

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  
32  
33  
34  
35  
36  
37  
38  
39  
40

# Perspectives on Codebook: sequence specificity of uncharacterized human transcription factors

Arttu Jolma<sup>1\*</sup>, Kaitlin U. Lavery<sup>1,2\*</sup>, Ali Fathi<sup>1,3\*</sup>, Ally W.H. Yang<sup>1\*</sup>, Isaac Yellan<sup>1,3\*</sup>, Ilya E. Vorontsov<sup>4\*</sup>, Sachi Inukai<sup>5,6</sup>, Judith F. Kribelbauer-Swietek<sup>5,6</sup>, Antoni J. Gralak<sup>5,6</sup>, Rozita Razavi<sup>1</sup>, Mihai Albu<sup>1</sup>, Alexander Brechalov<sup>1</sup>, Zain M. Patel<sup>13</sup>, Vladimir Nozdrin<sup>7</sup>, Georgy Meshcheryakov<sup>8</sup>, Ivan Kozin<sup>8</sup>, Sergey Abramov<sup>4,9</sup>, Alexandr Boytsov<sup>4,9</sup>, The Codebook Consortium, Oriol Fornes<sup>10</sup>, Vsevolod J. Makeev<sup>4,#</sup>, Jan Grau<sup>11</sup>, Ivo Grosse<sup>11</sup>, Philipp Bucher<sup>12</sup>, Bart Deplancke<sup>5,6\*\*</sup>, Ivan V. Kulakovskiy<sup>4,8\*\*</sup>, and Timothy R. Hughes<sup>1,3\*\*</sup>

<sup>1</sup>Donnelly Centre, University of Toronto, Toronto, ON M5S 3E1, Canada

<sup>2</sup>Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

<sup>3</sup>Department of Molecular Genetics, University of Toronto, Toronto, ON M5S 1A8, Canada

<sup>4</sup>Vavilov Institute of General Genetics, Russian Academy of Sciences, 119991, Moscow, Russia

<sup>5</sup>Laboratory of Systems Biology and Genetics, Institute of Bioengineering, School of Life Sciences, École Polytechnique Fédérale de Lausanne, 1015, Lausanne, Switzerland

<sup>6</sup>Swiss Institute of Bioinformatics, 1015, Lausanne, Switzerland

<sup>7</sup>Faculty of Bioengineering and Bioinformatics, Lomonosov Moscow State University, 119991, Moscow, Russia

<sup>8</sup>Institute of Protein Research, Russian Academy of Sciences, 142290, Pushchino, Russia

<sup>9</sup>Altius Institute for Biomedical Sciences, Seattle, WA 98121, USA

<sup>10</sup>Department of Medical Genetics, Centre for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, University of British Columbia, Vancouver, BC V5Z 4H4, Canada

<sup>11</sup>Institute of Computer Science, Martin Luther University Halle-Wittenberg, 06099, Halle, Germany

<sup>12</sup>Swiss Institute of Bioinformatics, 1015, Lausanne, Switzerland

<sup>#</sup>Present address: Cancer Research UK National Biomarker Centre, University of Manchester, Manchester, Manchester, M20 4BX, UK

\*These authors contributed equally

\*\* To whom correspondence should be addressed: [bart.deplancke@epfl.ch](mailto:bart.deplancke@epfl.ch), [ivan.kulakovskiy@gmail.com](mailto:ivan.kulakovskiy@gmail.com), [t.hughes@utoronto.ca](mailto:t.hughes@utoronto.ca)

41 **The Codebook Consortium**

42

43 **Principal investigators (steering committee)**

44 Philipp Bucher, Bart Deplancke, Oriol Fornes, Jan Grau, Ivo Grosse, Timothy R.

45 Hughes, Arttu Jolma, Fedor A. Kolpakov, Ivan V. Kulakovskiy, Vsevolod J. Makeev

46

47 **Analysis Centers:**

48 **University of Toronto (Data production and analysis):** Mihai Albu, Marjan

49 Barazandeh, Alexander Brechalov, Zhenfeng Deng, Ali Fathi, Arttu Jolma, Chun Hu,

50 Timothy R. Hughes, Samuel A. Lambert, Kaitlin U. Lavery, Zain M. Patel, Sara E. Pour,

51 Rozita Razavi, Mikhail Salnikov, Ally W.H. Yang, Isaac Yellan, Hong Zheng

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

53 Meshcheryakov

54 **EPFL, École polytechnique fédérale de Lausanne (Data production and analysis):**

55 Giovanna Ambrosini, Bart Deplancke, Antoni J. Gralak, Sachi Inukai, Judith F.

56 Kribelbauer-Swietek

57 **Martin Luther University Halle-Wittenberg (Data analysis):** Jan Grau, Ivo Grosse,

58 Marie-Luise Plescher

59 **Sirius University of Science and Technology (Data analysis):** Semyon Kolmykov,

60 Fedor Kolpakov

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

62 **Faculty of Bioengineering and Bioinformatics, Lomonosov Moscow State**

63 **University (Data analysis):** Nikita Gryzunov, Ivan Kozin, Mikhail Nikonov, Vladimir

64 Nozdrin, Arsenii Zinkevich

65 **Institute of Organic Chemistry and Biochemistry (Data analysis):** Katerina

66 Faltejskova

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

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

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

70 **Vavilov Institute of General Genetics (Data analysis):** Sergey Abramov, Alexandr

71 Boytsov, Vasilii Kamenets, Vsevolod J. Makeev, Dmitry Penzar, Anton Vlasov, Ilya E.

72 Vorontsov

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

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

75 Morris

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

77 **SUMMARY**

78 **We describe an effort (“Codebook”) to determine the sequence specificity of 332**  
79 **putative and largely uncharacterized human transcription factors (TFs), as well as**  
80 **61 control TFs. Nearly 5,000 independent experiments across multiple *in vitro* and**  
81 ***in vivo* assays produced motifs for just over half of the putative TFs analyzed**  
82 **(177, or 53%), of which most are unique to a single TF. The data highlight the**  
83 **extensive contribution of transposable elements to TF evolution, both in *cis* and**  
84 ***trans*, and identify tens of thousands of conserved, base-level binding sites in the**  
85 **human genome. The use of multiple assays provides an unprecedented**  
86 **opportunity to benchmark and analyze TF sequence specificity, function, and**  
87 **evolution, as further explored in accompanying manuscripts. 1,421 human TFs**  
88 **are now associated with a DNA binding motif. Extrapolation from the Codebook**  
89 **benchmarking, however, suggests that many of the currently known binding**  
90 **motifs for well-studied TFs may inaccurately describe the TF’s true sequence**  
91 **preferences.**

92 **KEYWORDS:** Transcription factor, TF, ChIP-seq, HT-SELEX, GHT-SELEX, SELEX,  
93 SMiLE-seq, Motif, DNA-binding specificity, PWM, PBM, Codebook

## 94 Introduction and motivations

95 The human genome encodes >1,600 putative transcription factors (TFs), defined as  
96 proteins that bind specific DNA sequence motifs and regulate gene expression<sup>1</sup>. These  
97 DNA binding motifs are most commonly modelled as a Position Weight Matrix (PWM)  
98 that describes the relative preference of the TF for each nucleotide base pair in the  
99 binding site<sup>2,3</sup>, and can be visualized as a sequence logo<sup>4</sup>. Several hundred putative  
100 human TFs still lack DNA binding motifs<sup>1</sup>, and even for well-characterized TFs, it  
101 remains controversial whether the reported motif model is accurate<sup>5,6</sup>, and to what  
102 degree the TF's sequence specificity contributes to binding site selection in living  
103 cells<sup>7,8</sup>. These uncertainties are due in part to the fact that different methods for  
104 measuring TF binding, and for deriving PWMs from these data, can have different  
105 inherent limitations and biases<sup>2</sup>. Such shortcomings represent fundamental hurdles for  
106 the analysis of gene regulation, as well as a myriad of related tasks in genome analysis,  
107 including the interpretation of conserved genomic elements and sequence variants, or  
108 genetic engineering such as synthetic enhancer design.

109 To address these issues, we analyzed a large majority of the as-yet uncharacterized  
110 human TFs<sup>1</sup>, as well as several dozen previously studied control TFs<sup>9,10</sup>, using a panel  
111 of assays that provide different perspectives on DNA sequence specificity. This  
112 unprecedented effort generated what we believe is the largest uniform data structure of  
113 its kind. We refer to this international collaborative project as the “Codebook/GRECO-  
114 BIT Collaboration”: the reagent set and laboratory experiments were initiated as the  
115 “Codebook Project”, alluding to the fact that TFs decode individual “words” in the  
116 genome, and the existing Gene REgulation COnsortium Benchmarking IniTiative,  
117 GRECO-BIT, was then engaged for much of the data analysis and benchmarking.

118 In this paper, we present an overview of the data collection and its analysis, the  
119 resulting data, several major outcomes and findings of the Codebook study, and  
120 examples of prevalent phenomena and applications. We also introduce web resources  
121 that can be used to access the primary and processed data, including the PWMs.  
122 Accompanying manuscripts provide greater depth regarding biological findings, new  
123 assays, and intriguing TF families, as well as methods for identifying binding patterns  
124 (i.e. PWM derivation) and PWM benchmarking (**Table S1**).

## 125 Codebook reagents, assays, and data structure

126 **Figure 1** provides a schematic of the Codebook project. We chose 332 putative TFs  
127 (i.e., “Codebook TFs”) (**Table S2**) for study by starting with a previously described list of  
128 427 hand-curated “likely” human TFs that lacked known motifs or any large-scale DNA  
129 binding data<sup>1</sup>. We removed 95 C2H2 zinc finger (C2H2-zf) proteins for which we were  
130 already aware of unpublished data (mainly from our prior collaboration with ENCODE<sup>11</sup>).  
131 As of June 2024, most of these putative TFs still lack motifs, outside of the Codebook  
132 study: of the 332, only 107 have PWMs on Factorbook<sup>12</sup> and/or HOCOMOCO<sup>13</sup>. Many  
133 of these motifs appear to be simple repeats, or common cofactor motifs (such as CTCF,  
134 REST, and CRE sites) (examples in **Figure S1**), but among the 107, 59 have at least  
135 one PWM that appears plausible for representing specificity of the TF (see below).

136 Among the 332 Codebook TFs, 180 contain C2H2-zf domains, while another 103  
137 contain another type of well-known DNA-binding domain (DBD). Forty-nine did not have  
138 an established DBD at the outset of the study; these were mainly identified as  
139 sequence-specific in studies of individual proteins or regulatory sites<sup>1</sup>. We  
140 simultaneously analyzed 61 control TFs, encompassing 29 well-characterized TFs  
141 representing diverse human DBD classes<sup>9</sup>, and an additional 32 C2H2-zf proteins for  
142 which published ChIP-seq data were available and had led to a binding motif<sup>10</sup>. For  
143 these controls, we incorporated the published SMiLE-seq and ChIP-seq data, rather  
144 than repeating the experiments.

145 To study the 332 Codebook proteins, we manually designed 716 protein-coding inserts,  
146 corresponding to full-length coding regions of the dominant isoform, and one or more  
147 DBDs (or subsets of C2H2-zf domain arrays), if there was a known DBD (**Table S3**).  
148 We employed up to three different expression vectors for each insert, as required for the  
149 different assays in **Figure 1**, resulting in a total of 1352 new distinct constructs (**Table**  
150 **S4**). One of the assays, GHT-SELEX (Genomic high-throughput SELEX), is a new  
151 variant of HT-SELEX which is performed with fragmented genomic DNA. As described  
152 in the accompanying manuscript<sup>14</sup>, GHT-SELEX yields peaks, analogous and often in  
153 agreement with ChIP-seq. GHT-SELEX thus provides a new perspective that bridges *in*  
154 *vitro* and *in vivo* DNA binding. HT-SELEX and GHT-SELEX were performed with  
155 multiple protein sources (mammalian cell extracts, and two different systems for *in vitro*  
156 transcription/translation) whereas SMiLE-seq and PBMs were performed with only one  
157 protein source. Multiple replicates were performed in many cases, for all assays.

158 The full Codebook data structure is composed of a total of 4,873 technically successful  
159 experiments (i.e. they produced data that could be analyzed by at least some  
160 subsequent processes) (**Table S5**). The Codebook data structure, experimental  
161 information, and PWMs (see below) are accessible at multiple sources (see **Data**  
162 **Availability**). Each experiment corresponds to one of the Codebook constructs (or one  
163 of the control constructs), analyzed using one of the assays, with one of the protein  
164 sources. Not every protein or every insert was analyzed in every assay, by design. For  
165 example, the ChIP-seq data only utilize full-length proteins, while Protein Binding  
166 Microarray data include only DBD constructs. Long human C2H2-zf domain arrays  
167 typically fail in PBMs, and such experiments were omitted. We note that, in general,  
168 experiments that are technically successful may not yield motifs that are specific to the  
169 TF assessed and supported by other data types (see below). For example, ChIP-seq  
170 can detect both indirect and non-sequence-specific DNA binding, as we explored  
171 separately<sup>15</sup>. We also emphasize that the *in vitro* assays described here were  
172 conducted with unmethylated DNA. We explored the sensitivity of a subset (79) of the  
173 Codebook TFs to DNA methylation in an accompanying study, however, which  
174 introduces the methylation-sensitive SMiLE-seq variant (meSMS)<sup>16</sup>. DNA binding  
175 interactions of 17 of the 79 were impacted by methylation, encompassing inhibition (10)  
176 and increased binding or alternative binding sites (7); these data were not incorporated  
177 in the analyses described herein.

178 **Motifs are obtained from most C2H2-zf proteins, and half of those containing**  
179 **other DBD classes, but only a few proteins with previously unknown DNA binding**  
180 **domains**

181 We next derived and examined motifs as PWMs for all the experiments in a semi-  
182 automated expert curation format, to identify “approved” experiments (i.e. experiments  
183 that contained clear enrichment of credible binding motifs (see **Methods**)). This effort is  
184 described in detail in a separate manuscript that describes motif benchmarking, data  
185 sets, and success measures<sup>17</sup>, and also introduces a web resource that makes all of the  
186 motifs available for browsing and download. Briefly, our primary approach was to ask  
187 whether similar motifs were obtained for the same protein from different assays and  
188 whether the PWMs scored highly by a panel of criteria, including predictive capacity in  
189 other data types (depicted schematically in **Figure 1, bottom left**), adapting a previously  
190 described motif benchmarking framework<sup>18</sup>. To increase our ability to derive motifs that  
191 would score highly across data sets, we employed ten motif discovery tools, ranging  
192 from the widely used MEME suite<sup>19</sup> to approaches based on machine learning or  
193 biophysical modeling, such as ExplaiNN<sup>20</sup> and ProBound<sup>21</sup>, thus producing hundreds of  
194 motifs per TF. In total, 177 Codebook TFs were associated with “approved” datasets  
195 (**Figure 1, bottom right**), and a total of 1,072 experiments associated with these 177  
196 TFs were approved (**Tables S2 and S5**). 59/61 controls were also approved, suggesting  
197 a low per-TF false-negative rate.

198 The 177 Codebook TFs for which there are approved experiments are dominated by the  
199 C2H2-zf domain class, for which 67% (121/180) had approved experiments. These  
200 proteins typically contain an array of C2H2-zf domains that bind DNA in tandem<sup>22</sup>.  
201 Some C2H2-zf domains can bind RNA, protein, or other ligands<sup>23-25</sup>. The Codebook  
202 outcome indicates that most C2H2-zf proteins are indeed DNA-binding, although it does  
203 not rule out their other activities. Experiments for roughly half (50/103, or 49%) of  
204 Codebook TFs in other established DBD classes were also successful. Lack of  
205 approved experiments for a putative TF could represent false negatives, which could  
206 arise from lack of an obligate binding partner, a requirement for epigenetically modified  
207 DNA, lack of requisite post-translational modification in our experiments, or limitations of  
208 the methods. Alternatively, they could represent true negatives which are not  
209 unexpected; some *bona fide* DBD classes are known to have subtypes that lack  
210 sequence specificity (e.g. HMG<sup>26</sup>). Among the Codebook proteins lacking a well-  
211 established DBD, only 6/49 (12%) yielded approved experiments (and thus motifs)  
212 (discussed in more detail below), suggesting that many of them may indeed lack  
213 sequence specificity.

214 We emphasize that our approval process was intentionally conservative, and many  
215 experiments were not approved despite being informative in some way (e.g. ChIP-seq  
216 yielding reproducible peaks, but no motif, which could indicate indirect association  
217 through other TFs or chromatin binding; these are explored in an accompanying  
218 manuscript<sup>15</sup>). We also note that our success criteria assume that the sequence  
219 preferences of TFs can be represented by PWMs. It is conceivable that uncharacterized  
220 TFs could instead recognize interspersed sequence patterns or other features of the  
221 DNA sequence that are not readily captured by PWM models or short k-mers.

## 222 Diversity and complexity among Codebook TF motifs

223 To gain an overview of the Codebook TF motifs, and to generate a representative PWM  
224 set, we next used expert curation to select a single PWM that is (i) high-performing  
225 among all “approved” experiments<sup>17</sup> (see **Methods**), (ii) representative of other high-  
226 performing PWMs for the same TF, (iii) consistent with expectation for the class of TF  
227 (e.g. the C2H2-zf “recognition code”<sup>27</sup>), and (iv) high information content (IC) (i.e. with a  
228 “tall” sequence logo), provided it does not compromise PWM performance. The PWM  
229 selected in this process is typically not the highest scoring by criterion (i) alone, as our  
230 extensive process typically generated dozens of high-performing PWMs from which to  
231 choose, for approved experiments<sup>17</sup>. **Table S6** shows sequence logos for these curated  
232 PWMs and their properties; the PWM IDs are given in **Table S2**, and all PWMs can be  
233 downloaded (see **Data Availability**). Notably, no data type or motif derivation method  
234 stood out as highly preferred by the curators, who were blinded to the source (i.e. data  
235 type and derivation method for the PWMs).

236 **Figure 2** shows an overview of similarity<sup>28</sup> among the curated PWMs. Small clusters  
237 along the diagonal mostly correspond to the handful of paralogs analyzed (e.g. TIGD4  
238 and 5, SP140 and SP140L, DACH1 and 2, CAMTA1 and 2, and ZXDA, B, and C). In the  
239 middle of **Figure 2** is a set of eight TFs that mainly bind CG dinucleotides, leading to  
240 similarity in DNA-binding, and in the lower right is a group of five AT-hook proteins that  
241 have similar preferences to A/T containing sequences. Most of the Codebook TF PWMs  
242 are unlike each other, however, and display a low similarity to any other known PWM<sup>17</sup>  
243 (examples are shown in **Figure 2**). This result is partly explained by the large number of  
244 C2H2-zf proteins, which are known to differ in their DNA-contacting “specificity  
245 residues”<sup>29</sup>. Regardless, a large majority of the Codebook TF motifs are apparently new,  
246 and all previous analyses in human regulatory genomics would have been unaware of  
247 the ~150 visibly distinct, curated motifs described here.

248 For dozens of TFs, the curated PWM had a degenerate appearance, i.e. there are few  
249 or no positions at which a specific base is absolutely required. Indeed, for fifty-two of  
250 them, no individual base at any position achieved a bit score of  $\geq 1.4$  in the curated  
251 PWM (equivalent to roughly >10% of aligned binding sites having a variant base at that  
252 position) (**Figure S2A**). Systematically increasing the information content (IC) (i.e.,  
253 “unflattening” the sequence logo, and increasing the specificity) of the low-IC curated  
254 PWMs almost universally reduced performance (**Figure S2B,C**), indicating that the  
255 degeneracy is required for accuracy. We also found that, overall, IC is not predictive of  
256 motif performance in the benchmarking effort<sup>17</sup>. It is counterintuitive that degeneracy  
257 (i.e. lower inherent specificity) would lead to better predictive capacity, but we note that  
258 similar findings by others support the validity of the result<sup>30-32</sup>.

259 We propose several explanations for this observation. First, lower IC tends to make  
260 affinity distributions across all possible k-mers less digital (i.e. it removes all-or-nothing  
261 dependence on specific base positions), which could facilitate the gradual evolution of  
262 *cis*-regulatory sequences. Second, homomeric binding (possibly via “avidity”<sup>33</sup>), which a  
263 body of evidence suggests is a widespread mechanism<sup>14,34</sup>, should reduce reliance on  
264 optimal specificity to a single binding site, and strong binding sites may evolve more

265 readily if weak binding sites tend to occur more frequently (and are selected). Third,  
266 motif degeneracy may be a consequence of forcing a single PWM to represent the  
267 specificity of TFs that, in reality, recognize multiple related motifs. For example, the  
268 dependency of binding energy on both enthalpy and entropy can lead to two distinct  
269 sequence optima<sup>35</sup>; in another example, different spacings of bZIP half-sites cannot be  
270 represented by a single PWM<sup>36</sup>. Consistent with this last possibility, the accompanying  
271 manuscript<sup>17</sup> finds that combining multiple PWMs (by Random Forests) typically  
272 produces models that are more accurate across platforms, relative to any single PWM.

273 The C2H2-zf proteins present a special case in which a single TF might be anticipated  
274 to require multiple PWMs, because long C2H2-zf domain arrays could utilize different  
275 segments of the array to bind to either overlapping or distinct sites<sup>37</sup>. Until now,  
276 however, examples were sparse and anecdotal. In an accompanying manuscript<sup>14</sup>, we  
277 present evidence that C2H2-zf proteins often bind multiple sequence motifs that  
278 correspond to different subsets of the extended motif predicted by the recognition code  
279 (i.e. protein-sequence-based computational prediction of C2H2-zf-domain specificities),  
280 consistent with varying usage of the C2H2-zf domains at different genomic binding sites  
281 being commonplace.

## 282 **Underappreciated DNA-binding domains**

283 The six Codebook proteins that were lacking canonical DBDs, yet yielded “approved”  
284 experiments and thus motifs (CGGBP1, NACC2, TCF20, PURB, DACH1, and DACH2),  
285 appear to represent cases of DBDs that were poorly described at the outset of the  
286 study. We and others have recently described CGGBP1 as the founding member of an  
287 extensive family of eukaryotic TFs derived from the DBDs of transposons<sup>38,39</sup>. NACC2  
288 contains a BEN domain, which over the last decade has been clearly established as a  
289 sequence-specific DBD<sup>40,41</sup>. TCF20 contains a potential AT-hook<sup>42</sup> (below the  
290 conventional Pfam scoring threshold), and yielded an AT-hook-like motif. PURB is  
291 composed largely of three copies of the PUR (Purine-rich-element binding) domain; it  
292 yielded a motif on four different PBM assays (resembling ACCnAC/GTnGGT), which is  
293 unlike its previously established binding site (CTTCCCTGGAAG)<sup>43</sup>. The sequence  
294 specificity of this protein thus remains enigmatic.

295 DACH1 and DACH2 are paralogs that yielded very similar motifs (**Figure 3A**). They  
296 contain a SKI/SNO/DAC domain, shared with their *Drosophila* counterpart Dachshund,  
297 from which their name is derived. A Forkhead-like motif (different from the one we  
298 obtained) was previously described for DACH1<sup>44</sup>, but to our knowledge, no other  
299 homolog has been reported as being sequence-specific. The SKI/SNO/DAC domain  
300 includes a helix-turn-helix (HTH), a feature found in many DBDs. Alphafold<sup>45</sup> predicts  
301 that the HTH inserts into the major groove precisely at the PWM-predicted binding site  
302 within an extended DNA sequence (**Figure 3A**). Interpro<sup>46</sup> lists over 7,000 proteins  
303 containing SKI/SNO/DAC domains, entirely in metazoans, with specific expansions in  
304 several fish lineages, particularly barbels and salmonids<sup>47</sup> (**Figure 3A**). SKI/SNO/DAC  
305 therefore may represent an expansive class of poorly-characterized DBDs.



306 In addition to these six examples, the sequence specificity of SLC2A4RG and ZNF395 –  
307 both C2H2-zf proteins – appears to reside in their C-clamp. The domain is also present  
308 in TCF7L and LEF proteins, where it is known to bind DNA alongside their HMG  
309 domains<sup>48</sup>. Alphafold3<sup>45</sup> predicts that the single C2H2-zf domains in SLC2A4RG and  
310 ZNF395 are not the main determinants of DNA-binding (although they may contact the  
311 major groove), but instead that a region corresponding to the C-clamp model on the  
312 SMART database of protein domains<sup>49</sup> binds the major groove precisely at the PWM-  
313 predicted binding site within an extended DNA sequence (**Figure 3B**). There is one  
314 additional human TF matching the C-clamp model, ZNF704, with a published PWM that  
315 is virtually identical to that of SLC2A4RG and ZNF395 (CCGGCCGG)<sup>50</sup> (**Figure 3B**).  
316 Like the SKI/SNO/DAC domain, the C-clamp is found broadly across animals<sup>46</sup>, and  
317 may therefore also represent a large class of unexplored DBDs.

### 318 **Widespread contribution of transposons to the human TF repertoire**

319 Sixteen of the Codebook TFs (and two controls) that yielded approved experiments  
320 possess a DBD that has been co-opted from a DNA transposon: CGGBP1<sup>39</sup>, five  
321 proteins containing BED-zf domains<sup>51</sup>, six with the related CENBP or Brinker domains<sup>52</sup>,  
322 two with transposon-derived Myb/SANT domains<sup>53</sup>, one with a MADF domain, and  
323 FLYWCH1<sup>54</sup>. The PWMs obtained for CENPB/Brinker TFs are often long (**Figure 3C**). A  
324 striking example is JRK, a TF that is derived from an ancient domesticated Tigger  
325 element DBD<sup>55</sup>, and is found broadly in mammals<sup>47</sup>. All DNA transposons, including  
326 Tigger, have been extinct in the human lineage for over 40 million years<sup>56</sup>. Remarkably,  
327 genomic binding of JRK is enriched for binding to a subset of Tigger elements, and the  
328 consensus sequence for these same elements has a PWM-predicted binding site for  
329 JRK in the terminal repeats of these elements (**Figure 3C**), consistent with its presumed  
330 ancestral role in transposition. We speculate that JRK may represent a case of co-  
331 option in which the same DNA transposon simultaneously introduced both a multitude of  
332 *cis*-regulatory elements, and the TF that binds them.

333 The Codebook data also underscore that many TFs bind preferentially and intrinsically  
334 to specific repeat classes. These interactions are explored in greater detail in the  
335 accompanying manuscripts<sup>14,15</sup>. Binding to endogenous retroelements is known to be a  
336 common property of the KRAB-domain-containing C2H2-zf (KZNF) subfamily *in vivo*<sup>27</sup>,  
337 but until now it has not been clear that the recruitment is defined almost entirely by the  
338 sequence specificity of the KZNFs alone. The combination of assays run here,  
339 particularly GHT-SELEX, extends earlier observations by pinpointing the exact binding  
340 sites, and demonstrating that these proteins typically have high specificity for these  
341 elements, because they bind preferentially to precisely the same elements *in vitro*.  
342 Binding preferentially to retroelements is not limited to KZNFs, but includes other C2H2-  
343 zf proteins and other classes of TFs. For example, binding sites for TIGD3, a  
344 transposon-derived TF which is closely related to JRK, are enriched for binding to L1s,  
345 SINEs, and DNA transposons<sup>15</sup>.

## 346 **Codebook PWMs predict TF binding in independent data and across cell types**

347 The Codebook project was conducted over a period of nearly six years, and during this  
348 time, several large-scale studies aimed at systematic ChIP-seq analysis of human TFs  
349 (e.g. ENCODE) were published<sup>11,57,58</sup>. Combined, the ENCODE data portal<sup>59</sup> and  
350 GTRD<sup>60</sup>, a compilation database, contain ChIP-seq and ChIP-exo peak data for 214 of  
351 the Codebook proteins, including 105 that were among the 166 with either “approved”  
352 Codebook ChIP-seq experiments (**Table S7**), or with ChIP-seq replicates that yielded  
353 reproducible peak sets<sup>15</sup>. We grouped both types of ChIP-seq data in our study and  
354 compared them to the external data. We first asked whether Codebook peak sets  
355 overlapped with these external peak sets for the same TF. Among the major ENCODE  
356 cell lines, the highest overlap values (Jaccard index) were found with experiments  
357 utilizing the same cell type (HEK293 cells) (**Figure S3A,B**). Slightly lower Jaccard  
358 values were obtained for experiments performed in HepG2 and other cell types, which  
359 would be expected given the altered chromatin profiles in different cell types, but over  
360 one-third were still clearly nonrandom (Jaccard > 0.1) (**Figure S3C**). Overlap scores  
361 with published K562 data, which dominate the external ChIP data due to a single large  
362 ChIP-exo study<sup>58</sup>, were much lower, overall (**Figure S3D**). We conclude from these  
363 analyses that the Codebook ChIP-seq data provide mainly new information.

364 We next asked how effectively the Codebook PWMs predict binding of TFs to peak sets  
365 in the published datasets. Consistent with the fact that the Codebook and external  
366 peaks often overlap, the Codebook PWMs had a median AUROC of 0.71 on the  
367 external HEK293 data, and were nearly as effective in predicting peak sets in other cell  
368 types (**Figure S3E**), illustrating that the Codebook PWMs are predictive across studies  
369 and cell types. We also asked how the predictive capacity of the Codebook PWMs  
370 compared to PWMs that appear in the latest versions of Factorbook<sup>12</sup>, JASPAR<sup>61</sup>, and  
371 HOCOMOCO<sup>13</sup> (**Table S8**). We identified 19 TFs with at least one successful Codebook  
372 ChIP-seq experiment and Codebook PWM, at least one external ChIP-seq experiment,  
373 and at least one PWM from an external database. In most cases, both the Codebook  
374 and external PWMs scored well on both Codebook and external peak sets (**Figure**  
375 **S3F,G**), supporting the validity of both PWMs and both peak sets. For seven proteins,  
376 low scores were obtained in at least some tests, however. For four of them, the  
377 independent Codebook *in vitro* data support the Codebook PWM; for two of the others,  
378 the external PWM scores poorly on Codebook peaks, while the Codebook PWM scores  
379 well on Codebook and external peak sets (**Figure S3H**). We conclude that the  
380 Codebook PWMs are generally more reliable than those published previously, likely  
381 because they are aided by confirmation of PWM performance across multiple data  
382 types that were not available in previous studies

## 383 **Codebook TF binding sites suggest functions for tens of thousands of conserved** 384 **elements**

385 Together, the Codebook assays and PWMs can be used to pinpoint genomic loci that  
386 are bound directly by each TF *in vivo* (i.e., in ChIP-seq), by identifying those that are  
387 also bound *in vitro* (i.e., GHT-SELEX), and that contain a PWM hit, thus allowing base-  
388 level resolution. We refer to these as “triple overlap” (TOP) sites, which are taken as the

389 overlap of the three sets (ChIP-seq, GHT-SELEX, and PWM hits) after applying  
390 optimized score thresholds for each (see **Methods** for details). This process produced a  
391 median of 455 TOP sites for 101 Codebook proteins, and a median of 3,014 TOP sites  
392 for 36 control TFs.

393 To gauge functionality of the TOP sites, we examined whether the pattern of per-  
394 nucleotide conservation<sup>13</sup> at each site is consistent with the TF's sequence preference  
395 driving local sequence constraint (see **Methods** for details). **Figure 4A** shows several  
396 examples illustrating that this approach readily detects apparent conservation of PWM  
397 hits, for both control and Codebook TFs. In total, 85/101 Codebook TFs (as well as  
398 33/36 controls) displayed conservation of at least one TOP site (FDR < 0.1), and in total  
399 we identified 121,785 such conserved TOP sites ("CTOP" sites) (83,621 for Codebook  
400 TFs and 38,164 for controls), encompassing 1,577,298 bases. These results,  
401 summarized in **Figure S4** and in greater detail in an accompanying manuscript<sup>15</sup>,  
402 provide strong support for the functional importance of Codebook TF binding sites in the  
403 genome.

404 Many of the CTOP sites were either overlapping or adjacent to CTOP sites for the same  
405 or other TFs. We grouped them into 50,375 clusters, based on proximity (allowing a  
406 maximum of 100 bases, to capture binding to different segments of what may be the  
407 same regulatory element). Codebook TFs with the largest number of CTOP sites were  
408 typically associated with CpG islands, which represented 37.5% of all the clusters  
409 (**Figure 4B**). The majority of protein-coding promoter CpG islands (58.7%,  
410 7,892/13,427) contained CTOP sites, with an average of 4.3 CTOP sites per CpG  
411 island. Moreover, 59/101 (58%) of all Codebook TFs had at least one CTOP site within  
412 a CpG island. An example CTOP that overlaps a CpG island is shown in **Figure 4C**.

413 The extent of specific, conserved, and intrinsic occupancy of CpG islands by many TFs  
414 of diverse classes is, to our knowledge, unexpected. The abundance of CG  
415 dinucleotides in CpG islands has been attributed primarily to their lack of methylation in  
416 the germline, rather than primary sequence constraint<sup>62</sup>. There is one class of TFs (the  
417 CXXC proteins) that is known to specifically recognize unmethylated CG dinucleotides  
418 and to modulate chromatin at promoters<sup>62</sup>, and we do observe this property for the  
419 CXXC proteins KDM2A, CXXC4, FBXL19, and TET3. Intriguingly, however, many of the  
420 Codebook TFs with CTOP sites in CpG islands recognize elaborate C/G rich motifs,  
421 rather than CG dinucleotides (**Figure 4C**).

422 CTOP clusters were also found in non-CpG island protein-coding promoters (**Figure**  
423 **4B**) (855/6,606 such promoters, defined as -1000 to +500 relative to TSS). These  
424 clusters are not dominated by any specific TFs, although some TFs are more prevalent  
425 than others (e.g. CTOPs for the controls ELF3 and CTCF, and Codebook TF ZBTB41,  
426 are each found in ~10% of all non-CpG promoters) (**Figure 4D**). **Figure 4E** shows an  
427 example of one such non-CpG promoter cluster, occurring early in the first intron of the  
428 TSPAN31 gene, which exhibits apparent conserved spacing and orientation of multiple  
429 Codebook TF binding sites. In contrast, CTOP clusters outside of promoters and CpG  
430 islands often contain just one or two CTOP sites (**Figure 4B**). One example is a very  
431 strongly conserved intergenic ZNF689 binding site found in an L1ME1 transposon; this

432 site is just over 100 bp from a predicted enhancer containing a CTCF binding site  
433 (**Figure 4F**).

434 A total of 42,200 distinct CTOP clusters (out of 50,375) overlapped catalogued  
435 conserved elements (UCSC PhastCons track), thus indicating a likely biochemical  
436 function for these elements. For the remaining 8,175, detection of functional elements  
437 from base-level scores is now augmented by the TF binding information. Relatively few  
438 CTOP clusters overlapped with known enhancers, however: only 4,768 are found in the  
439 extensive GeneHancer annotation set<sup>63</sup>, and 2,819 overlap with HEK293 enhancers (  
440 defined by ChromHMM<sup>15</sup>). This low overlap could be attributed to the relatively rapid  
441 evolution of enhancers<sup>64</sup>, or to lack of complete knowledge of enhancer identities. We  
442 also note that, even for well-studied TFs, most TOP sites were classified by our  
443 methods as not conserved, and that roughly half of the Codebook TFs had few or no  
444 conserved TOP sites (particularly the aforementioned retroelement-binding KZNFs)  
445 (**Figure S4**). Lack of conservation does not demonstrate that a sequence is not a  
446 functional binding site, however, as turnover in functional genomic binding sites of TFs  
447 is common<sup>65</sup>. This result is nonetheless consistent with the notion that many TF binding  
448 sites are coincidental, redundant, or serve(d) a purpose other than host genome  
449 regulation. In the accompanying manuscript<sup>15</sup>, we explore potential functions for  
450 proteins that frequently bind non-conserved sites in genomic “dark matter”.

#### 451 **Relationships between Codebook TFs, SNVs and chromatin**

452 Because the CTOP sites are evolutionarily constrained, we reasoned that they might  
453 also be less frequently associated with human sequence variation, and indeed, 92.6%  
454 of CTOPs lack SNPs and other common short variants, while only 82.1% of  
455 unconserved TOPs are variant-free. Both are depleted of common SNPs, however,  
456 when examined separately (Fisher’s exact test  $p \sim 2.4 \times 10^{-307}$  and odds ratio = 0.657,  $p$   
457  $\sim 0$  and ratio = 0.872, respectively). The CTOP SNPs also have a lower impact on PWM  
458 scores: on average, the relative PWM score for SNP-containing CTOP sequences  
459 declines by 0.027, while PWM scores for unconserved TOPs decline by 0.057 (median  
460 declines of 0.011 and 0.0285, respectively). CTOPs are furthermore depleted of  
461 common short indels (Fisher’s exact test,  $p \sim 1 \times 10^{-150}$ , ratio = 0.77), while unconserved  
462 TOPs (which often overlap with simple repeats) are enriched ( $p < 1 \times 10^{-150}$ , ratio =  
463 3.318), relative to genomic background. The depletion of common SNPs is consistent  
464 with ongoing purifying selection of CTOPs within recent human populations, and the  
465 association of SNPs with specific TFs should provide a ready means for directed study  
466 of the functionality of the encompassed SNPs.

467 We reasoned that the GHT-SELEX and ChIP-seq experiments would also allow direct  
468 assessment of allele-specific binding (ASB) of TFs, by quantifying allelic imbalance of  
469 read counts at SNVs. We note that the data were not initially intended for this purpose,  
470 and caveats included relatively low read counts, linked SNVs, and the fact that HEK293  
471 has an abnormal karyotype and was derived from a single individual. Nonetheless,  
472 there was sufficient coverage in the sequencing data to make 925,003 variant calls  
473 overlapping with dbSNP common SNPs (889,820 variant calls from 362 ChIP-seq  
474 experiments and 35,183 from 374 GHT-SELEX multi-cycle experiments), at 122,364

475 unique genomic locations (**Figure 5A, Figure S5A, Table S9**). 10,009 of these genomic  
476 locations were associated with 12,056 ASBs of 160 Codebook TFs and 46 positive  
477 controls in ChIP-seq (10,571 ASBs) or GHT-SELEX (1,485 ASBs) samples, i.e. there  
478 was a significant imbalance in the sequencing reads for the two alleles overlapping the  
479 respective SNPs. Among these ASBs, 3,569 also overlapped a PWM hit for the TF, and  
480 for 2,367 of them, the read count imbalance was concordant with the change in PWM  
481 scores, i.e. the allele with the higher read count also has a higher PWM score (**Figure**  
482 **S5A,B, Table S9**). (ASBs that do not overlap a PWM hit may be linked to a “causative”  
483 SNV, which may act indirectly). ASBs for control TFs were strongly enriched with  
484 previously-known ASBs of those TFs (ADAstra database, odds ratio of 5.7,  $p < 10^{-15}$ ,  
485 Fisher's exact test)<sup>66</sup>, and nearly three-quarters of ASBs coincided with eQTLs (GTEx  
486 database, odds ratio of 1.2,  $p < 10^{-15}$ , Fisher's exact test)<sup>67</sup> (**Figure S5C**), supporting the  
487 reliability of the detected ASBs as well as the validity of detected PWM hits.

488 Compared to whole-length peaks, TOP regions had an increased density of variant calls  
489 (~258 sufficiently covered variants per Mb in TOPs, versus 52 per Mb for peaks), and a  
490 larger fraction of ASB calls in SNVs (30%, compared to 9% for full peaks), presumably  
491 due to detection bias from higher ChIP-Seq or GHT-SELEX coverage at the TOPs.  
492 Nonetheless, variants in TOPs had a significantly higher predicted effect on protein  
493 binding (i.e. PWM score change) for both controls and Codebook TFs ( $p < 2.22 \times 10^{-5}$   
494 and  $p < 2.98 \times 10^{-12}$ , Mann-Whitney U test), relative to full peaks or non-ASB SNPs  
495 overlapping PWM hits (**Figure 5B**). Thus, the ASBs in TOPs are more likely to induce  
496 an effect than those elsewhere within peaks, presumably because they represent direct  
497 TF binding.

498 Among the mechanisms connecting TF binding to biological function are TF-mediated  
499 chromatin state changes. Hence, in heterozygotes, variant-dependent TF binding may  
500 co-occur with allele-specific chromatin accessibility variants (ASVs) (**Figure 5A**), which  
501 are SNVs with imbalanced read counts in ATAC-seq and/or DNase-seq experiments.  
502 To ask whether the Codebook TFs may be involved in control of ASVs, we utilized the  
503 UDACHA database, which contains ASVs from 577 ATAC-seq and 321 DNase-seq  
504 datasets from individual cell types<sup>68</sup> (**Table S9, Figure S5D**). Using a multi-tiered  
505 procedure (see **Methods**), we identified cases in which (1) ASVs in a specific cell type  
506 overlap significantly with PWM hits for a TF in the Codebook motif collection, (2) the  
507 change in the PWM score is concordant with the read imbalance in the ASVs, (i.e.  
508 stronger predicted binding is associated with more accessible chromatin), and (3) the  
509 concordance is significant across cell types detected in step (1). This procedure  
510 identified 53 TFs whose PWM hits were found often at, and concordant with, ASVs  
511 (**Figure S5E**). Twenty of these TFs were positive controls including well-known pioneers  
512 or activators (such as SOX2, GABPA, or JUN/FOS-family TFs), while 33 were  
513 previously unexplored Codebook TFs, including ZNF70, GRHL3, MYPOP, SP140(L),  
514 and DMTF1. An example ASV for ZNF70, in a region upstream of the PTMS gene that  
515 is annotated with multiple ENCODE enhancer elements is shown in **Figure 5C**.

516 For 34 of these 53 TFs, there was at least one ASV-overlapping TOP site (the non-TOP  
517 sites may represent sites that are not bound in HEK293). To assess whether ASVs in  
518 PWM hits have a greater effect at TOP sites than in other regions, we first removed

519 cases in which the TF does not appear to impact chromatin directly, by grouping the  
520 TFs into ASV-concordant (i.e. having overall concordance between ASVs and PWM hits  
521 in ChIP-seq or GHT-SELEX peaks; 18 TFs), and others (16 TFs). We separated the  
522 ASV-concordant group into Codebook and control TFs. For each of the groups, we then  
523 calculated the concordant-to-discordant ratio for loci that corresponded to PWM hits that  
524 are non-ASV for that TF, ASV, ASV in TF's peaks, and ASV in TOPs, and observed an  
525 overall monotonic increase in concordance (**Figure 5D**). Thus, the highest-confidence  
526 Codebook TF binding sites for these TFs are those most likely to impact the chromatin  
527 state. Moreover, the fraction of ASVs within PWM hits also increased monotonously as  
528 the ASV confidence increased, and the ASVs preferably occur at binding site positions  
529 that are most important for the PWM score (**Figure 5E, Figure S5F**), further supporting  
530 relevance of the TF sequence preferences.

531 Overall, the Codebook motifs provide a valuable resource for SNV interpretation,  
532 including identification of mechanisms that underpin variation in chromatin and  
533 transcription.

### 534 **Lessons from Codebook: prospects for a complete human TF motif collection**

535 Codebook yielded several clear outcomes, and guidance for future efforts. The high  
536 success rate is particularly striking. We obtained motifs for 177 previously  
537 uncharacterized human TFs, a number larger than the entire TF repertoire for many  
538 eukaryotes<sup>69</sup>. The selected PWMs for most of these TFs are unique, and unlike any  
539 previous TF motif. Most are from C2H2-zf proteins, and most C2H2-zf proteins analyzed  
540 were successful. Thus, a majority of putative and uncharacterized human TFs are *bona*  
541 *fide* TFs, and not annotation errors. We envision that the data produced will be broadly  
542 and immediately useful for a variety of applications. Motifs (especially as PWMs) are a  
543 standard component of the computational genomics toolkit, due to their utility in a range  
544 of tasks ranging from identification of key regulatory factors to building and interpreting  
545 models of gene expression<sup>70-73</sup>. For example, differential binding of TFs to noncoding  
546 SNVs (Single Nucleotide Variants) is thought to be a major mechanism by which these  
547 variants contribute to phenotypic differences<sup>74</sup>, and the Codebook data therefore  
548 provide vital new information for the analysis of *cis*-regulatory variation.

549 A key technical demonstration of the Codebook project is that the simultaneous  
550 application of multiple experimental strategies and multiple motif-derivation and motif-  
551 scoring strategies was highly beneficial. No single experiment type or data analysis  
552 approach dominated all others, or was universally successful, although specific assays  
553 were more or less advantageous for different classes of proteins (as evident in **Figure**  
554 **1**). For example, PBMs were uniquely successful with AT-hook proteins, while ChIP-seq  
555 and SELEX variants were most successful for C2H2-zf proteins. We caution that there  
556 are confounding variables limiting what conclusions can be drawn regarding the  
557 strengths and weaknesses of experimental platforms. The protein production and  
558 purification method can differentially impact success of specific DBD classes, even  
559 when the same assay is used, and the different assays we employed were tied to  
560 different affinity tags and expression systems. Data pre-processing (i.e. read filtering

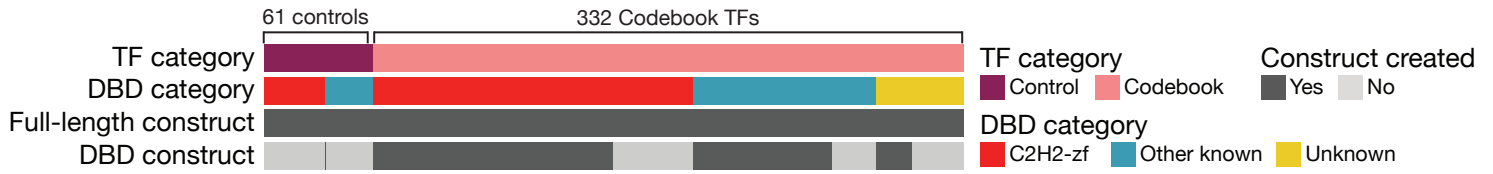
561 and background estimation) is an additional variable that we did not systematically  
562 explore, but is known to impact all of the assays used here.

563 As noted above, a subset of the Codebook TFs, as well as other poorly characterized  
564 TFs, have been analyzed by others since our study began. To evaluate the current  
565 scope of known human TF specificities, we surveyed JASPAR, HOCOMOCO, and  
566 Factorbook for PWMs for putative TFs that were not included in this study or not found  
567 among 177 Codebook successes. These databases reported PWMs for 107 proteins,  
568 63 of which we had tested, and 44 were among the 95 putative TFs not included in our  
569 experiments. We manually curated these external PWMs, using procedures similar to  
570 those we applied to our own data, to assess whether they are likely to represent the  
571 *bona fide* specificity of the TF analyzed. Many of them were comprised of simple  
572 repeats (which are common artifacts in virtually all assays) or appeared to correspond  
573 to indirect binding and/or recruitment by other TFs in ChIP-seq (See **Table S8** for  
574 annotations and classification, and **Figure S1** for examples of nonspecific, concordant,  
575 and likely correct PWMs in the external datasets).

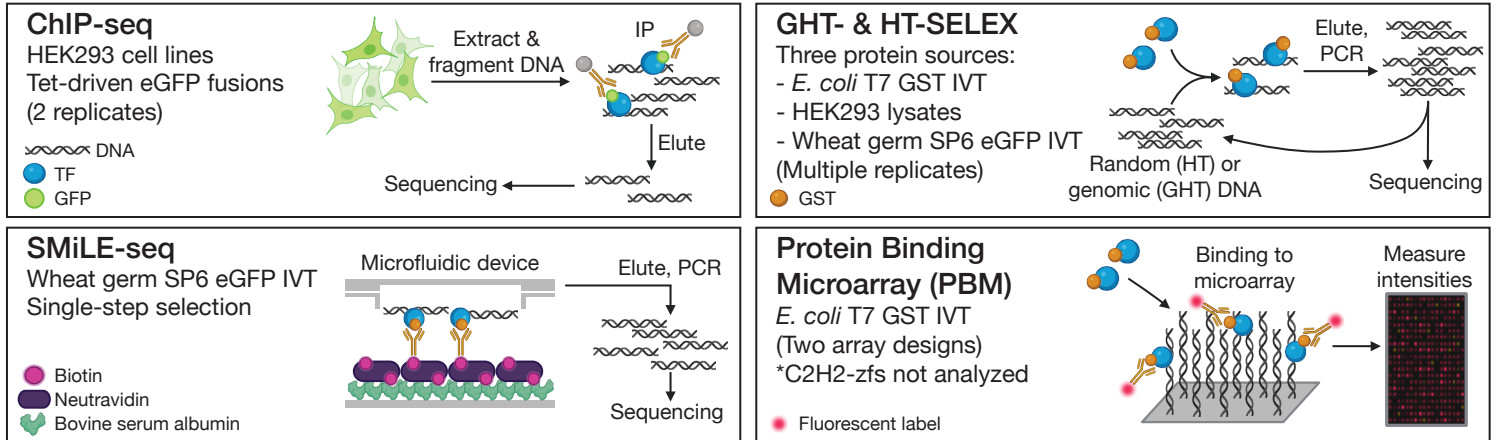
576 Based on this curation, 33 additional human TFs (i.e. beyond the 177 described here)  
577 have at least one plausible motif available in datasets that have been performed since  
578 our 2018 TF census<sup>1</sup>, leading to a total of 1,421 human TFs now with characterized  
579 sequence specificities (**Figure 6** and **Table S10**). Altogether, only 175 proteins with  
580 conventional DBDs now lack known sequence specificity. Not all proteins with such  
581 domains are necessarily TFs; for example, one systematic trend we observed is that  
582 almost all 36 proteins we tested with only a single C2H2-zf domain failed in every assay  
583 (**Figure 6**). At the same time, however, new DBD classes continue to appear, such as  
584 the aforementioned BEN, CGGBP, Dachshund, and C-clamp. Some TFs may bind only  
585 to methylated DNA, and ongoing advances in the prediction of protein and protein-DNA  
586 structures<sup>45</sup> have the potential to identify additional candidates for sequence-specific  
587 DNA binding. Thus, while completion of the objective to obtain a motif for every human  
588 TF now appears much closer, the list of likely human TFs continues to evolve.

589 Many of the Codebook TFs are now among the best characterized human DNA-binding  
590 proteins in terms of their sequence specificity. As illustrated in the accompanying  
591 papers (**Table S1**), and consistent with previous benchmarking efforts<sup>18,32</sup>, validation  
592 across platforms can lead to very different conclusions regarding PWM reliability.  
593 Moreover, obtaining *in vivo* and *in vitro* binding to the genome facilitates  
594 disentanglement of direct and indirect binding, as well as the contribution of the cellular  
595 environment. Obtaining *in vitro* binding data to both genomic-sequence and random-  
596 sequence DNA can provide insight into the importance of local sequence context. Only  
597 a small handful of the 1,000+ previously characterized TFs have such a combination of  
598 data types. A much better perspective on human gene regulation and genome function  
599 and evolution could presumably be obtained from generation of such data for all human  
600 TFs.

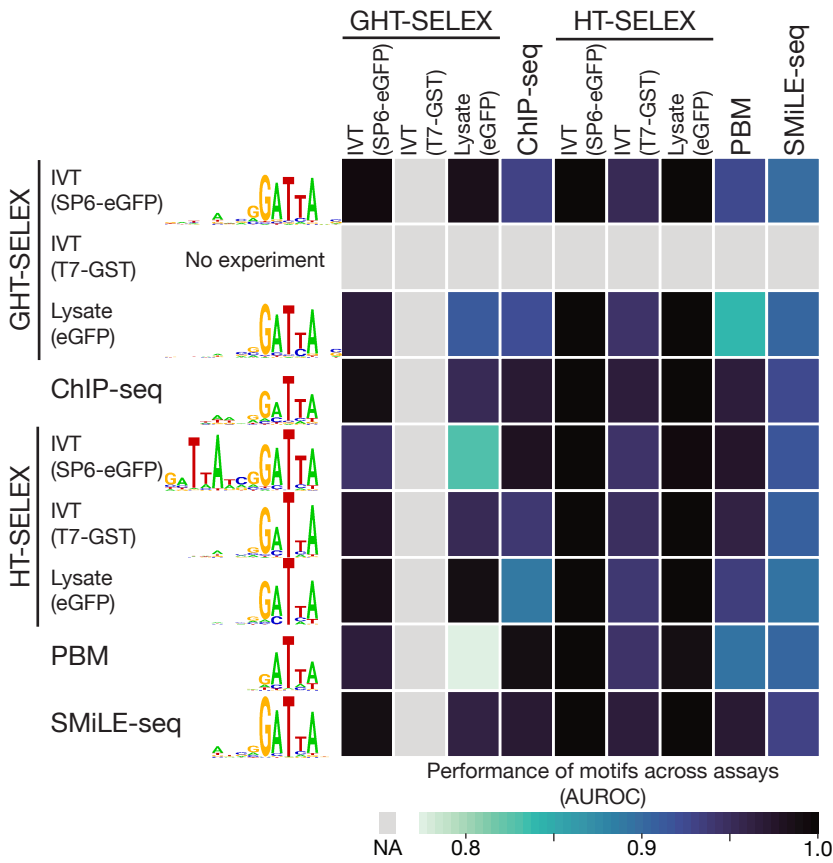
## TF selection & construct creation



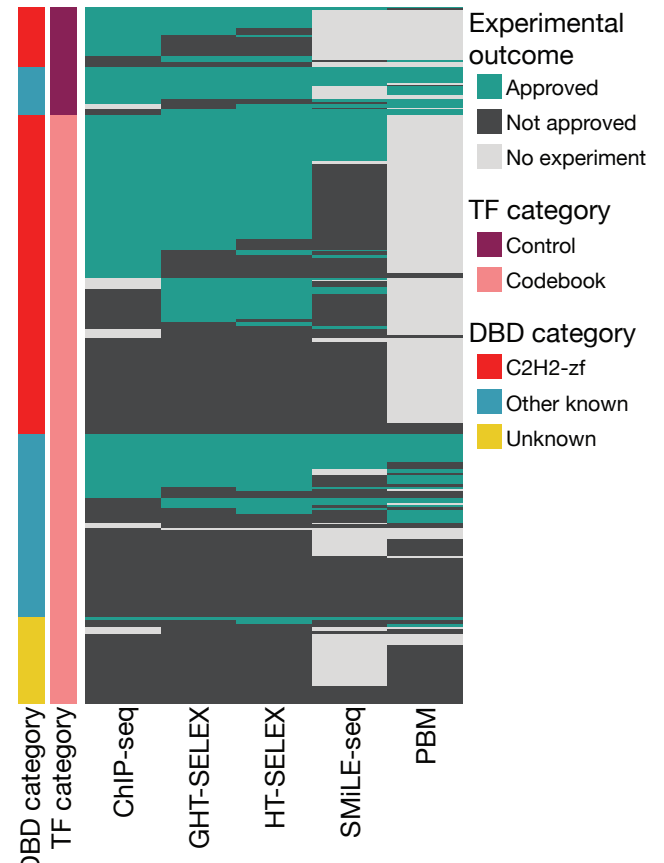
## Assays



## Motif derivation & benchmarking



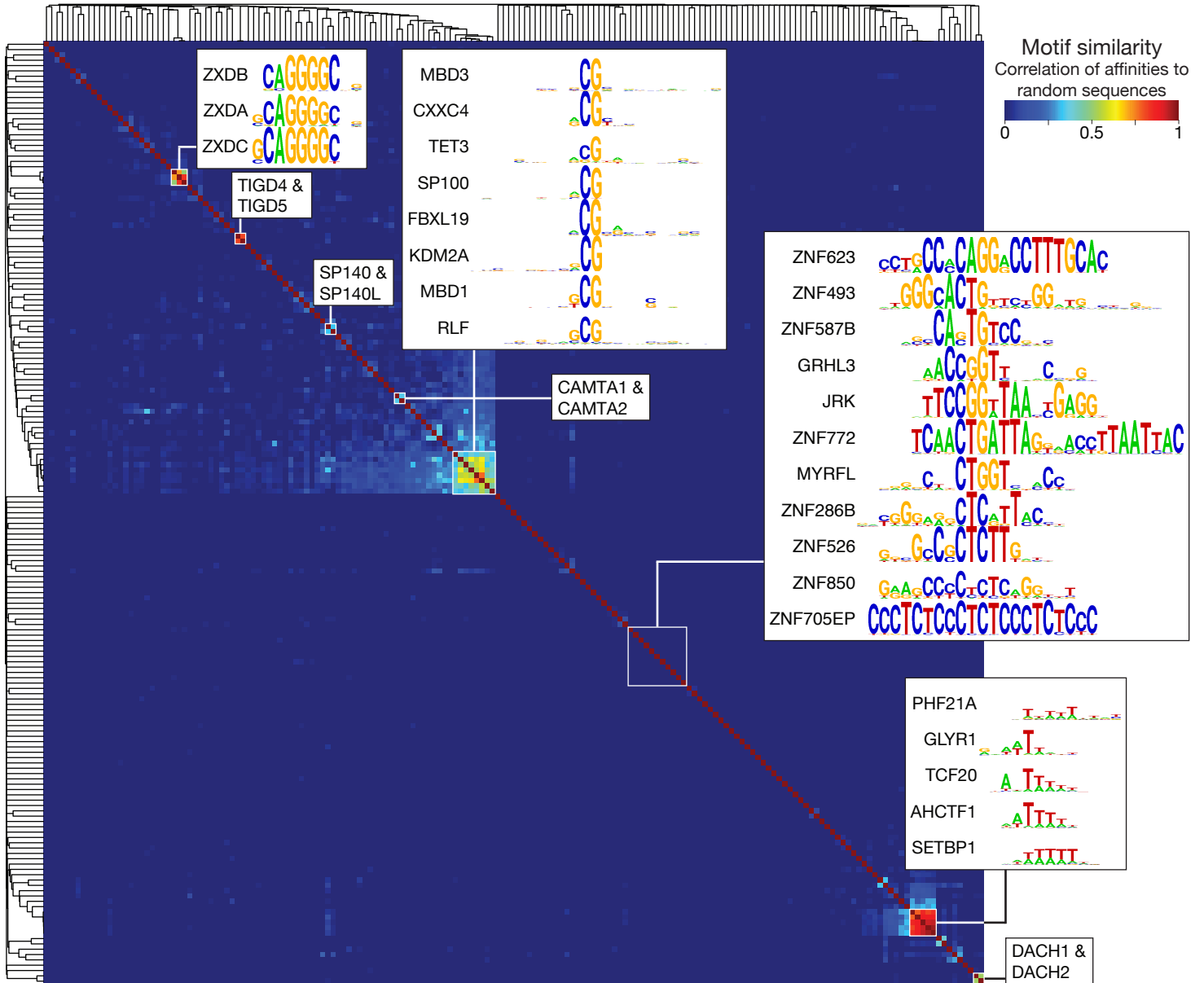
## Expert curation



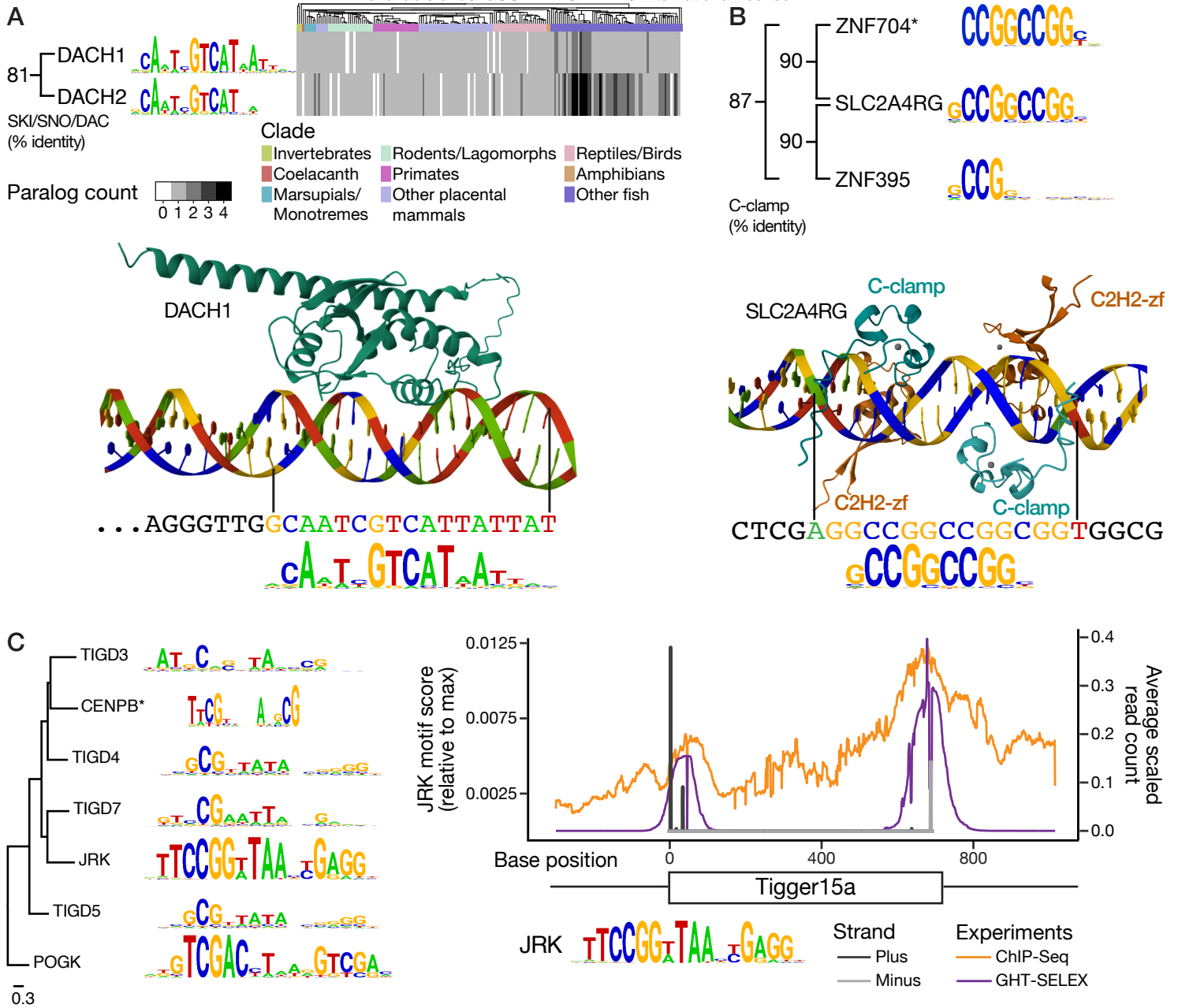
## Data analysis and exploration



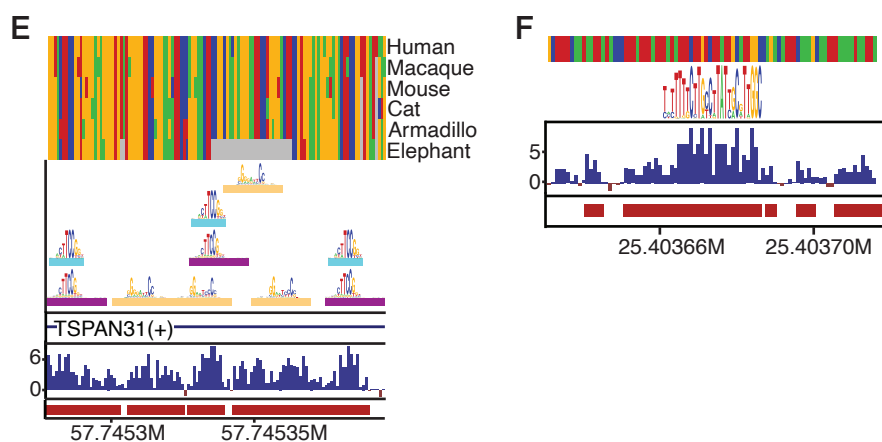
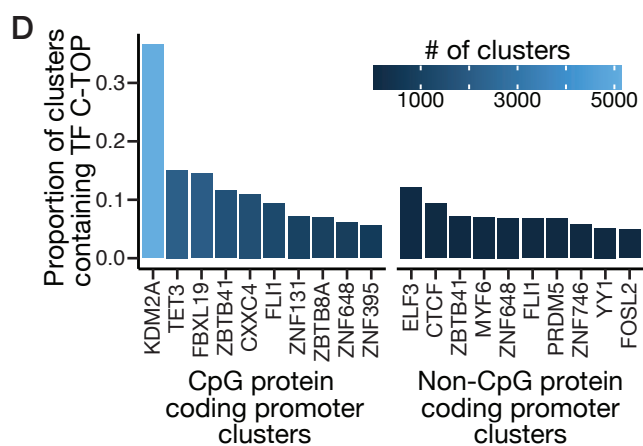
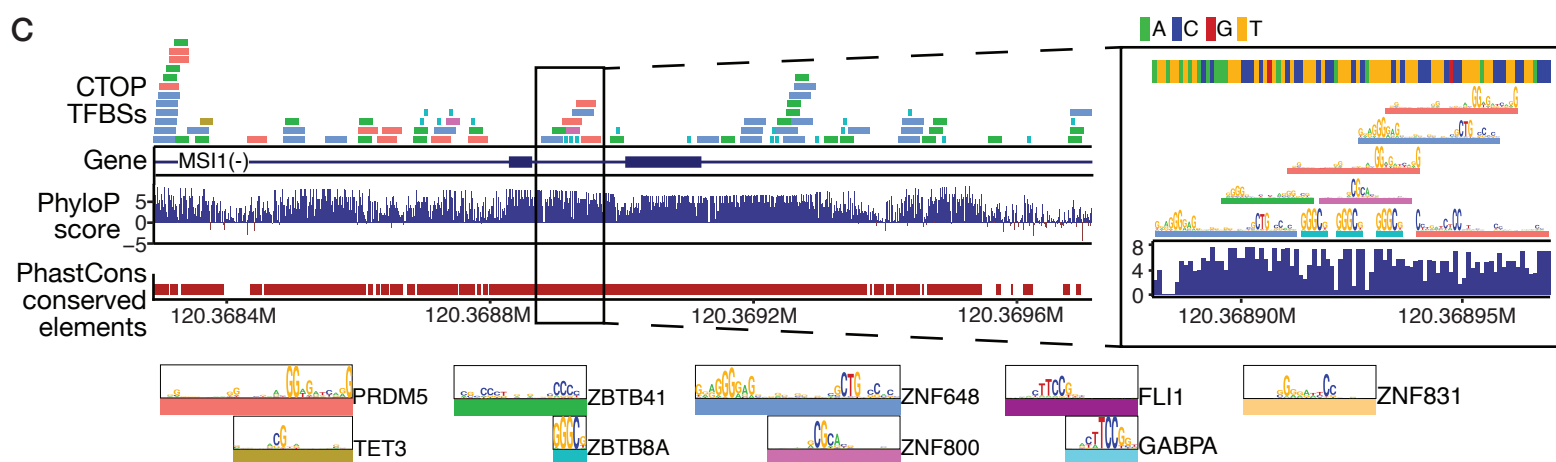
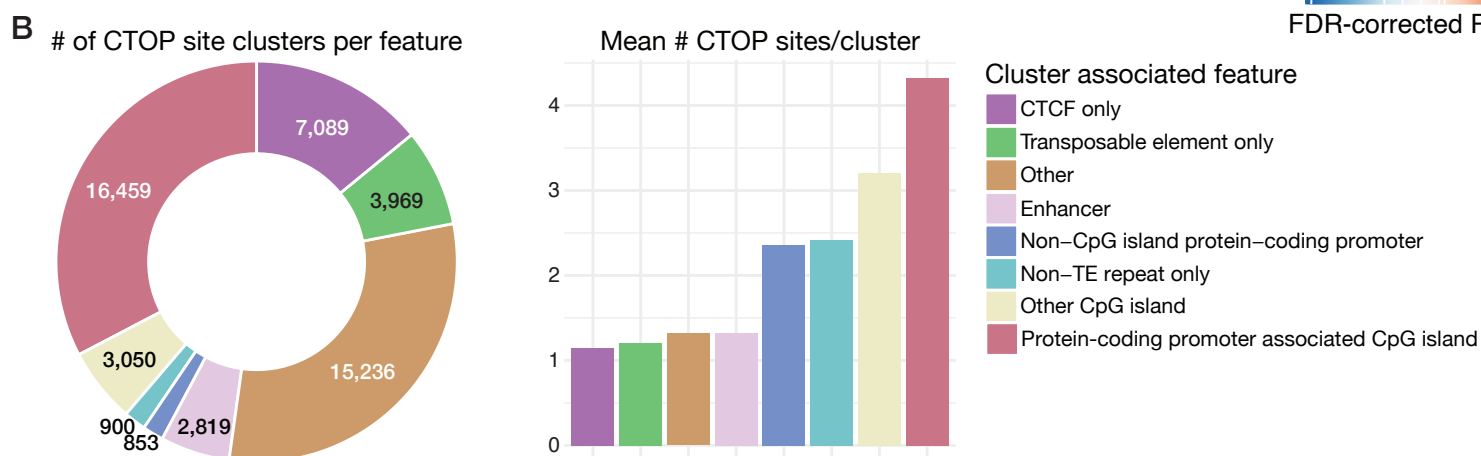
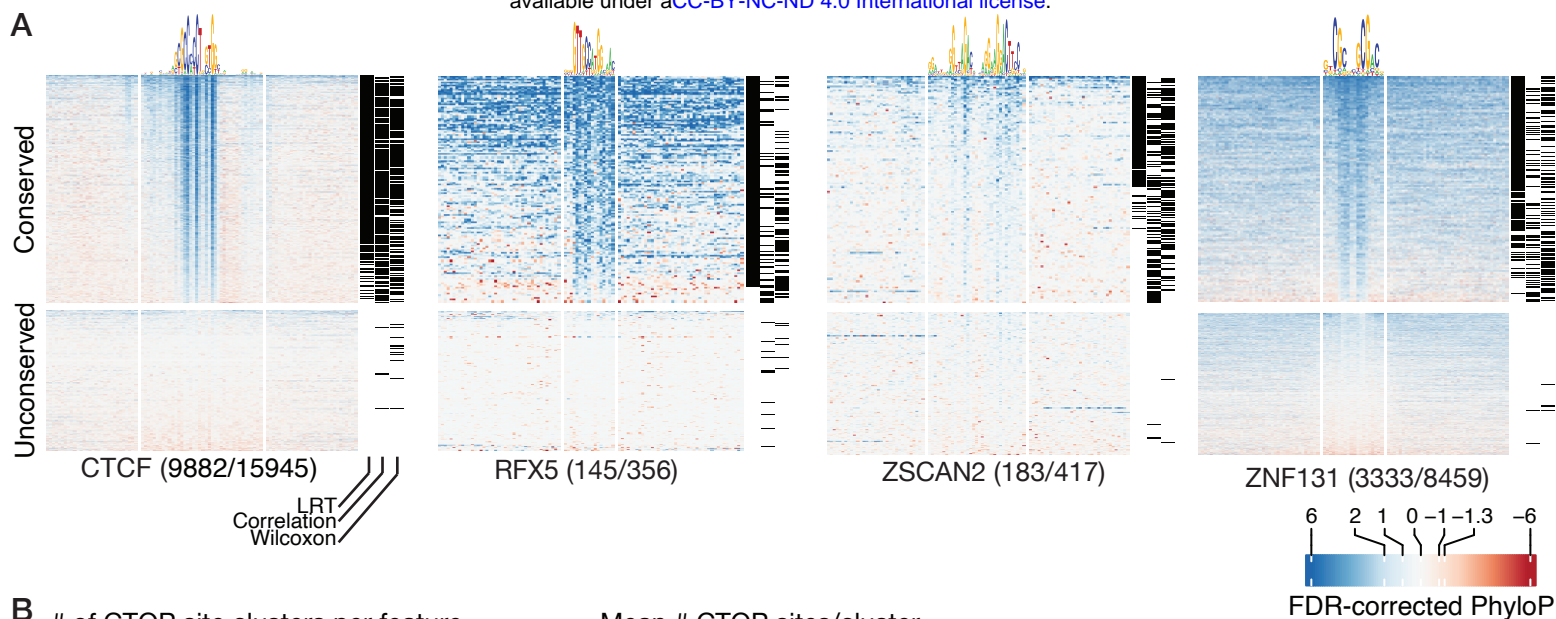
**Figure 1. Codebook project overview.** *Top*, Categories of 393 TFs assayed and their associated constructs. *Middle*, Graphical summary of assays employed. *Bottom left*, Example of performance (as AUROC) of the best performing PWM for TPRX1, for each combination of experiment type – one for motif derivation (rows), and one for motif testing (columns). *Bottom right*, Depiction of the approval process for each individual experiment, including comparison of motifs and/or binding sites between replicates, evaluation of motifs across experiments, and motif similarity between related TFs (see **Experiment evaluation by expert curation**). Heatmap shows approved experiments for all 393 TFs across all experiment types.



**Figure 2. Similarity of Codebook TF motifs.** Symmetric heatmap displaying the similarity between expert-curated PWMs for each pair of Codebook TFs, clustered by Pearson correlation with average linkage. The PWM similarity metric is the correlation between pairwise affinities to 200,000 random sequences of length 50, as calculated by MoSBAT<sup>28</sup>. Pullouts and labels illustrate specific points in the main text.

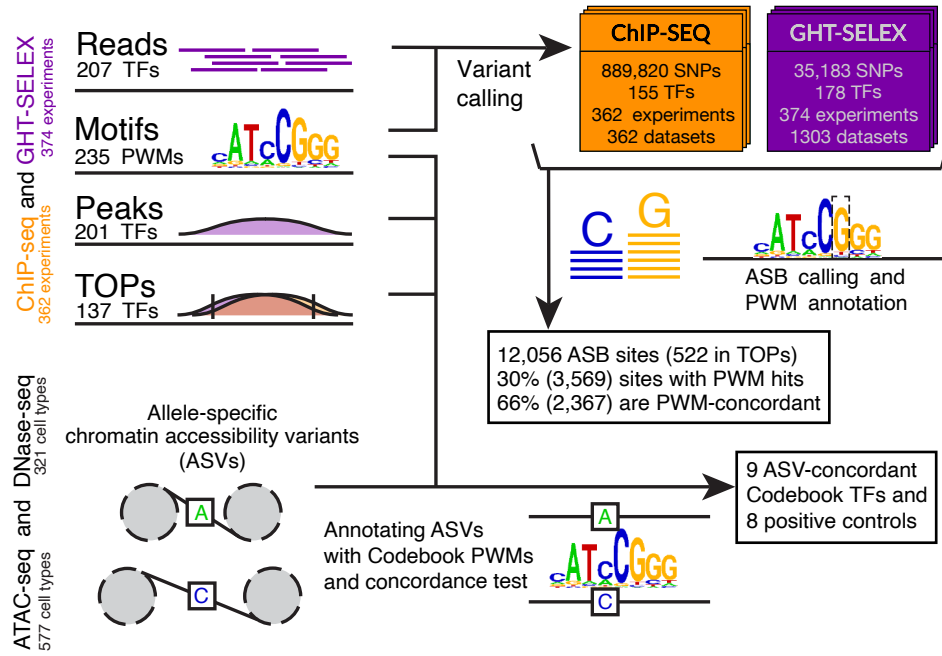


**Figure 3. Neglected DNA-binding domains.** Overview of new motifs for previously understudied TF families. **A**, *Top*, Number of DACH1 and DACH2 orthologs (union of one-to-one and one-to-many) across Ensembl v111 vertebrates and selected invertebrates. Species order reflects the Ensembl species tree. *Bottom*, AlphaFold3-predicted structure of the DACH1 SKI/SNO/DAC region (residues 130 – 390) bound to an HT-SELEX ligand sequence with a high-scoring PWM hit. **B**, *Top*, Sequence logos and sequence relationships of human C-Clamp domains (\*ZNF704 motif from <sup>50</sup>). *Bottom*, AlphaFold3-predicted structure of two full-length SLC2A4RG proteins bound to a CTOP sequence with flanking sequences (chr17:48,048,369-48,048,401), and four Zn<sup>2+</sup> ions (grey). The remainder of the proteins (beyond the C-clamp and C2H2-zf domains) are hidden, for visual simplicity. **C**. *Left*, Sequence logos of human TFs that are derived from the domestication of *Tigger* and *Pogo* DNA transposon DBDs elements and have known DNA binding motifs. Tree is a maximum-likelihood phylogram from FastTree<sup>92</sup>, using DBD sequence alignment with MAFFT L-INS-I<sup>93</sup>, rooted on POGK, which is derived from an older family of Tigger-like elements<sup>94,95</sup>. Sequence logos are Codebook-derived, except for CENPB<sup>96</sup>. *Right*, average per-base read count over Tigger15a TOPs in the human genome, for JRK ChIP-seq (orange) and GHT-SELEX (purple), with sequences aligned to the Tigger15a consensus sequence. JRK PWM scores at each base of the Tigger15a consensus sequence are shown in black (plus strand) and grey (minus strand).

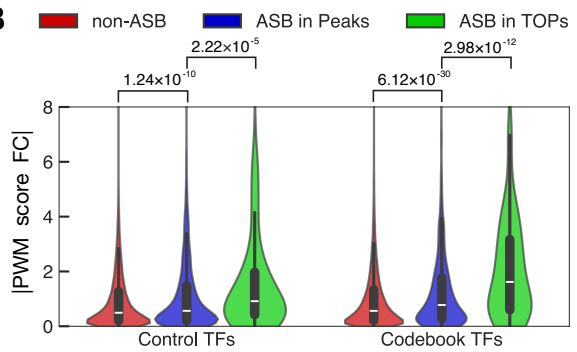


**Figure 4. Conservation of Codebook TF binding sites and association with genomic features.** **A**, Heatmaps of phyloP scores over the PWM hit and 50 bp flanking for TOP sites for four TFs (two controls and two Codebook TFs). Statistical test results (see main text and **Methods**) are indicated at right. **B**, *Left*, Donut plot displays the proportion and number of clusters of conserved TOP (CTOP) sites that overlap the genomic features indicated. *Middle*, Bar plot displays the mean # of individual CTOPs contained within clusters that overlap the examined genomic regions. **C**, A 1,420-base, CpG-island-overlapping CTOP cluster (chr12:120368293-120369713). Zoonomia 241-mammal phyloP scores and Multiz 471 Mammal alignment PhastCons Conserved Elements are shown. **D**, Bar plot of the frequency of TFs with CTOPs that occur most frequently in CTOP clusters that overlap CpG and non-CpG protein coding promoters, respectively. **E**, CTOP cluster overlapping the non-CpG promoter at chr12:57,745,278-57,745,396. **F**, CTOP site for the KRAB-C2H2-zf protein ZNF689, overlapping an L1ME4a located at chr16:25,403,631-25,403,717.

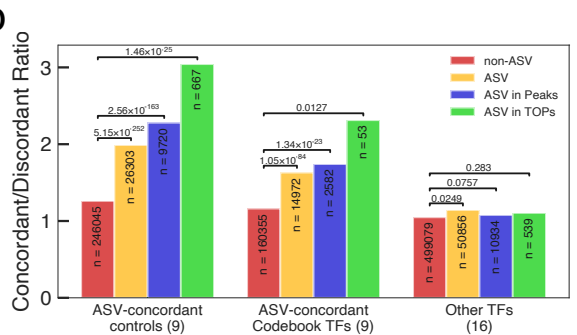
**A**



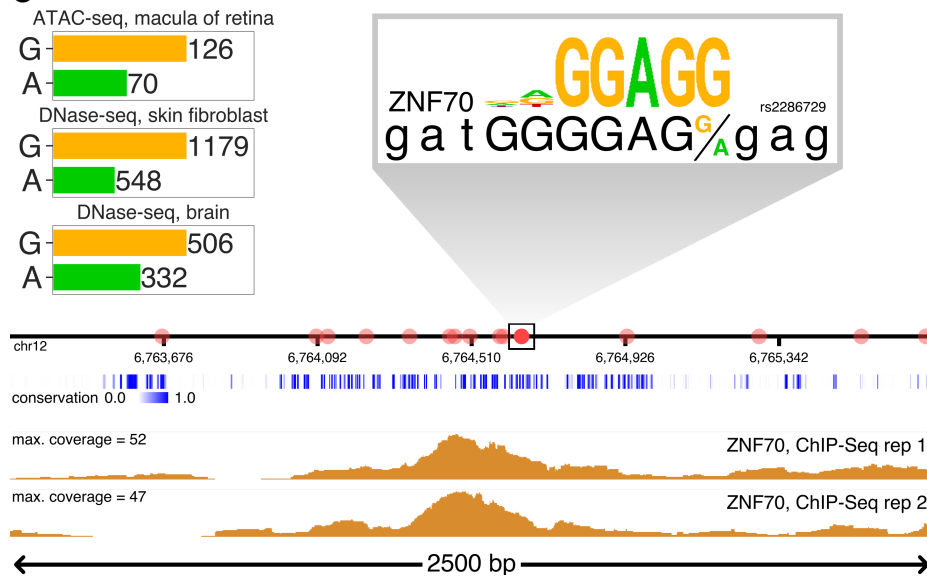
**B**



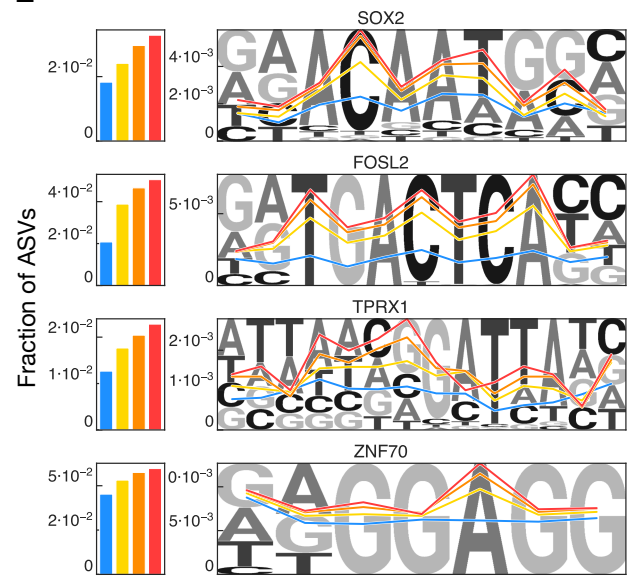
**D**



**C**



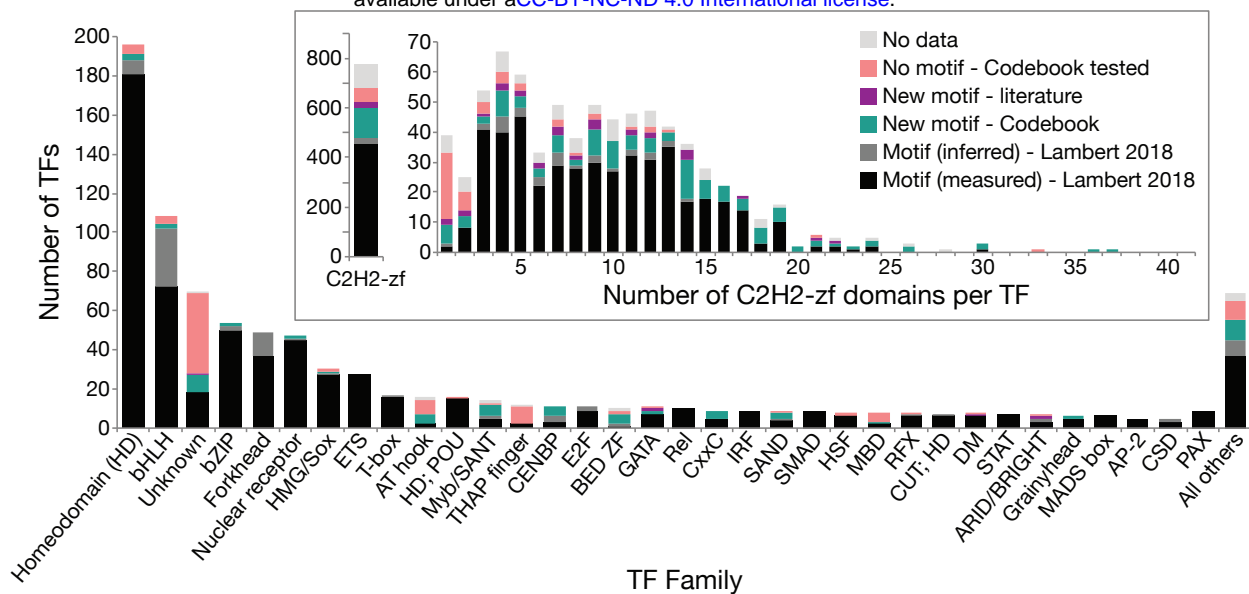
**E**





**Figure 5. Allele-specific transcription factor binding and chromatin accessibility.**

**A**, Scheme of the analysis: identification of allele-specific binding sites (ASBs) from Codebook ChIP-Seq and GHT-SELEX data and annotation of allele-specific chromatin accessibility variants (ASVs) with the Codebook motifs. **B**, Distribution of PWM score (log-odds) fold changes between alleles for non-ASB SNPs, ASBs in peaks, and ASBs in TOPs. *Left*, 32 positive control TFs, *Right*, 85 Codebook TFs. P-values: Mann-Whitney U test. **C**. An example ASV for ZNF70, in chr12:6,763,200-6,765,850, around 1kb upstream of the PTMS gene. Onset shows the exact location of the ASV (with A/G alleles) together with the corresponding PWM hit. Allelic read counts for three available ATAC- and DNase-seq samples are shown on the side. **D**. The ratio of concordant-to-discordant PWM hits for <SNP, TF> pairs for non-ASVs (red), all ASVs (yellow), ASVs overlapping with peaks (blue), and ASVs in TOPs (green). P-values: Fisher's exact test. **E**. *Left*, Fraction of ASVs overlapping with PWM hits for four example TFs, using 4 different thresholds on ASV significance: all SNPs (blue), 25% FDR ASVs (yellow), 10% FDR ASVs (orange), and 5% FDR ASVs (red). *Right*, Fraction of ASVs at each location within the genome-wide PWM hits of the representative TFs using four thresholds (same colors as in bar plots).



**Figure 6. Motif coverage of human TFs, by DBD family.** TFs are categorized into structural classes based on Lambert et al.<sup>1</sup>. See **Table S10** for underlying information.

## 601 METHODS

602 **Plasmids and inserts.** Sequences and accompanying information are given in **Table**  
603 **S3**, and the relationships between constructs, samples, and experiments are compiled  
604 in the information provided online at [codebook.ccb.utoronto.ca](http://codebook.ccb.utoronto.ca). Briefly, we selected  
605 Codebook TFs (and their DNA-binding domains catalogued) from information  
606 accompanying Lambert 2018<sup>1</sup>) and posted at <https://humantfs.ccb.utoronto.ca>. Inserts  
607 named with an “-FL” suffix correspond to the full-length ORF of a representative isoform  
608 of the protein. Those with a “-DBD” suffix contain all of the predicted DBDs in the protein  
609 flanked by either 50 amino-acids, or up to the N or C-terminus of the protein. Those with  
610 a “-DBD1”, “-DBD2” or “-DBD3” suffix contain a subset of the DBDs present in the  
611 proteins; these were designed manually, mainly for large C2H2-zf arrays. Inserts were  
612 obtained as recoded synthetic ORFs (BioBasic, US) flanked by *Ascl* and *Sbfl* sites, and  
613 subcloned into up to three plasmids: (i) pTH13195, a tetracycline-inducible, N-terminal  
614 eGFP-tagged expression vector with FLiP-in recombinase sites<sup>10</sup>; (ii) pTH6838, a T7-  
615 promoter driven, N-terminal GST-tagged bacterial expression vector<sup>75</sup>, and (iii)  
616 pTH16500 (pF3A-ResEnz-egfp), an SP6-promoter driven, N-terminal eGFP-tagged  
617 bacterial expression vector, modified from pF3A–eGFP<sup>9</sup> to contain the two restriction  
618 sites after the eGFP.

619 **Protein production.** Each experiment used a protein expressed from one of the  
620 following systems: (a) FLiP-in HEK293 cells (catalog number: R78007), induced with  
621 Doxycycline for 24 hours, used for inserts in pTH13195; (b) PURExpress T7  
622 recombinant IVT system (NEB Cat.#E6800L), for inserts in pTH6838; or (c) SP6-driven  
623 wheat germ extract-based IVT (Promega Cat#L3260), for inserts in pTH16500.

624 **DNA binding assays.** We followed previously-described methods for ChIP-seq<sup>10</sup>,  
625 PBMs<sup>32</sup>, and SMiLE-seq<sup>9</sup>. Detailed descriptions of GHT-SELEX, HT-SELEX, ChIP-seq,  
626 and SMiLE-seq data collection and initial analysis are found in the accompanying  
627 papers (**Table S1**). For PBMs, we analyzed proteins on two different PBM arrays (HK  
628 and ME), with differing probe sequences<sup>76</sup>.

629 **Data processing and motif derivation.** The accompanying paper<sup>17</sup> describes motif  
630 derivation and evaluation in detail. Briefly, after initial data processing steps, we  
631 obtained a set of 'true positive' (likely bound) sequences for each individual experiment.  
632 (721 / 4,873) experiments were removed at this step, due to a low number of peaks, or  
633 other technical issues, as documented in **Table S5**). We then applied a suite of tools to  
634 a training subset of the data from each experiment, and tested the resulting motifs on a  
635 test subset of the data from the same experiment, and also on the independent data for  
636 the same TF (i.e. the test sets from all other experiments done for the same TF). We  
637 employed a binary classification regime for all experiments and all motifs, and scored  
638 the motifs by a variety of criteria such as the areas under the receiver operating  
639 characteristic (AUROC) or the precision-recall curve (AUPRC).

640 **Experiment evaluation by expert curation.** To gauge the success of individual  
641 experiments, we employed an “expert curation” workflow with an initial voting scheme in  
642 which a committee of annotators gauged whether individual experiments should be

643 “approved”, i.e. included in subsequent analyses. All experiments were examined by at  
644 least three annotators. A subcommittee (AJ, IVK, and TRH) jointly resolved all cases of  
645 disagreement among initial annotators (~300 experiments), and then reviewed all  
646 approved experiments. Annotators had available an early version of the MEX portal  
647 (<https://mex.autosome.org>) containing results of all PWMs scored against all  
648 experiments, and were tasked with gauging whether the experiments yielded PWMs  
649 that were similar across experiments, or scored highly across experiments. Annotators  
650 also considered whether the motif was consistent with those for other members of their  
651 protein family (e.g. BHLHA9 yielded an E-box-like motif, CAnCTG), and/or similar  
652 between closely related paralogs (e.g. ZXDA, ZXDB, and ZXDC all yielded similar  
653 motifs). We also considered whether (and how many) “peaks” were obtained from ChIP-  
654 seq or GHT-SELEX, and whether these peaks were common to independent  
655 experiments (e.g. both ChIP-seq and GHT-SELEX). Annotators were further given a  
656 measure of similarity between Codebook PWMs and any PWMs in the public domain,  
657 as well as enrichment of known or suspected common contaminant motifs in any  
658 experiment.

659 **Post-evaluation peak processing.** After identification of “approved” experiments, we  
660 re-derived peaks sets for ChIP-seq and GHT-SELEX experiments in order to obtain a  
661 single peak set for each TF, as described in the accompanying papers<sup>14,15</sup>. Briefly, for  
662 ChIP-seq we repeated the peak calling using MACS2 and experiment-specific  
663 background sets, using a procedure previously described<sup>10</sup>, then merged the peak sets  
664 for replicates of the same TF with BEDTools merge<sup>77</sup> (see accompanying manuscript<sup>15</sup>:  
665 “ChIP peak replicate analysis and merging”). We derived GHT-SELEX peaks using a  
666 novel method that calculates enrichment of reads in each cycle, and treats different  
667 experiments as independent statistical samples in order to obtain a single enrichment  
668 coefficient per peak<sup>14</sup>.

669 **Expert motif curation.** For this study, to identify a single representative PWM for each  
670 TF, we first compiled a set of highest-scoring candidate PWMs for each TF (as  
671 summarized above and elsewhere<sup>17</sup>, then ran additional tests with them, utilizing the  
672 reprocessed peak data, and manually evaluated the outputs. We first took the union of  
673 three sets of 20 PWMs for each TF: the 20 PWMs with the highest AUROC (as  
674 calculated elsewhere<sup>17</sup>) on (i) any approved ChIP-seq experiment for the given TF, (ii)  
675 any approved GHT-SELEX experiment for the given TF, and (iii) any approved HT-  
676 SELEX experiment for the given TF. These PWMs were selected regardless of the data  
677 set from which they were derived. We then reassessed these PWMs against ChIP-seq  
678 and GHT-SELEX data with two parallel methodologies. First, we recalculated AUROC  
679 for each of the candidate top PWMs on the merged, thresholded sets of ChIP-seq  
680 peaks ( $P < 10^{-10}$ )<sup>15</sup> using AffiMX<sup>28</sup> to score each peak. We generated negative sets  
681 using BEDTools shuffle<sup>77</sup> with the *-noOverlapping* option to create sets of random  
682 genomic regions with the same number of peaks, and the same peak width distribution  
683 as the corresponding ChIP peak sets. We used the same technique to calculate  
684 AUROC values for GHT-SELEX, with thresholded peak sets (using a “Kneedle”<sup>78</sup>  
685 specificity value of 30 in the sorted enrichment values<sup>15</sup>). In parallel, we calculated the  
686 Jaccard index to measure the overlap between PWM hits (identified by MOODS<sup>79</sup> with -  
687 p 0.001) vs. the ChIP-seq peaks, and GHT-SELEX peaks, as two separate

688 measures. The overlap in each case was maximized by applying different thresholds on  
689 the peak sets and choosing the cutoff at which the Jaccard index was the highest<sup>14</sup>. We  
690 then applied expert curation (by a committee consisting of AJ, TRH, AF, KUL, RR, MA,  
691 and IY) to choose a single representative PWM with high performance on all compiled  
692 scores that, all else equal, also reflects reasonable expectation from the DBD class  
693 (including recognition-code predicted motifs, see accompanying manuscript<sup>14</sup>) and has  
694 high information content.

695 **Motif degeneracy analysis.** We adjusted the information content (IC) of PWMs on a  
696 per-base-pair basis, with all locations boosted equally, by incrementally scaling weights  
697 (e.g. probabilities in the PWM) until the PWM reached an adjusted to an average IC of 1  
698 bit per base pair. The script, “logo\_rescale.pl”, is available at  
699 <https://gitlab.sib.swiss/EPD/pwmscan>.

700 **Comparison to external peak sets and PWMs.** We downloaded comparison peak  
701 sets from GTRD<sup>60</sup> and ENCODE (4.12.2023)<sup>59</sup>, for all Codebook TFs. We then divided  
702 this date into four categories corresponding to cell type: HEK293/HEK293T, HepG2,  
703 K562, and other cells. Then, for each combination of TF and cell type category, we  
704 selected a single peak set. We preferentially selected the peak sets from GTRD,  
705 because it contains systematically derived peak sets; we also note that GTRD contains  
706 the majority of ENCODE consortium experiments, together with many non-ENCODE  
707 experiments. When multiple experiments were available for a TF in a cell type category,  
708 we selected the experiment with higher counts. If multiple computational methods had  
709 been used to derive peak sets for the selected experiment, we chose the peak set using  
710 a preferential order MACS, GEM, SISRIS, PICS and PEAKZILLA. See **Table S7** for  
711 identifiers and metadata of the reference datasets.

712 For PWM scoring, the external peak sets were used as downloaded, with the exception  
713 of peak sets that were generated with the GEM peak caller, which have a peak width of  
714 1, and were therefore expanded 250 bases in both directions. For Codebook data, we  
715 used the merged and thresholded Codebook ChIP peak sets as in “Expert motif  
716 curation”. We generated negative peak sets for each ChIP-seq peak set using  
717 BEDTools shuffle<sup>77</sup> with the *-noOverlapping* option to create sets of random genomic  
718 regions with the same number of peaks and the same peak width distribution as the  
719 corresponding ChIP peak sets. We downloaded PWMs for all Codebook TFs from  
720 JASPAR<sup>80</sup> (2024 version), HOCOMOCO<sup>13</sup> (Version 12) and Factorbook<sup>12</sup> (downloaded  
721 15.12.2023). We scanned Codebook and external peak sets (and corresponding  
722 negative sets) with the expert curated Codebook motifs as PWMs using AffiMX<sup>28</sup>, and  
723 calculated AUROC values. Additionally, for the 19 Codebook TFs with a successful  
724 Codebook ChIP-seq experiment, a Codebook PWM, an external ChIP-seq experiment,  
725 and an external PWM, we compared the performance of PWMs across the different  
726 peak sets as follows. We first selected a single external PWM for each of the 19 TFs by  
727 scanning each PWM for a given TF on each external peak set for the same TF and  
728 identifying the PWM that produced the highest AUROC. We then used these highest  
729 scoring PWMs to scan the corresponding Codebook data and calculate AUROC values.

730 **TOP (Triple Overlap) and CTOP (Conserved Triple Overlap) peak set analyses.** To  
731 obtain TOP sites, we first identified thresholds for ChIP-seq peaks, GHT-SELEX peaks,  
732 and PWM score “peaks” that maximize the three-way Jaccard metric (overlap/union) of  
733 the three sets, with the thresholds calculated for each TF independently. We converted  
734 PWM hits (derived from MOODS<sup>79</sup> using a p-value cut-off of 0.001) into peaks by  
735 merging neighboring matches with a distance less than 200bp and re-scoring them  
736 using the sum-of-affinities for clusters. We then identified TOPs were as peaks  
737 exceeding these thresholds in all three sets, and overlap in all three sets. To obtain  
738 CTOP sites, we then extracted PhyloP scores for each base at each TOP site (and 100  
739 flanking bases) from the Zoonomia consortium<sup>81</sup>, removed sites overlapping the  
740 ENCODE Blacklist<sup>82</sup> or protein coding sequences (due to the skew in phyloP scores  
741 caused by codons), and applied three different statistical tests for significance of phyloP  
742 scores over the PWM hit: two that test correlation between the IC and the phyloP value  
743 at each base position of the PWM (using either Pearson correlation or linear  
744 regression), and one that tests for higher phyloP scores over the PWM hit (Wilcoxon  
745 test). Greater detail on these specific operations is given in the accompanying  
746 manuscripts<sup>14,15</sup>.

747 **Intersection of TOPs/CTOPs and genomic features.** We first clustered all CTOPs  
748 using BEDTools merge<sup>77</sup>, with a max distance of 100 bp, then intersected with the  
749 following genomic feature sets: basic canonical protein coding promoters from  
750 GENCODE version 44<sup>83</sup>, defined as 1000 bp upstream and 500 bp downstream of the  
751 canonical TSS; the “Unmasked CpG Island” track, PhastCons Conserved Elements  
752 from the Multiz 470 Mammalian alignment, and RepeatMasker track from UCSC<sup>84</sup>;  
753 ChromHMM HEK293 enhancers<sup>15</sup>. We classified promoters as CpG island or non-CpG  
754 island based on the GENCODE basic TSS being within +/- 50 bp of a CpG island from  
755 the unmasked track. We classified the CTOP clusters as associated with a single type  
756 of genomic feature in the following order of priority: CpG island associated with a protein  
757 coding promoter; other CpG islands; a non-CpG island-associated protein-coding  
758 promoter; an enhancer; containing a CTCF binding site but not overlapping a CpG  
759 island, promoter or enhancer; overlapping a transposable element and none of the  
760 previous categories; overlapping a non-TE repeat and none of the prior categories; and  
761 “Other” for CTOP clusters not intersecting any examined features.

762 **SNV analyses.** *TOPs and CTOPs.* For analysis of common variants, we intersected  
763 TOPs with the common short variants from dbSNP version 53, defined as a minor allele  
764 frequency of  $\geq 1\%$  in the 1000 Genomes project<sup>85</sup>. We determined genomic overlap  
765 enrichment between CTOPs/unconserved TOPs and dbSNP variants using the Fisher's  
766 Exact Test implemented in BEDTools<sup>77</sup>.

767 **Variant calling for allele-specific binding analysis.** We performed variant calling on  
768 our GHT-SELEX and ChIP-seq datasets by mapping raw ChIP-Seq and pre-trimmed  
769 GHT-SELEX reads<sup>17</sup> for 207 TFs to the hg38 human genome assembly using *bwa-mem*  
770 (v.0.7.1) with default settings (workflow is shown in **Figure S5A**). Next, we used  
771 *filter\_reads.py* from *stampipes* ([https://github.com/StamLab/stampipes/tree/encode-](https://github.com/StamLab/stampipes/tree/encode-release/)  
772 [release/](https://github.com/StamLab/stampipes/tree/encode-release/), accessed Sept 2022) to filter out reads with  $>2$  mismatches and mapping  
773 quality  $<10$ . Then, we used a previously-described approach<sup>86</sup> for SNV calling and read

774 counting: (1) *samtools reheader* (v.1.16.1) was used to set the identical sample SM field  
775 in all alignment files; (2) SNP calling was performed using *bcftools mpileup* (v.1.10.2)  
776 with `--redo-BAQ --adjust-MQ 50 --gap-frac 0.05 --max-depth 10000` and *bcftools call*  
777 with `--keep-alts --multiallelic-caller`; (3) the resulting SNPs were split into biallelic records  
778 using *bcftools norm* with `--check-ref x -m` - followed by filtering with *bcftools filter -i*  
779 `"QUAL>=10 & FORMAT/GQ>=20 & FORMAT/DP>=10"` `--SnpGap 3 --IndelGap 10` and  
780 *bcftools view -m2 -M2 -v snps* leaving only biallelic SNPs covered by 10 or more reads;  
781 (4) SNPs were annotated using *bcftools annotate* with `--columns ID,CAF, TOPMED` and  
782 dbSNP (v.151)<sup>87</sup> (5) heterozygous variants located on the reference chromosomes with  
783 GQ  $\geq 20$ , depth  $\geq 10$ , and allelic counts  $\geq 5$  on each allele were filtered with *awk* (v.5.0.1),  
784 (6) WASP (v.0.3.4)<sup>88</sup> was used with *bwa mem* and *filter\_reads.py* to account for  
785 reference mapping bias, (7) *count\_tags\_pileup\_new.py* was used to obtain allelic read  
786 counts with *pysam* (v.0.20.0), (8) *recode\_vcf.py* was used to convert the resulting BED  
787 files to VCF. In total, we made 925,003 candidate variant calls supported by five reads  
788 for both alleles and listed in the dbSNP common subset<sup>87</sup>.

789 **ASB calling and analysis.** ASB calling was performed independently for GHT-SELEX  
790 and ChIP-seq data. To account for aneuploidy and copy-number variation, the profiles  
791 of relative background allelic dosage were reconstructed with BABACHI (v.2.0.26) using  
792 default settings<sup>89</sup>, Abstract O3). The allelic imbalance was estimated with MIXALIME  
793 (v.2.14.7)<sup>68</sup> starting with *mixalime create*. Next, we fitted a marginalized compound  
794 negative binomial model (MCNB) using *mixalime fit* specifying MCNB and setting `--`  
795 `window-size` to 1000 and 10000 for GHT-SELEX and ChIP-Seq, respectively, taking into  
796 account lower coverage and SNP counts in GHT-SELEX. Finally, we used *mixalime test*  
797 followed by TF-wise *mixalime combine* to obtain the TF-specific ASB calls (**Figure**  
798 **S5A**).

799 We then identified ASBs that overlap a PWM hit (P-value < 0.001) for the associated  
800 TF. For those ASBs, we calculated the PWM score for both alleles and estimated the P-  
801 value of those scores against a uniform background distribution for each allele using  
802 PERFECTOS-APE<sup>90</sup>. The fold-change between allele P-values (P1/P2) was then  
803 calculated with the P-value of the more abundant allele as the numerator. ASBs with a  
804  $\log_2(\text{fold-change}) \geq 1$  were labelled “strongly concordant”, i.e., the allele we observed  
805 to be bound more often is consistent with the PWM score (**Figure S5B**).

806 To assess the enrichment of Codebook ASBs within GTEx eQTLs<sup>67</sup> and ADAstra  
807 ASBs<sup>66</sup> we combined the ASB P-values from ChIP-Seq and GHT-SELEX data across all  
808 TFs and datasets (*logitp* method<sup>91</sup>) to generate a single P-value for each TF (**Figure**  
809 **S5C**).

810 **Analysis of allele-specific chromatin accessibility.** In this analysis, we relied on 321  
811 and 577 cell type-specific chromatin accessibility datasets derived from DNase- and  
812 ATAC-Seq experiments, respectively, and available in the UDACHA database (Release  
813 IceKing 1.0.3)<sup>68</sup>. We identified 4,048 instances in which ASVs in a specific cell type  
814 overlap significantly with PWM hits (P<0.0005) for a TF in the Codebook motif collection  
815 (236 PWMs) (Right-tailed Fisher’s exact test P < 0.05, and requiring 10 or more  
816 overlapping PWM hits) (**Figure S5D**). Then, for each ASV in each combination of TF



817 and cell type passing the PWM enrichment filter, we asked whether the change in the  
818 PWM score is concordant with the read imbalance in the ASVs, e.g. whether a higher  
819 PWM score at a given locus corresponds to a higher read count, and we assigned a P-  
820 value for each combination of TF and cell type, using a right-tailed Fisher's exact test,  
821 including only sites with at least two-fold change in PWM-predicted affinity. Finally, to  
822 obtain a single significance estimate per TF, we combined these P-values for each TF  
823 across the different cell types passing the first stage, i.e. for which the overlap between  
824 PWM hits and ASVs is significant (Fisher's method, considering DNase-Seq and ATAC-  
825 Seq data separately and FDR-adjusted). TFs passing  $FDR < 0.05$  in the final stage  
826 were considered ASV-concordant.

827 To further verify the concordance between ASVs and Codebook motifs, we selected 34  
828 (out of 53 TFs) with at least one TOP region overlapping ASVs, and re-evaluated the  
829 concordant-to-discordant ratio for ASVs within peaks and TOP regions (see **Results**  
830 and **Figure 5C**). For this analysis, for each TF, we picked the most significant ASV at  
831 each unique genomic position (SNP) from all available cell types, and performed a right-  
832 tailed Fisher's Exact Test (**Table S9**). At this stage, we considered SP140 and SP140L  
833 jointly they share short and highly similar DNA-binding motifs.

## 834 DATA AVAILABILITY

835 The sequencing raw data for the HT-SELEX and GHT-SELEX experiments have been  
836 deposited into the SRA database under identifiers PRJEB78913 (ChIP-seq),  
837 PRJEB76622 (GHT-SELEX), and PRJEB61115 (HT-SELEX). Genomic interval  
838 information generated for the GHT-SELEX and ChIP-seq have been deposited into  
839 GEO under accessions GSE280248 (ChIP-seq) and GSE278858 (GHT-SELEX). PWMs  
840 can be browsed at <https://mex.autosome.org> and downloaded at  
841 <https://doi.org/10.5281/ZENODO.8327372>. An updated list of human TFs is available at  
842 <https://humantfs.ccb.utoronto.ca>. Information on constructs, experiments, analyses,  
843 processed data, comparison tracks, and browsable pages with information and results  
844 for each TF is available at <https://codebook.ccb.utoronto.ca>.

## 845 ACKNOWLEDGEMENTS

846 We thank the IT Group of the Institute of Computer Science at Halle University for  
847 computational resources, Maximilian Biermann for valuable technical support, Gherman  
848 Novakovsky for providing feedback, Berat Dogan for testing earlier versions of  
849 RCADEEM, and Debashish Ray for assistance with database depositions.

850 This work was supported by the following:

- 851 • Canadian Institutes of Health Research (CIHR) grants FDN-148403, PJT-  
852 186136, PJT-191768, and PJT-191802, and NIH grant R21HG012258 to T.R.H
- 853 • CIHR grant PJT-191802 to T.R.H. and H.S.N.
- 854 • Natural Sciences and Engineering Research Council of Canada (NSERC) grant  
855 RGPIN-2018-05962 to H.S.N.
- 856 • A Russian Science Foundation grant [20-74-10075] to I.V.K.
- 857 • A Swiss National Science Foundation grant (no. 310030\_197082) to B.D.
- 858 • Marie Skłodowska-Curie (no. 895426) and EMBO long-term (1139-2019)  
859 fellowships to J.F.K.
- 860 • NIH grants R01HG013328 and U24HG013078 to M.T.W., T.R.H., and Q.D.M.
- 861 • NIH grants R01AR073228, P30AR070549, and R01AI173314 to M.T.W.
- 862 • NIH grant P30CA008748 partially supported Q.M.
- 863 • Canada Research Chairs funded by CIHR to T.R.H. and H.S.N.
- 864 • Ontario Graduate Scholarships to K.U.L and I.Y.
- 865 • A.J. was supported by Vetenskapsrådet (Swedish Research Council)  
866 Postdoctoral Fellowship (2016-00158)
- 867 • The Billes Chair of Medical Research at the University of Toronto to T.R.H.
- 868 • EPFL's Center for Imaging
- 869 • Resource allocations from Digital Research Alliance of Canada

## 870 DECLARATION OF COMPETING INTERESTS

871 O.F. is employed by Roche.

## 872 SUPPLEMENTARY TABLES

873 **Table S1. Accompanying manuscripts.** Table lists the 5 studies performed by the  
874 Codebook Consortium, providing basic information for each of the manuscripts,  
875 including title and author list.

876 **Table S2. TF list and assay success.** Table lists the Codebook proteins and positive  
877 control TFs that were analyzed in the Codebook studies and provides metadata and  
878 information on whether they showed sequence-specific DNA binding activities in  
879 different types of experiments, together with the ID of the representative PWM selected  
880 in this study, if any.

881 **Table S3. List of inserts used in this study.** Table provides the amino acid sequence  
882 and type (full-length or DBD) for the 716 inserts used in the Codebook studies.

883 **Table S4. List of plasmids used in this study.** Table lists the plasmid backbone and  
884 insert for each of the 1,387 plasmids used in the Codebook studies.

885 **Table S5. List of experiments performed in this study.** Table lists the 4,873  
886 experiments performed on Codebook and control TFs, along with 20 additional GFP  
887 control experiments. The experiment ID, experiment type, TF assayed, expert curation  
888 result, and plasmid ID are listed for each experiment. Each experiment is mapped to its  
889 ID in an accompanying manuscript<sup>17</sup>, and 9 additional experiments used only in an  
890 accompanying manuscript<sup>17</sup> are listed.

891 **Table S6. Representative PWMs.** Table shows logo representations for the PWMs  
892 that were selected as the representative for each of the TFs (i.e. the expert-curated  
893 motifs) and provides metadata describing the role of the TF in the study, DBD that it  
894 belongs to, source of the experimental data and motif derivation approach.

895 **Table S7. External peak datasets.** Table lists external peak location datasets obtained  
896 from GTRD database and ENCODE consortium, that were used in the comparisons  
897 carried out in this study.

898 **Table S8. External PWM datasets.** Table lists PWM identifiers, manual curation and  
899 other metadata for external motifs available from the databases Jaspar, HOCOMOCO  
900 and Factorbook.

901 **Table S9. ASE and ASV data.** Allele-specific binding sites detected in Codebook data  
902 and motif annotation of allele-specific chromatin accessibility events.

903 **Table S10. Updated census of human transcription factors and their motif  
904 coverage.** Table is modified from Lambert et al. to display an updated motif coverage of  
905 human TFs.

## 906 REFERENCES

- 907
- 908 1. Lambert, S.A. *et al.* The Human Transcription Factors. *Cell* **175**, 598-599 (2018).
  - 909 2. Stormo, G.D. & Zhao, Y. Determining the specificity of protein-DNA interactions.  
910 *Nat Rev Genet* **11**, 751-60 (2010).
  - 911 3. Stormo, G.D. Consensus patterns in DNA. *Methods Enzymol* **183**, 211-21  
912 (1990).
  - 913 4. Schneider, T.D. & Stephens, R.M. Sequence logos: a new way to display  
914 consensus sequences. *Nucleic Acids Res* **18**, 6097-100 (1990).
  - 915 5. Benos, P.V., Bulyk, M.L. & Stormo, G.D. Additivity in protein-DNA interactions:  
916 how good an approximation is it? *Nucleic Acids Res* **30**, 4442-51 (2002).
  - 917 6. Yan, J. *et al.* Systematic analysis of binding of transcription factors to noncoding  
918 variants. *Nature* **591**, 147-151 (2021).
  - 919 7. Wasserman, W.W. & Sandelin, A. Applied bioinformatics for the identification of  
920 regulatory elements. *Nat Rev Genet* **5**, 276-87 (2004).
  - 921 8. Srivastava, D. & Mahony, S. Sequence and chromatin determinants of  
922 transcription factor binding and the establishment of cell type-specific binding  
923 patterns. *Biochim Biophys Acta Gene Regul Mech* **1863**, 194443 (2020).
  - 924 9. Isakova, A. *et al.* SMiLE-seq identifies binding motifs of single and dimeric  
925 transcription factors. *Nat Methods* **14**, 316-322 (2017).
  - 926 10. Schmitges, F.W. *et al.* Multiparameter functional diversity of human C2H2 zinc  
927 finger proteins. *Genome Res* **26**, 1742-1752 (2016).
  - 928 11. Consortium, E.P. *et al.* Perspectives on ENCODE. *Nature* **583**, 693-698 (2020).
  - 929 12. Pratt, H.E. *et al.* Factorbook: an updated catalog of transcription factor motifs and  
930 candidate regulatory motif sites. *Nucleic Acids Res* **50**, D141-D149 (2022).
  - 931 13. Vorontsov, I.E. *et al.* HOCOMOCO in 2024: a rebuild of the curated collection of  
932 binding models for human and mouse transcription factors. *Nucleic Acids Res*  
933 **52**, D154-D163 (2024).
  - 934 14. Jolma, A. *et al.* GHT-SELEX demonstrates unexpectedly high intrinsic sequence  
935 specificity and complex DNA binding of many human transcription factors.  
936 *bioRxiv*, 2024.11.11.618478 (2024).
  - 937 15. Razavi, R. *et al.* Extensive binding of uncharacterized human transcription  
938 factors to genomic dark matter. *bioRxiv*, 2024.11.11.622123 (2024).
  - 939 16. Gralak, A. *et al.* Identification of methylation-sensitive human transcription factors  
940 using meSMiLE-seq. *bioRxiv*, 2024.11.11.619598 (2024).
  - 941 17. Vorontsov, I.E. *et al.* Cross-platform DNA motif discovery and benchmarking to  
942 explore binding specificities of poorly studied human transcription factors.  
943 *bioRxiv*, 2024.11.11.619379 (2024).
  - 944 18. Ambrosini, G. *et al.* Insights gained from a comprehensive all-against-all  
945 transcription factor binding motif benchmarking study. *Genome Biol* **21**, 114  
946 (2020).
  - 947 19. Bailey, T.L., Johnson, J., Grant, C.E. & Noble, W.S. The MEME Suite. *Nucleic*  
948 *Acids Res* **43**, W39-49 (2015).
  - 949 20. Novakovskiy, G., Fornes, O., Saraswat, M., Mostafavi, S. & Wasserman, W.W.  
950 ExplainNN: interpretable and transparent neural networks for genomics. *Genome*  
951 *Biol* **24**, 154 (2023).

- 952 21. Rube, H.T. *et al.* Prediction of protein-ligand binding affinity from sequencing  
953 data with interpretable machine learning. *Nat Biotechnol* (2022).
- 954 22. Wolfe, S.A., Nekludova, L. & Pabo, C.O. DNA recognition by Cys2His2 zinc  
955 finger proteins. *Annu Rev Biophys Biomol Struct* **29**, 183-212 (2000).
- 956 23. Brayer, K.J., Kulshreshtha, S. & Segal, D.J. The protein-binding potential of  
957 C2H2 zinc finger domains. *Cell Biochem Biophys* **51**, 9-19 (2008).
- 958 24. Bird, A.J., Gordon, M., Eide, D.J. & Winge, D.R. Repression of ADH1 and ADH3  
959 during zinc deficiency by Zap1-induced intergenic RNA transcripts. *EMBO J* **25**,  
960 5726-34 (2006).
- 961 25. Font, J. & Mackay, J.P. Beyond DNA: zinc finger domains as RNA-binding  
962 modules. *Methods Mol Biol* **649**, 479-91 (2010).
- 963 26. Stros, M., Launholt, D. & Grasser, K.D. The HMG-box: a versatile protein domain  
964 occurring in a wide variety of DNA-binding proteins. *Cell Mol Life Sci* **64**, 2590-  
965 606 (2007).
- 966 27. Najafabadi, H.S. *et al.* C2H2 zinc finger proteins greatly expand the human  
967 regulatory lexicon. *Nat Biotechnol* (2015).
- 968 28. Lambert, S.A., Albu, M., Hughes, T.R. & Najafabadi, H.S. Motif comparison  
969 based on similarity of binding affinity profiles. *Bioinformatics* **32**, 3504-3506  
970 (2016).
- 971 29. Emerson, R.O. & Thomas, J.H. Adaptive evolution in zinc finger transcription  
972 factors. *PLoS Genet* **5**, e1000325 (2009).
- 973 30. Zhao, Y. & Stormo, G.D. Quantitative analysis demonstrates most transcription  
974 factors require only simple models of specificity. *Nat Biotechnol* **29**, 480-3 (2011).
- 975 31. Ruan, S., Swamidass, S.J. & Stormo, G.D. BEESEM: estimation of binding  
976 energy models using HT-SELEX data. *Bioinformatics* **33**, 2288-2295 (2017).
- 977 32. Weirauch, M.T. *et al.* Evaluation of methods for modeling transcription factor  
978 sequence specificity. *Nat Biotechnol* **31**, 126-34 (2013).
- 979 33. Kuznetsov, V.A. Mathematical Modeling of Avidity Distribution and Estimating  
980 General Binding Properties of Transcription Factors from Genome-Wide Binding  
981 Profiles. *Methods Mol Biol* **1613**, 193-276 (2017).
- 982 34. Horton, C.A. *et al.* Short tandem repeats bind transcription factors to tune  
983 eukaryotic gene expression. *Science* **381**, eadd1250 (2023).
- 984 35. Morgunova, E. *et al.* Two distinct DNA sequences recognized by transcription  
985 factors represent enthalpy and entropy optima. *Elife* **7**(2018).
- 986 36. Siggers, T. & Gordan, R. Protein-DNA binding: complexities and multi-protein  
987 codes. *Nucleic Acids Res* **42**, 2099-111 (2014).
- 988 37. Iuchi, S. Three classes of C2H2 zinc finger proteins. *Cell Mol Life Sci* **58**, 625-35  
989 (2001).
- 990 38. Yellan, I., Yang, A.W.H. & Hughes, T.R. Diverse Eukaryotic CGG-Binding  
991 Proteins Produced by Independent Domestications of hAT Transposons. *Mol Biol*  
992 *Evol* **38**, 2070-2075 (2021).
- 993 39. Singh, U. & Westermark, B. CGGBP1--an indispensable protein with ubiquitous  
994 cytoprotective functions. *Ups J Med Sci* **120**, 219-32 (2015).
- 995 40. Aoki, T., Sarkeshik, A., Yates, J. & Schedl, P. Elba, a novel developmentally  
996 regulated chromatin boundary factor is a hetero-tripartite DNA binding complex.  
997 *Elife* **1**, e00171 (2012).

- 998 41. Dai, Q. *et al.* The BEN domain is a novel sequence-specific DNA-binding domain  
999 conserved in neural transcriptional repressors. *Genes Dev* **27**, 602-14 (2013).
- 1000 42. Vetrini, F. *et al.* De novo and inherited TCF20 pathogenic variants are associated  
1001 with intellectual disability, dysmorphic features, hypotonia, and neurological  
1002 impairments with similarities to Smith-Magenis syndrome. *Genome Med* **11**, 12  
1003 (2019).
- 1004 43. Gupta, M., Zak, R., Libermann, T.A. & Gupta, M.P. Tissue-restricted expression  
1005 of the cardiac alpha-myosin heavy chain gene is controlled by a downstream  
1006 repressor element containing a palindrome of two ets-binding sites. *Mol Cell Biol*  
1007 **18**, 7243-58 (1998).
- 1008 44. Zhou, J. *et al.* Attenuation of Forkhead signaling by the retinal determination  
1009 factor DACH1. *Proc Natl Acad Sci U S A* **107**, 6864-9 (2010).
- 1010 45. Abramson, J. *et al.* Accurate structure prediction of biomolecular interactions with  
1011 AlphaFold 3. *Nature* (2024).
- 1012 46. Mitchell, A.L. *et al.* InterPro in 2019: improving coverage, classification and  
1013 access to protein sequence annotations. *Nucleic Acids Res* **47**, D351-D360  
1014 (2019).
- 1015 47. Harrison, P.W. *et al.* Ensembl 2024. *Nucleic Acids Res* **52**, D891-D899 (2024).
- 1016 48. Hoverter, N.P. *et al.* The TCF C-clamp DNA binding domain expands the Wnt  
1017 transcriptome via alternative target recognition. *Nucleic Acids Res* **42**, 13615-32  
1018 (2014).
- 1019 49. Letunic, I., Khedkar, S. & Bork, P. SMART: recent updates, new developments  
1020 and status in 2020. *Nucleic Acids Res* **49**, D458-D460 (2021).
- 1021 50. Yin, Y. *et al.* Impact of cytosine methylation on DNA binding specificities of  
1022 human transcription factors. *Science* **356**(2017).
- 1023 51. Hayward, A., Ghazal, A., Andersson, G., Andersson, L. & Jern, P. ZBED  
1024 evolution: repeated utilization of DNA transposons as regulators of diverse host  
1025 functions. *PLoS One* **8**, e59940 (2013).
- 1026 52. Smit, A.F. & Riggs, A.D. Tiggers and DNA transposon fossils in the human  
1027 genome. *Proc Natl Acad Sci U S A* **93**, 1443-8 (1996).
- 1028 53. Etchegaray, E., Baas, D., Naville, M., Haftek-Terreau, Z. & Volff, J.N. The  
1029 neurodevelopmental gene MSANTD2 belongs to a gene family formed by  
1030 recurrent molecular domestication of Harbinger transposons at the base of  
1031 vertebrates. *Mol Biol Evol* **39**(2022).
- 1032 54. Marquez, C.P. & Pritham, E.J. Phantom, a new subclass of Mutator DNA  
1033 transposons found in insect viruses and widely distributed in animals. *Genetics*  
1034 **185**, 1507-17 (2010).
- 1035 55. Toth, M., Grimsby, J., Buzsaki, G. & Donovan, G.P. Epileptic seizures caused by  
1036 inactivation of a novel gene, jerky, related to centromere binding protein-B in  
1037 transgenic mice. *Nat Genet* **11**, 71-5 (1995).
- 1038 56. Pace, J.K., 2nd & Feschotte, C. The evolutionary history of human DNA  
1039 transposons: evidence for intense activity in the primate lineage. *Genome Res*  
1040 **17**, 422-32 (2007).
- 1041 57. Partridge, E.C. *et al.* Occupancy maps of 208 chromatin-associated proteins in  
1042 one human cell type. *Nature* **583**, 720-728 (2020).

- 1043 58. Lai, W.K.M. *et al.* A ChIP-exo screen of 887 Protein Capture Reagents Program  
1044 transcription factor antibodies in human cells. *Genome Res* **31**, 1663-1679  
1045 (2021).
- 1046 59. Luo, Y. *et al.* New developments on the Encyclopedia of DNA Elements  
1047 (ENCODE) data portal. *Nucleic Acids Res* **48**, D882-D889 (2020).
- 1048 60. Kolmykov, S. *et al.* GTRD: an integrated view of transcription regulation. *Nucleic*  
1049 *Acids Res* **49**, D104-D111 (2021).
- 1050 61. Castro-Mondragon, J.A. *et al.* JASPAR 2022: the 9th release of the open-access  
1051 database of transcription factor binding profiles. *Nucleic Acids Res* **50**, D165-  
1052 D173 (2022).
- 1053 62. Cohen, N.M., Kenigsberg, E. & Tanay, A. Primate CpG islands are maintained by  
1054 heterogeneous evolutionary regimes involving minimal selection. *Cell* **145**, 773-  
1055 86 (2011).
- 1056 63. Fishilevich, S. *et al.* GeneHancer: genome-wide integration of enhancers and  
1057 target genes in GeneCards. *Database (Oxford)* **2017**(2017).
- 1058 64. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-  
1059 66 (2015).
- 1060 65. Weirauch, M.T. & Hughes, T.R. Conserved expression without conserved  
1061 regulatory sequence: the more things change, the more they stay the same.  
1062 *Trends Genet* **26**, 66-74 (2010).
- 1063 66. Abramov, S. *et al.* Landscape of allele-specific transcription factor binding in the  
1064 human genome. *Nat Commun* **12**, 2751 (2021).
- 1065 67. Consortium, G.T. The Genotype-Tissue Expression (GTEx) project. *Nat Genet*  
1066 **45**, 580-5 (2013).
- 1067 68. Buyan, A. *et al.* Statistical framework for calling allelic imbalance in high-  
1068 throughput sequencing data. *bioRxiv*, 2023.11.07.565968 (2023).
- 1069 69. Lambert, S.A. *et al.* Similarity regression predicts evolution of transcription factor  
1070 sequence specificity. *Nat Genet* **51**, 981-989 (2019).
- 1071 70. Avsec, Z. *et al.* Effective gene expression prediction from sequence by  
1072 integrating long-range interactions. *Nat Methods* **18**, 1196-1203 (2021).
- 1073 71. de Boer, C.G. & Taipale, J. Hold out the genome: a roadmap to solving the cis-  
1074 regulatory code. *Nature* **625**, 41-50 (2024).
- 1075 72. Wang, Y. *et al.* SNP rs17079281 decreases lung cancer risk through creating an  
1076 YY1-binding site to suppress DCBLD1 expression. *Oncogene* **39**, 4092-4102  
1077 (2020).
- 1078 73. Degtyareva, A.O., Antontseva, E.V. & Merkulova, T.I. Regulatory SNPs: Altered  
1079 Transcription Factor Binding Sites Implicated in Complex Traits and Diseases. *Int*  
1080 *J Mol Sci* **22**(2021).
- 1081 74. Deplancke, B., Alpern, D. & Gardeux, V. The Genetics of Transcription Factor  
1082 DNA Binding Variation. *Cell* **166**, 538-554 (2016).
- 1083 75. Weirauch, M.T. *et al.* Determination and inference of eukaryotic transcription  
1084 factor sequence specificity. *Cell* **158**, 1431-43 (2014).
- 1085 76. Narasimhan, K. *et al.* Mapping and analysis of *Caenorhabditis elegans*  
1086 transcription factor sequence specificities. *Elife* **4**(2015).
- 1087 77. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing  
1088 genomic features. *Bioinformatics* **26**, 841-2 (2010).

- 1089 78. Satopaa, V., Albrecht, J., Irwin, D. & Raghavan, B. Finding a "kneedle" in a  
1090 haystack: Detecting knee points in system behavior. in *2011 31st international*  
1091 *conference on distributed computing systems workshops* 166-171 (IEEE, 2011).
- 1092 79. Korhonen, J., Martinmaki, P., Pizzi, C., Rastas, P. & Ukkonen, E. MOODS: fast  
1093 search for position weight matrix matches in DNA sequences. *Bioinformatics* **25**,  
1094 3181-2 (2009).
- 1095 80. Rauluseviciute, I. *et al.* JASPAR 2024: 20th anniversary of the open-access  
1096 database of transcription factor binding profiles. *Nucleic Acids Res* **52**, D174-  
1097 D182 (2024).
- 1098 81. Armstrong, J. *et al.* Progressive Cactus is a multiple-genome aligner for the  
1099 thousand-genome era. *Nature* **587**, 246-251 (2020).
- 1100 82. Amemiya, H.M., Kundaje, A. & Boyle, A.P. The ENCODE Blacklist: Identification  
1101 of Problematic Regions of the Genome. *Sci Rep* **9**, 9354 (2019).
- 1102 83. Frankish, A. *et al.* GENCODE: reference annotation for the human and mouse  
1103 genomes in 2023. *Nucleic Acids Res* **51**, D942-D949 (2023).
- 1104 84. Nassar, L.R. *et al.* The UCSC Genome Browser database: 2023 update. *Nucleic*  
1105 *Acids Res* **51**, D1188-D1195 (2023).
- 1106 85. Sayers, E.W. *et al.* Database resources of the National Center for Biotechnology  
1107 Information. *Nucleic Acids Res* **47**, D23-D28 (2019).
- 1108 86. Vierstra, J. *et al.* Global reference mapping of human transcription factor  
1109 footprints. *Nature* **583**, 729-736 (2020).
- 1110 87. Sherry, S.T. *et al.* dbSNP: the NCBI database of genetic variation. *Nucleic Acids*  
1111 *Res* **29**, 308-11 (2001).
- 1112 88. van de Geijn, B., McVicker, G., Gilad, Y. & Pritchard, J.K. WASP: allele-specific  
1113 software for robust molecular quantitative trait locus discovery. *Nat Methods* **12**,  
1114 1061-3 (2015).
- 1115 89. Selected abstracts of Bioinformatics: from Algorithms to Applications 2021  
1116 Conference. *BMC Bioinformatics* **22**, 591 (2021).
- 1117 90. Kulakovskiy, I., Vorontsov, I. & Makeev, V. *PERFECTOS-APE – predicting*  
1118 *regulatory functional effect of SNPs by approximate P-value estimation*, (2015).
- 1119 91. George, E.O. & Mudholkar, G.S. On the convolution of logistic random variables.  
1120 *Metrika* **30**, 1-13 (1983).
- 1121 92. Price, M.N., Dehal, P.S. & Arkin, A.P. FastTree 2--approximately maximum-  
1122 likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 1123 93. Katoh, K., Kuma, K., Toh, H. & Miyata, T. MAFFT version 5: improvement in  
1124 accuracy of multiple sequence alignment. *Nucleic Acids Res* **33**, 511-8 (2005).
- 1125 94. Dupeyron, M., Baril, T., Bass, C. & Hayward, A. Phylogenetic analysis of the  
1126 Tc1/mariner superfamily reveals the unexplored diversity of pogo-like elements.  
1127 *Mob DNA* **11**, 21 (2020).
- 1128 95. Gao, B. *et al.* Evolution of pogo, a separate superfamily of IS630-Tc1-mariner  
1129 transposons, revealing recurrent domestication events in vertebrates. *Mob DNA*  
1130 **11**, 25 (2020).
- 1131 96. Jolma, A. *et al.* DNA-Binding Specificities of Human Transcription Factors. *Cell*  
1132 **152**, 327-39 (2013).
- 1133 97. Worsley Hunt, R. & Wasserman, W.W. Non-targeted transcription factors motifs  
1134 are a systemic component of ChIP-seq datasets. *Genome Biol* **15**, 412 (2014).



**A** **ARID1** Factorbook ENCSR491EBY  
YVSYGCCMYCTGSTG



**KDM5B** Factorbook ENCSR000AQA  
GCCGCCATCTY



**DRAP1** Factorbook ENCSR765MKZ  
YKSYSATTGGYYSN



**GATAD2A** Factorbook ENCSR160QYK  
RTKRTGCAAYM



**GATAD2A** Factorbook ENCSR925BFV  
HWRWGYAAACA



**MBD1** Factorbook ENCSR396QWK  
CGCTGTCCRYGGTGCTGAA



**B** **GATAD2A** Factorbook ENCSR160QYK  
WGATAAGV



**ZBED2** Jaspar MA1971.1



**ZNF623** Jaspar UN0210.1



**C** **ZNF592** Factorbook ENCSR701AQS  
AGYRACTCCATCTTG



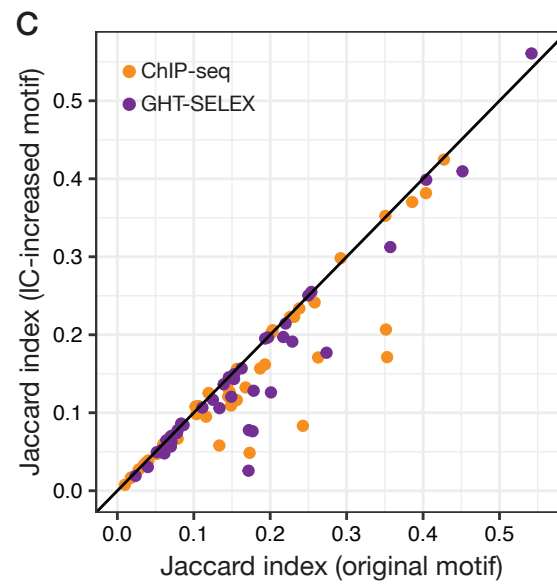
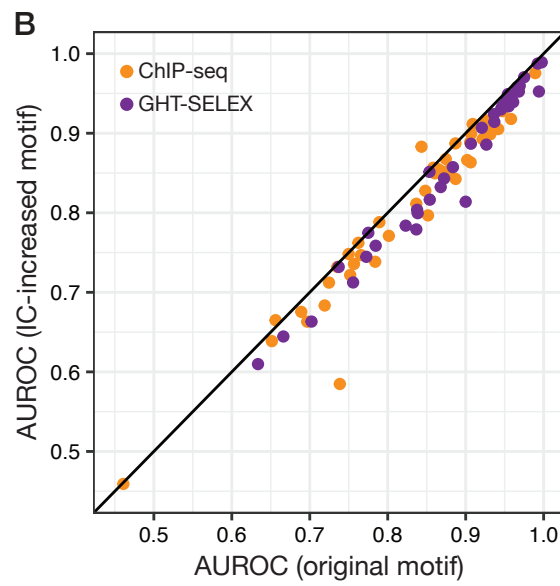
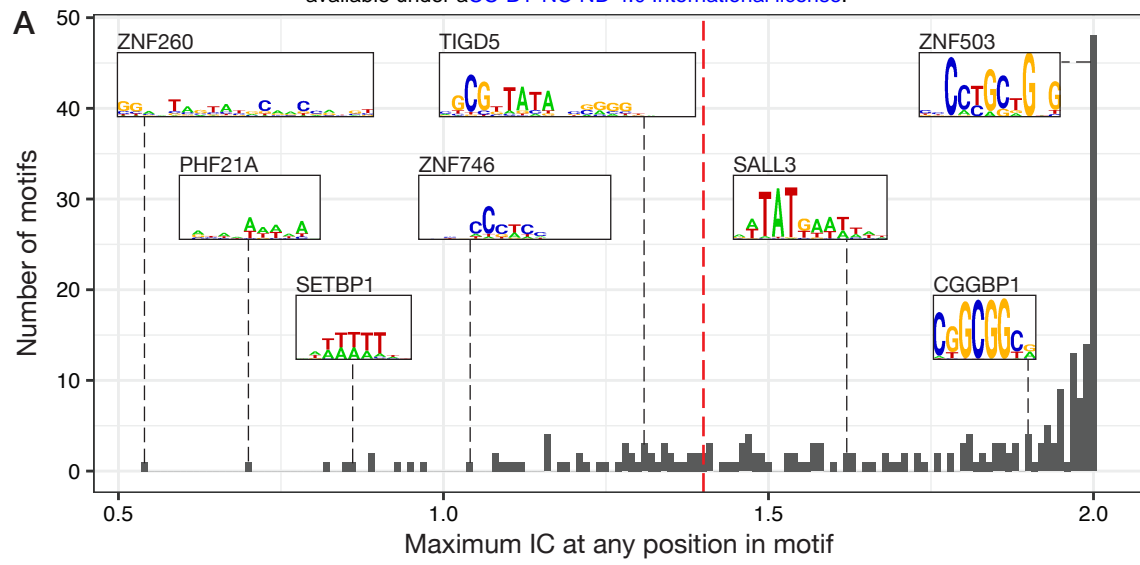
**ZNF577** HOCOMOCO  
ZN577.H12CORE.0.P.B



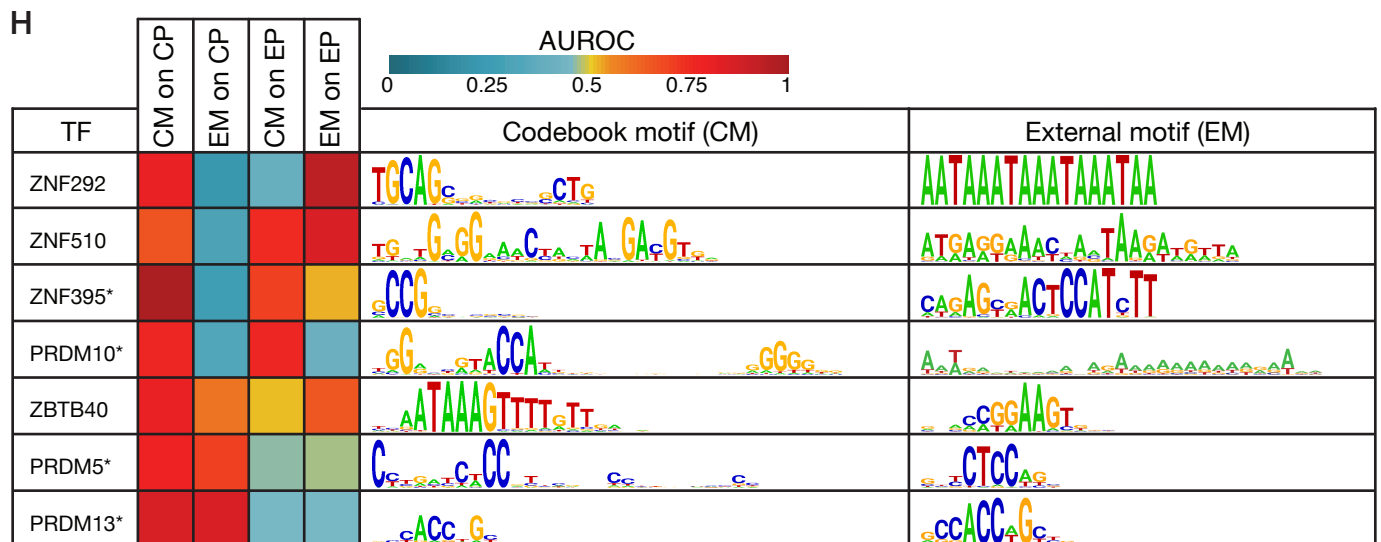
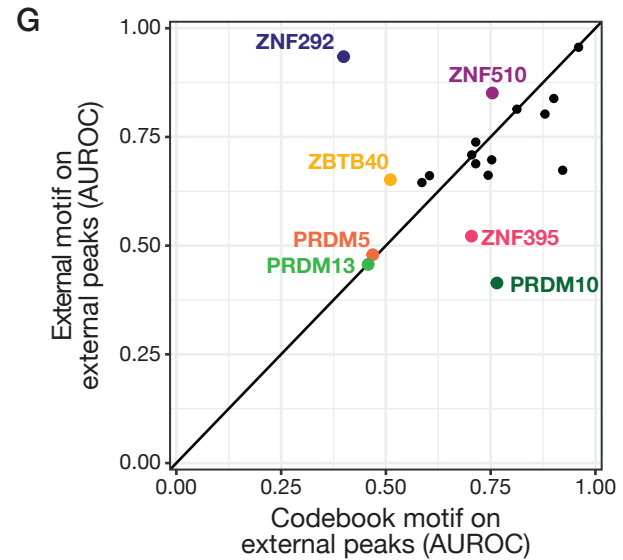
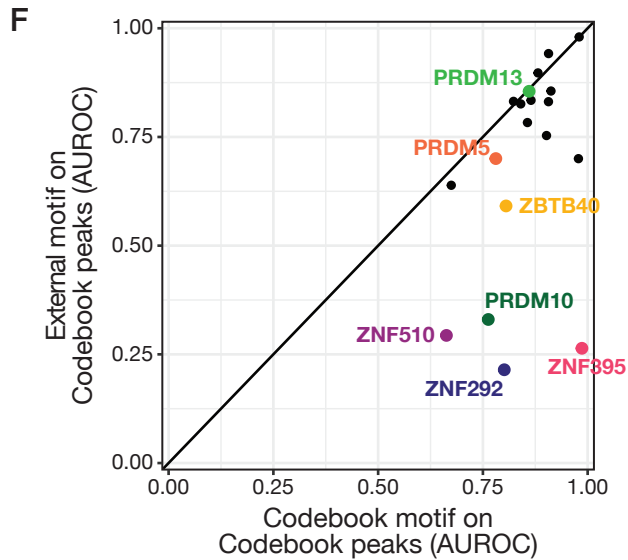
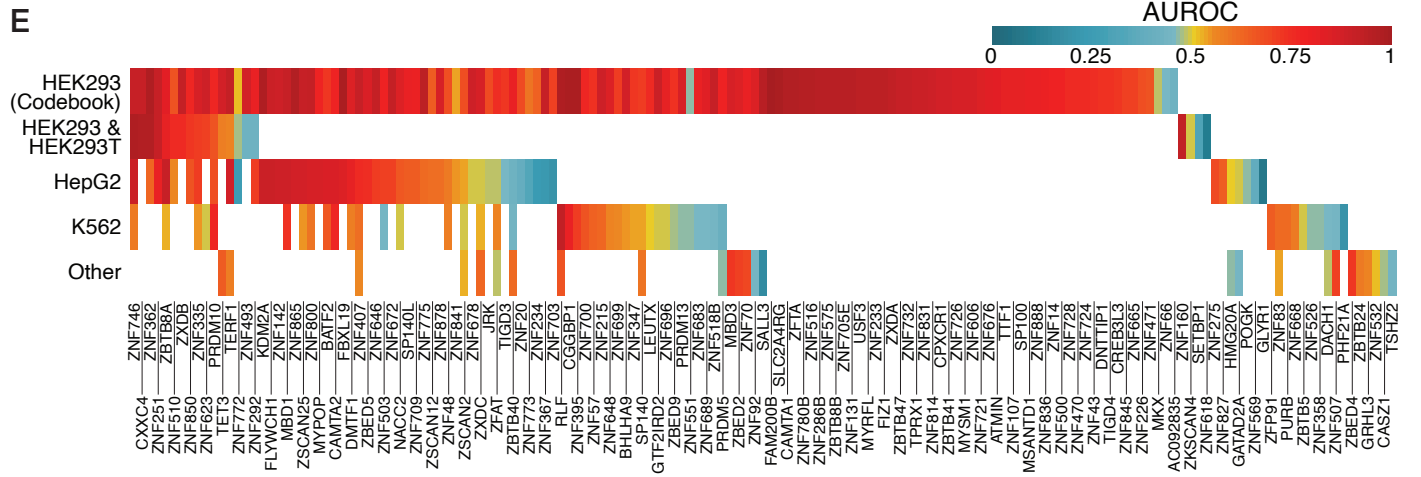
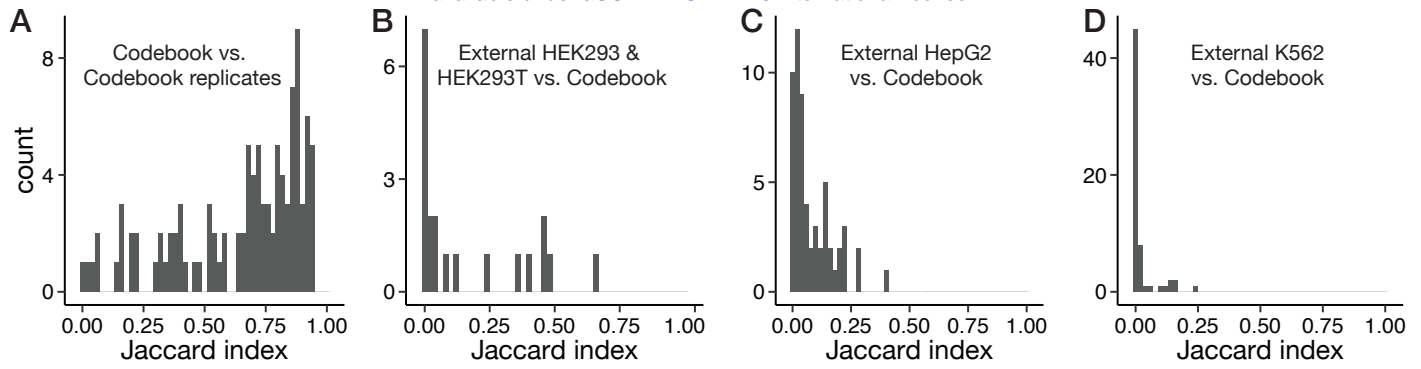
**GATAD2B** Factorbook ENCSR547LKC  
HNNDNWN YCTTATCTVYHNHY



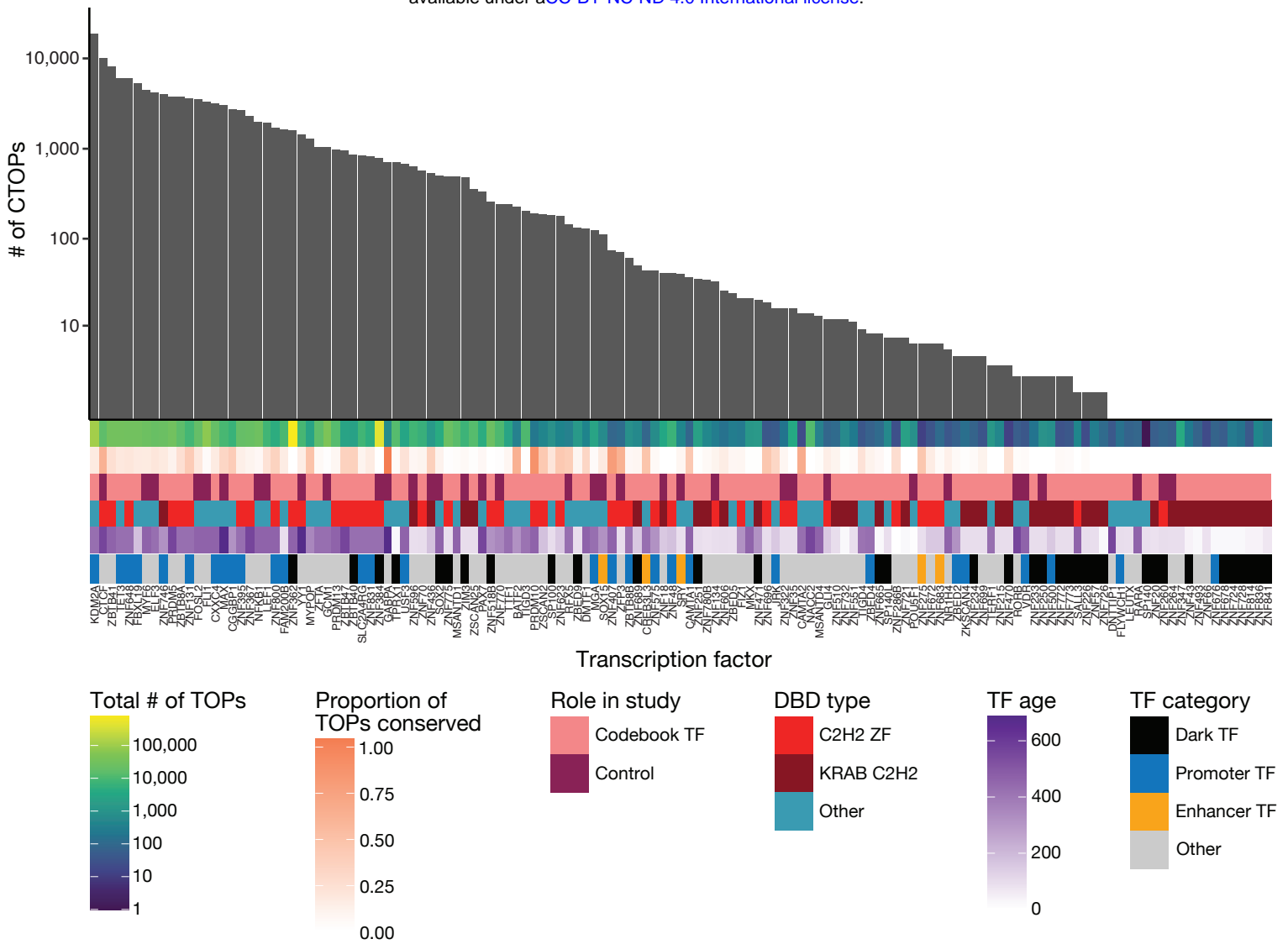
**Figure S1. Examples of evaluation of external PWMs.** **A**, Cases in which the external PWM matches that of a well-studied TF that is a frequent “contaminant” motif in CHIP-seq<sup>97</sup>. In each example, the top sequence logo represents the external PWM, and the bottom sequence logo represents a highly-similar CisBP PWM. **B**, Cases in which the external PWM (top in each example) is consistent with the Codebook PWM for the same TF (bottom in each example). **C**, External PWM sequence logos that cannot be explained as known contaminants or artifacts, some of which are supported by multiple lines of evidence, and thus appear accurate.



**Figure S2. Motif degeneracy analysis.** **A**, Histogram displays the maximum information content (IC) for any position within the expert-curated PWM for all Codebook and control TFs. Logos are shown for TFs at various maximum positional IC values, for illustration. Red dashed line indicates an IC of 1.4. **B**, and **C**, comparison of original PWMs to IC-increased PWMs for the 52 TF PWMs for which no base position exceeded an IC of 1.4. **B**, AUROC scores for original vs. IC-increased PWMs, discriminating ChIP-seq or GHT-SELEX peaks vs. random genomic background loci. **C**, Maximum Jaccard index for ChIP-seq or GHT-SELEX peak sets; using the approach described for optimized TOPs in **Methods**, for original vs. IC-increased PWMs.



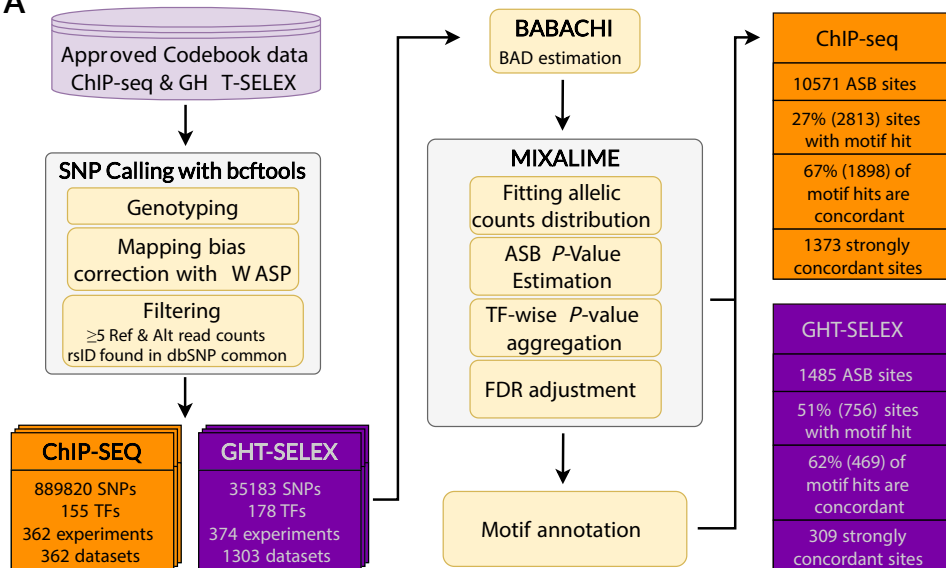
**Figure S3. Comparison to external ChIP-seq datasets and PWMs.** **A-D**, Histograms of Jaccard indices measuring the overlap between two ChIP-seq peak sets for the same TF: **A**, Codebook ChIP-seq replicates; **B**, **C**, **D**: Codebook ChIP-seq vs. external ChIP-seq performed in HEK293 cells (**B**), HepG2 cells (**C**), or K562 cells (**D**). **E**, AUROC scores for expert curated Codebook PWMs (columns), discriminating ChIP-seq peaks vs. random genomic background loci. Rows show different cell types. **F**, **G**, comparison of Codebook and external PWMs at the task of discriminating ChIP-seq peak sets from random sequences (as in **E**), for the 19 TFs that have a Codebook peak set (CP), a Codebook motif (CM), an external peak set (EP), and an external motif (EM), for Codebook ChIP-seq data (**F**) and external ChIP-seq data (**G**). The seven TFs with an AUROC of  $< 0.55$  on either axis of either plot are highlighted. **H**, Sequence logos for the seven TFs highlighted in **F** and **G**. All Codebook PWMs shown are supported by ChIP-seq, GHT-SELEX, and HT-SELEX. Asterisk indicates that the Codebook PWM is additionally supported by SMiLE-seq data.



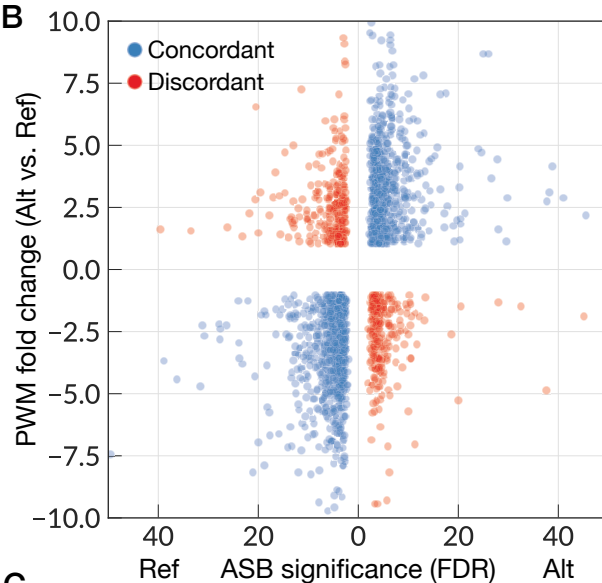
**Figure S4. Number of CTOP sites per TF.** Bar graph displays the number of individual CTOP sites obtained for each TF. Heatmap and annotations below indicate other properties of each TF and its TOP sites.



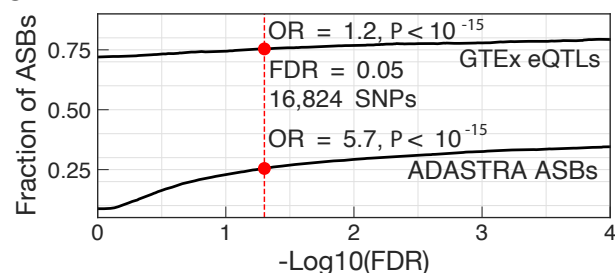
**A**



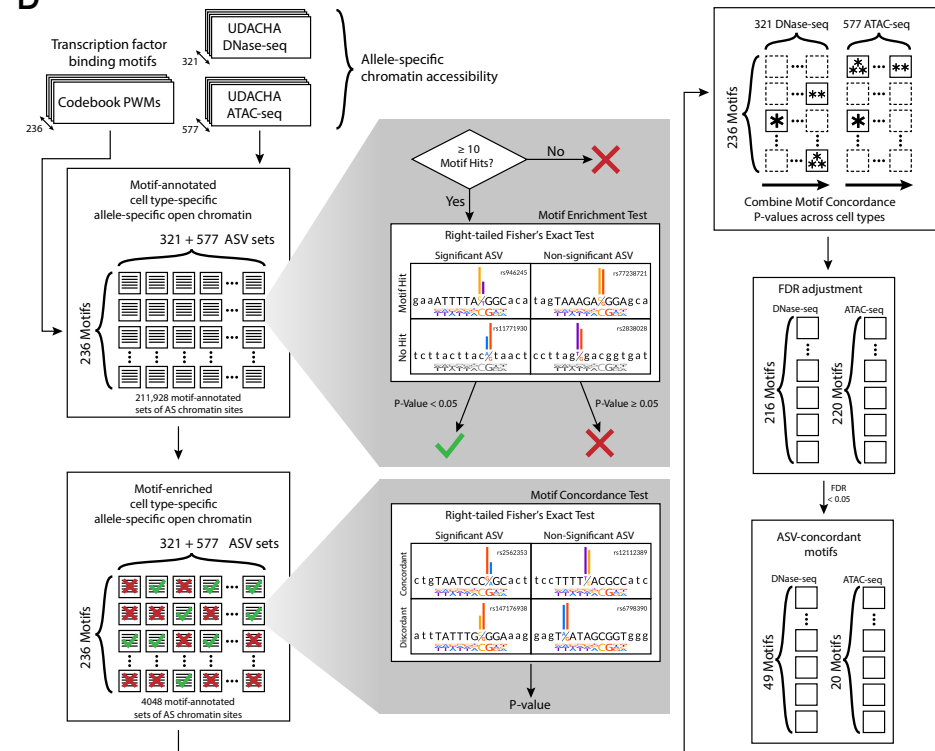
**B**



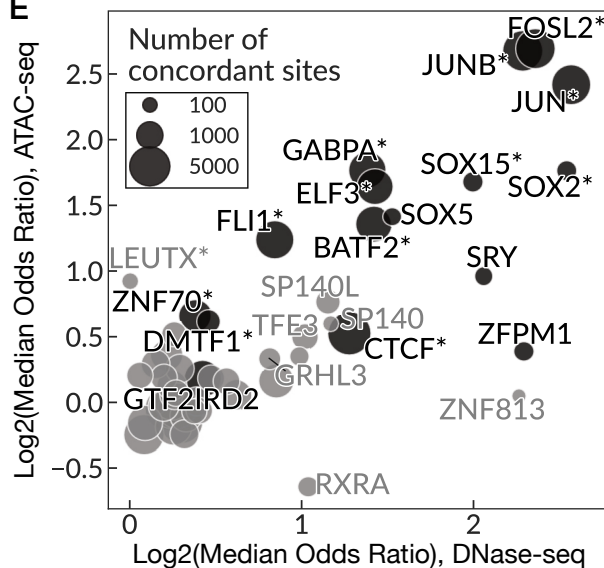
**C**



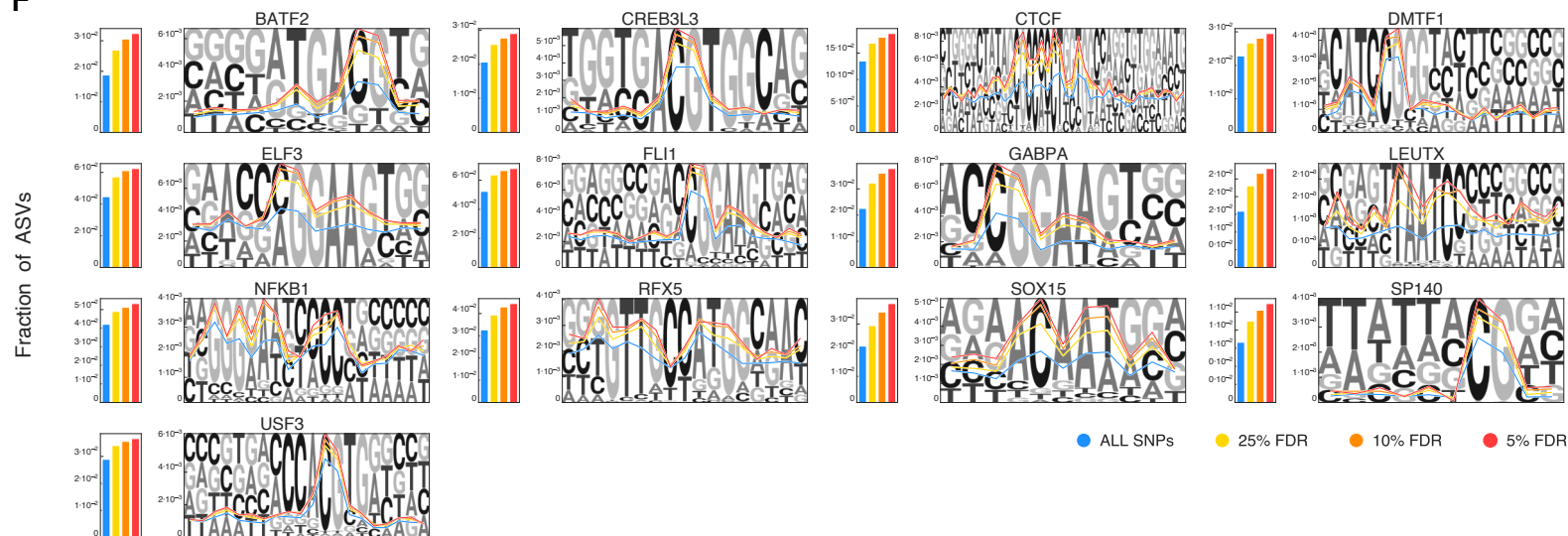
**D**



**E**



**F**



**Figure S5. Identifying allele-specific TF binding in HEK293 cells and analyzing allele-specific chromatin accessibility events using Codebook motifs.**

**A**, Codebook ASB calling workflow: SNP calling with bcftools, mapping bias correction with WASP, background allelic dosage reconstruction with BABACHI, statistical scoring of the allelic imbalance with MIXALIME, and motif annotation with PERFECTOS-APE. **B**, Motif concordance of Codebook ASBs. X-axis: ASB significance (i.e., allelic preference;  $\log_{10}$  FDR, *minus side*: preference for Ref, *plus side*: preference for Alt). Y-axis:  $\log_2$  PWM score fold-change between Alt vs. Ref. The plot shows only strongly concordant and strongly discordant sites with  $|\log_2(\text{Fold Change})| \geq 1$ . **C**, Fraction of Codebook ASBs (combined) coinciding with GTEx eQTLs and ADAstra known ASBs at different FDR thresholds for ASB calling. Fisher's exact test odds ratios (OR) and P-values for ASBs at 5% FDR (covering 16,724 SNPs, dashed line) are labeled on the plot. **D**, Workflow for detection of TFs involved in allele-specific chromatin accessibility. UDACHA DNase-seq and ATAC-seq ASVs across different cell types were annotated with Codebook motifs, followed by motif enrichment and motif concordance analysis, combining the resulting P-values across the cell types, and FDR correction for multiple tested motifs. Central call-outs: details of the motif enrichment and motif concordance test using SP140 motif for illustration. SNPs (rs946245, rs77238721, rs11771930, rs2838028, rs2562353, rs12112389, rs147176938, rs6798390) illustrating the cells of the 2x2 contingency tables are actual UDACHA ASVs with or without motif hits of selected TFs. **E**, Scatterplot of Median Odds Ratios of PWM scores within the ASVs enriched in and concordant with the PWM matches. Motifs significant for both DNase-seq and ATAC-seq (black), or just one assay (gray). The asterisk denotes TFs that exhibit significant enrichment considering peaks-supported PWM hits only. **F**, *Bar plots*: Fraction of ASVs overlapping with PWM hits for 13 TFs, using 4 different thresholds on ASV significance: all SNPs (blue), 25% FDR ASVs (yellow), 10% FDR ASVs (orange), and 5% FDR ASVs (red). *Line plots*: Fraction of ASVs at each location within the genome-wide PWM hits of the representative TFs (P-value<0.001) using four thresholds (the same colors as in bar plots). SNP: single-nucleotide polymorphism, ASB: allele-specific binding, ASV: allele-specific chromatin accessibility variant.