

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31

Extensive binding of uncharacterized human transcription factors to genomic dark matter

Rozita Razavi^{1*}, Ali Fathi^{1*}, Isaac Yellan^{1*}, Alexander Brechalov^{1*}, Kaitlin U. Lavery^{1,2}, Arttu Jolma¹, Aldo Hernandez-Corchado³, Hong Zheng¹, Ally W.H. Yang¹, Mihai Albu¹, Marjan Barazandeh¹, Chun Hu¹, Ilya E. Vorontsov⁴, Zain M. Patel¹, The Codebook Consortium, Ivan V. Kulakovskiy⁵, Philipp Bucher⁶, Quaid Morris², Hamed S. Najafabadi^{3,7}, and Timothy R. Hughes^{1**}

¹Donnelly Centre and Department of Molecular Genetics, 160 College Street, Toronto, ON M5S 3E1, Canada

²Memorial Sloan Kettering Cancer Center, Rockefeller Research Laboratories, New York, NY 10065, USA

³Victor P. Dahdaleh Institute of Genomic Medicine, 740 Dr. Penfield Avenue, Room 7202, Montréal, Québec, H3A 0G1, Canada

⁴Vavilov Institute of General Genetics, Russian Academy of Sciences, 119991, Moscow, Russia

⁵Institute of Protein Research, Russian Academy of Sciences, 142290, Pushchino, Russia

⁶Swiss Institute of Bioinformatics, 1015, Lausanne, Switzerland

⁷Department of Human Genetics, McGill University, Montréal, Québec, H3A 0C7, Canada

*These authors contributed equally

** To whom correspondence should be addressed: t.hughes@utoronto.ca

32 **The Codebook Consortium**

33 **Principal investigators (steering committee)**

34 Philipp Bucher, Bart Deplancke, Oriol Fornes, Jan Grau, Ivo Grosse, Timothy R.
35 Hughes, Arttu Jolma, Fedor A. Kolpakov, Ivan V. Kulakovskiy, Vsevolod J. Makeev

36 **Analysis Centers:**

37 **University of Toronto (Data production and analysis):** Mihai Albu, Marjan
38 Barazandeh, Alexander Brechalov, Zhenfeng Deng, Ali Fathi, Arttu Jolma, Chun Hu,
39 Timothy R. Hughes, Samuel A. Lambert, Kaitlin U. Lavery, Zain M. Patel, Sara E. Pour,
40 Rozita Razavi, Mikhail Salnikov, Ally W.H. Yang, Isaac Yellan, Hong Zheng

41 **Institute of Protein Research (Data analysis):** Ivan V. Kulakovskiy, Georgy
42 Meshcheryakov

43 **EPFL, École polytechnique fédérale de Lausanne (Data production and analysis):**
44 Giovanna Ambrosini, Bart Deplancke, Antoni J. Gralak, Sachi Inukai, Judith F.
45 Kribelbauer-Swietek

46 **Martin Luther University Halle-Wittenberg (Data analysis):** Jan Grau, Ivo Grosse,
47 Marie-Luise Plescher

48 **Sirius University of Science and Technology (Data analysis):** Semyon Kolmykov,
49 Fedor Kolpakov

50 **Biosoft.Ru (Data analysis):** Ivan Yevshin

51 **Faculty of Bioengineering and Bioinformatics, Lomonosov Moscow State
52 University (Data analysis):** Nikita Gryzunov, Ivan Kozin, Mikhail Nikonov, Vladimir
53 Nozdrin, Arsenii Zinkevich

54 **Institute of Organic Chemistry and Biochemistry (Data analysis):** Katerina
55 Faltejskova

56 **Max Planck Institute of Biochemistry (Data analysis):** Pavel Kravchenko

57 **Swiss Institute for Bioinformatics (Data analysis):** Philipp Bucher

58 **University of British Columbia (Data analysis):** Oriol Fornes

59 **Vavilov Institute of General Genetics (Data analysis):** Sergey Abramov, Alexandr
60 Boytsov, Vasilii Kamenets, Vsevolod J. Makeev, Dmitry Penzar, Anton Vlasov, Ilya E.
61 Vorontsov

62 **McGill University (Data analysis):** Aldo Hernandez-Corchado, Hamed S. Najafabadi

63 **Memorial Sloan Kettering (Data production and analysis):** Kaitlin U. Lavery, Quaid
64 Morris

65 **Cincinnati Children's Hospital (Data analysis):** Xiaoting Chen, Matthew T. Weirauch

66 **SUMMARY**

67 **Most of the human genome is thought to be non-functional, and includes large**
68 **segments often referred to as “dark matter” DNA. The genome also encodes**
69 **hundreds of putative and poorly characterized transcription factors (TFs). We**
70 **determined genomic binding locations of 166 uncharacterized human TFs in**
71 **living cells. Nearly half of them associated strongly with known regulatory**
72 **regions such as promoters and enhancers, often at conserved motif matches and**
73 **co-localizing with each other. Surprisingly, the other half often associated with**
74 **genomic dark matter, at largely unique sites, via intrinsic sequence recognition.**
75 **Dozens of these, which we term “Dark TFs”, mainly bind within regions of closed**
76 **chromatin. Dark TF binding sites are enriched for transposable elements, and are**
77 **rarely under purifying selection. Some Dark TFs are KZNFs, which contain the**
78 **repressive KRAB domain, but many are not: the Dark TFs also include known or**
79 **potential pioneer TFs. Compiled literature information supports that the Dark TFs**
80 **exert diverse functions ranging from early development to tumor suppression.**
81 **Thus, our results sheds light on a large fraction of previously uncharacterized**
82 **human TFs and their unappreciated activities within the dark matter genome.**

83 **KEYWORDS:** Transcription factor (TF), ChIP-seq, SELEX, GHT-SELEX, PWM, Gene
84 regulation, KRAB zinc finger protein, C2H2, Codebook

85

86 INTRODUCTION

87 Deciphering *cis* and *trans* gene regulation is a long-standing challenge in molecular
88 biology and computational genomics. Transcription factors (TFs) are the sequence-
89 specific DNA binding proteins that regulate gene expression, typically by associating
90 with promoters and/or enhancers^{1,2}. The human genome encodes over 1,600 apparent
91 TFs, but hundreds of them have been identified as such only on the basis of conserved
92 protein domain structures, and are otherwise poorly characterized, with no known DNA
93 binding motif³. The function of much of the conserved (and presumably functional)
94 noncoding DNA in human is also largely unknown, although at least some of it is
95 involved gene regulation⁴. Indeed, phylogenetic footprinting, which identifies islands of
96 conserved DNA sequence, is a long-established approach to identify functional
97 regulatory elements^{5,6}, despite being limited by the frequent turnover of TF binding
98 sites⁷.

99 The functionality of most of the human genome is similarly ambiguous and uncertain, as
100 most of it is both non-genic and unconserved^{4,8}. Roughly half is composed of
101 transposable elements (TEs), especially endogenous retroelements (EREs)⁹. Broadly,
102 the human genome can be divided into gene-rich regions and gene deserts¹⁰; some of
103 the latter contain transcriptional enhancers¹¹, while others appear to be dispensable¹².
104 Genome-wide chromatin contact maps also reveal two major genomic compartments
105 that partially mirror the dichotomy in gene density: the “A” compartment is broadly
106 associated with the presence of genes and accessible chromatin, while the “B”
107 compartment is enriched for LINE-1 (L1) elements and constitutive heterochromatin¹³.
108 The functions of the compartments and the mechanisms that create them are not well
109 understood i.e. it is not clear whether they are defined by and/or recruit specific TFs,
110 and how often (if at all) they recruit factors that contribute to regulation and/or structure
111 of the genome.

112 It is known, however, that many human TFs bind to specific classes of TEs. TEs can
113 rewire regulatory circuitry by introducing transcription factor binding sites, thus spawning
114 novel cis-regulatory elements¹⁴⁻¹⁶. In some cases, these elements derive from the
115 promoter of the TE (e.g. endogenous retrovirus long terminal repeats (ERV LTRs))^{17,18}.
116 Other cases may represent inadvertent matches to a host TF motif within a TE¹⁹. In
117 addition, mammalian genomes encode a large family of KRAB-containing C2H2 zinc-
118 finger (KZNF) TFs (~350 members in human), which evolves rapidly in parallel to ERE
119 classes bound by its members²⁰. KZNFs silence EREs via direct recruitment of
120 KAP1/TRIM28, which associates physically with both readers (e.g. HP1/CBX proteins)
121 and writers (SETDB1)^{21,22} of the H3K9me3 mark that defines constitutive
122 heterochromatin^{23,22}. The KZNFs are also known for the potential to have very long
123 binding sites, enabled by their long C2H2-zf domain arrays²⁴. In general, the DNA
124 binding preferences of C2H2-zf proteins, KRAB or otherwise, have proven difficult to
125 characterize precisely, due to a lack of antibodies for ChIP-seq and low apparent
126 functionality in biochemical assays; such that they are depleted from systematic studies
127 of human TF motifs (e.g.²⁵). The largest collections of binding data for C2H2-zf proteins
128 have come from studies using ChIP-seq with epitope-tagged proteins in cultured
129 cells^{26,27}. These data reveal what families of EREs are bound by C2H2-zf proteins, but it

130 is difficult to accurately determine their precise sequence specificity (and thus the exact
131 binding sites) because the repeat elements are related by common descent, which
132 confounds motif discovery²⁸. The lack of accurate knowledge of DNA sequence
133 specificity of TFs complicates interpretation of ChIP-seq data, in general, because
134 ChIP-seq readily detects indirect recruitment and nonspecific binding^{29,30}.

135 Regulatory DNA represents an expanding frontier in genetics, and it is critical that we
136 gain a complete picture of human TF-DNA binding. As part of an international initiative
137 termed the “Codebook consortium”, aimed at obtaining binding motifs for all human
138 TFs³¹, we analyzed 315 uncharacterized human TFs by ChIP-seq in HEK293 cells,
139 together with 58 controls. We evaluated the data in conjunction with other data from the
140 Codebook project, which allowed base-level identification of direct binding sites.
141 Previous ChIP-seq analyses have focused mainly on preferential association of TFs
142 with promoters vs. enhancers¹, and indeed, many of the 217 TFs yielding reliable data
143 in our study (i.e. reproducible and/or enriched for TF motif matches) bound
144 predominantly and directly to such sites, consistent with conventional roles in gene
145 regulation, and providing a likely explanation for the fact that their binding sites are
146 frequently conserved. Surprisingly, however, roughly half of the uncharacterized
147 Codebook TFs, including most KZNFs as well as other TF families, bound to apparently
148 unique sites that are located in regions depleted from activating epigenetic marks. A
149 subset of these TFs also bind mainly to closed chromatin; we refer to this subset as
150 “Dark TFs”. Multiple lines of evidence suggest diverse biochemical, cellular, and
151 physiological functions of the Dark TFs, and by extension, the dark matter genome.

152 **RESULTS**

153 **Generation of ChIP-seq data for hundreds of putative TFs**

154 We surveyed the genomic binding sites of 314 poorly characterized, putative human
155 TFs (the “Codebook” set, derived from Lambert 2018³², and described in detail
156 elsewhere³¹), and 58 previously characterized TFs as controls (selected from Isakova
157 2017³³ and Schmitges 2016²⁷), using ChIP-seq in HEK293 cells (**Figure 1A, Table S1**).
158 We used an inducible eGFP-tagged transgene system (**Figure 1B**) previously employed
159 for ChIP-seq and to identify protein-protein interactions^{27,34,35}. Using this system, we
160 have shown that KZNFs bind to specific classes of retroelements, and that their binding
161 sites are often depleted for open chromatin, indicating that the transgene system can
162 readily assess binding to inactive or repressed regions of the genome^{27,34}. The present
163 study includes biological replicates performed by different experimentalists, such that
164 the resulting dataset includes 678 ChIP-seq experiments for Codebook TFs, and 112
165 experiments for control TFs (**Table S1**). A full list of experiments is given in **Table S2**.
166 Representative motifs obtained from control TFs of various families illustrate that the
167 assay recovers known sequence-binding preferences, as expected (**Figure 1C**).

168 We used two criteria to determine which experiments were successful. First, the data
169 were analyzed as part of a larger Codebook benchmarking effort, which is described in
170 more detail in accompanying manuscripts^{31,36}. The Codebook benchmarking included
171 expert curation that relied mainly on obtaining similar motifs for the same TF from

172 different data types (ChIP-seq, Protein Binding Microarrays³⁷, SMiLE-seq³³, and several
173 variants of HT-SELEX³⁸) as evidence of direct, sequence-specific DNA binding. This
174 Codebook motif benchmarking identified 130 Codebook TFs and 49 controls with
175 “approved” ChIP-seq data, meaning that sequence-specific DNA binding is observed in
176 ChIP-seq, and it is supported in almost all cases by *in vitro* experimental data.

177 Second, we identified experiments in which the peak overlaps of biological replicates
178 exceeded what is expected at random (i.e. with TF identities permuted). ChIP-seq
179 experiments that were classified as “approved” based on the motif similarity analysis
180 described above displayed a higher overlap between TF replicates relative to
181 mismatched TFs (median Kulczynski II coefficient of 0.57 vs. 0.034; **Figure S1A**; this
182 statistic is a modified Jaccard value that compensates for class imbalance). A
183 Kulczynski II coefficient threshold of 0.4 captures 78% of approved experiments with
184 replicates, and 90% of controls, but eliminates 94.5% of mismatched experiments.
185 Among proteins for which there was no “approved” experiment but for which there were
186 biological ChIP-seq replicates, 36 putative TFs (and two controls) displayed peak
187 overlaps between replicates that exceeded a Kulczynski II coefficient of 0.4. In these
188 cases, the DNA binding may be indirect, i.e. these proteins may be DNA-associated
189 chromatin factors, rather than TFs. Alternatively, they may recognize properties of the
190 DNA sequence that are not captured by common motif models, or the constructs used
191 may be inactive for direct DNA binding in HEK293 cells, but competent for association
192 with chromatin. We included these 38 putative TFs in subsequent analyses, and we
193 refer to the entire set of 217 successful proteins (130+49 “approved”, and 36+2 with
194 matching replicates; **Table S3**) as “TFs”, for simplicity, although we caution that the
195 subsets that are not “approved” may instead be chromatin factors.

196 We merged the peaks from TF replicates among the 489 ChIP-seq experiments
197 deemed successful to produce a dataset for downstream analysis (**Table S3**). This
198 dataset encompasses 217 proteins (166 Codebook putative TFs, and 51 controls), with
199 a median of 12,681 peaks per protein (range 76-163,602) (using a MACS threshold of
200 $P < 10^{-10}$; see **Methods** for explanation of threshold choice)³⁹.

201 **Overview of ChIP-seq data illustrates that half of the Codebook TFs bind genomic** 202 **dark matter**

203 To begin characterizing the ChIP-seq data, we surveyed for preferential association of
204 the putative TFs with promoters and/or enhancers. **Figure 2A** shows the fraction of
205 peaks from each protein that overlaps with protein-coding promoters (defined by
206 RefSeq⁴⁰) and enhancers (defined by HEK293 chromatin state⁴¹, and corresponding
207 mainly to H3K4me1 signal; see **Methods**). Indeed, many proteins are highly associated
208 with promoters, and a smaller number with enhancers, although a large number of
209 proteins did not associate with either promoters or enhancers. On average, the
210 HEK293-derived enhancer set yielded higher overlaps than the larger, more universal
211 “GeneHancer” set⁴² (**Figure S2A**), indicating that this lower number of enhancer-
212 favouring TFs (relative to promoter-favouring) is not due to incomplete enhancer
213 annotations.

214 To gain further insight into the properties of the ChIP-seq binding sites, we compared
215 the peak sets for each of the proteins to those of all other proteins in the dataset, and to
216 a panel of genome annotations. **Figure 2B** (bottom) shows a symmetric heatmap of
217 Jaccard similarity (intersection/overlap) between all 217 ChIP-seq datasets, providing
218 an overview of the overlap between all pairs of TF peak sets. The heatmap at the top of
219 **Figure 2B** shows the fraction of each corresponding TF peak set that overlaps with
220 each type of genome annotation. The chromatin states were derived mainly from public-
221 domain data for unperturbed HEK293 cells; we therefore expect them to reflect the state
222 of the chromosomes prior to induction of the tagged TF transgene. This state could be
223 involved in recruiting the TF, but it could also result from endogenous expression of the
224 native, untagged protein, as most of the studied TFs are already expressed in HEK293
225 cells²⁷ (**Figure S3**).

226 **Figure 2B** reflects and expands upon trends observed in **Figure 2A**. The large bright
227 square in the lower right quadrant of the bottom heatmap corresponds to TFs that
228 associate primarily with open chromatin (ATAC-seq) and H3K4me3. These TFs also
229 often associate with many promoters (median 6,140 coding gene promoters; see
230 below), leading to high overlap between the peak sets. The observation that many TFs
231 co-bind promoters and/or enhancers is prevalent in the literature (e.g.¹); we note,
232 however, that the Codebook proteins were considered uncharacterized TFs at the
233 outset of this study, and therefore it appears that even well-known regulatory sequences
234 often contain previously unidentified TF binding sites.

235 A second main feature of **Figure 2B** is the diagonal line in the upper left quadrant of the
236 heatmap at the bottom. These are proteins for which there is very little overlap in peaks
237 with any other Codebook TF. In addition, these peaks often do not overlap with any
238 peak from any other protein in the Codebook ChIP-seq dataset (**Figure S1B**). The
239 unique binding profiles are not explained by experimental error or random events: there
240 is strong overlap between replicates of the same protein (**Figure S1C**), and these
241 proteins often bind to the same unique sequences in ChIP-seq and *in vitro* (see below).
242 The ChIP-seq peaks for these proteins also tend to be outside open chromatin, and
243 outside of either promoters or apparent enhancers in HEK293 cells (**Figure 2B, top**).
244 Instead, roughly half of these TFs' peak sets are either enriched for marks that
245 characterize heterochromatin, or lack any of the diagnostic marks of promoter or
246 enhancer activity (i.e. "empty" ChromHMM regions, **Figure S2B**). Peaks for a subset of
247 these TFs are mainly associated with the Hi-C "B" compartment, and many associate
248 with specific classes of repeats (**Figure 2B** and see below). Roughly half (50/94) of
249 these TFs are KZNFs, which would be expected to display these properties³⁴ (bar in
250 middle of **Figure 2B**).

251 In subsequent analyses, we sought to gain a better understanding of the properties and
252 characteristics of TFs (and their binding sites) that represent the major patterns shown
253 in **Figure 2**. For simplicity, we defined four mutually exclusive groups (see **Table S4** for
254 TF labels). One group we named "Promoter binders" (55 proteins); for these proteins,
255 more than 37% of peaks overlap with promoters (this threshold captures the most
256 prominent features in **Figure 2A**). Another group was designated "Enhancer binders" (9
257 proteins); for these proteins, >35% of peaks overlap with enhancers (this threshold

258 corresponds to the visual separation of data points on the vertical axis in **Figure 2A**). A
259 third group we named “Dark TFs”, after the genomic dark matter (54 proteins); for these,
260 most peaks lie within either the “empty”, “constitutive heterochromatin”, or “facultative
261 heterochromatin” states (i.e. in HEK293 cells, they are outside of the states that
262 represent promoters, enhancers, insulators, or gene bodies), and fewer than half of the
263 peaks overlap with ATAC-seq peaks in unperturbed HEK293 cells. These thresholds
264 exclude some TFs that may associate significantly with specific regions, but also bind
265 many other locations; for example, the control TF YY1 bound 53% of all promoters in
266 human, but it also had 32,046 additional binding sites outside promoters, which
267 represent 77% of all YY1 peaks. Similarly, a subset of KZNFs were not classified as
268 Dark TFs because they had many binding sites within open chromatin. The remaining
269 ~40% of TFs we labeled as “Other” (i.e. not Promoter, Enhancer, or Dark TFs); they
270 include 32 control TFs and 64 uncharacterized proteins which display a diversity of
271 attributes and patterns in the data. The “Other” TFs thus present a rich landscape for
272 further exploration, but we did not attempt to further subclassify them here.

273 **Direct DNA binding by Codebook TFs to specific types of genomic elements**

274 We next asked whether preferential association of the TFs with different types of
275 genomic regions and chromatin states could be accounted for by intrinsic sequence
276 recognition of the individual TF analyzed. We mainly compared the Promoter TFs and
277 Dark TFs, which are large groups that contrast in many ways. ChIP-seq can detect both
278 direct binding (i.e. the TF intrinsically recognizes the bound DNA sequences) and
279 indirect binding (e.g. recruitment by another factor)²⁹. ChIP-seq also readily detects non-
280 specific DNA-binding (e.g. by histones), and is biased towards open chromatin since the
281 sonication step preferentially releases these regions³⁰. Therefore, to accurately identify
282 direct binding sites in the ChIP-seq data, we used two independent sources of
283 information that were available as part of the larger Codebook initiative. First, we
284 employed data from a novel assay, GHT-SELEX (Genomic HT-SELEX; described in
285 detail in ⁴³), which surveys binding of synthetic TFs to fragmented, purified, and
286 unmodified genomic DNA *in vitro*; GHT-SELEX yields peaks that resemble those from
287 ChIP-seq, but with greater resolution due to the smaller DNA fragment lengths (~64
288 bases). Second, for each TF, we obtained genomic matches to its DNA binding motif,
289 modeled as Position Weight Matrices (PWMs) with an associated PWM score. PWM
290 derivation and benchmarking are described in more detail in accompanying
291 manuscripts^{31,36}.

292 To make a conservative assessment of direct binding, we considered a ChIP-seq peak
293 to be bound directly by a TF if the peak overlapped with a GHT-SELEX peak for the
294 same TF, and also contained at least one motif match. In addition, the significance
295 thresholds for all three (ChIP-seq and GHT-SELEX peaks, and PWM hits treated as
296 peaks) were adjusted to maximize the Jaccard value (intersection/union) between all
297 three peak sets; we refer to these as “triple overlap” or TOP sites (see **Methods** for
298 details). In this approach, false negatives will arise due to any experimental error or
299 inaccuracy of motif models, as well as widespread non-specific DNA binding, which will
300 tend to raise the thresholds for sequence-specific binding in these procedures. Thus,
301 the number of direct binding sites obtained are underestimates. In addition, 37% of the

302 217 proteins lacked GHT-SELEX data, and/or did not have motifs; therefore, this
303 analysis could be conducted only on 137 TFs (101 Codebook TFs and 36 controls; see
304 **Table S4**). For these 137, the fraction of peaks that could be accounted for by direct
305 binding ranged from 0.04% (for SP140, which binds a short motif composed mainly of a
306 CG dinucleotide) to 65.7% (for ZNF728, which has a unique 21-base motif), with a
307 median of 10%. The fraction of ChIP-seq peaks that are due to apparent direct binding
308 (i.e. % of all ChIP-seq peaks that are TOP sites) is similar for the Promoter TFs (9.5%)
309 and Dark TFs (9.3%) (**Figure 3A**), and the absolute number of direct binding sites is
310 similarly high for a subset of both groups (**Figure 3B**). We conclude that there is no
311 systematic difference between Promoter TFs and Dark TFs in direct binding
312 characteristics, and that many of the observed TF binding sites are direct.

313 To ask whether the relative preference of TFs for different types of genomic regions and
314 chromatin states in ChIP-seq is intrinsic, we examined the fraction of GHT-SELEX and
315 ChIP-seq peaks for each TF that are found within genomic regions corresponding to
316 each type of genome annotation. To avoid circularity, we used universal peak
317 thresholds (i.e. the same cutoff across all experiments, see **Methods**) which lowers the
318 overlap between GHT-SELEX and ChIP-seq peak sets. The fraction of intrinsic (i.e.
319 GHT-SELEX) and cellular (i.e. ChIP-seq) sites for each TF that overlap with protein-
320 coding gene promoters, repeat sequences (of any kind), and the combination of the
321 “empty” and “heterochromatin” states are shown in **Figures 3C, D, and E**, respectively.
322 In each case, there is preferential binding *in vitro* which corresponds to that observed in
323 cells, with Promoter TFs having much higher intrinsic preference for promoter DNA, and
324 Dark TFs having higher preference for repeats and empty/heterochromatin. We note
325 that many Promoter and Enhancer TFs have a greater tendency to bind “empty” and
326 “heterochromatin” state DNA *in vitro* than in cells, which could be due to a *bona fide*
327 preference for open chromatin, functional binding at these loci in other cell types (but
328 not HEK293), or preferential extraction of these proteins at open chromatin in ChIP-seq
329 experiments.

330 The Promoter TFs also displayed intrinsic preference for the regions that overlap or are
331 just upstream of transcription start sites (TSS), by multiple measures (**Figure 3F**),
332 similar to that described for characterized TFs in a variety of genomes⁴⁴⁻⁴⁶. This
333 observation is consistent with functional roles for these uncharacterized TFs in promoter
334 definition, delineation of TSS location, and/or gene regulation.

335 **Distinct conservation patterns of Dark TF vs Promoter TF binding sites**

336 To further query functionality of direct binding sites (i.e. TOPs), we examined
337 conservation of the TOP sites, producing an estimate of whether each site is under
338 purifying selection. In essence, for many TFs, the TOP sites in aggregate display
339 conservation patterns that mimic the selectivity of each base position in the TF’s PWM.
340 **Figure 4A** shows a graphic demonstration: when TOP sites are aligned to the PWM hit,
341 and displayed as heatmaps that show base-level conservation scores (here, phyloP⁴⁷),
342 there are often vertical blue lines. These lines represent positions in the PWM hits that
343 are preferentially conserved across many TOP sites. Similar to previous observations
344 made with well-characterized TFs⁴⁷, the positions with highest conservation often

345 correspond to tall letters in the sequence logo (i.e. high information content), indicating
346 selection on the binding site to match the sequence preferences of the TF.

347 We developed three heuristics to discriminate conserved vs. unconserved TOP sites
348 (see **Methods**). Two of them test for a relationship between the information content at
349 each base position of the PWM and the conservation score, while the third tests for
350 higher overall conservation at the PWM hit than in immediate flanking sequence. As
351 shown in **Figure 4A**, and in similar diagrams for all 137 TFs for which these tests could
352 be run (**Document S1**), these tests together detected sites that appear plausible by
353 visual inspection (i.e. apparent conservation signal relative to flanks). We considered a
354 site to be conserved if any of the three criteria were met, and at least one nucleotide in
355 the PWM hit had an FDR-corrected PhyloP score ≥ 1 . By these criteria, conservation of
356 TOPs is observed for both Promoter TFs and Dark TFs (**Figure 4A**), but the fraction
357 and absolute number of conserved TOP sites for Promoter TFs is much higher (**Figure**
358 **4B,C** and **Table S4**). This outcome suggests that many Promoter TF binding sites are
359 functional, and that the corresponding TFs have conserved functions at promoters. An
360 individual conserved TOP site (hereafter, “CTOP” site), for the Promoter TF ZNF407, is
361 shown in **Figure 4D**; like many Promoter TF TOPs and CTOPs (**Figure 3F**), it overlaps
362 with a transcription start site. CTOPs are also often found adjacent to other CTOPs
363 (explored in greater detail in an accompanying manuscript³¹); an example of multiple
364 sites for ZNF131 and YY1 is shown in **Figure 4E**.

365 Despite their lower numbers, there are still thousands of CTOPs for Dark TFs: in
366 aggregate, the criteria used here yielded 6,086. They tend to be distant from promoters,
367 or each other (e.g., 2,916 are > 1000 bp away from any other CTOP), and they tend to
368 have lower PhyloP scores than CTOPs for Promoter TFs. The Dark TF ZBTB40
369 recognized nearly 1,000 CTOP sites, the vast majority of which correspond to remnants
370 of *hAT/Charlie* DNA transposons (**Figure 4A,F**). Its most strongly conserved CTOP falls
371 outside of a transposon, however, and instead is within the PRKACA 3' UTR (**Figure**
372 **4G**), which may be relevant to its known function (see below). ZNF689, in contrast, is an
373 example of Dark TF that has a much smaller number of CTOP sites, and is enriched for
374 binding L1M5 elements across its TOPs (**Figure 4A**; example CTOP shown in **Figure**
375 **4H**). Overall, these analyses indicate that Dark TFs occupy a unique and expansive
376 fraction of the genome, and thousands of their direct binding sites show indications of
377 conserved function. The interactions of TFs with TEs, and the known and potential
378 functions of these and other Dark TFs, are explored in the next sections.

379 **Widespread and specific binding of Codebook TFs to transposable elements**

380 We reasoned that the generally low conservation in direct binding sites for Dark TFs
381 could be due to domination by TEs, which are typically under neutral selection. In
382 addition, TEs are only present in a subset of species that have retained an ancestral
383 insertion, limiting power to detect purifying selection. Indeed, 92.0% of the Dark TF
384 binding sites overlap with repeats (aggregated TOPs vs. Repeatmasker track) (vs.
385 25.3% for Promoter TFs, 39.9% for Enhancer TFs, and 46.8% for Other TFs).

386 The combination of ChIP-seq, GHT-SELEX and Codebook PWM data enables us to
387 circumvent previous challenges in analysis of repeat sequences, and to examine
388 binding of TFs to TEs with unprecedented precision, including detection of direct, base-
389 level binding. **Figure 5A** provides an overview of high-confidence ($P < 10^{-8}$, Fisher's
390 Exact Test) interactions between the Codebook TFs and specific TE classes, with ChIP-
391 seq and GHT-SELEX peak sets calculated separately (**Figure S4** shows an expanded
392 version with all rows labelled). A first observation that emerges from this analysis is that
393 Dark TFs are much more likely than Promoter TFs to significantly bind a specific class
394 of TEs (36% vs. 8%, respectively), but binding to a specific TE class is not a universal
395 or discriminating property of either the Dark TFs or KZNFs. Among the 42 TFs that
396 passed the cutoff, 20 are Dark TFs, 22 are KZNFs, 14 are both, and 14 are neither.

397 A second observation is that the TE enrichments in the GHT-SELEX data are virtually
398 identical to those in the ChIP-seq data (shown adjacent to each other in **Figure 5A**),
399 illustrating that specific binding to these elements is an intrinsic property of individual
400 TFs. To our knowledge, this is the first experimental demonstration that KZNFs
401 independently possess sufficient sequence specificity to discriminate ERE subfamilies
402 from the rest of the genome: previous motif models derived from ChIP-seq data were
403 unable to specify individual elements as precisely⁴⁸. A third observation is that TEs of all
404 major classes (LINE, SINE, LTR/ERV, and DNA transposons) are recognized by
405 specific TFs. Moreover, for all four major classes of TEs, there are cases in which
406 greater than 10% of a TF's TOP sites overlap one type of TE, and are conserved
407 (**Figure S5**), consistent with a function for the host genome

408 A fourth observation is that the encompassed TEs span a very wide age range, from
409 human-specific AluY elements, to L2, L3, and MIR, which pre-date eutherian mammals.
410 These associations can provide insight into the evolution and molecular function(s) of
411 the TFs. For example, ZNF836 and ZNF841 (which are both Dark TFs and KZNFs) are
412 paralogs that arose from a pre-simian duplication event⁴⁹ and bind to distinct subtypes
413 of the closely related, simian-specific MaLR LTR elements. They bind distinct motifs that
414 specify the differing base identities at homologous positions within the diverged LTR,
415 suggesting neofunctionalization and retention to maintain silencing of both LTR
416 subtypes (**Figure 5B**). There are also cases in which the TFs and TEs they bind are
417 grossly mismatched in age. For example, ZNF286B is a human-specific duplicate of
418 ZNF286A which has lost its KRAB domain^{50,51}, but its binding sites are enriched for
419 LINE-3 (L3), an ancient element found across all mammals, suggestive of coincidental
420 adaptation (**Figure 5A**). In contrast, ZNF362 and ZNF384 (both non-KZNFs) are
421 products of a duplication ~429 MYA (the duplication is found across bony vertebrates),
422 but the binding sites for both proteins are enriched for the much younger, primate-
423 specific Alu elements, as well as poly-A repeats, consistent with their DNA binding
424 motifs (**Figure 5C**). These proteins have the largest number of TOP sites within the
425 Codebook dataset, and it is possible that the recently-expanded target range of these
426 proteins is a coincidental liability, as rearrangements of both ZNF362 and ZNF384
427 genes (most commonly as fusions to activating TFs and cofactors) are found frequently
428 in leukemia^{52,53}.

429 Four of the Promoter TFs bind to specific TEs in this analysis (**Figure 5A**), potentially
430 providing direct links between TE insertions and regulation of host genes. One of them,
431 ZNF676, is a KZNF that was previously shown to associate with LTR12, and to repress
432 “transpochimeric” gene transcripts⁵⁴, which are generated during human early
433 embryogenesis. The GHT-SELEX data and Codebook motif pinpoint its exact binding
434 site in LTR12 (**Figure 5D**). ZNF676 may also have other roles at promoters: at the TOP
435 site upstream of QSER1 (**Figure 5D**), the LTR12 element is in the opposite orientation
436 from the gene. Another Promoter TF, JRK, preferentially associates with the DNA
437 transposon Tigger15a. JRK is itself derived from a Tigger element, and the Tigger15a
438 consensus sequence contains binding sites for JRK at its terminus³¹. Thus, Tigger15a
439 may have simultaneously contributed both to the rise of JRK protein and a set of JRK
440 binding sites that are still utilized; this hypothesis is supported by the taxonomic
441 distribution of JRK and Tigger15a to therian mammals (dating to 160 MYA).

442 **Older TFs tend to bind older DNA**

443 In addition to insertions such as TEs, new TF binding sites can emerge from random
444 mutations in pre-existing sequences. This mechanism is thought to be dominant for
445 traditional enhancer-binding TFs⁵⁵. To determine whether binding sites for Codebook
446 TFs evolved from ancestral DNA, we estimated the age of each TOP site for each TF,
447 gauged as that of the oldest ancestral genome that contains the entire site (i.e. a
448 gapless alignment, even if the base identities are different) in the Zoonomia mammalian
449 reconstructions⁵⁶ (**Figure 6A**). This is a simple heuristic, but we obtained a qualitatively
450 similar conclusions using other approaches to estimate binding site age (**Figure S6**). As
451 expected, TOP sites for the Dark TFs (most of which correspond to TEs) are estimated
452 to be younger on average than those of Promoter TFs (median ages of 46 and 72 MYA,
453 respectively (**Figure 6B, Table S4**), but there is a large overlap of age distributions
454 between the two TF groups. Both groups contain TFs with binding sites at both
455 extremes (i.e. very old or very young binding sites). Thus, average age of the binding
456 site is not a discriminating characteristic of Promoter vs. Dark TFs.

457 We also estimated the ages of the TFs, using catalogued ortholog and paralog
458 relationships⁵⁷ and species divergence times⁵⁸ (**Figure 6C, Table S4**). Overall, TFs in
459 both classes tend to be older than the sites they bind: Promoter TFs have a median age
460 of 429 MYA, while Dark TFs have a median age of 97 MYA, which is still older than a
461 typical binding site even for Promoter TFs. These results are consistent with the
462 established phenomenon of TF binding site turnover⁵⁵. They are also consistent with
463 previous observations with KZNFs, which concluded that the correlation between age of
464 binding sites and age of the KZNF is weaker than expected if they evolve only to silence
465 TEs²⁶. Together with the retention of many KZNFs that bind extinct TEs, this finding
466 supports the notion that KZNFs must frequently take on additional regulatory roles, e.g.
467 in regulation of host genes.

468 **Functions of Dark TFs**

469 Finally, we examined existing literature and databases to survey known and potential
470 functions for the Dark TFs, and related it to the data we collected (**Figure 7, Table S5**).

471 Most Dark TFs have apparent roles in repression of transcription. 35 out of 54 are
472 KZNFs, and for 14 of them (and one non-KRAB TF, ZNF888), physical association with
473 KAP1 has been verified^{27,59-62}. The KZNFs may also have repressive functions beyond
474 the recruitment of KAP1^{27,59}. Five of these Dark TF KZNFs also interact with TRIM39,
475 which itself interacts with numerous ubiquitin conjugating enzymes, H3K4 demethylase
476 KDM1A, and dozens of other KZNFs^{60,63}. One additional Dark TF KZNFs (and two other
477 Dark TFs) interact with TRIM33, a member of the TIF1Y complex that specifically
478 suppresses TGF β -responsive genes by directly interacting with the histone subunits as
479 well as E3 ubiquitin ligase⁶⁴.

480 Ten of the twelve non-KRAB Dark C2H2-zf proteins also appear to contribute to the
481 formation and maintenance of heterochromatin, by association with chromatin proteins
482 (CBX/HP1) directly, or via recruitment of other C2H2-zf proteins. One of them,
483 ZNF518B, was identified as a partner of both H3K27 methylase EZH2 and H3K9/H3K27
484 methylase G9A, and to promote H3K9me2⁶⁵. ZNF518B and ZNF280D both associate
485 with multiple CBX/HP1 proteins^{60,63}. ZNF518B binds many primate-specific L1
486 elements, but its most conserved binding sites are in its own promoter, suggesting a
487 critical negative feedback mechanism (**Figure S7**). In another example, ZNF516
488 associates with the multifunctional CTBP1/KDM1A/RCOR1 corepressor complex, and
489 its repressive function was shown in reporter assays⁶⁶. Intriguingly, 21 of the 47 C2H2-
490 zf proteins, including both KRAB and non-KRAB C2H2-zf proteins (as well as
491 transposon-derived ZBED9) interact with other C2H2-zf proteins, often extensively⁶⁰,
492 suggesting a potentially widespread role in organization of chromosome topology.

493 Four additional Dark TFs have other potential roles in repression of transcription. Three
494 of them are the paralogous nuclear speckle proteins SP100, SP140, and SP140L. Each
495 contains a SAND domain, which we confirmed binds to unmethylated CG dinucleotides
496 *in vitro*⁶⁷, and CG-containing motifs are enriched in their ChIP-seq peaks³⁶. These
497 proteins also contain PHD and BRD domains, which typically function as epigenetic
498 readers⁶⁸. In our ChIP-seq data, they are enriched at sites of H3K27me3 methylation
499 (**Figure 7, Table S5**). SP140 is an exceptional TF among the Codebook data set; its
500 ChIP-seq sites predominantly overlap with “GeneHancer” loci – i.e. these sites are
501 catalogued as enhancers in other cell types, but not HEK293 enhancers (**Figure S2A**),
502 suggesting that these loci may be actively silenced in HEK293. The fourth protein is
503 SCML4, a polycomb group protein that was included in our study because it contains an
504 AT hook, but we did not obtain evidence for its sequence-specific DNA binding. Thus, it
505 may be more properly described as a chromatin protein. SCML4 is reported to
506 associate with H3K4 demethylase KDM5C⁶⁹ as well as ubiquitination factors FBXO11
507 and UBR1⁷⁰.

508 Five of the non-KZNF Dark TFs may have roles other than repression. One of them,
509 SOX2, is a well-known pioneer factor that can bind to motif matches within unmodified
510 closed chromatin, but is inhibited to some extent by H3K9me3⁷¹. Indeed, in the ChIP-
511 seq data reported here, most (66%) of its TOP sites are in “empty” chromatin, and only
512 5% overlap with ATAC-seq peaks in unperturbed HEK293 cells, consistent with its
513 pioneer function. Less than 2% of SOX2 peaks overlap with heterochromatin (defined
514 by ChromHMM mainly by H3K9me3 and H327me3), consistent with H3K9me3 being

515 refractory to SOX2 binding. Two additional Dark TFs may also represent pioneers:
516 TPRX1 has recently been described as a master regulator in zygotic genome
517 activation⁷², while SALL3 controls the differentiation of hiPSCs into cardiomyocytes vs
518 neural cells⁷³. In contrast, two other Dark TFs have been described as impacting DNA
519 metabolism. ZNF384, which we find binds many Alu and Poly-A repeats, as described
520 above, is also known to bind Ku and recruit NHEJ factors to double-strand breaks⁷⁴.
521 ZNF146 binds L1 elements, and its depletion slows the replication fork⁷⁵.

522 The distinct binding sites and diversity of apparent effector mechanisms and cellular
523 roles of the Dark TFs suggest that they may each regulate specific biological functions,
524 and that they may also be multifunctional. Indeed, physiological consequences that
525 have been reported for perturbation of the Dark TFs vary widely (**Figure 7**, right column;
526 **Table S5** provides the values and sources), ranging from basic cellular processes to
527 development. For example, the KZNF ZNF689, which we show above binds ~50
528 conserved sites enriched for L1M5 (**Figure 4G**), also binds promoters of various L1
529 subtypes, preventing genomic instability conferred by L1 retrotransposition⁷⁶. ZBTB40,
530 which we observe almost exclusively at DNA *hAT/Charlie* transposons (**Figure 4E**), and
531 which is one of the oldest Dark TFs (**Figure 7**), was recently shown to bind telomeric
532 dsDNA breaks and maintain telomeric length⁷⁷. In mouse, *Zbtb40* deficiency impacts
533 spermatogenesis through disrupted telomeric lengthening and maintenance in
534 spermatocytes⁷⁸. *hAT/Charlie* transposons are enriched in telomeric regions of human
535 DNA⁷⁹, suggesting that this function may be conserved. The most conserved ZBTB40
536 binding site, however, is within the 3'UTR of PRKACA (**Figure 4F**). PRKACA encodes
537 the catalytic subunit α of protein kinase A, whose deficiency is associated with fertility
538 defects in male mice and humans⁸⁰. This site is also less than 1 kb from the TSS of the
539 chromatin regulator SAMD1, which impacts sperm cells⁸¹.

540 DISCUSSION

541 The Codebook ChIP-seq data provide cellular binding sites for 130 putative TFs,
542 defined as previously lacking PWMs or other models of sequence specificity. It
543 represents a valuable resource for studying TF function and evolution in the context of
544 regulatory genomics. For a large majority of the proteins assayed, we have also now
545 identified a binding motif which is supported by independent assays³⁶. Thirty-six of the
546 proteins did not produce a motif and may not be *bona fide* TFs. Their ChIP-seq profiles
547 are nonetheless informative: enrichment of ChIP-seq peaks at different types of
548 genomic features (e.g. promoters, repeats) or chromatin states, as well as co-
549 occurrence with peaks for other proteins (e.g. TFs), can yield clues as to potential
550 function.

551 Previous large-scale ChIP-seq analyses have mainly focused on the established roles
552 of TFs in binding to promoters and enhancers (e.g.¹). A major exception has been
553 studies of KZNFs, which focus on binding to TEs, and specifically EREs^{26,27,34}. The
554 analysis scheme described here considers these models of TF function as hypotheses
555 with equal weight. The known categories are clearly present, including preferences for
556 promoters and enhancers, as well as the strong tendency for KZNFs to bind specific
557 classes of EREs, and within constitutive heterochromatin. Overall, however, TF

558 behaviour with respect to chromatin states and genomic landmarks appears more
559 varied than a simple categorization scheme would imply. We did not systematically
560 explore the “Other” category, used here as a catch-all. Like the Dark TFs, it appears to
561 encompass proteins that satisfy some expectations of “pioneers”, given that they bind
562 both *in vitro* and *in vivo* to many regions that are labeled as inactive and/or closed
563 chromatin in HEK293 cells prior to induction of the TF. There are many intriguing TFs in
564 the “Other” category: one example is the non-KRAB protein ZSCAN2, which is involved
565 in spermatogenesis and fertility in mice (and perhaps human)⁸². We catalogued 183
566 CTOPs for ZSCAN2, and found that its binding sites are enriched for mammal-wide L3
567 elements.

568 Establishing functions, if any, for individual TF binding sites is a long-standing and
569 difficult problem in regulatory genomics. The level of binding site turnover observed on
570 evolutionary timescales requires that binding sites arise at random, many of which are
571 presumably irrelevant for gene regulation or reproductive fitness, at least initially. By this
572 reasoning, we expect that many biochemically verified, direct TF binding sites should be
573 non-functional, and indeed we find that, overall, most TOP sites are not conserved,
574 even for Promoter TFs. Phylogenetic footprinting does not discriminate between false
575 negatives due to binding site turnover or redundancy, and *bona fide* non-functional
576 sites, and therefore lack of conservation cannot be taken as lack of biological purpose.
577 Nonetheless, conserved TOP sites would seem most likely to yield interpretable results
578 in targeted laboratory studies. Sites overlapping TSS may be particularly fruitful, given
579 apparent constraint on both sequence and position of the binding site. More generally,
580 the Codebook TOP catalogue will provide a rich resource for future efforts in examining
581 genome function.

582 The low primary sequence conservation of Dark TF binding sites, especially relative to
583 those of Promoter TFs, could have several explanations. One possibility is that very few
584 of the binding sites are functional; in theory, only a single binding site that confers
585 modest selective advantage would be sufficient to drive retention of both the site and
586 the TF, with all other sites arising at random (and under neutral evolutionary pressure,
587 provided they are not detrimental). Another possibility is that the exact positioning of the
588 sites is not critical to their function, unlike Promoter TFs, which by definition must be
589 close to TSSs. Dark TF functions could simply require that binding sites are distributed
590 widely across non-functional DNA, and thus be highly redundant over large sequence
591 windows (e.g. TADs). Such functions would not preclude a small subset of sites being
592 co-opted for regulation of host genes, which would become constrained (i.e.
593 conserved). Regardless of what biochemical, cellular, and physiological functions are
594 revealed, the Dark TFs represent a new contribution to the decades-old odyssey into
595 the function and significance of the dark matter genome.

596

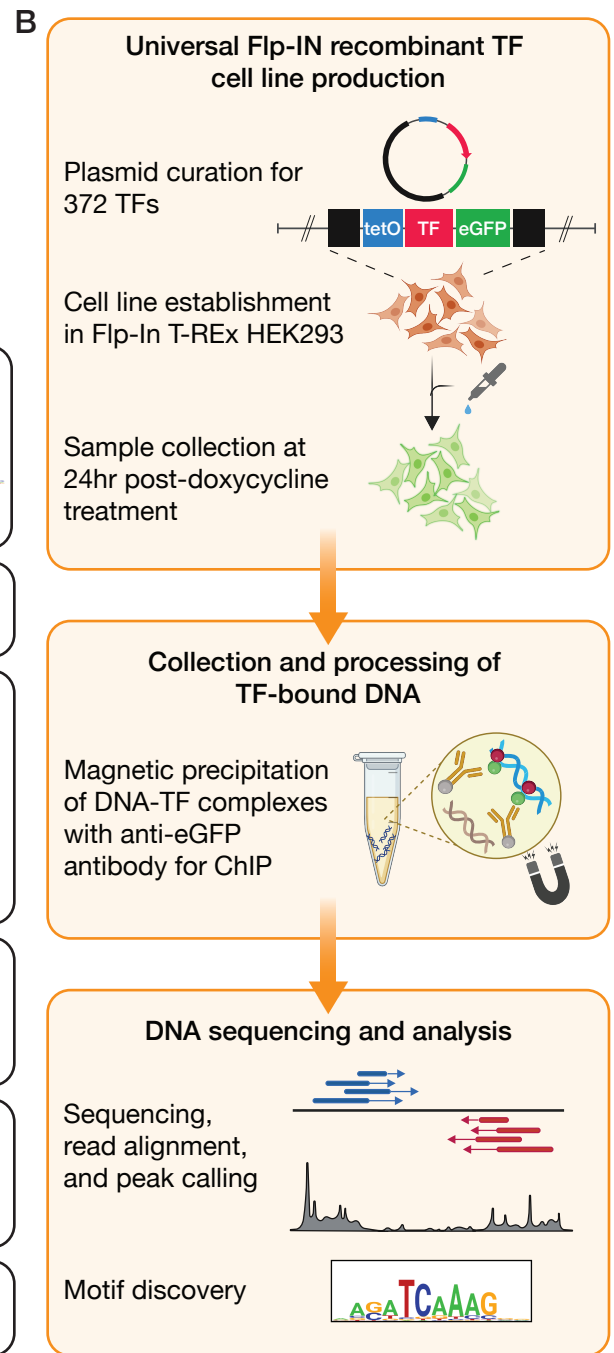
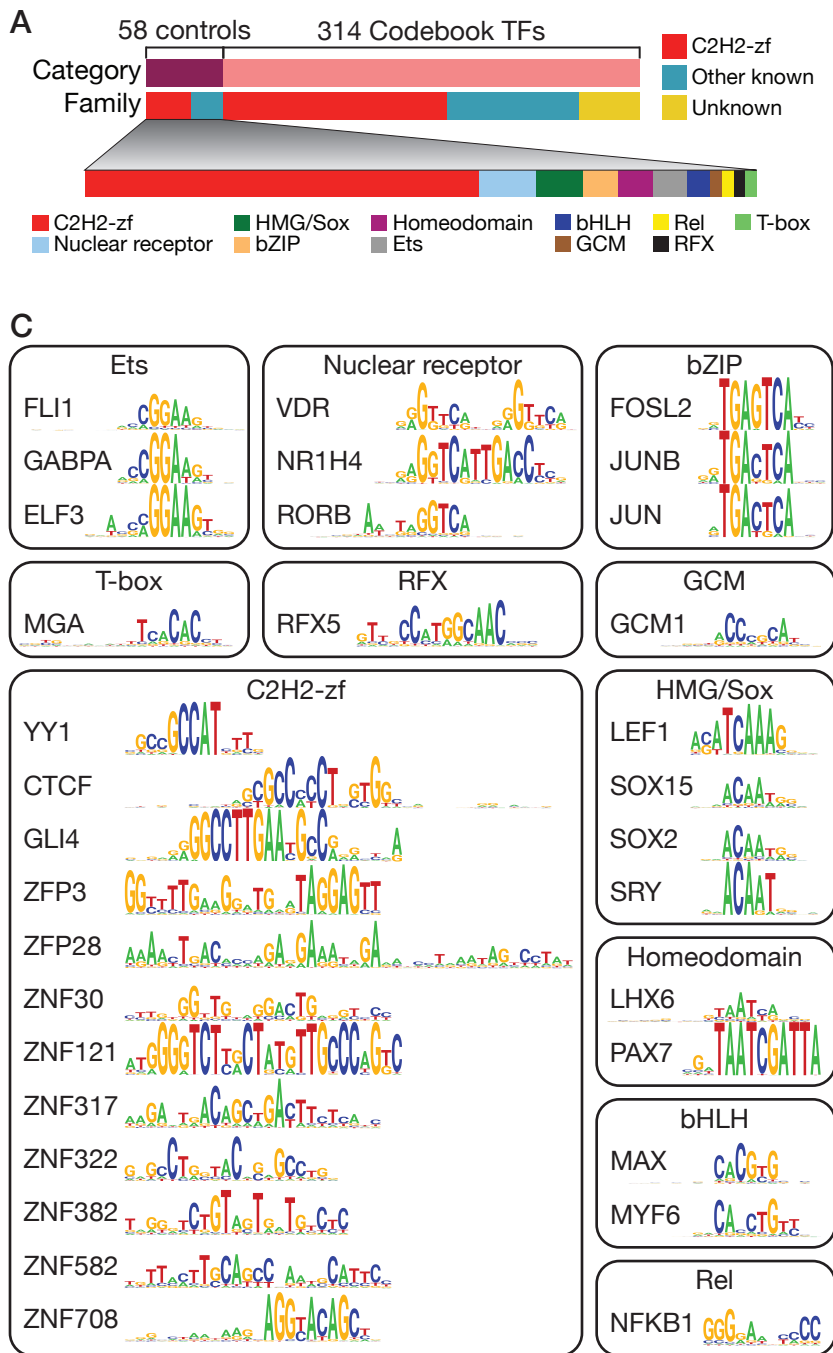


Figure 1. Project overview. (A) Overview of the TF categories assayed in this study. (B) A schematic of the experimental pipeline for production of 372 inducible EGFP-labelled TF cell lines used in ChIP experiments and deriving TF binding sites. (C) Samples of representative motifs obtained from different families of control TFs.

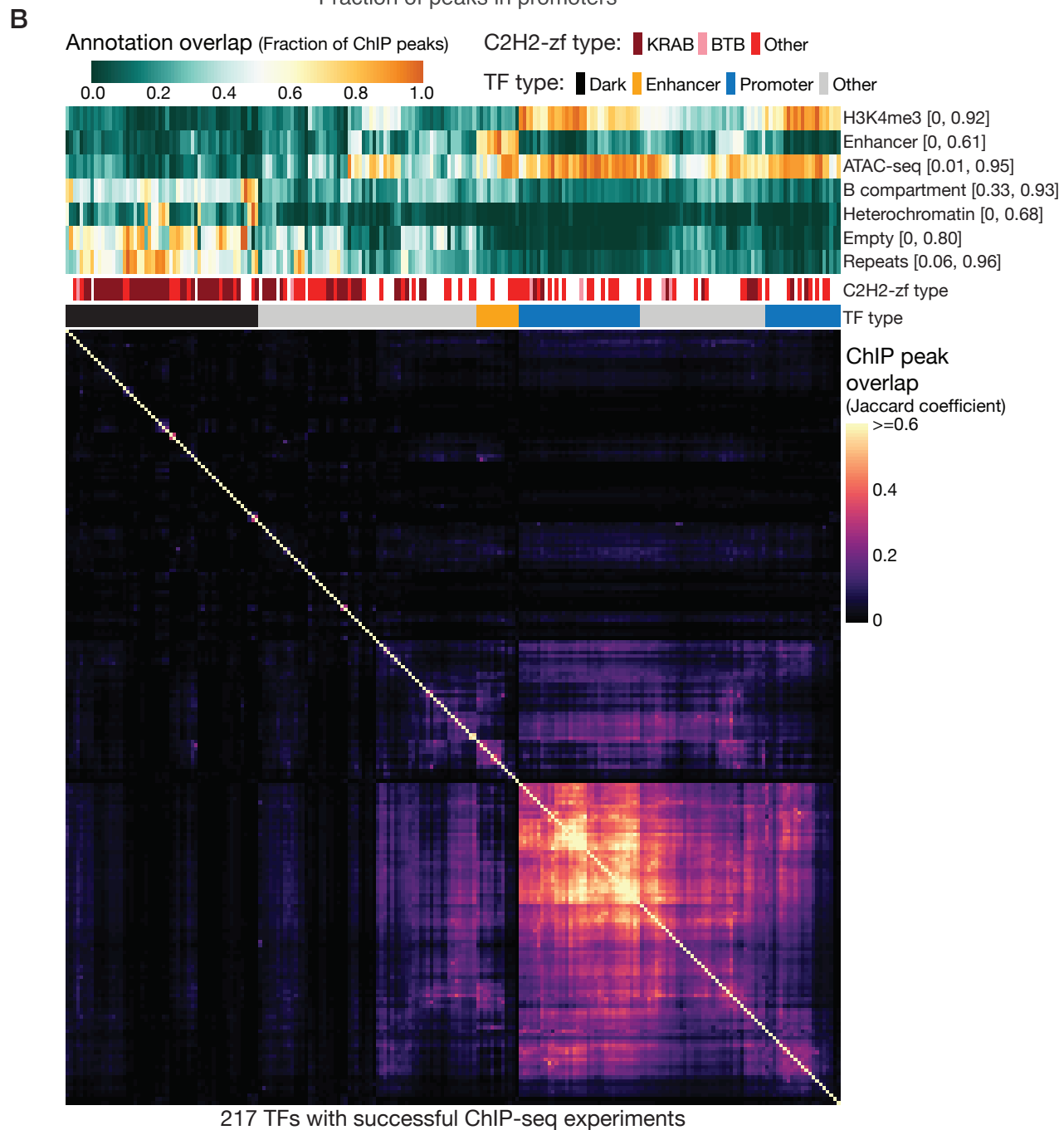
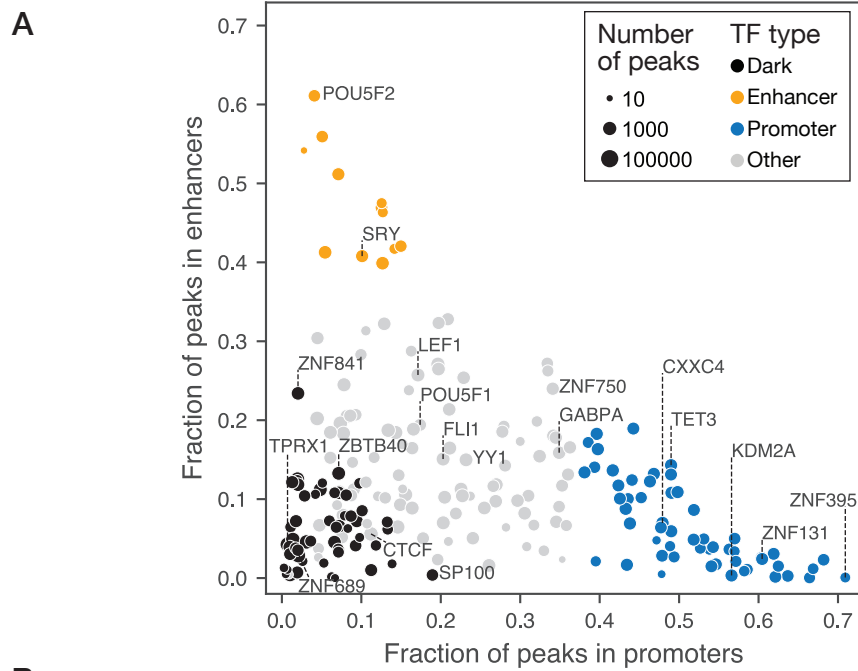


Figure 2. Overlapping *in vivo* binding sites of 217 TFs with each other and with various genomic regions. (A) Fraction of ChIP-seq peaks in protein-coding promoters (x-axis) and HEK293 enhancers (y-axis). Point sizes are proportional to the number of peaks for each TF (log scale). (B) *Bottom (square) heatmap*: Jaccard similarity coefficient between ChIP-seq peaks of all TF pairs. *Top heatmap*: Fraction of ChIP-seq peaks falling within genomic regions, as indicated, and other properties of the TFs. Fractions are scaled to fit in [min, max] range across the TFs for better visualization, as indicated in the right. TF ordering is determined by hierarchical clustering with Ward linkage and Euclidean distance, using the tracks 'H3K4me3', 'ATAC-seq', 'B compartment', 'Empty' + 'Heterochromatin', 'Repeats', 'CpG', 'Protein-coding promoters', 'H3K27ac' (the last three not shown), along with the one-hot encoded 'TF type' to aid in illustration.

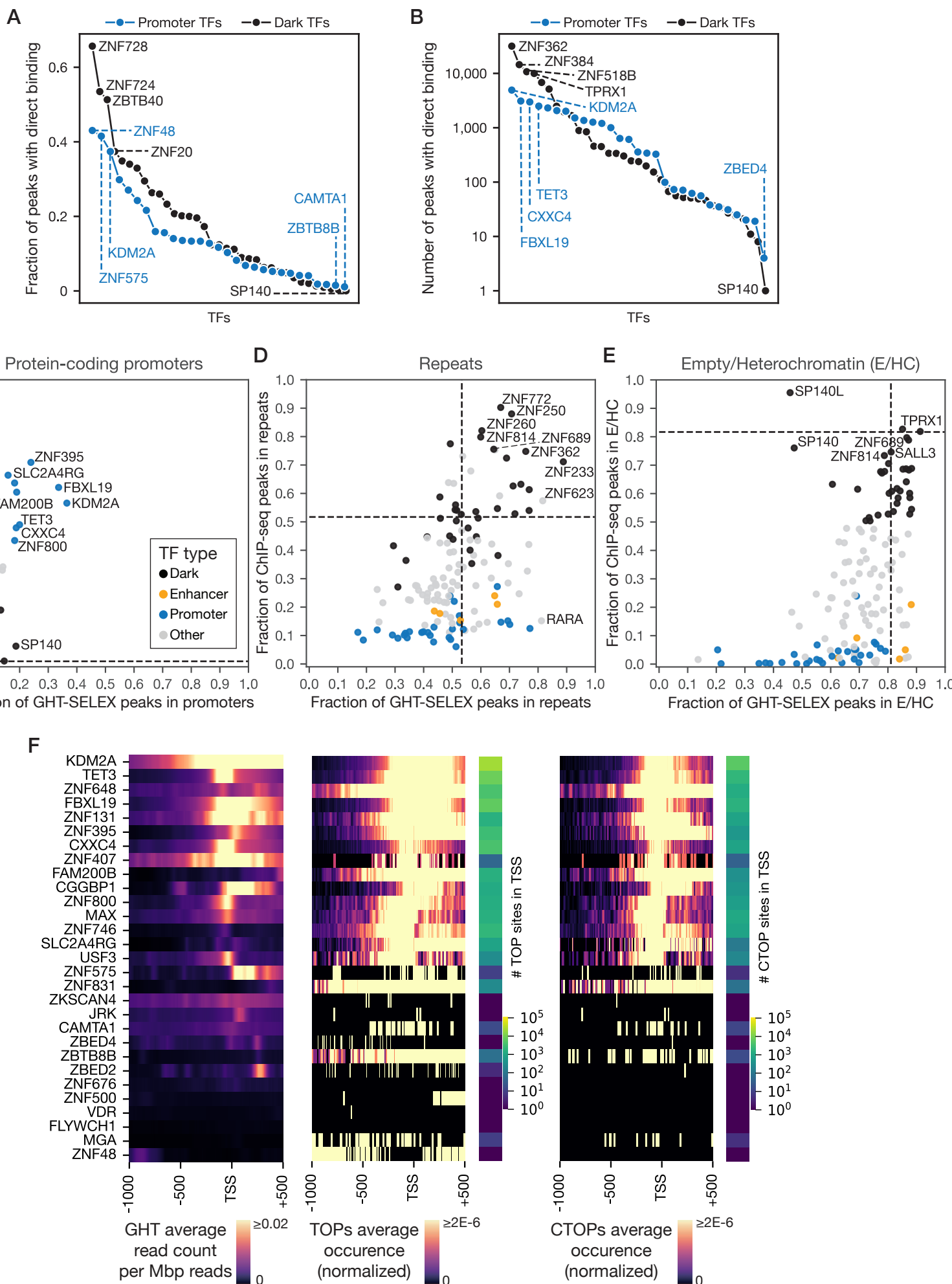


Figure 3. Characteristics of Promoter TFs, Enhancer TFs, and Dark TFs interaction with specific genomic sites. Fraction (**A**) and absolute number (**B**) of peaks with direct binding (i.e. TOP sites) for Promoter TFs and Dark TFs. TFs are sorted to compare distributions. The denominator for (A) is the total number of ChIP peaks at the same optimized threshold. (**C, D, E**) Fraction of GHT-SELEX (x-axis) and ChIP-seq (y-axis) peaks falling in the specified genomic regions (protein-coding promoters, repeats, and empty or heterochromatin), using the peaks at the universal threshold. Dashed lines show the expected fraction if peaks were distributed at random. (**F**) Density of GHT-SELEX signal (left), TOP sites (middle), and CTOP sites (right) by position relative to TSS of protein-coding promoters, for 29 Promoter TFs that have available GHT-SELEX data. Intensity of heatmaps for TOPs (middle) and CTOPs (right) have been normalized by the total number of PWM hits (of TOPs and CTOPs, respectively) in promoters (shown at the right of each heatmap).

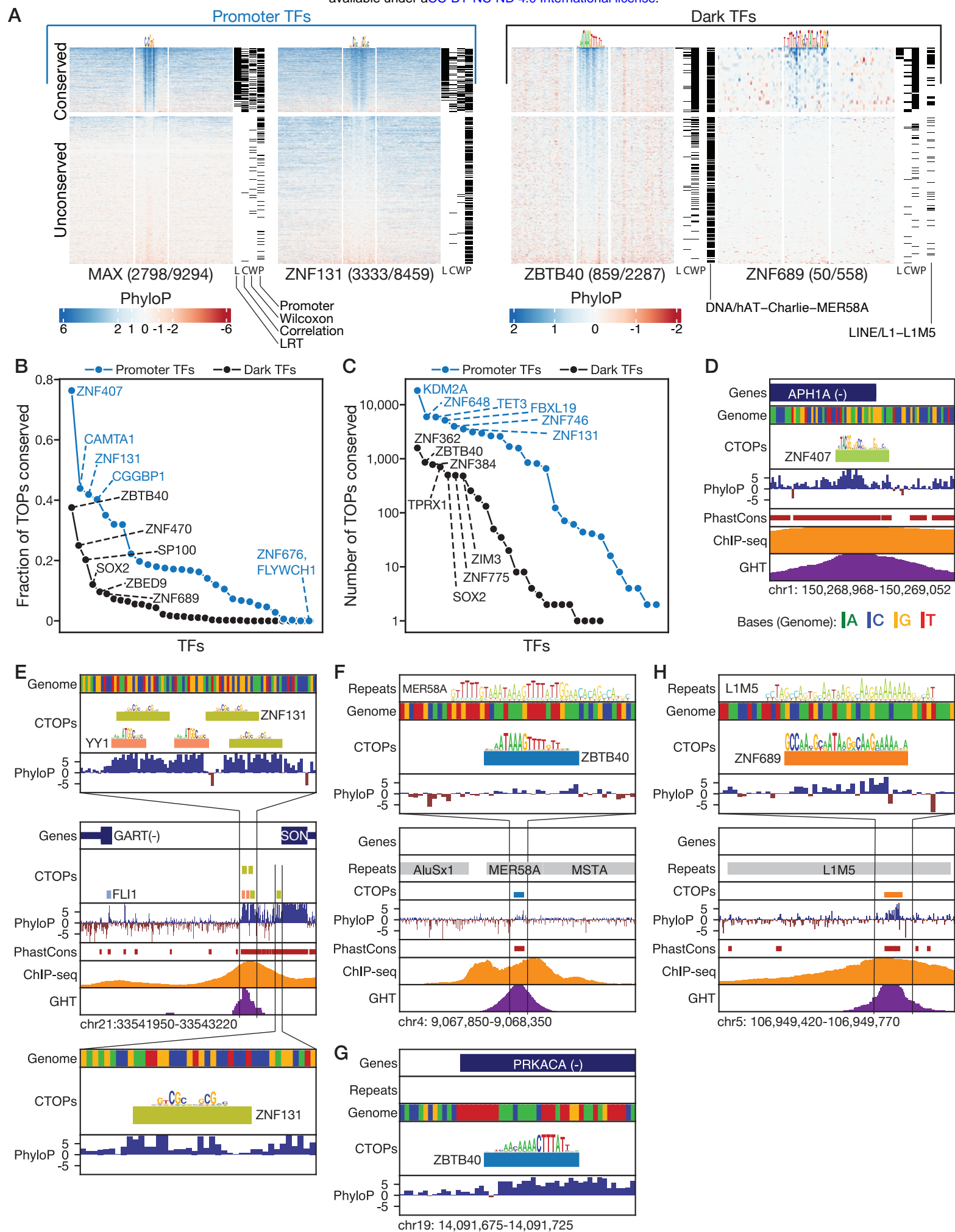
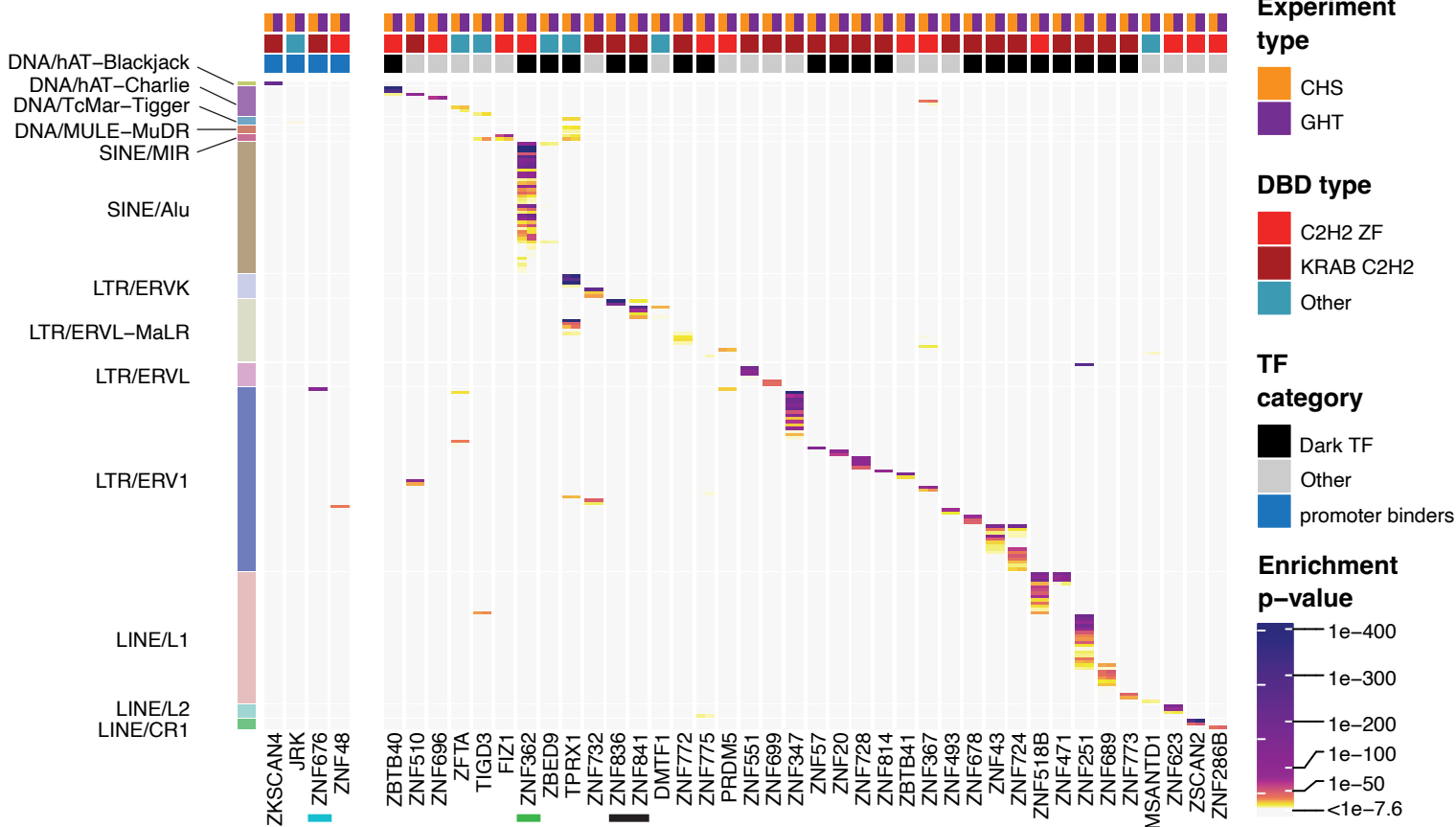


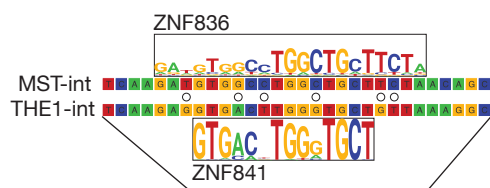
Figure 4. Conservation patterns of sequence-dependent TFs' target sites (TOPs).

(A) Heatmaps of FDR-corrected phyloP scores across the TOP sites (rows), split into top and bottom segments that contain conserved and unconserved sites. Bars to the right indicate which tests of conservation are satisfied (Likelihood-ratio, Correlation, Wilcoxon), along with overlaps with promoters (P) and specific repeat families if applicable. 100 bp segments are shown with the PWM hit in the middle. Blue/positive phyloP indicates purifying selection, and red/negative phyloP values represent diversifying selection. **(B, C)** Fraction **(B)** and absolute number **(C)** of TOPs that are conserved, for Promoter TFs and Dark TFs, sorted to compare distributions. **(D, E, F, G, H)** Genome track displays of CTOP sites for ZNF407 **(D)**, ZNF131 and YY1 **(E)**, ZBTB40 at a *hAT/Charlie* (MER58A) element **(F)** and its most-conserved TOP (at the PRKACA promoter) **(G)**, ZNF689 at an L1M5 element **(H)**. The Dfam¹⁰³ repeat model sequence logo is also shown for MER58A **(F)** and L1M1 **(H)**.

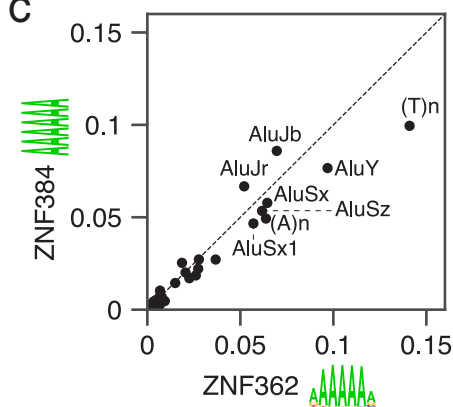
A



B



C



D

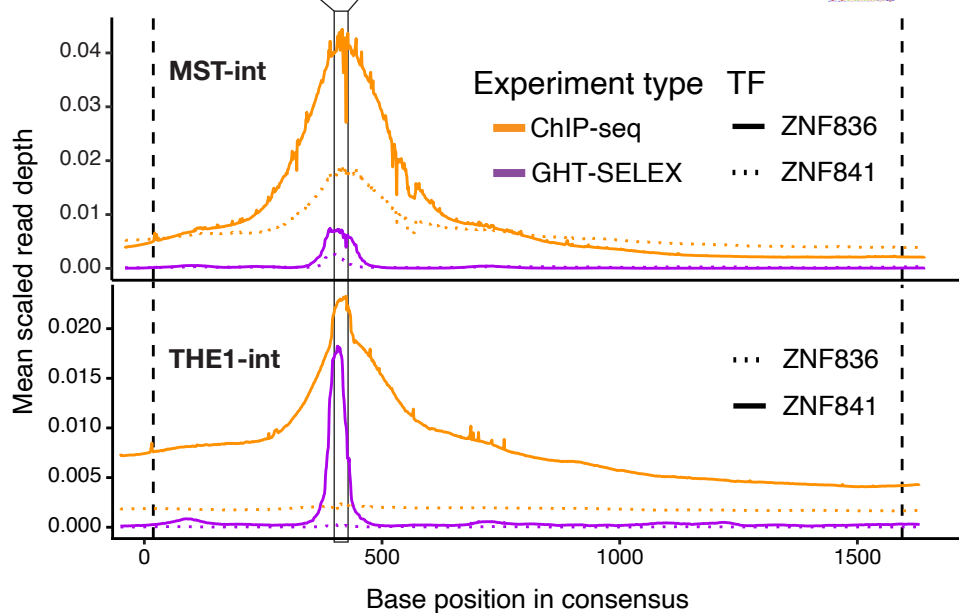
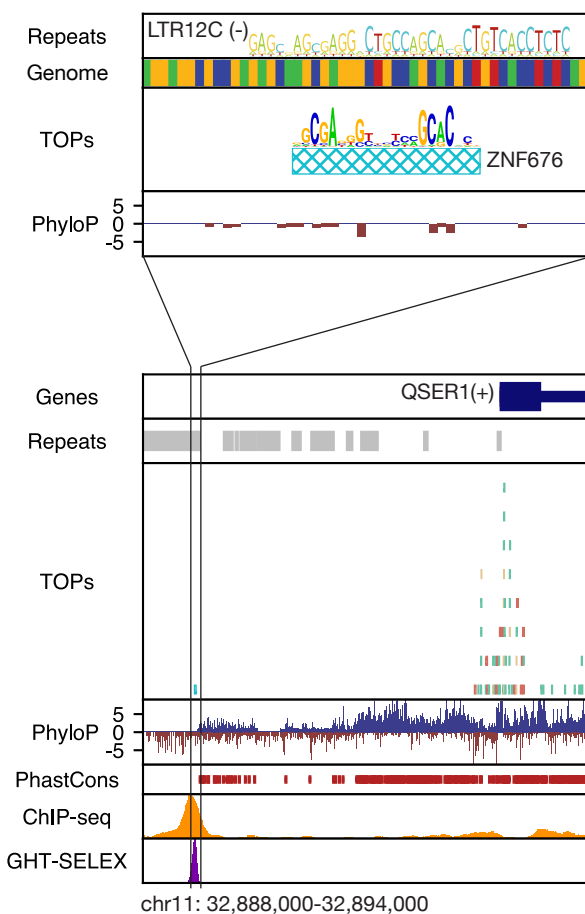


Figure 5. Enrichment pattern of transposable elements in TFs' TOPs. (A) Heatmap of $-\log_{10}$ p-values for TFs (x-axis) that are enriched for binding specific TE families (y-axis). Labels show superfamily/family. (B) Binding of paralogous TFs, ZNF836 and ZNF841, to a homologous region in the two related LTR families, MSTA-int and THE1-int. Bottom plot shows the average ChIP-seq and GHT-SELEX signal (i.e. read count) across all the instances of MST-int and THE1-int aligned to their consensus. (C) Fraction of TOP sites in various repeat elements for two poly-A binding TFs ZNF362 and ZNF384. (D) An example of the Promoter TF ZNF676 binding site targeting an unconserved LTR12C sequence.

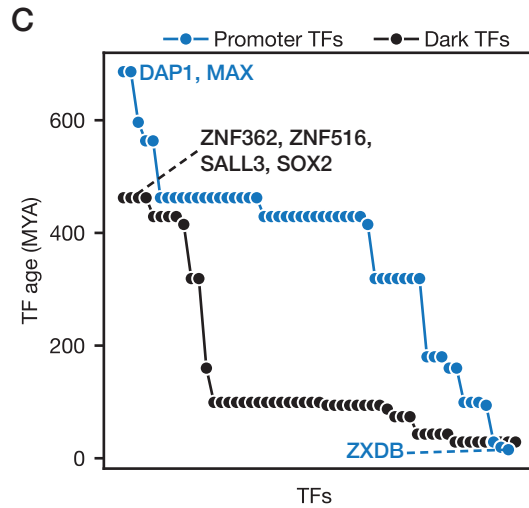
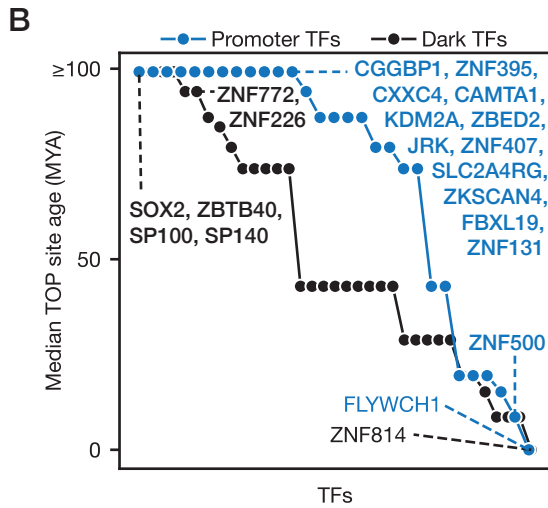
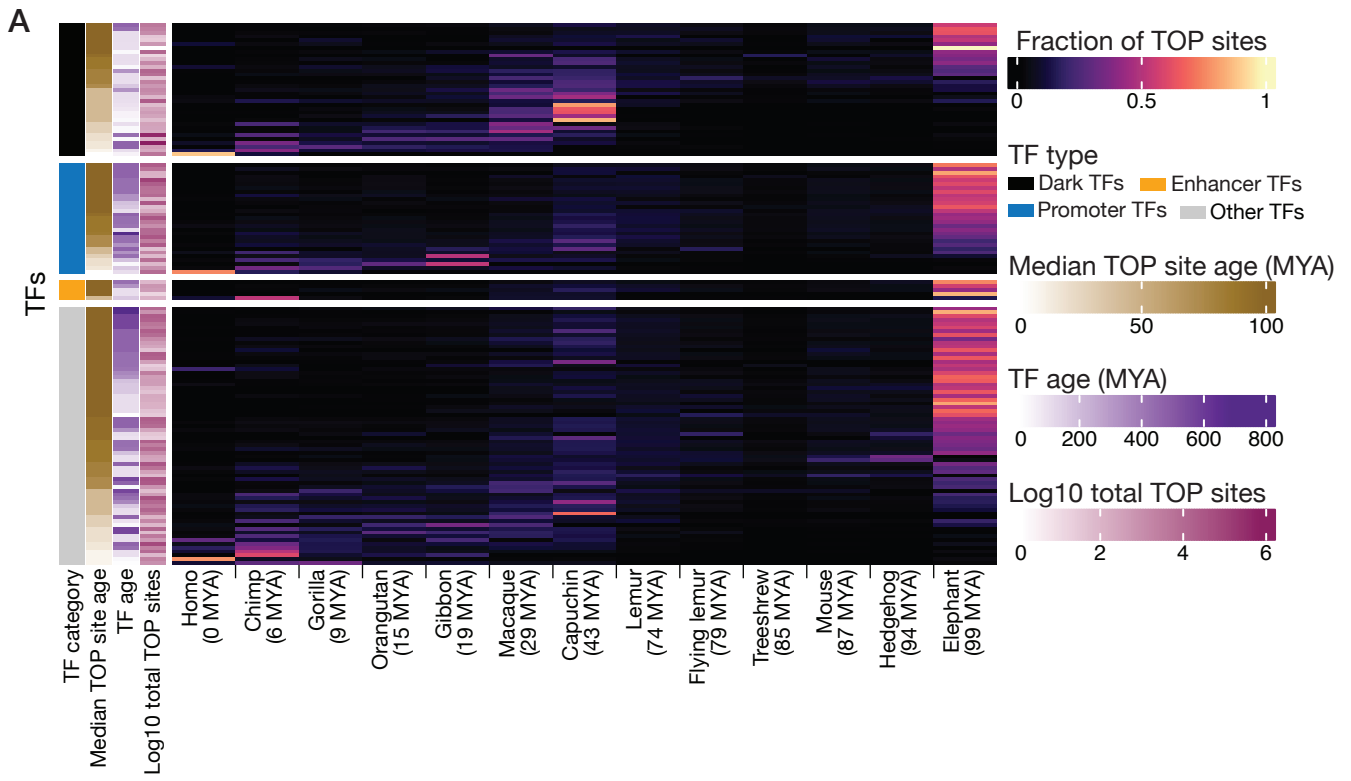


Figure 6. Age distribution of TOPs and their corresponding TFs. (A) Heatmap showing the fraction of TOP sites for each TF dating to different mammalian clades in the human lineage, along with information about the TF category, median age of TOP sites and TFs (million years ago, MYA), and log₁₀ of total TOP sites. (B, C) Sorted median age of the TOP sites (B) and the age of the TFs (C) are compared for Dark TFs and Promoter TFs.

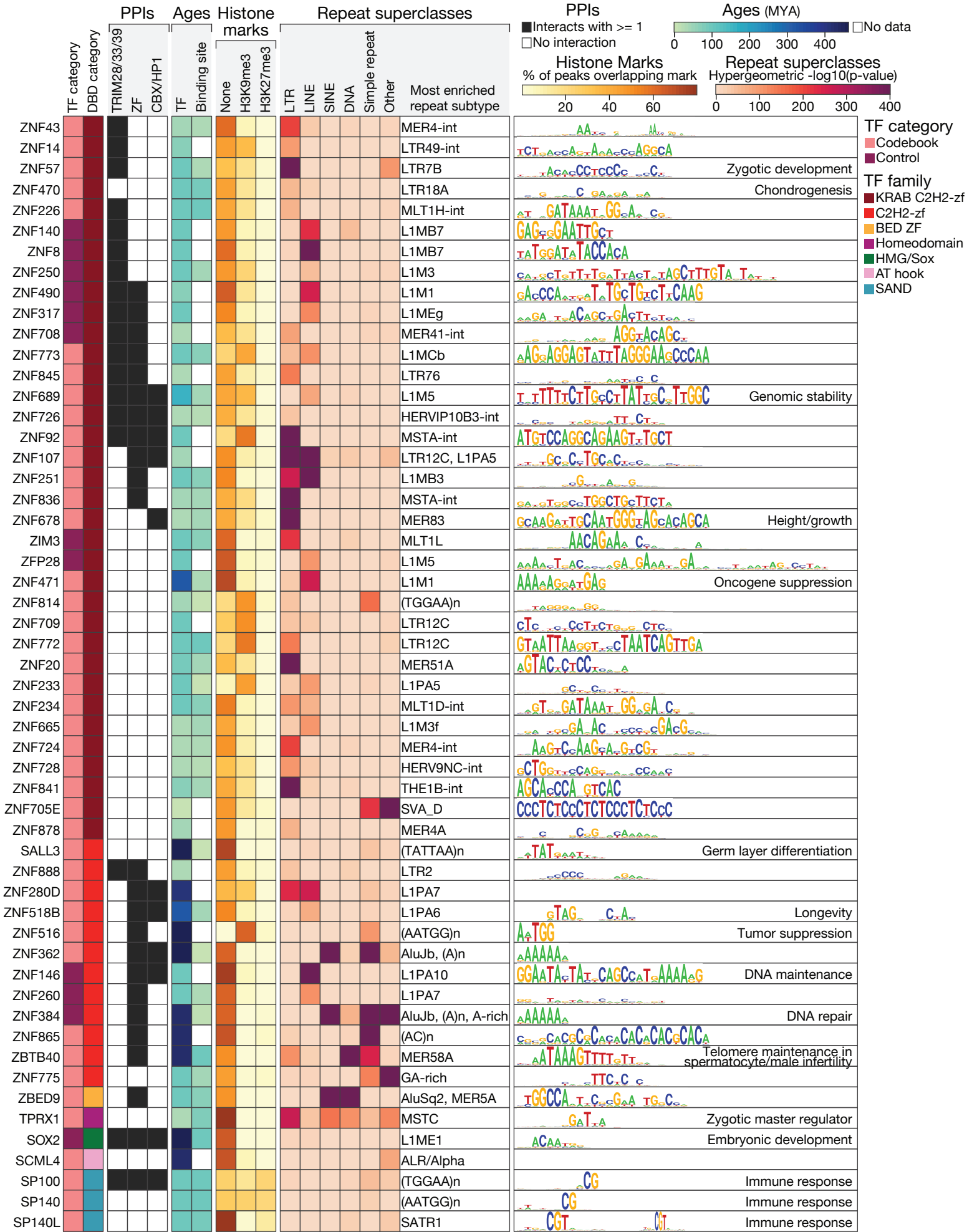


Figure 7. Consolidated functional information for Dark TFs. Compiled protein-protein interactions (PPIs)⁶³ mostly supported by two independent lines of support and grouped into three categories of TRIM28/33/39 interactions, zinc-finger (ZF) protein interactions, and CBX/HP1 interactions are shown at left. Median binding site age was calculated for TOP sites, only for the TFs with available GHT-SELEX data, shown along with the age of the TF. The fraction of ChIP-seq peaks (using the universal threshold) overlapping with H3K9me3 and H3K27me3 histone marks and with the ChromHMM “empty” state (None) are shown in the middle. For the repeat, in each superclass, the enrichment score (-log(p-value) hypergeometric test) for the most enriched repeat element within that superclass is plotted as a heatmap, and the most enriched repeat subtype across all the superfamilies is mentioned beside. The expert-curated sequence logos are displayed to the right (except for ZNF280D and SCML4 which did not produce any approved PWM), along with the corresponding phenotype for any TF with known biological function through literature review (in the same block).

597 METHODS

598 **Plasmids.** The Codebook project design is described elsewhere³¹. Putative TFs were
599 those from Lambert et al.³, with TFs we had already attempted as part of ENCODE
600 removed. We attempted ChIP-seq analysis of all Codebook TFs. We designed full-
601 length ORFs for synthesis (BioBasic.ca) and used conventional restriction cloning to
602 insert them into Flp-In destination plasmid pTH13195 (a modified pDEST
603 pcDNA5/FRT/TO-eGFP vector), which places ORFs under the control of a tet-on, CMV-
604 driven promoter, with an N-terminally EGFP tag. We obtained the 58 controls
605 independently^{27,33} and cloned them into pTH13195 by the same process. See
606 accompanying manuscript³¹ for the sequences and other information about the inserts.

607 **Cell line production.** We used a previously established protocol for creating individual
608 cell lines for each TF³⁴. In brief, we cultured and maintained HEK293 Flp-In T-REx cells
609 (Invitrogen) in Dulbecco's modified Eagle's medium with 10% fetal bovine serum and
610 antibiotics. We created individual cell lines for each TF by flp recombination, in which
611 individual destination plasmids were co-transfected into Flp-In T-REx 293 cells together
612 with the pOG44 Flp recombinase expression plasmid using FuGene (Roche,
613 11814443001). We then selected cells for FRT site-specific recombination into the
614 genome using selection media containing hygromycin (200µg/ml) for 1 to 4 weeks. For
615 each TF cell line, we confirmed expression of EGFP by fluorescent microscopy after 24
616 hours of Doxycycline treatment (1ug/ml), at which point 10M-20M cells were used for
617 downstream experiments.

618 **Chromatin immunoprecipitation.** We fixed ~20M cells on 15cm plates using 1%
619 paraformaldehyde for 10 min on ice followed by 10 min quenching with 0.125 M glycine.
620 We washed fixed cells twice with cold PBS, scrape collected, pelleted, flash froze, and
621 stored the cells at -80C. Upon completion of cell collection for a panel of TFs, we
622 thawed cell pellets on ice, lysed them as previously described⁸³, and sonicated them
623 using a BioRuptor to shear chromatin. We then trapped the protein-DNA complexes on
624 Dynabeads using an anti-EGFP antibody (Ab290, Abcam). Following wash, crosslink
625 reversal, and elution, we assessed the size and concentration of DNA fragments with an
626 Agilent bioanalyzer and Qubit, prior to sequencing.

627 **Library preparation and sequencing.** DNA library preparation and sequencing was
628 performed at three different facilities (Memorial Sloan Kettering Cancer Center, SickKids
629 Hospital in Toronto, and the Donnelly Centre at the University of Toronto) over a period
630 of four years. The facilities prepared DNA libraries using the NEBNext® Ultra™ II DNA
631 Library Prep Kit for Illumina. Samples were paired end sequenced (50-150bp) with a
632 target depth of 20M reads per sample.

633 **ChIP-seq data processing steps.** Read mapping: We mapped raw ChIP-seq reads to
634 the human genome build hg38 with *bowtie2*⁸⁴ (options: *--very-sensitive*, and *--no-unal*).
635 We used Samtools^{85,86} (options: *-q 30*, and *-F 1548*) to remove reads that were
636 unmapped, failed platform/vendor quality checks, were PCR duplicates, or had a
637 mapping quality <30. Peak calling: We created sample-specific background models
638 following a procedure established previously²⁷, with minor modifications. Specifically, for

639 each pull-down experiment, we pooled reads from different control experiments together
640 in a manner that maximizes the similarity of the pooled dataset and the background
641 signal observed in that specific pull-down experiment, while ensuring high coverage. To
642 do so, we first identified genome regions that show high background signal in at least
643 one control experiment, by performing peak calling directly on control experiments with
644 MACS2 (options: p-value < 0.001, and *--nomodel*)^{39,87,88}. We pooled these “background
645 hotspots” from all control experiments, and merged those whose summits were within
646 50 bp of each other, to create a unified set of hotspots. Then, in each control
647 experiment, we calculated the number of reads overlapping each of the hotspots from
648 the unified set, resulting in a read count matrix (with hotspots as rows and control
649 experiments as columns). Similarly, for each pull-down experiment, we calculated the
650 number of reads overlapping each of the hotspots, which we then used as the response
651 variable in a non-negative Poisson regression, with the matrix described above as the
652 set of independent variables. This regression results in a set of non-negative
653 coefficients, representing a weighted mix of the control experiments that reconstructs
654 the read count profile of the pull-down experiment across the hotspot regions as closely
655 as possible. We then pooled the BAM files from the control experiments, by sampling a
656 number of reads from each file that is proportional to this experiment-specific coefficient,
657 to create a pulldown-specific background file, which we subsequently used for peak
658 calling on the pull-down dataset using MACS2^{39,87,88}.

659 **ChIP peak replicate analysis and merging.** For each TF with one or more replicate,
660 we calculated the Kulczynski II similarity metric for each pair of replicates (**Figure S1**).
661 We used the Kulczynski II metric in place of Jaccard as it is less affected by the uneven
662 size of the peak sets. We additionally calculated the Kulczynski II similarity metric for
663 each pair of mismatched replicates (i.e., with TF identities permuted). Based on the
664 distributions of “approved” experiment replicates and mismatch replicates, we defined a
665 Kulczynski II threshold of 0.4 as the separating value for those two distributions (**Figure**
666 **S1**). For TFs with “not approved” experiments (i.e., two ChIP-seq experiments did
667 produce a reliable motif) we retained 36 (plus two controls) that achieved a Kulczynski II
668 value >0.4 for inclusion in downstream analyses.
669 To generate a single peak set for each transcription factor, we merged the peak data
670 from all successful experiments by merging overlapping peaks from one or more
671 replicates using BEDTools⁸⁹ merge to generate new, wider peaks, with the sum of
672 component peak -log(p-values) assigned as the new peak score, and center of mass of
673 summits as the new peak summit. By default, we employed a peak cutoff of $P < 10^{-10}$
674 (MACS2). Modulation of thresholds is described below.

675 **ATAC-seq experiment and data analysis.** We performed ATAC-seq in HEK293 cells,
676 in four replicates, as described⁹⁰. Briefly, 50,000 viable HEK293 cells were pelleted (500
677 RCF at 4°C for 5 min). After removing the supernatant, the cells were lysed in 50 µl of
678 cold ATAC–resuspension buffer (RSB) containing 0.1% NP40, 0.1% Tween 20, and
679 0.01% digitonin by pipetting up and down three times followed by 3 min incubation on
680 ice. The lysate was then washed out with 1 ml of cold ATAC-RSB containing 0.1%
681 Tween 20 and the nuclei were pelleted at 500 RCF for 10 min at 4°C. 50µl of
682 transposition mixture (25 µl of 2× TD buffer, 2.5 µl of transposase, 16.5 µl of PBS, 0.5 µl
683 of 1% digitonin, 0.5 µl of 10% Tween 20, and 5 µl of H₂O) was added to each pellet,

684 mixed well by pipetting up and down, and incubated for 30 min at 37C. Tagment DNA
685 TDE1 Enzyme and Buffer Kit (Illumina) was used for this step. The tagmented DNA was
686 then purified with DNA Clean and Concentrator kit (Zymo Research) in 21 ul of elution
687 buffer. DNA amplification and barcoding were performed using Nextera DNA Library
688 Prep kit (Illumina) and NEB barcoding oligos. Subsequent sequencing was performed at
689 the Donnelly Center sequencing facility using 100bp paired end sequencing at 60M
690 reads per sample. Adapter sequences were first trimmed using *cutadapt*. The resulting
691 reads were then mapped to the human genome (hg38) using *bowtie2*⁸⁴, followed by the
692 creation of BAM files using *samtools*^{85,86} *view*, and sorting with *samtools sort*. Peak
693 calling was performed on the sorted BAM reads by running *macs2*^{39,87,88} *callpeak* with
694 the options *-f BAMPE*, *-g hs*, *-B*, and *-q 0.01*. Finally, to generate a single peak file for
695 HEK293 open chromatin, all the peak sets were merged using *bedtools*⁸⁹ *merge*.

696 **Chromatin state analysis.** We obtained chromatin states by training a ChromHMM⁴¹
697 with ten states (see **Figure S2B**) on marks H3K4me1, H3H4me3, H3K36me3,
698 H3K27ac, H3K9me3, and K3K27me3, collected in HEK293 cells by ENCODE⁹¹, plus
699 ATAC-seq and CTCF peaks from HEK293 cells generated as part of this project. We
700 also employed promoter regions derived from the GENCODE annotation (release 44)⁹²
701 for our analyses (-1000 to +500). For Hi-C B compartment annotations, we labeled
702 genomic regions with a Hi-C first eigenvector value less than 0.4 in ENCODE data for
703 HAP1 cells⁹³, comprising 65% of the genome.

704 **Overlap of ChIP-seq peaks between all pairs of TFs.** Jaccard similarity is taken as $O / (N1+N2-O)$
705 where O is the number of intersecting peaks and $N1$ and $N2$ are the size of
706 each set. We utilized BEDTools⁸⁹ to calculate overlaps. To prevent miscounting of the
707 cases in which one peak in one set overlaps with multiple peaks in another set, we used
708 the average of overlapping peaks ($O = (O1+O2)/2$ where $O1$ is the number of peaks in
709 set 1 overlapping with any peak in set 2 and vice versa) to calculate the intersection in
710 Jaccard. The same methodology was used to calculate the overlap of ChIP-seq peaks
711 with the chromatin tracks.

712 **Selection of universal ChIP-seq and GHT-SELEX thresholds.** We calculated the
713 Jaccard similarity from all 137 pairs of TFs with ChIP-seq and GHT-SELEX data, using
714 the merged ChIP-seq peaks. We performed a grid search for all TFs simultaneously,
715 sampling ChIP-seq P-value and GHT-SELEX cutoffs (determined by selecting different
716 “knee” values⁹⁴ in the graph of sorted enrichment coefficients⁴³), to identify a pair of
717 thresholds that maximize median Jaccard. Two ChIP-seq cutoff (10^{-10} and 10^{-20}) yielded
718 an almost identical maximum; we chose 10^{-10} as it includes a larger number of peaks. A
719 corresponding knee threshold of 30 emerged for the GHT-SELEX knee-based cutoff.

720 **Derivation of TOP sites.** To define binding sites supported by ChIP-seq, GHT-SELEX,
721 and PWM hits, we optimized the cut-offs of all three to maximize the overlap between all
722 three data types. We first sorted the peaks based on their statistical scores, i.e., merged
723 p-values for ChIP-seq peaks, cycle enrichment coefficient for GHT-SELEX peaks (see
724 accompanying manuscript⁴³), and sum-of-affinities for clusters of PWM hits with a p-
725 value < 0.001 (from MOODS⁹⁵), merged with neighboring hits in the case of having a
726 distance less than 200 bp. Then, for different values of N , we took the top N peaks and

727 calculated the overlap (measured as the Jaccard index; intersection of all three divided
728 by the union of all) using the top N ChIP-seq peaks, top N GHT-SELEX peaks, and top
729 N merged PWM hits. The N that maximizes the Jaccard overlap was taken as the
730 optimized threshold, and the overlap of all three sets at this threshold (N) is referred to
731 as triple overlap or “TOP” sites.

732 **Analysis of purifying selection and classification as conserved and unconserved**
733 **binding sites.** We extracted phyloP scores⁴⁷ for each PWM hit, and for flanking regions
734 of equal length (for a total of 100 bp including the PWM hit and its flanks) from the 241
735 eutherian mammal Zoonomia alignment⁹⁶ using DeepTools⁹⁷. We excluded PWM hits
736 overlapping with ENCODE Blacklist sites⁹⁸ or protein coding sequences, due to the
737 skew in phyloP scores caused by codons. All phyloP scores reported here are FDR-
738 corrected. We conducted three tests to classify PWM hits as ‘conserved’ or
739 ‘unconserved’:

740 1, LRT (Likelihood-Ratio Test): This test scores the likelihood that the phyloP
741 scores are driven by the PWM information content (IC) at each base position in the
742 PWM. For each TF, we created a scoring model that represents the relationship
743 between the phyloP scores at a PWM hit, and the information content at each base
744 position of the PWM. This model is an $l \times 1$ vector, where l is the length of the motif. To
745 derive this vector, we first took the correlation of phyloP scores at each base position
746 within the PWM hit to the IC at that position, for each PWM hit in the TOP dataset. We
747 then selected the 100 PWM hits with the highest correlation and calculated the standard
748 deviation ($\bar{\sigma}$) of the phyloP score at each position of these 100 PWM hits. If a position
749 has an invariant phyloP score (i.e. $\bar{\sigma} = 0$), the $\bar{\sigma}$ at this position was replaced with a 1.
750 As a null model, an IC value of 0 was assumed at each position, and the same $\bar{\sigma}$ values
751 as the phyloP model. The LRT statistic for each PWM hit \bar{m} was then taken according to
752 the equation:

753 **Equation 1:**
$$L(m) = -2 \times \left(\sum_{i=1}^l \frac{P_i^2}{\sigma_i^2} - \sum_{i=1}^l \frac{(P_i - I_i \times C)^2}{\sigma_i^2} \right)$$

754 Where \bar{P}_i represents the phyloP score of position \bar{i} in PWM hits \bar{m} , $\bar{\sigma}_i$ is the standard
755 deviation of the phyloP model at position \bar{i} , and \bar{I}_i is the IC of the PWM at position \bar{i} . The
756 IC value is first multiplied by a coefficient \bar{C} , which is the linear regression coefficient
757 describing the relationship between the position-wise phyloP means of all of a TF’s
758 TOPs and position-wise PWM IC. It therefore has units phyloP/bits and converts \bar{I}_i to a
759 phyloP score. Based on manual inspection of phyloP patterns across TOPs at different
760 test statistic thresholds, we selected a threshold of $L < -10$ to be considered “conserved”
761 according to this test, which manual inspection indicated is conservative.

762 2, Correlation Test: For each TF, we permuted the position-wise IC of the PWM
763 using the *permute* R package, up to a maximum of 1,000 unique permutations (not
764 every PWM has 1,000 unique permutations). We then took the Pearson correlation of
765 each of these permuted PWM IC vectors using the phyloP scores of 1,500 randomly
766 selected PWM hits from the unfiltered PWM scan results (or fewer, if there are <1,500
767 total hits). This resulted in a maximum of 1,500,000 correlations per TF, dependent
768 upon the number of unique PWM IC permutations and number of PWM hits. We used
769 these correlation values as a null distribution, converted to Z scores, and determined a

770 threshold correlation value corresponding to an alpha of 0.05; this threshold was chosen
771 manually. PWM hits with a Pearson correlation to the unpermuted motif IC values
772 greater than this threshold were considered to have passed this test.

773 **3, Wilcoxon Test:** For each TOP site, we performed two Wilcoxon tests, one
774 comparing values in the PWM hits to those in the 25bp downstream flank, and the same
775 for the 25bp upstream flank. All p-values were FDR corrected, and an FDR-corrected p-
776 value less than a threshold of 0.1 for both flanks was considered a positive.

777 A TOP PWM hit was considered conserved if it passed one of the three
778 conservation tests above, and had at least one site with an FDR-corrected phyloP score
779 > 1 (i.e. corresponding to an FDR-corrected p-value < 0.1).

780 **Determination of binding site ages.** We used `halLiftover`⁹⁹ to map TOP PWM hits to
781 syntenic loci in all other genomes in the Zoonomia 241-mammal alignment⁹⁶, including 9
782 reconstructed genomes ancestral to human, and calculated the alignment's % identity to
783 the human sequence. We then assigned an age using multiple criteria. For the method
784 used in the main text, we identified the oldest ancestral genome with a gapless
785 alignment of any % identity to the human PWM hit. **Supplementary Figure S6** shows
786 alternative schemes for age inference, including the oldest extant *species* with a
787 gapless alignment to a human TOP (at various threshold % identities), or the oldest
788 *clade* wherein 60% of species have a gapless alignment to the human TOP (at various
789 threshold % identities). We acquired the age of each clade from `TimeTree`⁵⁸.

790 **Determination of TF ages.** To infer TF ages, we acquired all vertebrate ortholog
791 annotations and ortholog quality statistics for each TF from `Ensembl`¹⁰⁰, and ages of
792 each pair of species from `TimeTree`⁵⁸. The age of a TF was taken as the oldest ortholog
793 annotated as having a 1-1 relationship with human, or having a gene order conservation
794 (GOC, a metric of synteny) score ≥ 50 and classification as a high-confidence ortholog
795 by `Ensembl`.

796 **Repeat enrichment.** To calculate enrichment for each TF and repeat pair, we identified
797 the intersections of the peak summits from ChIP-seq peaks (or TOPs) and the middle
798 position of GHT-SELEX TOPs with the 2022-10-18 version of the UCSC Genome
799 Browser RepeatMasker track¹⁰¹. The enrichment significance between GHT-SELEX and
800 ChIP-seq TOPs and each repeat family was calculated using Fisher's Exact Test
801 implemented in `SciPy`¹⁰². The contingency table took the form of:
802

	TF y summit +	TF y summit -
Repeat x +	# of repeat x and TF y summit intersections	# of repeat x bases with no TF y summit intersections
Repeat x -	# of TF y summits with no repeat x intersections	# of repeat-annotated bases, excluding # of repeat x bases and TF y summits

803
804

805

806 DATA AVAILABILITY

807 The sequencing raw data for the experiments have been deposited into the SRA
808 database under identifiers PRJEB78913 (ChIP-seq), PRJEB76622 (GHT-SELEX), and
809 PRJEB61115 (HT-SELEX). Genomic interval information generated for the ChIP-seq
810 and GHT-SELEX have been deposited into GEO under accessions GSE280248 and
811 GSE278858, respectively. Information on constructs, experiments, analyses, processed
812 data, comparison tracks, with many accessory files and browsable results is available at
813 <https://codebook.ccb.utoronto.ca>. Larger collection of motifs generated for these
814 experiments in an accompanying study³⁶ can be browsed at <https://mex.autosome.org>
815 and downloaded at <https://doi.org/10.5281/ZENODO.8327372>.

816 ACKNOWLEDGEMENTS

817 We thank the IT Group of the Institute of Computer Science at Halle University for
818 computational resources and Maximilian Biermann for valuable technical support.
819 This work was supported by the following:

- 820 • Canadian Institutes of Health Research (CIHR) grants FDN-148403, PJT-
821 186136, PJT-191768, and PJT-191802, and NIH grant R21HG012258 to T.R.H.
- 822 • CIHR grant PJT-191802 to T.R.H. and H.S.N.
- 823 • Natural Sciences and Engineering Research Council of Canada (NSERC) grant
824 RGPIN-2018-05962 to H.S.N.
- 825 • Russian Science Foundation grant 20-74-10075 to I.V.K.
- 826 • Russian Science Foundation grant 24-14-20031 to F.A.K.
- 827 • Swiss National Science Foundation grant (no. 310030_197082) to B.D.
- 828 • Marie Skłodowska-Curie (no. 895426) and EMBO long-term (1139-2019)
829 fellowships to J.F.K.
- 830 • NIH grants R01HG013328 and U24HG013078 to M.T.W., T.R.H., and Q.M.
- 831 • NIH grants R01AR073228, P30AR070549, and R01AI173314 to M.T.W.
- 832 • NIH grant P30CA008748 partially supported Q.M.
- 833 • Canada Research Chairs funded by CIHR to T.R.H. and H.S.N.
- 834 • Ontario Graduate Scholarships to K.U.L and I.Y.
- 835 • A.J. was supported by Vetenskapsrådet (Swedish Research Council)
836 Postdoctoral Fellowship (2016-00158)
- 837 • The Billes Chair of Medical Research at the University of Toronto to T.R.H.
- 838 • EPFL Center for Imaging
- 839 • Institutional funding from EPFL
- 840 • Resource allocations from the Digital Research Alliance of Canada

841

842 **SUPPLEMENTAL TABLES AND DOCUMENTS**

843 **Table S1. Overview of the tested TFs.** This table lists the TFs that were tested in this
844 study using ChIP-seq.

845 **Table S2. List of all the ChIP-seq experiments.** This table lists the ChIP experiments,
846 their approval status, and related produced files.

847 **Table S3. List of 217 TFs.** This table lists the TFs with either “approved” ChIP-seq
848 experiments or significant overlap between replicates, together with the list of ChIP-seq
849 samples used in “merged” peaks for each TF.

850 **Table S4. Binding category of the TFs (i.e. Promoter TFs, Enhancer TF, Dark TF,
851 and Others) for 217 TFs with successful ChIP-seq experiments.** For the TFs with
852 available GHT-SELEX data (hence TOP sites, the number of optimized ChIP-seq peaks
853 (i.e. Triple peaks), number of TOP ChIP-seq peaks, fraction of direct binding sites (i.e.
854 #TOP peaks divided by #Triple peaks), number of TOPs, number of CTOPs, fraction of
855 conserved TOPs (i.e. #CTOPs / #TOPs), and the median age of the TOP sites are also
856 included. Note that the number of TOP ChIP-seq peaks might be different (less) than
857 TOPs (referring to triple-overlap PWM hits), since each peak might comprise multiple
858 PWM hits.

859 **Table S5. Consolidated functional information for Dark TFs.** This table provides the
860 data underlying **Figure 7** including the references in the literature.

861 **Document S1. Heatmaps of conservation/phyloP score across TOPs for 137 TFs.**
862 This document provides the same analysis as **Figure 4A**, for all TFs of the study,
863 heatmaps of phyloP scores in PWM hits (middle column) and flanking sequences of
864 tops are displayed. Bars to the right indicate which tests of conservation are satisfied
865 (Likelihood-ratio, Correlation, Wilcoxon), along with overlaps with promoters (P) and
866 specific repeat families if applicable.

867

868 REFERENCES

- 869 1. Partridge, E.C. *et al.* Occupancy maps of 208 chromatin-associated proteins in one
870 human cell type. *Nature* **583**, 720-728 (2020).
- 871 2. Long, H.K., Prescott, S.L. & Wysocka, J. Ever-Changing Landscapes: Transcriptional
872 Enhancers in Development and Evolution. *Cell* **167**, 1170-1187 (2016).
- 873 3. Lambert, S.A. *et al.* The Human Transcription Factors. *Cell* **175**, 598-599 (2018).
- 874 4. Sullivan, P.F. *et al.* Leveraging base-pair mammalian constraint to understand genetic
875 variation and human disease. *Science* **380**, eabn2937 (2023).
- 876 5. Lenhard, B. *et al.* Identification of conserved regulatory elements by comparative
877 genome analysis. *J Biol* **2**, 13 (2003).
- 878 6. Gumucio, D.L. *et al.* Phylogenetic footprinting reveals a nuclear protein which binds to
879 silencer sequences in the human gamma and epsilon globin genes. *Mol Cell Biol* **12**,
880 4919-29 (1992).
- 881 7. Dermitzakis, E.T. & Clark, A.G. Evolution of transcription factor binding sites in
882 Mammalian gene regulatory regions: conservation and turnover. *Mol Biol Evol* **19**, 1114-
883 21 (2002).
- 884 8. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29
885 mammals. *Nature* **478**, 476-82 (2011).
- 886 9. Lander, E.S. *et al.* Initial sequencing and analysis of the human genome. *Nature* **409**,
887 860-921 (2001).
- 888 10. Venter, J.C. *et al.* The sequence of the human genome. *Science* **291**, 1304-51 (2001).
- 889 11. Nobrega, M.A., Ovcharenko, I., Afzal, V. & Rubin, E.M. Scanning human gene deserts for
890 long-range enhancers. *Science* **302**, 413 (2003).
- 891 12. Nobrega, M.A., Zhu, Y., Plajzer-Frick, I., Afzal, V. & Rubin, E.M. Megabase deletions of
892 gene deserts result in viable mice. *Nature* **431**, 988-93 (2004).
- 893 13. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals
894 folding principles of the human genome. *Science* **326**, 289-93 (2009).
- 895 14. Cosby, R.L. *et al.* Recurrent evolution of vertebrate transcription factors by transposase
896 capture. *Science* **371**(2021).
- 897 15. Schmidt, D. *et al.* Waves of retrotransposon expansion remodel genome organization
898 and CTCF binding in multiple mammalian lineages. *Cell* **148**, 335-48 (2012).
- 899 16. Chuong, E.B., Elde, N.C. & Feschotte, C. Regulatory evolution of innate immunity
900 through co-option of endogenous retroviruses. *Science* **351**, 1083-7 (2016).
- 901 17. Cohen, C.J., Lock, W.M. & Mager, D.L. Endogenous retroviral LTRs as promoters for
902 human genes: a critical assessment. *Gene* **448**, 105-14 (2009).
- 903 18. Ernst, J. *et al.* Genome-scale high-resolution mapping of activating and repressive
904 nucleotides in regulatory regions. *Nat Biotechnol* **34**, 1180-1190 (2016).
- 905 19. Kunarso, G. *et al.* Transposable elements have rewired the core regulatory network of
906 human embryonic stem cells. *Nat Genet* **42**, 631-4 (2010).
- 907 20. Jacobs, F.M. *et al.* An evolutionary arms race between KRAB zinc-finger genes ZNF91/93
908 and SVA/L1 retrotransposons. *Nature* **516**, 242-5 (2014).
- 909 21. Bannister, A.J. *et al.* Selective recognition of methylated lysine 9 on histone H3 by the
910 HP1 chromo domain. *Nature* **410**, 120-4 (2001).

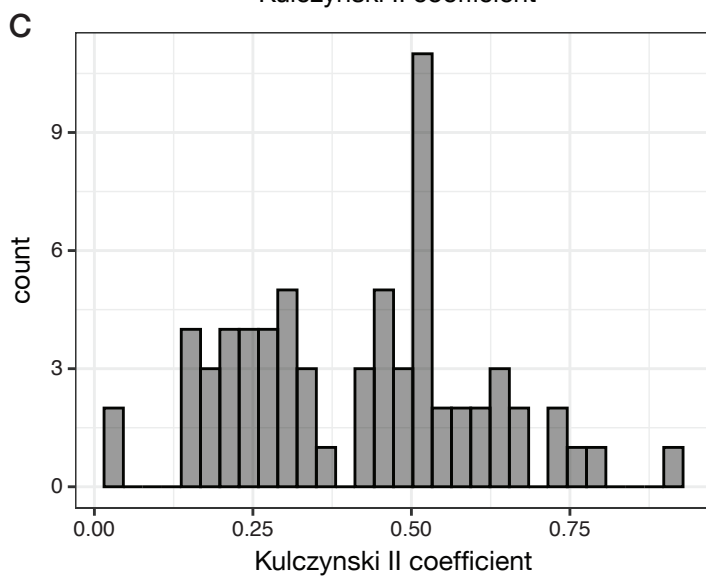
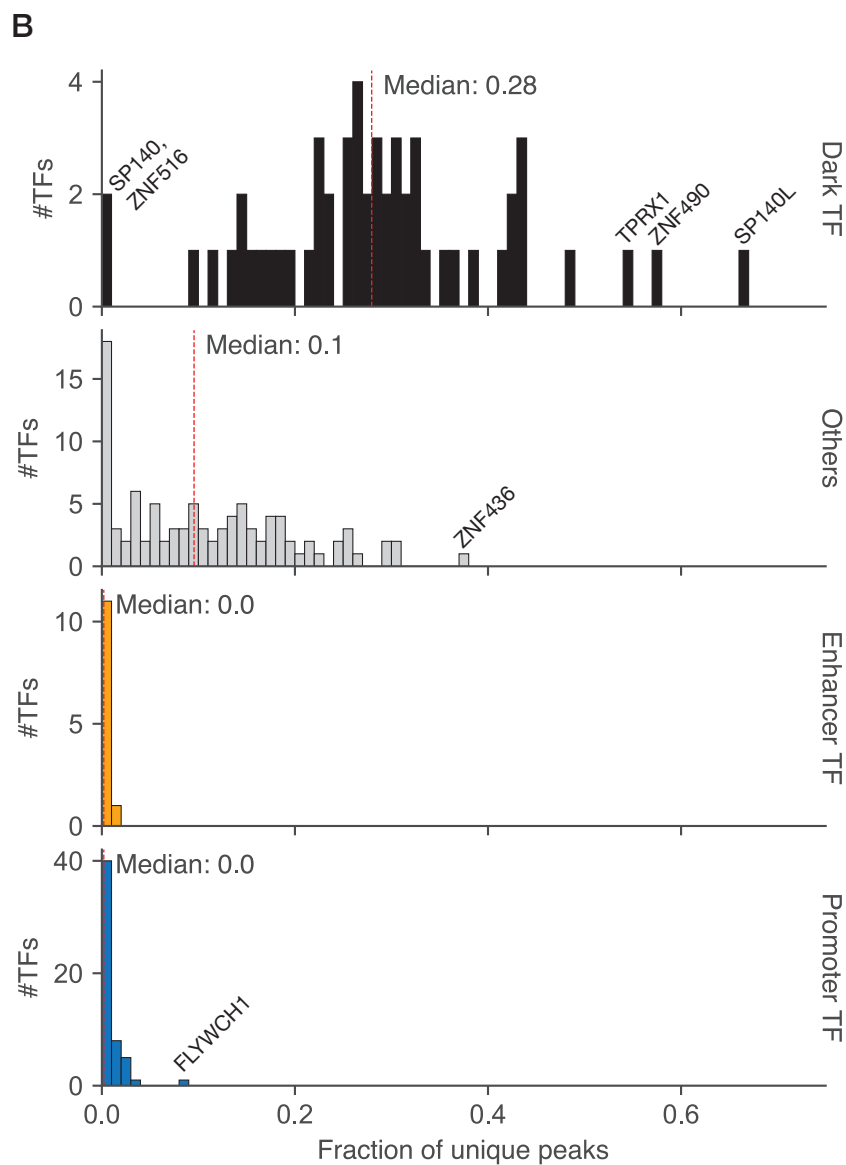
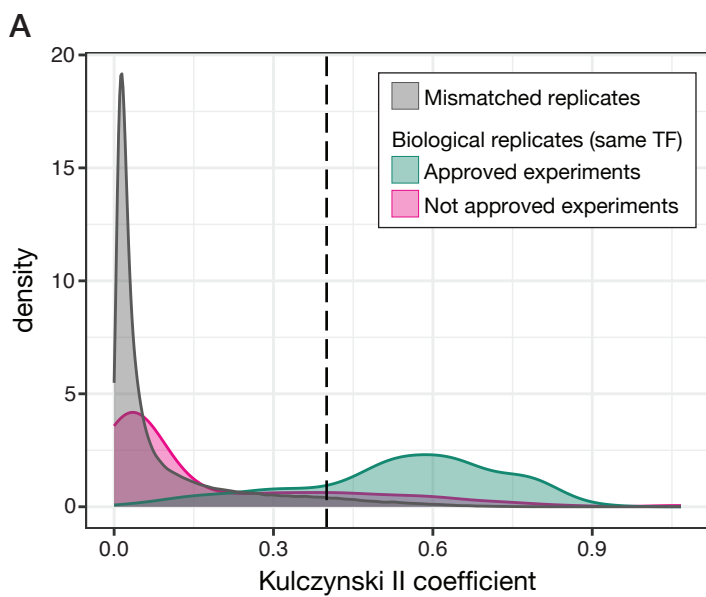
- 911 22. Lachner, M., O'Carroll, D., Rea, S., Mechtler, K. & Jenuwein, T. Methylation of histone H3
912 lysine 9 creates a binding site for HP1 proteins. *Nature* **410**, 116-20 (2001).
- 913 23. Ayyanathan, K. *et al.* Regulated recruitment of HP1 to a euchromatic gene induces
914 mitotically heritable, epigenetic gene silencing: a mammalian cell culture model of gene
915 variegation. *Genes Dev* **17**, 1855-69 (2003).
- 916 24. Stubbs, L., Sun, Y. & Caetano-Anolles, D. Function and Evolution of C2H2 Zinc Finger
917 Arrays. *Subcell Biochem* **52**, 75-94 (2011).
- 918 25. Jolma, A. *et al.* DNA-Binding Specificities of Human Transcription Factors. *Cell* **152**, 327-
919 39 (2013).
- 920 26. Imbeault, M., Helleboid, P.Y. & Trono, D. KRAB zinc-finger proteins contribute to the
921 evolution of gene regulatory networks. *Nature* **543**, 550-554 (2017).
- 922 27. Schmitges, F.W. *et al.* Multiparameter functional diversity of human C2H2 zinc finger
923 proteins. *Genome Res* **26**, 1742-1752 (2016).
- 924 28. Najafabadi, H.S., Albu, M. & Hughes, T.R. Identification of C2H2-ZF binding preferences
925 from ChIP-seq data using RCADE. *Bioinformatics* **31**, 2879-81 (2015).
- 926 29. Worsley Hunt, R. & Wasserman, W.W. Non-targeted transcription factors motifs are a
927 systemic component of ChIP-seq datasets. *Genome Biol* **15**, 412 (2014).
- 928 30. Auerbach, R.K. *et al.* Mapping accessible chromatin regions using Sono-Seq. *Proc Natl*
929 *Acad Sci U S A* **106**, 14926-31 (2009).
- 930 31. Jolma, A. *et al.* Perspectives on Codebook: sequence specificity of uncharacterized
931 human transcription factors. *bioRxiv*, 2024.11.11.622097 (2024).
- 932 32. Lambert, S.A. *et al.* The Human Transcription Factors. *Cell* **172**, 650-665 (2018).
- 933 33. Isakova, A. *et al.* SMiLE-seq identifies binding motifs of single and dimeric transcription
934 factors. *Nat Methods* **14**, 316-322 (2017).
- 935 34. Najafabadi, H.S. *et al.* C2H2 zinc finger proteins greatly expand the human regulatory
936 lexicon. *Nat Biotechnol* (2015).
- 937 35. Kean, M.J. *et al.* Structure-function analysis of core STRIPAK Proteins: a signaling
938 complex implicated in Golgi polarization. *J Biol Chem* **286**, 25065-75 (2011).
- 939 36. Vorontsov, I.E. *et al.* Cross-platform DNA motif discovery and benchmarking to explore
940 binding specificities of poorly studied human transcription factors. *bioRxiv*,
941 2024.11.11.619379 (2024).
- 942 37. Berger, M.F. *et al.* Compact, universal DNA microarrays to comprehensively determine
943 transcription-factor binding site specificities. *Nat Biotechnol* **24**, 1429-35 (2006).
- 944 38. Jolma, A. *et al.* Multiplexed massively parallel SELEX for characterization of human
945 transcription factor binding specificities. *Genome Res* **20**, 861-73 (2010).
- 946 39. Zhang, Y. *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).
- 947 40. O'Leary, N.A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,
948 taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733-45 (2016).
- 949 41. Ernst, J. & Kellis, M. ChromHMM: automating chromatin-state discovery and
950 characterization. *Nat Methods* **9**, 215-6 (2012).
- 951 42. Fishilevich, S. *et al.* GeneHancer: genome-wide integration of enhancers and target
952 genes in GeneCards. *Database (Oxford)* **2017**(2017).

- 953 43. Jolma, A. *et al.* GHT-SELEX demonstrates unexpectedly high intrinsic sequence specificity
954 and complex DNA binding of many human transcription factors. *bioRxiv*,
955 2024.11.11.618478 (2024).
- 956 44. Badis, G. *et al.* A library of yeast transcription factor motifs reveals a widespread
957 function for Rsc3 in targeting nucleosome exclusion at promoters. *Mol Cell* **32**, 878-87
958 (2008).
- 959 45. Core, L.J. *et al.* Analysis of nascent RNA identifies a unified architecture of initiation
960 regions at mammalian promoters and enhancers. *Nat Genet* **46**, 1311-20 (2014).
- 961 46. FitzGerald, P.C., Shlyakhtenko, A., Mir, A.A. & Vinson, C. Clustering of DNA sequences in
962 human promoters. *Genome Res* **14**, 1562-74 (2004).
- 963 47. Pollard, K.S., Hubisz, M.J., Rosenbloom, K.R. & Siepel, A. Detection of nonneutral
964 substitution rates on mammalian phylogenies. *Genome Res* **20**, 110-21 (2010).
- 965 48. Barazandeh, M., Lambert, S.A., Albu, M. & Hughes, T.R. Comparison of ChIP-Seq Data
966 and a Reference Motif Set for Human KRAB C2H2 Zinc Finger Proteins. *G3 (Bethesda)* **8**,
967 219-229 (2018).
- 968 49. Thomas, J.H. & Schneider, S. Coevolution of retroelements and tandem zinc finger
969 genes. *Genome Res* **21**, 1800-12 (2011).
- 970 50. Nowick, K. *et al.* Gain, loss and divergence in primate zinc-finger genes: a rich resource
971 for evolution of gene regulatory differences between species. *PLoS One* **6**, e21553
972 (2011).
- 973 51. Nowick, K., Hamilton, A.T., Zhang, H. & Stubbs, L. Rapid sequence and expression
974 divergence suggest selection for novel function in primate-specific KRAB-ZNF genes. *Mol*
975 *Biol Evol* **27**, 2606-17 (2010).
- 976 52. Alexander, T.B. *et al.* The genetic basis and cell of origin of mixed phenotype acute
977 leukaemia. *Nature* **562**, 373-379 (2018).
- 978 53. Arber, D.A. *et al.* International Consensus Classification of Myeloid Neoplasms and Acute
979 Leukemias: integrating morphologic, clinical, and genomic data. *Blood* **140**, 1200-1228
980 (2022).
- 981 54. Iouranova, A. *et al.* KRAB zinc finger protein ZNF676 controls the transcriptional
982 influence of LTR12-related endogenous retrovirus sequences. *Mob DNA* **13**, 4 (2022).
- 983 55. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-66
984 (2015).
- 985 56. Christmas, M.J. *et al.* Evolutionary constraint and innovation across hundreds of
986 placental mammals. *Science* **380**, eabn3943 (2023).
- 987 57. Harrison, P.W. *et al.* Ensembl 2024. *Nucleic Acids Res* **52**, D891-D899 (2024).
- 988 58. Kumar, S. *et al.* TimeTree 5: An Expanded Resource for Species Divergence Times. *Mol*
989 *Biol Evol* **39**(2022).
- 990 59. Helleboid, P.Y. *et al.* The interactome of KRAB zinc finger proteins reveals the
991 evolutionary history of their functional diversification. *EMBO J* **38**, e101220 (2019).
- 992 60. Huttlin, E.L. *et al.* Dual proteome-scale networks reveal cell-specific remodeling of the
993 human interactome. *Cell* **184**, 3022-3040 e28 (2021).
- 994 61. Silva, F.P., Hamamoto, R., Furukawa, Y. & Nakamura, Y. TIPUH1 encodes a novel KRAB
995 zinc-finger protein highly expressed in human hepatocellular carcinomas. *Oncogene* **25**,
996 5063-70 (2006).

- 997 62. Kim, J.J. *et al.* Systematic bromodomain protein screens identify homologous
998 recombination and R-loop suppression pathways involved in genome integrity. *Genes*
999 *Dev* **33**, 1751-1774 (2019).
- 1000 63. Oughtred, R. *et al.* The BioGRID database: A comprehensive biomedical resource of
1001 curated protein, genetic, and chemical interactions. *Protein Sci* **30**, 187-200 (2021).
- 1002 64. Agricola, E., Randall, R.A., Gaarenstroom, T., Dupont, S. & Hill, C.S. Recruitment of
1003 TIF1gamma to chromatin via its PHD finger-bromodomain activates its ubiquitin ligase
1004 and transcriptional repressor activities. *Mol Cell* **43**, 85-96 (2011).
- 1005 65. Maier, V.K. *et al.* Functional Proteomic Analysis of Repressive Histone Methyltransferase
1006 Complexes Reveals ZNF518B as a G9A Regulator. *Mol Cell Proteomics* **14**, 1435-46
1007 (2015).
- 1008 66. Li, L. *et al.* ZNF516 suppresses EGFR by targeting the CtBP/LSD1/CoREST complex to
1009 chromatin. *Nat Commun* **8**, 691 (2017).
- 1010 67. Huggenvik, J.I. *et al.* Characterization of a nuclear deformed epidermal autoregulatory
1011 factor-1 (DEAF-1)-related (NUDR) transcriptional regulator protein. *Mol Endocrinol* **12**,
1012 1619-39 (1998).
- 1013 68. Fraschilla, I. & Jeffrey, K.L. The Speckled Protein (SP) Family: Immunity's Chromatin
1014 Readers. *Trends Immunol* **41**, 572-585 (2020).
- 1015 69. Tumber, A. *et al.* Potent and Selective KDM5 Inhibitor Stops Cellular Demethylation of
1016 H3K4me3 at Transcription Start Sites and Proliferation of MM1S Myeloma Cells. *Cell*
1017 *Chem Biol* **24**, 371-380 (2017).
- 1018 70. Huttlin, E.L. *et al.* The BioPlex Network: A Systematic Exploration of the Human
1019 Interactome. *Cell* **162**, 425-440 (2015).
- 1020 71. Soufi, A., Donahue, G. & Zaret, K.S. Facilitators and impediments of the pluripotency
1021 reprogramming factors' initial engagement with the genome. *Cell* **151**, 994-1004 (2012).
- 1022 72. Zou, Z. *et al.* Translatome and transcriptome co-profiling reveals a role of TPRXs in
1023 human zygotic genome activation. *Science* **378**, abo7923 (2022).
- 1024 73. Kuroda, T. *et al.* SALL3 expression balance underlies lineage biases in human induced
1025 pluripotent stem cell differentiation. *Nat Commun* **10**, 2175 (2019).
- 1026 74. Singh, J.K. *et al.* Zinc finger protein ZNF384 is an adaptor of Ku to DNA during classical
1027 non-homologous end-joining. *Nat Commun* **12**, 6560 (2021).
- 1028 75. Feu, S. *et al.* OZF is a Claspin-interacting protein essential to maintain the replication
1029 fork progression rate under replication stress. *FASEB J* **34**, 6907-6919 (2020).
- 1030 76. Ge, L.P. *et al.* ZNF689 deficiency promotes intratumor heterogeneity and
1031 immunotherapy resistance in triple-negative breast cancer. *Cell Res* **34**, 58-75 (2024).
- 1032 77. Zhou, M. *et al.* ZBTB40 is a telomere-associated protein and protects telomeres in
1033 human ALT cells. *J Biol Chem* **299**, 105053 (2023).
- 1034 78. Cui, Y., Zhou, M., He, Q. & He, Z. Zbtb40 Deficiency Leads to Morphological and
1035 Phenotypic Abnormalities of Spermatocytes and Spermatozoa and Causes Male
1036 Infertility. *Cells* **12**(2023).
- 1037 79. Campos-Sanchez, R., Kapusta, A., Feschotte, C., Chiaromonte, F. & Makova, K.D.
1038 Genomic landscape of human, bat, and ex vivo DNA transposon integrations. *Mol Biol*
1039 *Evol* **31**, 1816-32 (2014).

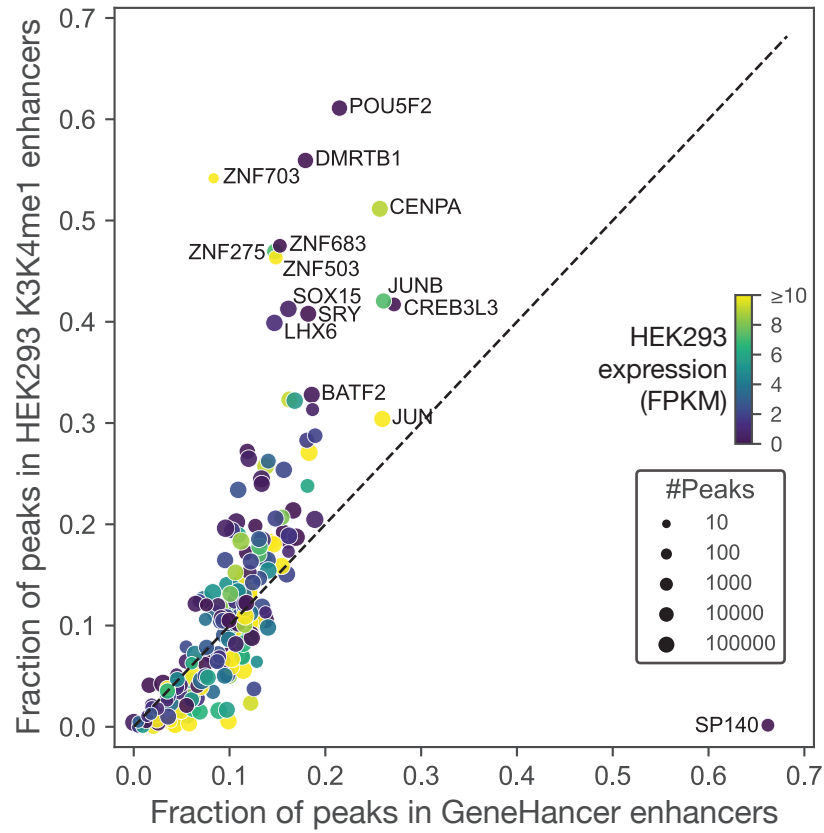
- 1040 80. Burton, K.A. *et al.* Haploinsufficiency at the protein kinase A RI alpha gene locus leads to
1041 fertility defects in male mice and men. *Mol Endocrinol* **20**, 2504-13 (2006).
- 1042 81. Stielow, B. *et al.* The SAM domain-containing protein 1 (SAMD1) acts as a repressive
1043 chromatin regulator at unmethylated CpG islands. *Sci Adv* **7**(2021).
- 1044 82. Lord, T. *et al.* A novel high throughput screen to identify candidate molecular networks
1045 that regulate spermatogenic stem cell functionsdagger. *Biol Reprod* **106**, 1175-1190
1046 (2022).
- 1047 83. Schmidt, D. *et al.* ChIP-seq: using high-throughput sequencing to discover protein-DNA
1048 interactions. *Methods* **48**, 240-8 (2009).
- 1049 84. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods*
1050 **9**, 357-9 (2012).
- 1051 85. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
1052 2078-9 (2009).
- 1053 86. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**(2021).
- 1054 87. Feng, J., Liu, T., Qin, B., Zhang, Y. & Liu, X.S. Identifying ChIP-seq enrichment using
1055 MACS. *Nat Protoc* **7**, 1728-40 (2012).
- 1056 88. Liu, T. Use model-based Analysis of ChIP-Seq (MACS) to analyze short reads generated
1057 by sequencing protein-DNA interactions in embryonic stem cells. *Methods Mol Biol*
1058 **1150**, 81-95 (2014).
- 1059 89. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic
1060 features. *Bioinformatics* **26**, 841-2 (2010).
- 1061 90. Grandi, F.C., Modi, H., Kampman, L. & Corces, M.R. Chromatin accessibility profiling by
1062 ATAC-seq. *Nat Protoc* **17**, 1518-1552 (2022).
- 1063 91. Consortium, E.P. *et al.* Perspectives on ENCODE. *Nature* **583**, 693-698 (2020).
- 1064 92. Frankish, A. *et al.* GENCODE: reference annotation for the human and mouse genomes
1065 in 2023. *Nucleic Acids Res* **51**, D942-D949 (2023).
- 1066 93. Dixon, J.R. *et al.* Topological domains in mammalian genomes identified by analysis of
1067 chromatin interactions. *Nature* **485**, 376-80 (2012).
- 1068 94. Satopaa, V., Albrecht, J., Irwin, D. & Raghavan, B. Finding a" kneedle" in a haystack:
1069 Detecting knee points in system behavior. in *2011 31st international conference on*
1070 *distributed computing systems workshops* 166-171 (IEEE, 2011).
- 1071 95. Korhonen, J., Martinmaki, P., Pizzi, C., Rastas, P. & Ukkonen, E. MOODS: fast search for
1072 position weight matrix matches in DNA sequences. *Bioinformatics* **25**, 3181-2 (2009).
- 1073 96. Armstrong, J. *et al.* Progressive Cactus is a multiple-genome aligner for the thousand-
1074 genome era. *Nature* **587**, 246-251 (2020).
- 1075 97. Ramirez, F., Dundar, F., Diehl, S., Gruning, B.A. & Manke, T. deepTools: a flexible
1076 platform for exploring deep-sequencing data. *Nucleic Acids Res* **42**, W187-91 (2014).
- 1077 98. Amemiya, H.M., Kundaje, A. & Boyle, A.P. The ENCODE Blacklist: Identification of
1078 Problematic Regions of the Genome. *Sci Rep* **9**, 9354 (2019).
- 1079 99. Hickey, G., Paten, B., Earl, D., Zerbino, D. & Haussler, D. HAL: a hierarchical format for
1080 storing and analyzing multiple genome alignments. *Bioinformatics* **29**, 1341-2 (2013).
- 1081 100. Martin, F.J. *et al.* Ensembl 2023. *Nucleic Acids Res* **51**, D933-D941 (2023).
- 1082 101. Jurka, J. Repbase update: a database and an electronic journal of repetitive elements.
1083 *Trends Genet* **16**, 418-20 (2000).

- 1084 102. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python.
1085 *Nat Methods* **17**, 261-272 (2020).
- 1086 103. Storer, J., Hubley, R., Rosen, J., Wheeler, T.J. & Smit, A.F. The Dfam community resource
1087 of transposable element families, sequence models, and genome annotations. *Mob DNA*
1088 **12**, 2 (2021).
- 1089 104. Lambert, S.A., Albu, M., Hughes, T.R. & Najafabadi, H.S. Motif comparison based on
1090 similarity of binding affinity profiles. *Bioinformatics* **32**, 3504-3506 (2016).
1091

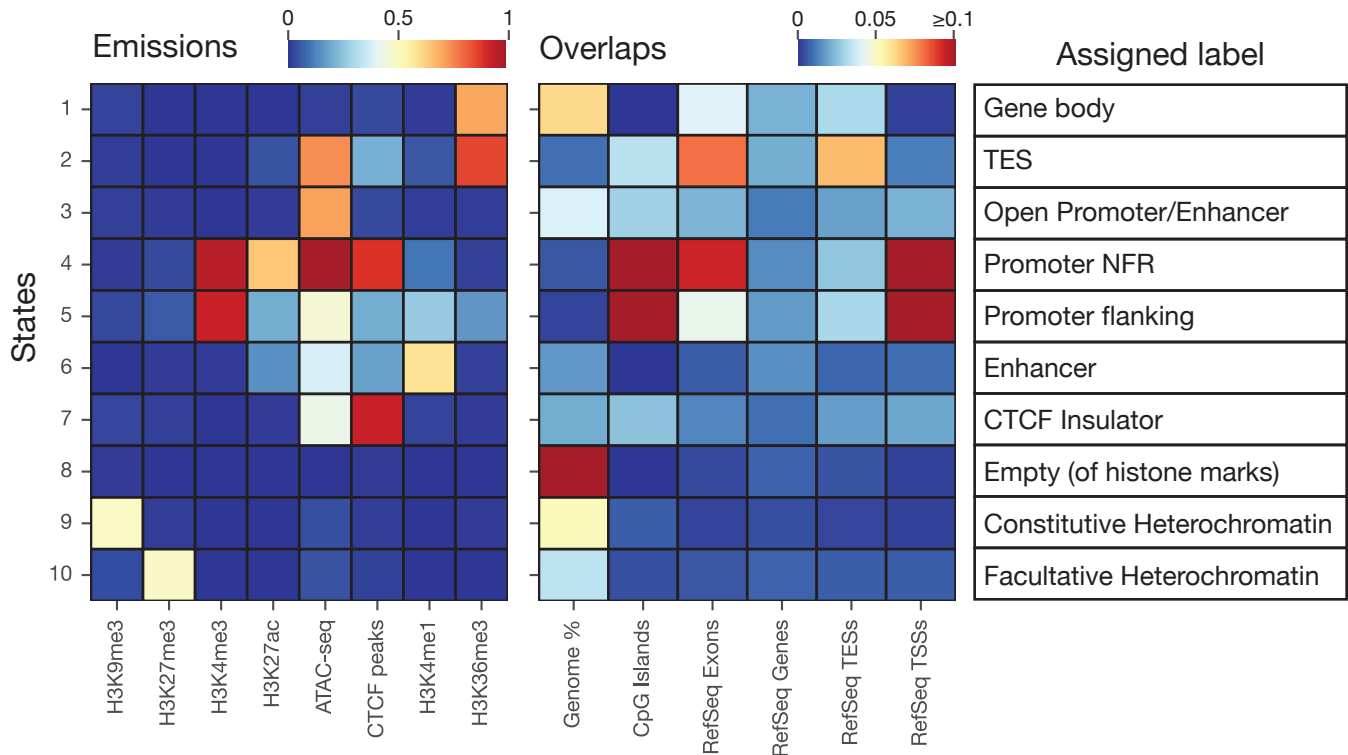


Supplementary Figure S1. Evaluation of ChIP-seq success by peak overlap assessment between experimental replicates. (A) Distribution of peak overlap between ChIP-seq replicates, for approved experiments (i.e., produced a motif), not approved experiments (i.e., did not produce a motif), and mismatch replicates (i.e., TF identities permuted), calculated by Kulczynski II similarity metric (i.e. average of overlaps). The dotted line indicates the threshold at which pairs of not approved experiments were considered successful and thus could be included in downstream analyses. (B) Distribution of the uniqueness of peaks for different categories of TFs, measured as the fraction of ChIP-seq peaks (at the universal threshold) not overlapping with *any* peak from any other TF in this study. (C) Distribution of Kulczynski II similarity metric between ChIP-seq replicates (as in (A)), restricted to the TFs that have a low peak overlap with other TFs (specifically, the 94 TFs in the upper left darker region of the square matrix in **Figure 2B**).

A



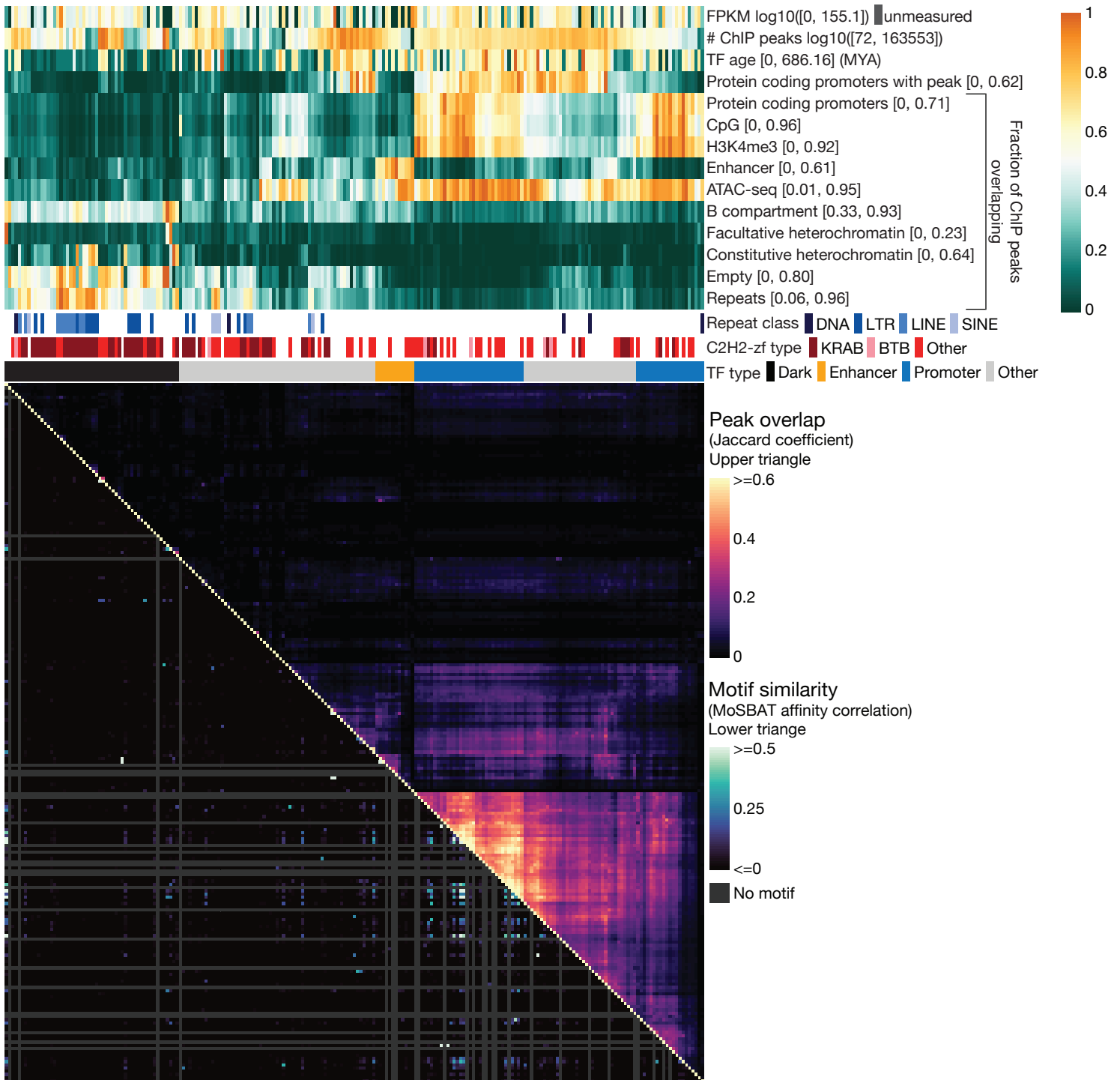
B



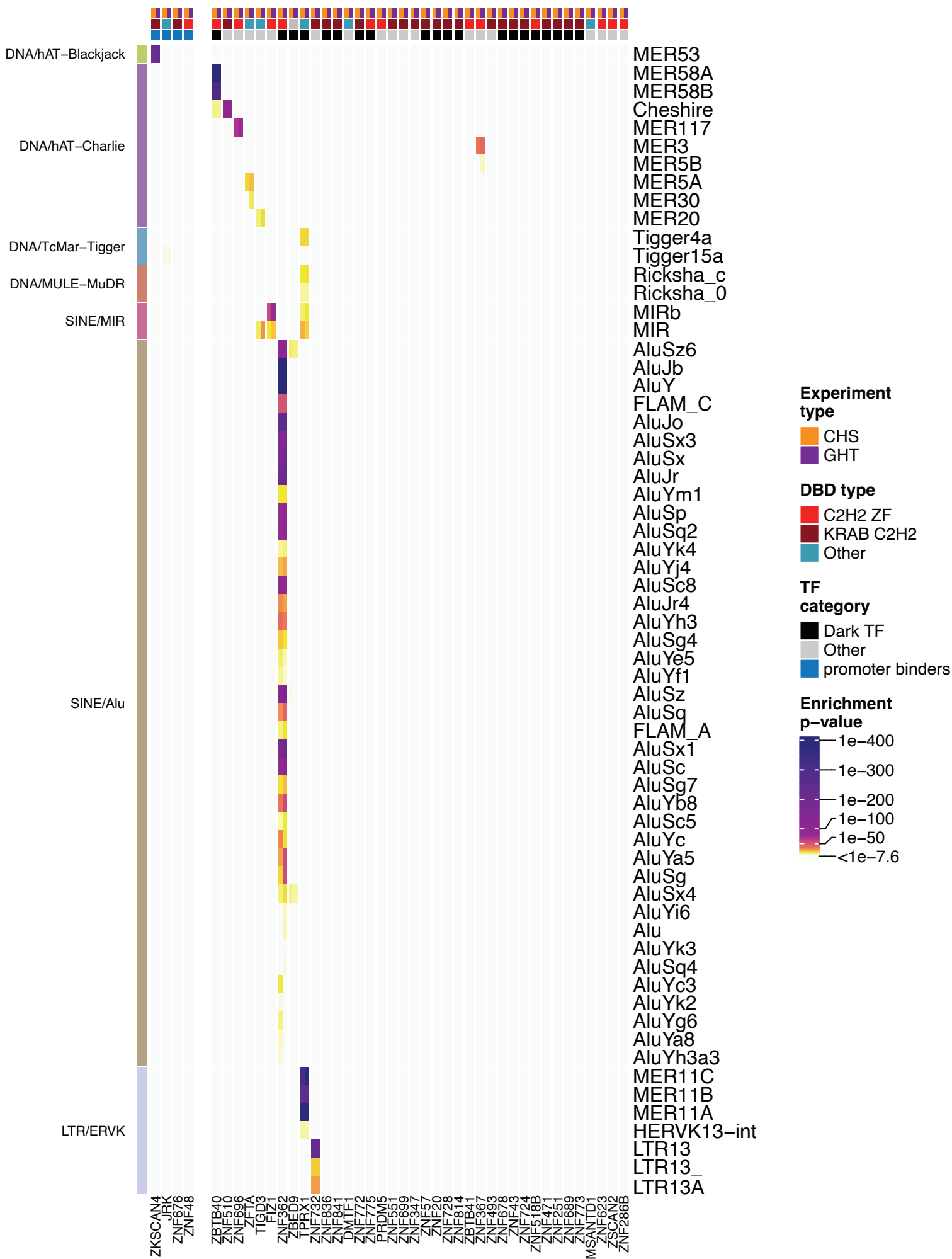
Supplementary Figure S2. Overlap of ChIP-seq peaks with different enhancer sets and ChromHMM tracks.

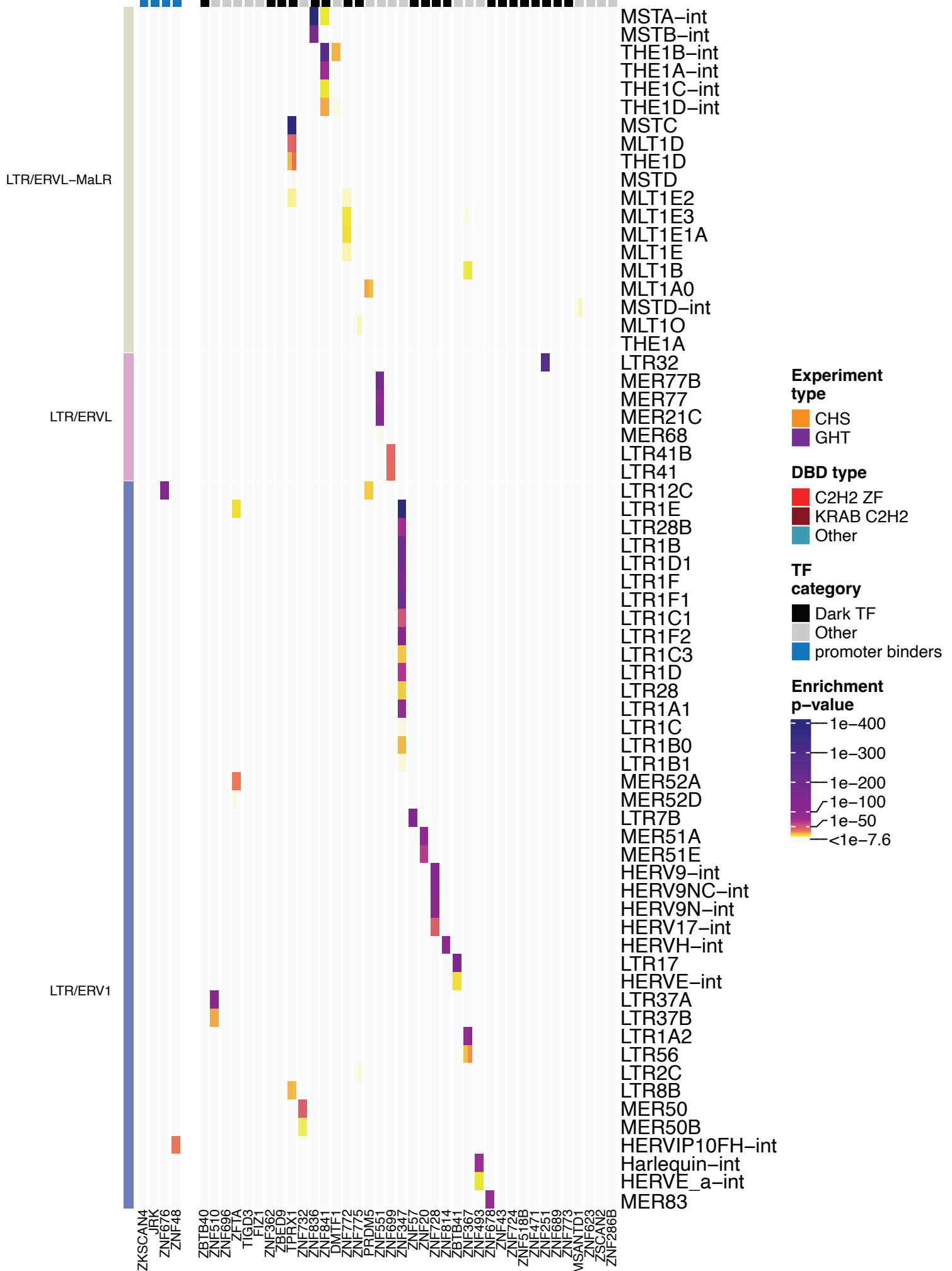
(A) Fraction of ChIP-seq peaks overlapping with GeneHancer annotated enhancers (x-axis) and HEK293 enhancers (defined by H3K4me1-positive regions from ChromHMM; y-axis). Points (TFs) are scaled based on their number of peaks. Colors also display the expression of TFs in HEK293 cells²⁷. **(B)**

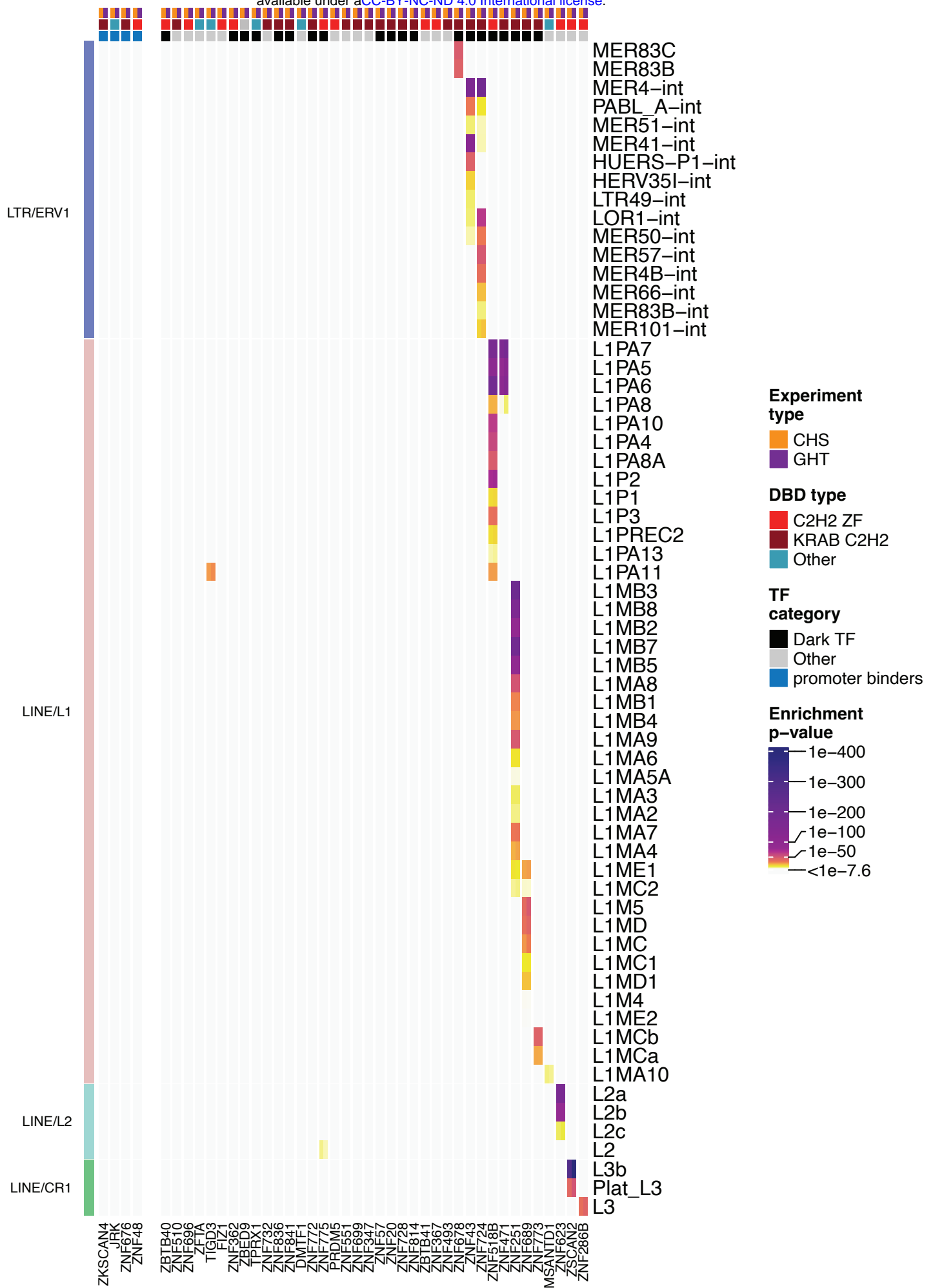
Characterization of the states of a ChromHMM model with 10 states trained on various HEK293 chromatin data (i.e., H3K9me3, H3K27me3, H3K4me1, H3K4me3, H3K36me3, and H3K27ac from ENCODE, and ATAC-seq and CTCF peaks from this study). Based on the correspondence between emissions and the chromatin marks and genome annotations, the states were assigned to Gene body, TES, Open Promoter/Enhancer, Promoter NFR (nucleosome-free regions), Promoter flanking, Enhancer, CTCF Insulator, Empty (of histone marks), Constitutive Heterochromatin, and Facultative Heterochromatin.



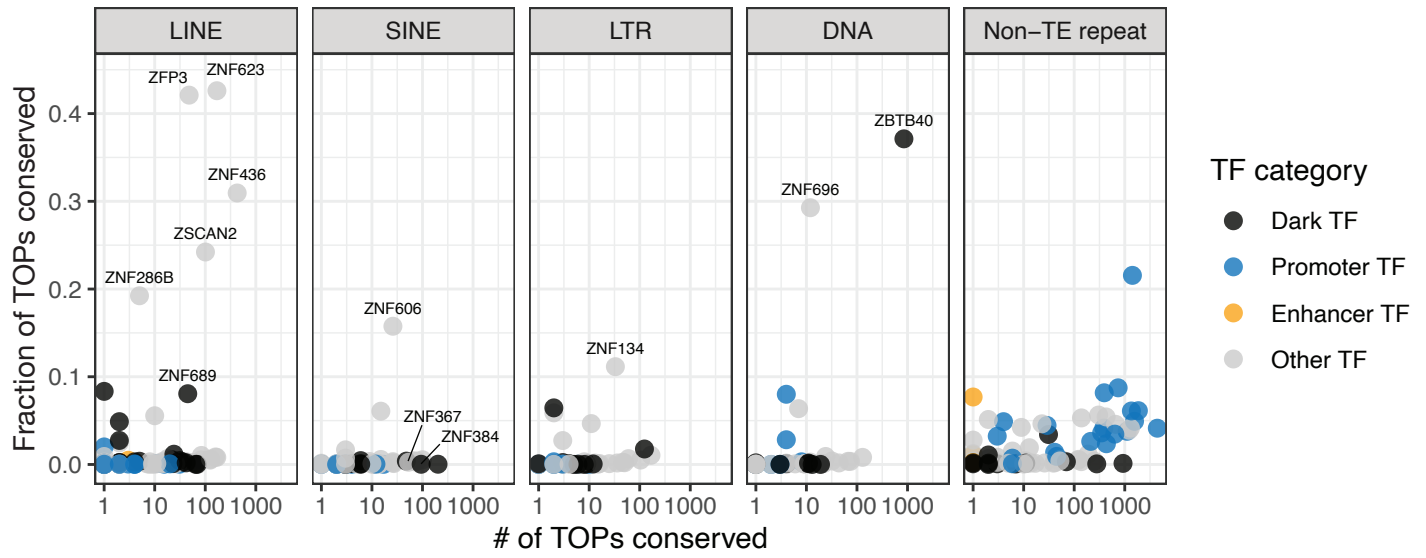
Supplementary Figure S3. A detailed version of **Figure 2** including additional tracks, such as gene expression in HEK293 cells (FPKM)²⁷, number of total ChIP-seq peaks (at the universal threshold of MACS2 P-value $\leq 10^{-10}$), TF age, fraction of human protein-coding promoters (out of 20,052) covered by TF peaks, fraction of ChIP-seq peaks falling within: CpG islands, H3K4me3-positive regions, facultative heterochromatin, and constitutive heterochromatin, with the main repeat class bound by the TFs included. The upper triangle in the bottom square is the same as **Figure 2**, however, the lower triangle here is the similarity between PWMs for each pair of TFs, calculated by MoSBAT¹⁰⁴. Gray stripes correspond to the TFs without a selected PWM in the Codebook set.





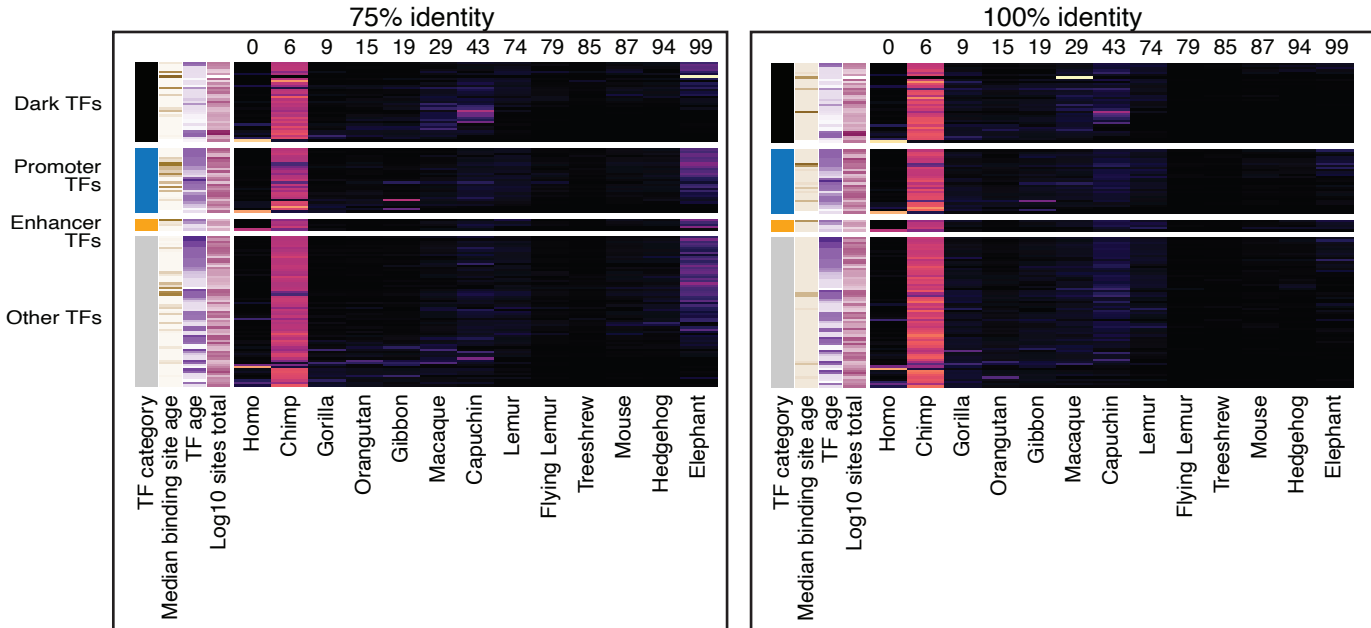


Supplementary Figure S4. Enrichment of transposable elements in TOPs, with expanded TE family classification. Heatmap is from **Figure 6**, with expanded labels for specific elements enriched in TOPs of each TF.

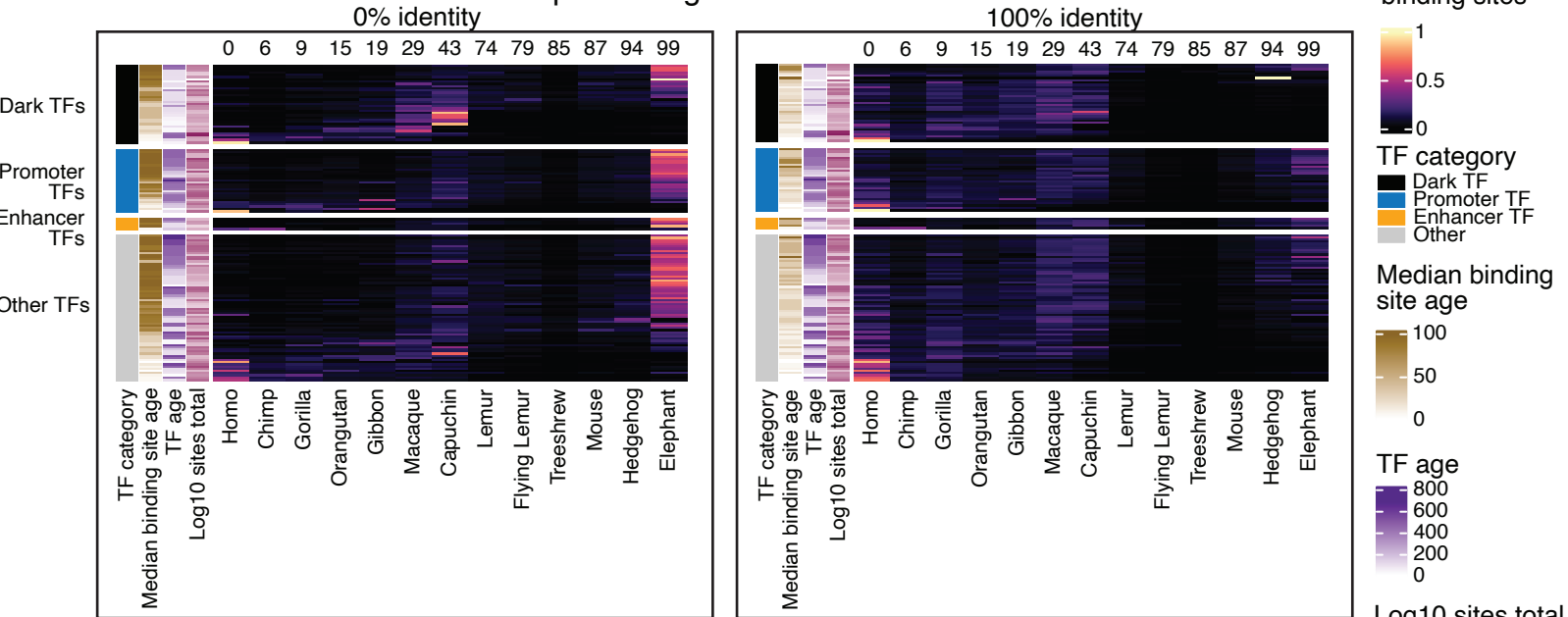


Supplementary Figure S5. Conservation of the binding sites for repeat-binding TFs. Plots showing the fraction of each TF's TOPs that are conserved (i.e. 'CTOPs') and overlap a major class of transposable elements or non-TE repeats. The proportion of TOPs that are conserved and overlap a repeat class is shown on the y-axis, and the log₁₀ count of these sites is shown on the x-axis. Each TF is coloured according to its classification as a Dark TF, Promoter TF, Enhancer and Other TFs. Only proteins with a fraction greater than 0.1 of conserved TOPs that fall in a repeat class are labeled. TFs discussed in the main text are also labeled.

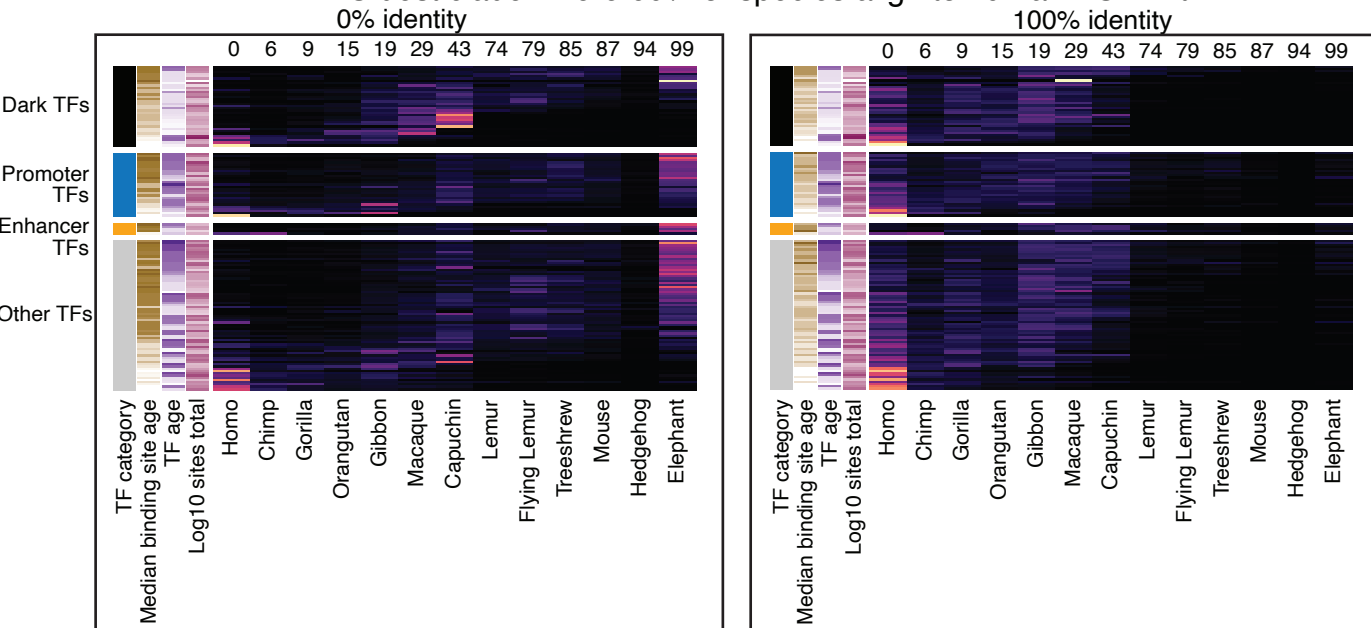
Oldest ancestral genome aligned to human TOP with:



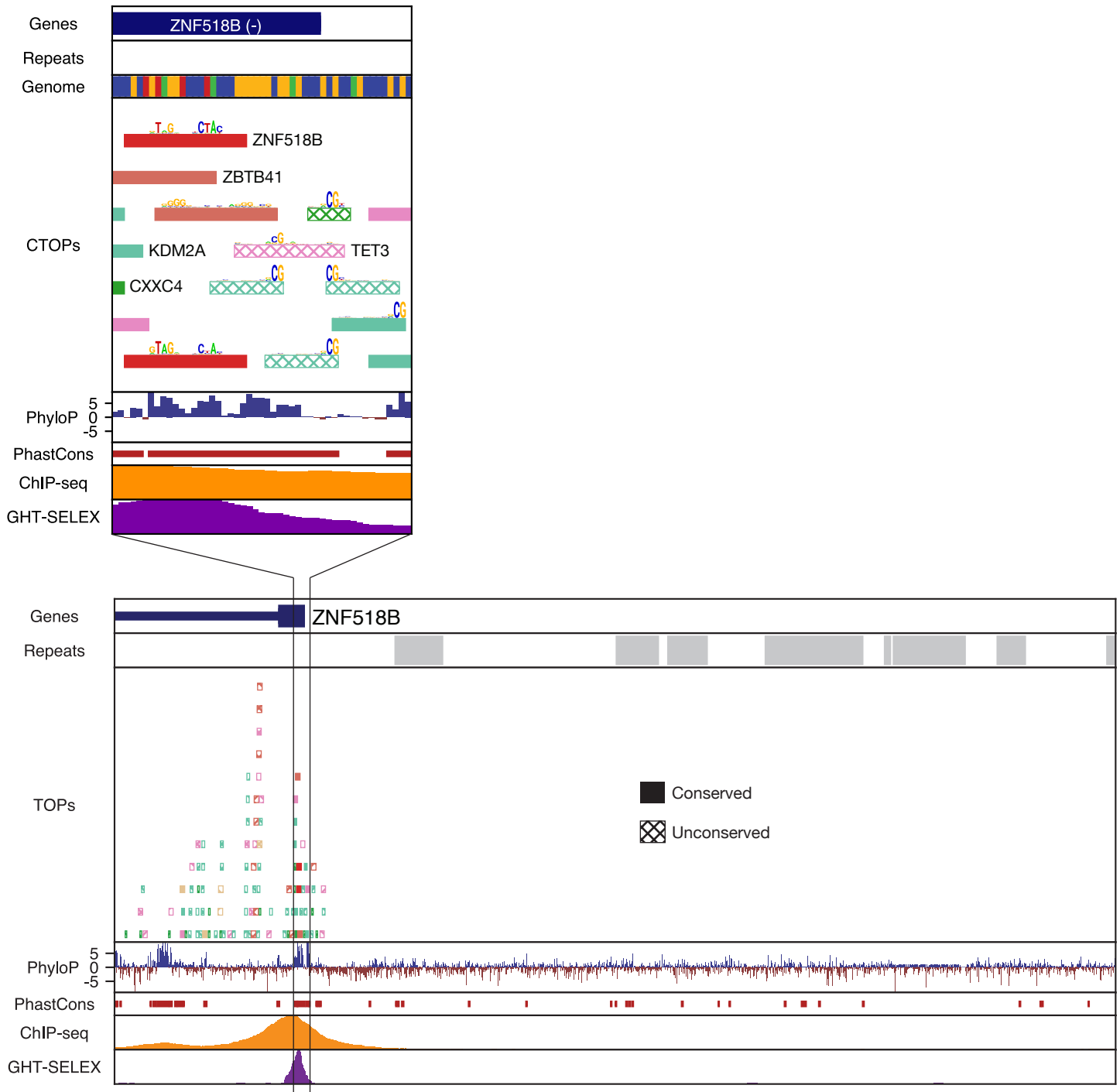
Oldest species aligned to human TOP with:



Oldest clade where 60% of species align to human TOP with:



Supplementary Figure S6. Estimation of binding site age using three different test and two different thresholds. Heatmaps show the proportion of each TF's TOPs (rows) inferred to be a certain age, as in Figure 5, but with each panel utilizing a different scheme. *Top row:* Age of each TOP site inferred as that of oldest ancestral genome with a gapless alignment to the human TOP site and minimum 75% identity (left) or 100% identity (right). (**Figure 5** shows this same analysis with a 0% identity threshold). *Middle row:* Age of each TOP site inferred as that of oldest *species* with a gapless alignment to the human TOP site and minimum 0% identity (left) or 100% identity (right). *Bottom row:* Age of each TOP site inferred as that of the oldest *clade* where 60% of the species have a gapless alignment to the human TOP with a minimum 0% identity (left) or 100% identity (right).



Supplementary Figure S7. A conserved binding site of ZNF518B as a potential self-regulator. Conserved binding sites for ZNF518B (red) located in the promoter of ZNF518B itself, and in a predicted enhancer-region ~4kb upstream of its promoter. Binding sites for ZBTB41, KDM2A, TET3, and CXXC4 are also present in this region.