of the TSS. To account for the issue of

714    possible duplicated counts of adjacent TSS loci, we repeated the analysis after

715    excluding all the TSSs separated by less than 2 kb for to prevent duplicated counts.

716    We also repeated the analysis for datasets under three different minimum DAFs (0.01,

717    0.1 and 0.5).

718

719    **Data availability**

720    All the analyses in this study were based on published datasets. A table of data source

721    links is given in **Supplementary Table 5**. A table containing the defined TSS

722    groups/subgroups in this study is provided in **Supplementary Table 6.** All other data

723    are available from the authors upon reasonable request.

724

725

## Acknowledgments

## Author Contributions

C.L. conceived the project, with considerable discussion with N.M.L. C.L. performed the analyses and drafted the manuscript; N.M.L. supervised the project and contributed extensively to the writing and revising of the manuscript. B. L. provided important advice and contributed to the writing and revising of the manuscript.

## Competing financial interests

The authors declare no competing financial interests.

## References

1       Clark, M. B. *et al.* The reality of pervasive transcription. *PLoS biology* **9**, e1000625; discussion e1001102, doi:10.1371/journal.pbio.1000625 (2011).
2       Wade, J. T. & Grainger, D. C. Pervasive transcription: illuminating the dark matter of bacterial transcriptomes. *Nature reviews. Microbiology* **12**, 647-653, doi:10.1038/nrmicro3316 (2014).
3       Djebali, S. *et al.* Landscape of transcription in human cells. *Nature* **489**, 101-108, doi:10.1038/nature11233 (2012).
4       FANTOM Consortium *et al.* A promoter-level mammalian expression atlas. *Nature* **507**, 462-470, doi:10.1038/nature13182 (2014).
5       Core, L. J. *et al.* Analysis of nascent RNA identifies a unified architecture of initiation regions at mammalian promoters and enhancers. *Nature genetics* **46**, 1311-1320, doi:10.1038/ng.3142 (2014).
6       Core, L. J., Waterfall, J. J. & Lis, J. T. Nascent RNA sequencing reveals widespread pausing and divergent initiation at human promoters. *Science* **322**, 1845-1848, doi:10.1126/science.1162228 (2008).

760    7    Kim, T. K. & Shiekhattar, R. Architectural and Functional Commonalities
761         between Enhancers and Promoters. *Cell* **162**, 948-959,
762         doi:10.1016/j.cell.2015.08.008 (2015).
763    8    Li, W., Notani, D. & Rosenfeld, M. G. Enhancers as non-coding RNA
764         transcription units: recent insights and future perspectives. *Nature reviews.*
765         *Genetics* **17**, 207-223, doi:10.1038/nrg.2016.4 (2016).
766    9    Andersson, R. *et al.* An atlas of active enhancers across human cell types and
767         tissues. *Nature* **507**, 455-461, doi:10.1038/nature12787 (2014).
768    10   Frith, M. C. *et al.* Evolutionary turnover of mammalian transcription start sites.
769         *Genome research* **16**, 713-722, doi:10.1101/gr.5031006 (2006).
770    11   Main, B. J., Smith, A. D., Jang, H. & Nuzhdin, S. V. Transcription start site
771         evolution in Drosophila. *Molecular biology and evolution* **30**, 1966-1974,
772         doi:10.1093/molbev/mst085 (2013).
773    12   Yokoyama, K. D., Thorne, J. L. & Wray, G. A. Coordinated genome-wide
774         modifications within proximal promoter cis-regulatory elements during
775         vertebrate evolution. *Genome biology and evolution* **3**, 66-74,
776         doi:10.1093/gbe/evq078 (2011).
777    13   Young, R. S. *et al.* The frequent evolutionary birth and death of functional
778         promoters in mouse and human. *Genome research* **25**, 1546-1557,
779         doi:10.1101/gr.190546.115 (2015).
780    14   Taylor, M. S. *et al.* Heterotachy in mammalian promoter evolution. *PLoS*
781         *genetics* **2**, e30, doi:10.1371/journal.pgen.0020030 (2006).
782    15   Ward, L. D. & Kellis, M. Evidence of abundant purifying selection in humans
783         for recently acquired regulatory functions. *Science* **337**, 1675-1678,
784         doi:10.1126/science.1225057 (2012).
785    16   Scala, G., Affinito, O., Miele, G., Monticelli, A. & Cocozza, S. Evidence for
786         evolutionary and nonevolutionary forces shaping the distribution of human
787         genetic variants near transcription start sites. *PloS one* **9**, e114432,
788         doi:10.1371/journal.pone.0114432 (2014).
789    17   Schor, I. E. *et al.* Promoter shape varies across populations and affects
790         promoter evolution and expression noise. *Nature genetics* **49**, 550-558,
791         doi:10.1038/ng.3791 (2017).
792    18   Faulkner, G. J. *et al.* The regulated retrotransposon transcriptome of
793         mammalian cells. *Nature genetics* **41**, 563-571, doi:10.1038/ng.368 (2009).
794    19   Nguyen, T. A. *et al.* High-throughput functional comparison of promoter and
795         enhancer activities. *Genome research* **26**, 1023-1033,
796         doi:10.1101/gr.204834.116 (2016).
797    20   van Arensbergen, J. *et al.* Genome-wide mapping of autonomous promoter
798         activity in human cells. *Nature biotechnology* **35**, 145-153,
799         doi:10.1038/nbt.3754 (2017).
800    21   Lenhard, B., Sandelin, A. & Carninci, P. Metazoan promoters: emerging
801         characteristics and insights into transcriptional regulation. *Nature reviews.*
802         *Genetics* **13**, 233-245, doi:10.1038/nrg3163 (2012).
803    22   Sawaya, S. *et al.* Microsatellite tandem repeats are abundant in human
804         promoters and are associated with regulatory elements. *PloS one* **8**, e54710,
805         doi:10.1371/journal.pone.0054710 (2013).
806    23   Bilgin Sonay, T. *et al.* Tandem repeat variation in human and great ape
807         populations and its impact on gene expression divergence. *Genome research*
808         **25**, 1591-1599, doi:10.1101/gr.190868.115 (2015).
809    24   Tarailo-Graovac, M. & Chen, N. Using RepeatMasker to identify repetitive
810         elements in genomic sequences. *Current protocols in bioinformatics* **Chapter**
811         **4**, Unit 4 10, doi:10.1002/0471250953.bi0410s25 (2009).

812  25  Benson, G. Tandem repeats finder: a program to analyze DNA sequences.
813      *Nucleic acids research* **27**, 573-580 (1999).
814  26  Willems, T. *et al.* The landscape of human STR variation. *Genome research*
815      **24**, 1894-1904, doi:10.1101/gr.177774.114 (2014).
816  27  Kanamori-Katayama, M. *et al.* Unamplified cap analysis of gene expression on
817      a single-molecule sequencer. *Genome research* **21**, 1150-1159,
818      doi:10.1101/gr.115469.110 (2011).
819  28  Hon, C. C. *et al.* An atlas of human long non-coding RNAs with accurate 5'
820      ends. *Nature* **543**, 199-204, doi:10.1038/nature21374 (2017).
821  29  Ohadi, M. *et al.* Core promoter short tandem repeats as evolutionary switch
822      codes for primate speciation. *American journal of primatology* **77**, 34-43,
823      doi:10.1002/ajp.22308 (2015).
824  30  Gymrek, M. *et al.* Abundant contribution of short tandem repeats to gene
825      expression variation in humans. *Nature genetics* **48**, 22-29,
826      doi:10.1038/ng.3461 (2016).
827  31  Slotkin, R. K. & Martienssen, R. Transposable elements and the epigenetic
828      regulation of the genome. *Nature reviews. Genetics* **8**, 272-285,
829      doi:10.1038/nrg2072 (2007).
830  32  Merkenschlager, M. & Odom, D. T. CTCF and cohesin: linking gene regulatory
831      elements with their targets. *Cell* **152**, 1285-1297,
832      doi:10.1016/j.cell.2013.02.029 (2013).
833  33  Pratto, F. *et al.* Recombination initiation maps of individual human genomes.
834      *Science* **346**, 1256442, doi:10.1126/science.1256442 (2014).
835  34  Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of
836      transposable elements: from conflicts to benefits. *Nature reviews. Genetics*
837      **18**, 71-86, doi:10.1038/nrg.2016.139 (2017).
838  35  Baudat, F., Imai, Y. & de Massy, B. Meiotic recombination in mammals:
839      localization and regulation. *Nature reviews. Genetics* **14**, 794-806,
840      doi:10.1038/nrg3573 (2013).
841  36  Eckert, K. A. & Hile, S. E. Every microsatellite is different: Intrinsic DNA
842      features dictate mutagenesis of common microsatellites present in the human
843      genome. *Molecular carcinogenesis* **48**, 379-388, doi:10.1002/mc.20499
844      (2009).
845  37  Kelkar, Y. D., Tyekucheva, S., Chiaromonte, F. & Makova, K. D. The genome-
846      wide determinants of human and chimpanzee microsatellite evolution.
847      *Genome research* **18**, 30-38, doi:10.1101/gr.7113408 (2008).
848  38  Tang, Z. *et al.* CTCF-Mediated Human 3D Genome Architecture Reveals
849      Chromatin Topology for Transcription. *Cell* **163**, 1611-1627,
850      doi:10.1016/j.cell.2015.11.024 (2015).
851  39  Albert, F. W. & Kruglyak, L. The role of regulatory variation in complex traits
852      and disease. *Nature reviews. Genetics* **16**, 197-212, doi:10.1038/nrg3891
853      (2015).
854  40  Cheng, C. *et al.* Understanding transcriptional regulation by integrative
855      analysis of transcription factor binding data. *Genome research* **22**, 1658-1667,
856      doi:10.1101/gr.136838.111 (2012).
857  41  GTEx Consortium. Human genomics. The Genotype-Tissue Expression (GTEx)
858      pilot analysis: multitissue gene regulation in humans. *Science* **348**, 648-660,
859      doi:10.1126/science.1262110 (2015).
860  42  Attig, J. *et al.* Splicing repression allows the gradual emergence of new Alu-
861      exons in primate evolution. *eLife* **5**, doi:10.7554/eLife.19545 (2016).
862  43  Wagner, A. *Robustness and Evolvability in Living Systems*. (Princeton
863      University Press, 2007).

864    44    Tyner, C. *et al.* The UCSC Genome Browser database: 2017 update. *Nucleic*
865          *acids research* **45**, D626-D634, doi:10.1093/nar/gkw1134 (2017).
866    45    Li, W. & Freudenberg, J. Characterizing regions in the human genome
867          unmappable by next-generation-sequencing at the read length of 1000 bases.
868          *Computational biology and chemistry* **53 Pt A**, 108-117,
869          doi:10.1016/j.compbiolchem.2014.08.015 (2014).
870    46    Pickrell, J. K., Gaffney, D. J., Gilad, Y. & Pritchard, J. K. False positive peaks
871          in ChIP-seq and other sequencing-based functional assays caused by
872          unannotated high copy number regions. *Bioinformatics* **27**, 2144-2146,
873          doi:10.1093/bioinformatics/btr354 (2011).
874    47    Zhao, X., Valen, E., Parker, B. J. & Sandelin, A. Systematic clustering of
875          transcription start site landscapes. *PloS one* **6**, e23409,
876          doi:10.1371/journal.pone.0023409 (2011).
877    48    Cohen, N. M., Kenigsberg, E. & Tanay, A. Primate CpG islands are maintained
878          by heterogeneous evolutionary regimes involving minimal selection. *Cell* **145**,
879          773-786, doi:10.1016/j.cell.2011.04.024 (2011).
880    49    Loytynoja, A. & Goldman, N. Phylogeny-aware gap placement prevents errors
881          in sequence alignment and evolutionary analysis. *Science* **320**, 1632-1635,
882          doi:10.1126/science.1158395 (2008).
883    50    Ashkenazy, H. *et al.* FastML: a web server for probabilistic reconstruction of
884          ancestral sequences. *Nucleic acids research* **40**, W580-584,
885          doi:10.1093/nar/gks498 (2012).
886    51    Guo, F. *et al.* The Transcriptome and DNA Methylome Landscapes of Human
887          Primordial Germ Cells. *Cell* **161**, 1437-1452, doi:10.1016/j.cell.2015.05.015
888          (2015).
889    52    International HapMap Consortium *et al.* A second generation human
890          haplotype map of over 3.1 million SNPs. *Nature* **449**, 851-861,
891          doi:10.1038/nature06258 (2007).
892    53    Pereira, V. Automated paleontology of repetitive DNA with REANNOTATE.
893          *BMC genomics* **9**, 614, doi:10.1186/1471-2164-9-614 (2008).
894    54    Stempor, P. & Ahringer, J. SeqPlots - Interactive software for exploratory
895          data analyses, pattern discovery and visualization in genomics. *Wellcome*
896          *open research* **1**, 14, doi:10.12688/wellcomeopenres.10004.1 (2016).
897    55    Ramirez, F., Dundar, F., Diehl, S., Gruning, B. A. & Manke, T. deepTools: a
898          flexible platform for exploring deep-sequencing data. *Nucleic acids research*
899          **42**, W187-191, doi:10.1093/nar/gku365 (2014).
900    56    Grubert, F. *et al.* Genetic Control of Chromatin States in Humans Involves
901          Local and Distal Chromosomal Interactions. *Cell* **162**, 1051-1065,
902          doi:10.1016/j.cell.2015.07.048 (2015).
903    57    Maurano, M. T. *et al.* Large-scale identification of sequence variants
904          influencing human transcription factor occupancy in vivo. *Nature genetics* **47**,
905          1393-1401, doi:10.1038/ng.3432 (2015).
906    58    Schultz, M. D. *et al.* Human body epigenome maps reveal noncanonical DNA
907          methylation variation. *Nature* **523**, 212-216, doi:10.1038/nature14465 (2015).
908    59    Tehranchi, A. K. *et al.* Pooled ChIP-Seq Links Variation in Transcription Factor
909          Binding to Complex Disease Risk. *Cell* **165**, 730-741,
910          doi:10.1016/j.cell.2016.03.041 (2016).
911
912

913

**Figure legends**

**Fig. 1 Classification of human TSSs by evolutionary age.** (**a**) Statistics of four TSS groups defined by sequence age using genomic alignments. At the bottom is the phylogeny with colors indicating the corresponding period of each TSS group. (**b**) An example gene locus shows two 'mammalian' TSSs (red shade) and one 'OWA' TSS (cyan shade). An LTR element overlapping the young TSS can be seen at the bottom of the alignment. CAGE tag counts and transcript isoforms shown at the top were from FANTOM CAT annotation (part of FANTOM 5). Genome alignments represented by grey blocks and lines were generated using UCSC genome browser (hg19). (**c**) Composition of transcription type in each TSS group. Transcript types are derived from FANTOM CAT annotation. (**d**) Violin and box plots for TSS peak widths of each TSS group. (**e**) Proportions of TATA-box containing and TATA-less TSSs. (**f**) Proportions of CGI-associated and non-CGI-associated TSSs. Statistical significances in panel **d** were calculated by one-tailed Wilcoxon rank sum tests; statistical significances in panels **e** and **f** by Fisher's exact tests; "**", $p < 0.01$; "***", $p < 0.001$.

**Fig. 2 Intrinsic and extrinsic factors contributing to the origin of new TSSs.** (**a**) Composition of major repeat families in four TSS groups. To obtain a non-redundant assignment, we considered the nearest repeat element within TSS±100 bp. (**b**) Distribution of young TSSs along the LTR/THE1B elements, with a bin size of 2% of its full-length consensus sequence. In the middle is the THE1B structure, which includes the original TSS, U3, R and U5 regions for the transposable element. (**c**) Distribution of young TSSs along the LINE/L1 elements, with a bin size of 2% of full-length consensus sequences. In the middle is the L1 structure, which indicates the sense and antisense L1 TSSs at 5'end. (**d**) Comparison of distances of TSS-associated and non-TSS-associated LTRs to the closest old TSSs. The distances of random intervals (generated by "bedtools shuffle" with TSS-associated LTRs as input) to the closest old TSSs are also provided for comparison. (**e**) Comparison of distances of TSS-associated and non-TSS-associated LTRs to the closest CTCF or RAD21 ChIA-PET peaks (GM12878). Random intervals used here is the same as that in panel **d**. (**f**) Exponential approximation for the number of genes with a certain number of TSSs

29

945 and number of TSSs per gene, based on data of all TSSs. $R^2$ is the coefficient of

946 determination for the linear regression in the figure. (**g**) Exponential approximation

947 for the number of genes and number of newly gained TSSs per gene, based on data of

948 newly emerged TSSs in three periods. $R^2$ is the coefficient of determination.

949 Statistical significances in panels **d** and **e** were calculated by one-tailed Wilcoxon

950 rank sum tests; "***", $p < 0.001$.

951 **Fig. 3 Rapid sequence evolution of young TSSs.** (**a**) Left, a phylogeny of genomes

952 used for evolutionary rate analysis, with arrows indicating the two evolutionary

953 periods considered for calculating rates. Right, relative substitution rates (normalized

954 by genomic average) inferred from genomic alignments for three TSS groups, using

955 40 bins along TSS±1 kb for calculating the average rate in each bin. Best-fit curves

956 were estimated by 'loess'. (**b**) Violin and box plots for germline DNA methylation

957 levels (a male germline dataset from Guo et al. 2015) for different TSS subgroups

958 defined by the retrotransposon context. For each TSS, the average methylation level

959 of CpGs was calculated for TSS±1 kb. (**c**) Frequencies of nucleotide substitution

960 types in different TSS groups, based on the variants and ancestral alleles from the

961 1000 genomes project. (**d**) Comparison of recombination rates among TSSs

962 associated with different types of transposable elements and genomic background

963 ('random'). The recombination rate of each TSS was defined as the average rate for

964 TSS±1 kb. Background recombination rates were generated for randomly selected 2-

965 kb windows in human genome. (**e**) The fraction of solitary LTRs in four TSS groups.

966 (**f**) Distribution of tandem repeat (TR) lengths in four TSS groups. (**g**) An example

967 plot depicting a possible TSS death event around an LTR. Statistical significances in

968 panels **b**, **d** and **f** were calculated by one-tailed Wilcoxon rank sum tests. "*", $p < 0.05$;

969 "**", $p < 0.01$; "***", $p < 0.001$; N.S., not significant.

970 **Fig. 4 Distinct functional signatures in different TSS groups.** (**a**) Meta-profiles of

971 DHS signals for four TSS groups using a 20bp bin size (same bin sizes for other

972 panels). (**b**) Meta-profiles of H3K4me3 signals. (**c**) Meta-profiles of CpG methylation

973 levels. (**d**) Meta-profiles of coverage ratio by TF ChIP-seq peaks. Previously called

974 peaks of 88 TF ChIP-seq datasets from ENCODE were merged together, and for each

975 bin of each TSS locus we calculated how much is covered by merged peaks. (**e**) Meta-

976 profiles of coverage ratio by RNAP II ChIA-PET peaks. (**f**) Meta-profiles of RNAP II

977 ChIP-seq signals. (**g**) Distribution of maximum expression levels of TSSs across

978 primary cell samples, based on the expression data of FANTOM CAT annotation. (**h-**
979 **n**) Produced using the same methods as for panels **a-g**, but specifically for the OWA
980 TSSs which were divided into subgroups of different transcript types. All functional
981 genomic data except the expression data are for the GM12878 cell line.

982 **Fig. 5 Temporal and spatial constraints on the regulatory evolution of young**
983 **TSSs.** (**a**) Top, proportion of TSSs harboring regulatory variants associated with
984 allele-specific DHS within TSS±1 kb for each TSS group; above the bars are the
985 numbers of TSSs with regulatory variants. Bottom, proportion of TSSs harboring
986 regulatory variants in different TSS subgroups, defined by transcript type. (**b**)
987 Proportion of TSSs harboring variants associated with allele-specific methylation
988 within TSS±1 kb. (**c**) Proportion of TSSs harboring H3K4me3 QTLs within TSS±1
989 kb. Data generated from lymphoblastoid cell lines (LCLs). (**d**) Proportion of TSSs
990 harboring NF-kb binding QTLs within TSS±1 kb. Data generated from LCLs. (**e**) A
991 schematic illustration depicting different possible evolutionary paths for young TSSs.
992 (**f**) A young TSS *cis*-proximal to old TSSs. Top, FANTOM CAT transcript models
993 (red for forward-strand, blue for reverse-strand); genome alignments and TE
994 annotations obtained from the UCSC genome browser. Bottom, enlarged region of an
995 'OWA' TSS inside a LINE element. Below the alignments are the common SNPs
996 (allele frequency ≥0.01) from the dbSNP database and SNPs associated with
997 regulatory variation within this region. (**g**) A young TSS *trans*-proximal to old TSSs.
998 Top, similar to panel **f** but with additional CTCF and RNAP II ChIA-PET interaction
999 data for GM12878 cell line. Bottom, enlarged region of the young TSS inside a LTR
1000 element. Below the alignments are the common SNPs (allele frequency ≥0.01) from
1001 dbSNP database and the SNPs associated with regulatory variation within this region.

1002 **Fig. 6. Proposed evolution model for young TSSs**. The origin of new TSSs is
1003 promoted by sequence-intrinsic and extrinsic factors. A typical intrinsic factor is the
1004 promoter element in newly inserted retrotransposons. An important extrinsic factor is
1005 the proximity to established regulatory modules as the proximity of a 'proto-TSS' to
1006 established regulatory elements provides easier access to transcription machinery.
1007 Newly emerged TSSs tend to be highly repressed and have limited regulatory capacity.
1008 In the early phase, young TSSs undergo rapid sequence evolution allow genomic
1009 conflicts associated with repeats to be resolved. Targeted mutational mechanisms
1010 enable this rapid evolution, including DNA hypermethylation (methylated C to T

1011 mutations), recombination and tandem repeat instability. The accumulated changes

1012 around young TSSs can reduce or eliminate the transpositional capacity of associated

1013 TEs and stabilize associated tandem repeats. They may also lead to deaths of some

1014 young TSSs. In the later phases, surviving TSSs gradually gain mutations in

1015 surrounding regions which could increase their regulatory capacity (e.g. TF binding,

1016 chromatin accessibility or transcription-associated histone modifications) and are

1017 exapted by the host for transcriptional regulation. At the mature phase, TSSs tend to

1018 have more permissive chromatin environments, enhanced spatial connectivity and

1019 higher expression.

1020

# Figure 1



**a**

Mammalian (n=141,117) 92.9%
Primate (n=6,668) 4.4%
OWA (n=3,318) 2.2%
Hominid (n=799) 0.5%

**b** *BAAT* locus

old TSS    young TSS

CAGE total counts

FANTOM transcripts

chr9:104,144,337-104,148,781

Human
Chimp
Gorilla
Orangutan
Gibbon
Rhesus
Baboon
Marmoset
Bushbaby
Mouse
Rat
Pig
Cow
Horse
Dog

LTR/MER11A

**c**
- mRNA
- proximal lncRNA
- intergenic lncRNA
- other RNA

% of TSSs

Mammalian  Primate  OWA  Hominid

**d**

*** ** p = 0.25

TSS peak width (bp)

Mammalian  Primate  OWA  Hominid

**e**
■ TATA-box  □ TATA-less

*** ** p = 0.63

% of TSSs

Mammalian  Primate  OWA  Hominid

**f**
■ CGI  □ non−CGI

*** *** p = 0.07

% of TSSs

Mammalian  Primate  OWA  Hominid

# Figure 2

# Figure 3

Figure 4

**a** DHS (GM12878)
**b** H3K4me3 (GM12878)
**c** Methylation (GM12878)
**d** TF ChIP-seq peaks (GM12878)
**e** RNAPII ChIA-PET peaks (GM12878)
**f** RNAPII ChIP-seq (GM12878)
**g** Expression level

All TSSs

Mammalian
Primate
OWA
Hominid

**h** **i** **j** **k** **l** **m** **n**

OWA TSSs (transcript type)

mRNA
proximal lncRNA
intergenic lncRNA

# Figure 5



**a** allele-specific DHS (Maurano et al. 2015) (multiple cell types)

**b** allele-specific methy. (Schultz et al. 2015) (multiple cell types)

**c** H3K4me3 QTL (Grubert et al. 2015) (LCLs)

**d** NF-kB binding QTL (Tehranchi et al. 2016) (LCLs)

**e** → observed evolutionary path  ⇢ other possible paths

Mammalian / Primate / OWA / Hominid

mRNA / proximal lncRNA / intergenic lncRNA

**f** *RNFT2* locus

**g** *TAGAP* locus

# Figure 6

## Evolutionary trajectories of young TSSs

Repression ('Conflicts')
Exaptation ('Benefits')

Time

**Birth**
**Origin of new TSSs**
Intrinsic factors
Extrinsic environment
Highly repressed

**Early periods**
**Resolving conflicts**
Accelerated evolution
Transposition capacity ↓
Sequence mutability ↓

**Later periods**
**Regulatory exaptation**
Regulatory mutations
Permissive chromatin
Enhanced connectivity

## Intrinsic factors

Promoter elements in repeats
(e.g. LTRs of ERVs)

LTR          LTR

OR

Death?

DNA methylation (mC > T)
Truncation by recombination

Regulatory mutations affecting
- TF binding
- chromatin accessibility
- histone marks, etc.

## Extrinsic environment

Proximal to established regulatory modules (in *cis* or in *trans*)

Enhancer
Gene
Pol II
Gene

Enhancer
Gene
Pol II
Gene

Enhancer
Gene
Pol II
Gene

### Legend

- Young TSS
- Repeat element
- Transcript
- Nucleosome
- DNA methylation
- Repressive histone marks
- Activating histone marks
- Mutations
- Regulatory mutations

- Repeat element
- Old TSS
- Young TSS
- Unexpressed proto-TSS
- CTCF
- Cohesin
- TFs
- // Mutations

**Supplementary Information**

**Supplementary Table 1 Species and genome assemblies used for estimating sequence ages of TSSs.**

| Species | Assembly version | Taxa | | | |
|---|---|---|---|---|---|
| Human | hg19 | Hominids | Old anthropoids | Primates | Mammals |
| Chimp | panTro4 | | | | |
| Gorilla | gorGor3 | | | | |
| Orangutan | ponAbe2 | | | | |
| Gibbon | nomLeu3 | | | | |
| Rhesus | rheMac3 | | World | | |
| Baboon | papHam1 | | | | |
| Marmoset | calJac3 | | | | |
| Tarsier | tarSyr1 | | | | |
| Mouse lemur | micMur1 | | | | |
| Bushbaby | otoGar1 | | | | |
| Mouse | mm10 | | | | |
| Rat | rn6 | | | | |
| Pig | susScr3 | | | | |
| Cow | bosTau7 | | | | |
| Horse | equCab2 | | | | |
| Dog | canFam3 | | | | |

**Supplementary Table 2 Statistics of grouping results with different sets of cutoffs for liftOver, after filtering the TSSs overlapping blacklist regions.**

| Min mapped % of TSS peaks | Min mapped % of TSS peak±100 bp | Min chain size | Mammalian | Primate | OWA | Hominid | Used in final analyses? |
|---|---|---|---|---|---|---|---|
| 0.9 | 0.5 | 10kb | 141,117 | 6,668 | 3,318 | 799 | Yes |
| 0.8 | 0.5 | 10kb | 142,782 | 5,531 | 2,902 | 687 | No |
| 0.5 | 0.5 | 10kb | 144,121 | 4,652 | 2,532 | 597 | No |
| 0.9 | 0.3 | 10kb | 141,288 | 6,559 | 3,264 | 791 | No |
| 0.9 | 0.7 | 10kb | 139,652 | 7,505 | 3,840 | 905 | No |
| 0.9 | 0.5 | 5kb | 141,328 | 6,716 | 3,109 | 749 | No |
| 0.9 | 0.5 | 20kb | 140,913 | 6,591 | 3,525 | 873 | No |

**Supplementary Table 3 Lists of blacklist genomic regions used for filtering TSSs.**

| File | Source |
|---|---|
| | |

| | |
|---|---|
| wgEncodeDukeMapabilityRegionsExcludable.bed | ENCODE |
| wgEncodeDacMapabilityConsensusExcludable.bed | ENCODE |
| seq.cov1.ONHG19.bed | Pickrell et al. 2011 |
| UM1K0M50BP.bed | Li and Freudenberg 2014 |

**Supplementary Table 4 Statistics of TSS subgroups defined by transcript types and the nearest retrotransposon elements.**

| Mammalian | | | Primate | | | Old World Anthropoid | | | Hominid | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| mRNA | SINE | 3427 | mRNA | SINE | 271 | mRNA | SINE | 100 | mRNA | SINE | 16 |
| | LINE | 3470 | | LINE | 299 | | LINE | 118 | | LINE | 43 |
| | LTR | 830 | | LTR | 433 | | LTR | 306 | | LTR | 61 |
| | Others | 75301 | | Others | 1569 | | Others | 555 | | Others | 145 |
| proximal lncRNA | SINE | 2019 | proximal lncRNA | SINE | 204 | proximal lncRNA | SINE | 89 | proximal lncRNA | SINE | 15 |
| | LINE | 2311 | | LINE | 266 | | LINE | 164 | | LINE | 34 |
| | LTR | 827 | | LTR | 465 | | LTR | 309 | | LTR | 67 |
| | Others | 35202 | | Others | 1071 | | Others | 426 | | Others | 87 |
| intergenic lncRNA | SINE | 966 | intergenic lncRNA | SINE | 84 | intergenic lncRNA | SINE | 43 | intergenic lncRNA | SINE | 8 |
| | LINE | 1232 | | LINE | 219 | | LINE | 173 | | LINE | 53 |
| | LTR | 1192 | | LTR | 799 | | LTR | 524 | | LTR | 146 |
| | Others | 9106 | | Others | 516 | | Others | 269 | | Others | 58 |
| other RNA | SINE | 324 | other RNA | SINE | 32 | other RNA | SINE | 11 | other RNA | SINE | 0 |
| | LINE | 368 | | LINE | 78 | | LINE | 39 | | LINE | 17 |
| | LTR | 147 | | LTR | 99 | | LTR | 67 | | LTR | 17 |
| | Others | 4395 | | Others | 263 | | Others | 125 | | Others | 32 |

**Supplementary Table 5 URL links of main published datasets used in this study.**

| | Download links |
|---|---|
| FANTOM TSSs | http://fantom.gsc.riken.jp/5/suppl/Hon_et_al_2016/data/ |
| liftOver chain files | http://hgdownload.cse.ucsc.edu/goldenPath/hg19/liftOver/ |
| RepeatMasker annotation | http://www.repeatmasker.org/genomes/hg19/RepeatMasker-rm405-db20140131/hg19.fa.out.gz |
| TRF | http://hgdownload.cse.ucsc.edu/goldenpath/hg19/database/simpleRepeat.txt.gz |
| STRcat | http://strcat.teamerlich.org/download |
| MULTIZ alignments | http://hgdownload.cse.ucsc.edu/goldenPath/hg19/multiz100way/maf/ |
| Germline methylation | https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63818 |
| Variants from | ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/ |

| 1000 genomes project | |
|---|---|
| ENCODE functional datasets | ftp://ftp.ebi.ac.uk/pub/databases/ensembl/encode/integration_data_jan2011/ |
| ChIA-PET data | https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE62742<br>https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE72816 |
| AS or QTL data | http://www.nature.com/ng/journal/v47/n12/extref/ng.3432-S5.txt<br>https://www.nature.com/nature/journal/v523/n7559/extref/nature14465-s2.zip<br>http://mitra.stanford.edu/kundaje/portal/chromovar3d/index.html<br>http://www.cell.com/cms/attachment/2062331538/2064077614/mmc2.xlsx |

**Supplementary Table 6 A table containing the defined TSS groups/subgroups used in analyses (in a separate file).**

**Supplementary Figure 1 Comparison of GC content and CpG content between four groups.**

**Supplementary Figure 2 Heatmap for repeat content, GC content, CpG content and TATA-box in four TSS groups.** Each TSS group is subdivided into subgroups based on transcript type. Within each subgroup, rows are sorted by the distance from the TSS to the nearest TE element (priority order: SINE > LINE > LTR > Others). For the TSSs in the 'Others' category, rows are sorted based on their distances to the nearest tandem repeat elements. The color gradients for repeat elements (SINE, LINE, LTR, tandem repeat (TR)) represent the repeat coverage in 10 bp bins. Regions shown in the TATA-box columns are TSS±200 bp. Because of the large number of TSSs in the 'mammalian' group (panel **d**), only the data of 5000 randomly selected TSSs are shown.

**Supplementary Figure 3 Distribution of young TSSs along LTR subfamilies.** These nine subfamilies are among the top 10 LTR subfamilies which harbor most young TSSs. The tenth, THEIB, has already been shown in **Fig. 2b**. Number of young TSSs for each subfamily is given in the bracket.



**Supplementary Figure 4 Percentages of TSSs associated with different retrotransposons which contain a TATA-box motif starting at 25-35 bp upstream regions of the dominant TSSs.**

**Supplementary Figure 5 Distribution of young TSSs along L1 subfamilies.** The top 10 L1 subfamilies, which harbor most young TSSs, show considerable heterogeneity regarding the positions of young TSSs within the consensus sequences. The gray barplots are background positional distributions of sequences from the corresponding subfamilies in the human genome. Number of young TSSs for each subfamily is given in the bracket.

7

**Supplementary Figure 6 Putative false positives associated with pseudogenes and poly(dA:dT) tracts in FANTOM 5 TSSs. (a)** Percentages of FANTOM 5 TSSs of GM12878 found in GRO-cap defined TSSs of GM12878 (from Core et al. 2014), based on the FANTOM TSSs found only in primate lineages. A FANTOM TSS is considered to be found in the GRO-cap dataset if it is within 100 bp of a GRO-cap TSS. **(b)** Distribution of FANTOM 5 TSSs along the Alu consensus element before and after filtering the suspicious TSSs.



**Supplementary Figure 7 TSSs associated with tandem repeats (TRs) but not associated with TEs. (a)** Comparison of TSSs with non-TE-associated TRs and TSSs with TE-associated TRs regarding sequence conservation scores among placental mammals and GC content. **(b-d)** Examples of non-TE-associated TR expansions which contribute to new TSSs in (b) 'hominid', (c) 'OWA' and (d) 'primate' groups respectively. In each panel, at the top are the CAGE total tag counts and transcripts from FANTOM (red, forward strand; blue, reverse strand) and genomic alignment

8

blocks from UCSC genome browser; at the bottom is the enlarged region of the young TSS, with grey shade indicating the TR expansion.



**Supplementary Figure 8 The exponential relationship between number of genes with a specific number of TSSs and number of TSSs per gene is independent of the gene lengths.** (**a**) Boxplots of transcript lengths for genes with different numbers of TSSs. For each gene, we used the length of its longest transcript. Although genes that have more TSSs tend to have longer transcripts, there are also many long genes that have small numbers of TSSs. Inspecting genes within specific length ranges (panels **b-d**), still reveals a clear exponential relationship. $R^2$ is the coefficient of determination for the linear regression in the figure.

**Supplementary Figure 9 Relationship between the number of old ('mammalian') TSSs per gene and the number of newly gained TSSs in primate lineages, on a log10 scale.** The red line is derived from linear regression based on the data points. Pearson's r and the corresponding p-value are also shown in the figure.



**Supplementary Figure 10 Relative deletion and insertion rates (normalized by genomic average) inferred from genomic alignments for three TSS groups.** Average rate were calculated for 40 bins along TSS±1kb. We estimated

insertion/deletion rates of two periods for 'mammalian' and 'primate' groups, but only one for the 'OWA' group so as to focus on the evolutionary rates after TSS loci emerged in the genome. Fitting curves were estimated by 'loess' method.



**Supplementary Figure 11 Single nucleotide polymorphism (SNP, panel a) and insertion/deletion (INDEL, panel b) densities around TSSs (40 bins along TSS±1kb), based on variants of the 1000 genomes project phase 3 release.** Only the biallelic variants with a minor allele frequency of ≥ 0.01 were considered. Because the genotype files in 1000 genomes project lack the ancestral allele information for insertion/deletion variants, the insertion and deletion variants were merged together for this analysis. Fitting curves were estimated by the 'loess' method.



**Supplementary Figure 12 DNA methylation in TSS loci in the germline.** Violin and box plots for germline CpG methylation levels (data from Guo et al. 2015) in different TSS subgroups defined by the types of associated retrotransposons. For each TSS, average methylation level of CpGs in the 2 kb around the TSS was calculated. The TSSs in the "Others" group are mostly non-TE-associated TSSs, except for a few that are associated with DNA transposons. Statistical significance was calculated using the one tailed Wilcoxon rank sum test ("*", $p < 0.05$; "**", $p < 0.01$; "***", $p < 0.001$; N.S., not significant).
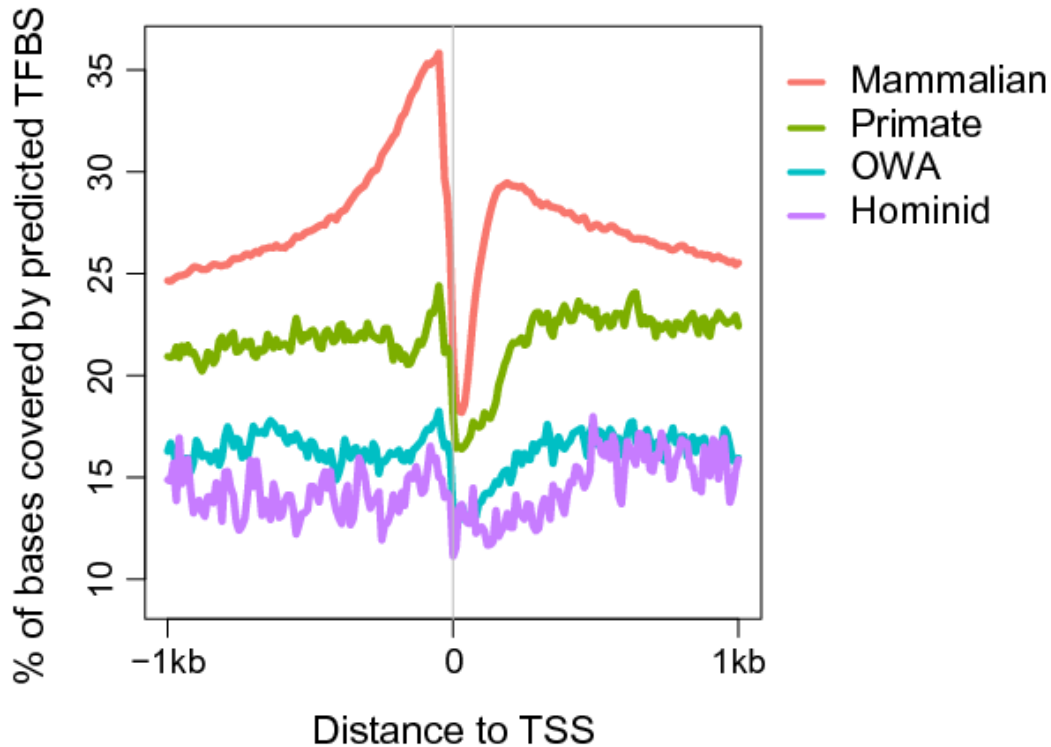
**Supplementary Figure 13** (**a**) Meta-profiles of functional signatures in H1-hESC cell line in different TSS groups. (**b**) Meta-profiles of functional signatures in K562 cell line in different TSS groups. Global hypomethylation in the K562 cell line has been previously reported, so the similar pattern of DNA methylation meta-profiles in K562 across TSS groups is not surprising. For the TFBS analysis, we merged the called peaks of TF ChIP-seq datasets and calculated how many bases around TSSs are covered by the peaks. The figure for RNAP II ChIA-PET in H1-hESC is missing because of lack of publicly available data.



**Supplementary Figure 14 Meta-profiles for histone modifications in GM12878, supplementary to that shown in Fig. 4.** All the data was obtained from ENCODE project.
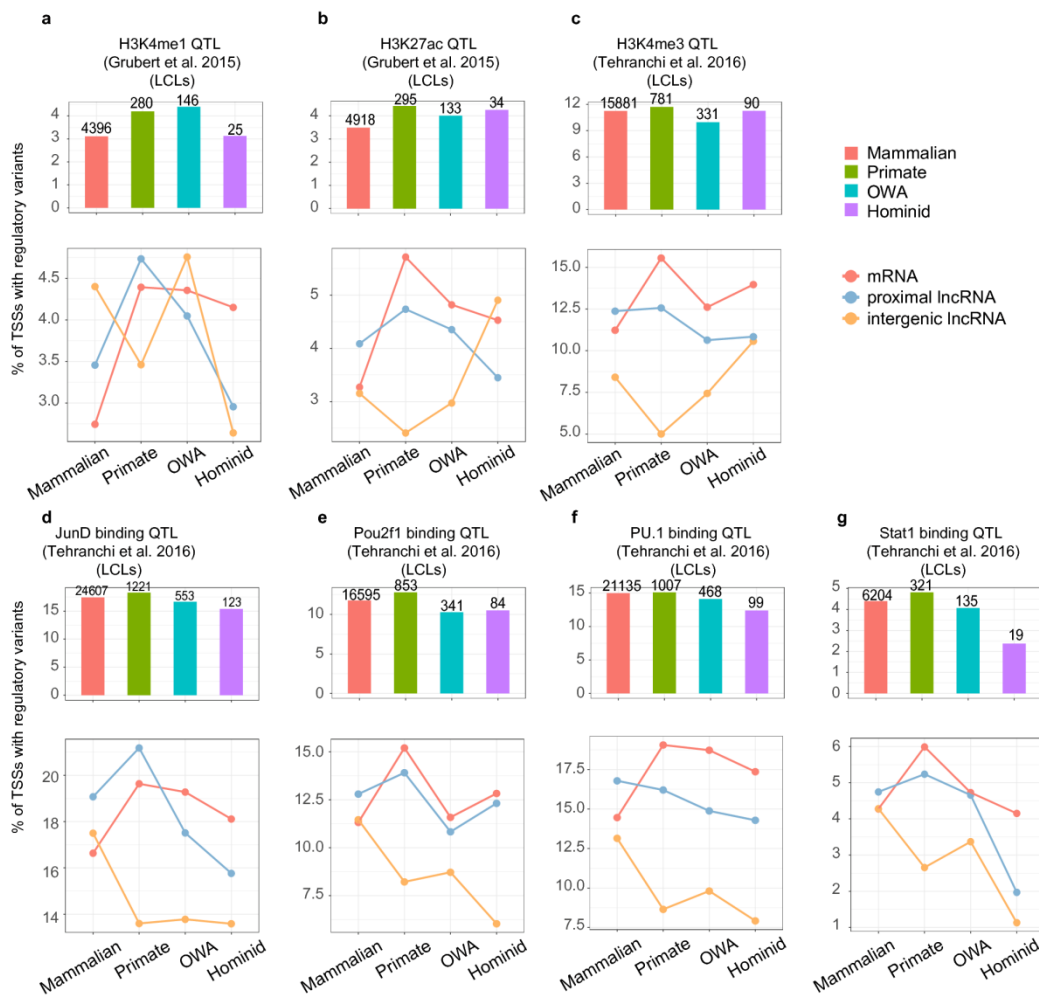
**Supplementary Figure 15 Meta-profiles for TF ChIP-seq signals in GM12878 cell line in different TSS groups.** All the data was obtained from ENCODE project.
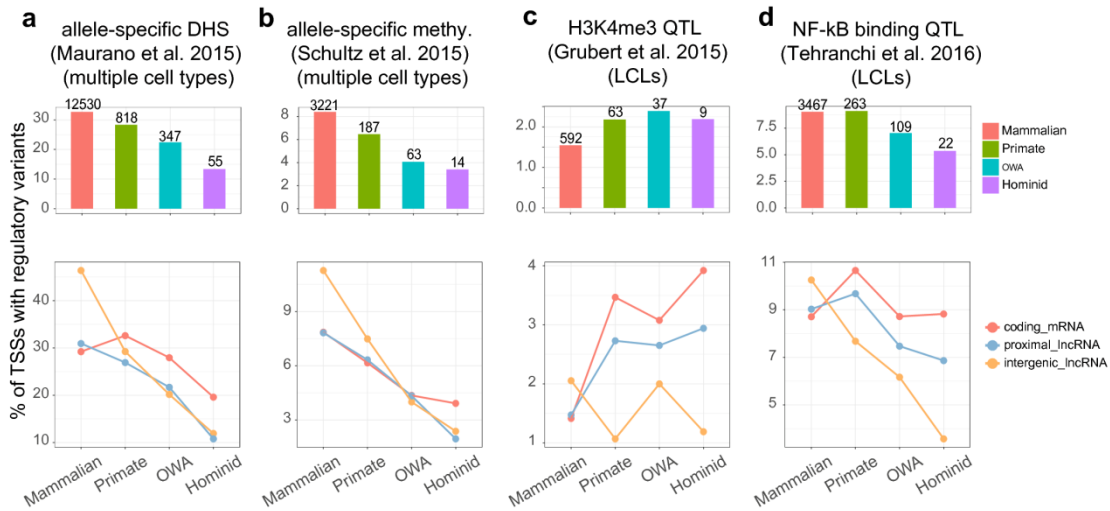
**Supplementary Figure 16 Comparison of the coverage by computationally predicted TFBSs between four TSS groups.** The computationally predicted TFBSs in human genome were from ENCODE project (http://compbio.mit.edu/encode-motifs/). Note that the TFBSs predicted by computational methods are based on binding motifs, usually smaller than the called peaks in the TF ChIP-seq data that were used for generating **Fig. 4**.
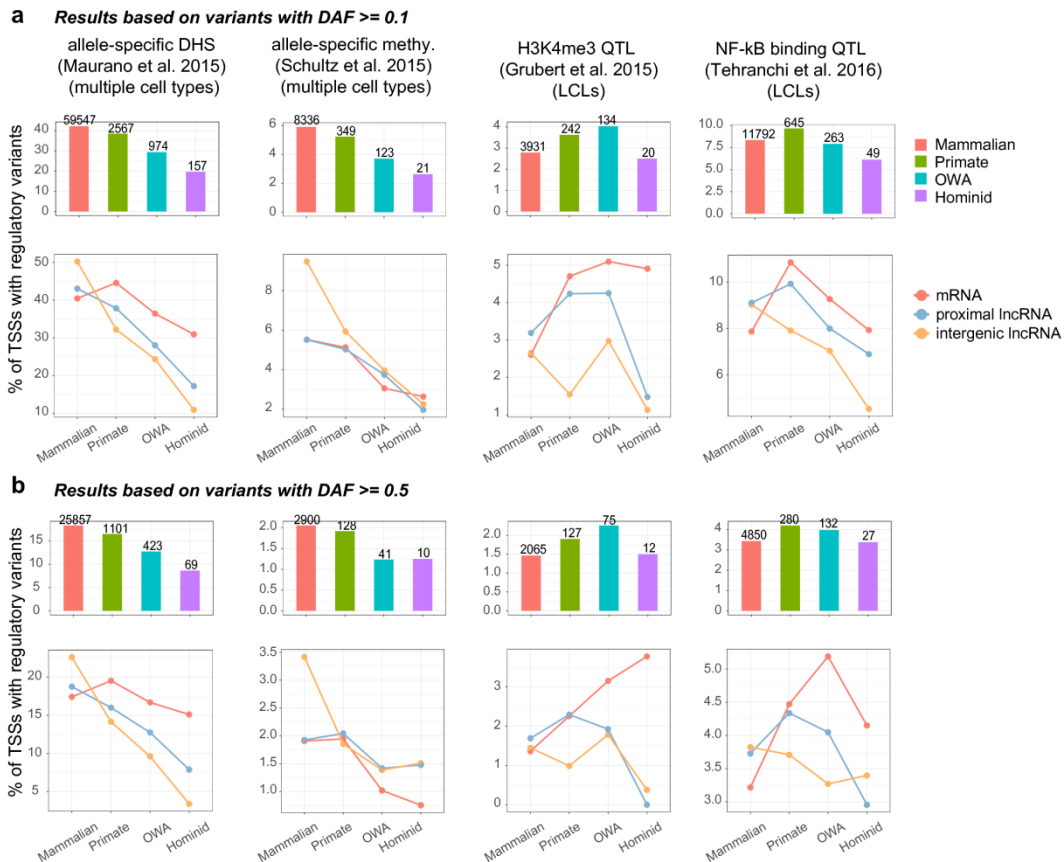


**Supplementary Figure 17 Meta-profiles of functional signatures in GM12878 cell line for different TSS subgroups, defined by transcript types**. (**a**) For 'primate' TSS subgroups. (**b**) For 'hominid' TSS subgroups. Statistical significance was calculated using the one-tailed Wilcoxon rank sum tests ("*", $p < 0.05$; "**", $p < 0.01$; "***", $p < 0.001$).

**Supplementary Figure 18 Proportions of TSSs harboring regulatory variants within TSS±1kb in different TSS groups in additional datasets.** The results in this figure were based on regulatory variants with derived allele frequency (DAF) $\geq 0.01$. Above the bars are the numbers of TSSs with regulatory variants. Note that for the H3K4me3 QTL dataset from Grubert et al. (2015), the numbers of regulatory variants found in the TSS groups/subgroups are very small, so the trends shown in the panels **a-b** for different transcript types may not accurately reflect actual trends. LCLs, lymphoblastoid cell lines.
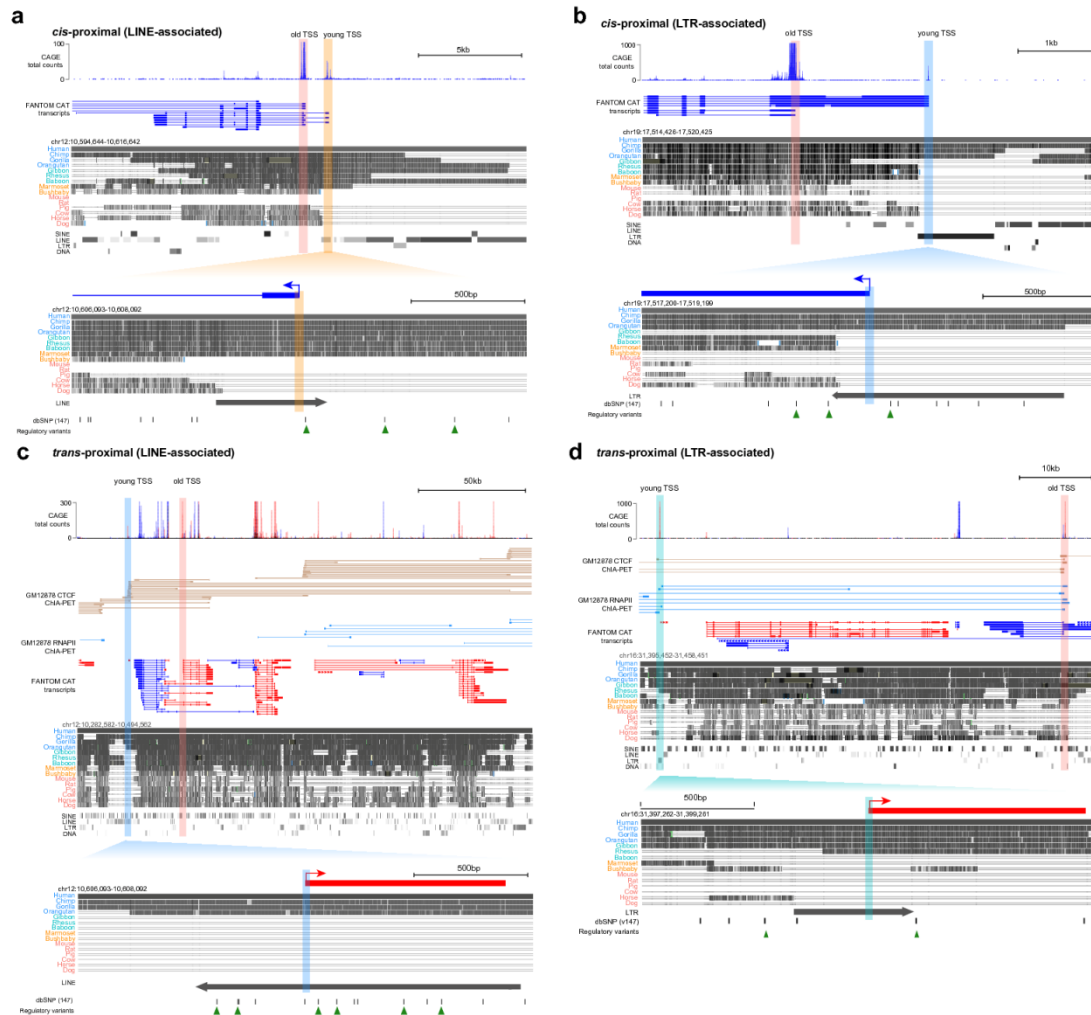
**Supplementary Figure 19 Proportions of TSSs harboring regulatory variants within TSS±1kb in different TSS groups excluding the TSSs separated by < 2 kb.** The shown results are based on variants with DAF ≥ 0.01. Above the bars are the numbers of TSSs with regulatory variants. Note that for the H3K4me3 QTL dataset from Grubert et al. (2015), the numbers of regulatory variants found in the TSS groups/subgroups are very small, so the changing trends shown in the panel **c** for different transcript types may not accurately reflect actual trends. LCLs, lymphoblastoid cell lines.



**Supplementary Figure 20 Proportions of TSSs harboring regulatory variants within TSS±1kb in different TSS groups, based on variants with higher**

**thresholds of derived allele frequency (DAF).** (**a**) Results based on variants with DAF $\geq 0.1$. (**b**) Results based on variants with DAF $\geq 0.5$. Above the bars are the numbers of TSSs with regulatory variants. Note that for the results based on variants with DAF $\geq 0.5$, the numbers of regulatory variants found in the TSS groups/subgroups are very small, so the changing trends shown in the some panels for different transcript types may not accurately reflect actual trends. LCLs, lymphoblastoid cell lines.



**Supplementary Figure 21 Additional examples for *cis*-proximal and *trans*-proximal young TSSs.** In each panel, from top to bottom: 1) CAGE total tag counts from FANTOM; 2) CTCF and RNAP II ChIA-PET interactions (only for *trans*-proximal examples); 3) FANTOM CAT transcript models, red for forward-strand and blue for reverse-strand transcripts; 4) genome alignments represented by grey blocks and transposable elements within this region, generated from UCSC genome browser; 5) the enlarged region of the young TSS. The old and young TSSs are indicated with shades of different colors (red, "mammalian"; green, "primate"; cyan, "OWA"; blue, "hominid"). The positions of regulatory variants are shown with small triangles in the enlarged figures.

17