

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

GHT-SELEX demonstrates unexpectedly high intrinsic sequence specificity and complex DNA binding of many human transcription factors

Arttu Jolma^{1,*}, Aldo Hernandez-Corchado^{2,3*}, Ally W.H. Yang^{1,*}, Ali Fathi^{4,*}, Kaitlin U. Lavery^{1,5*}, Alexander Brechalov¹, Rozita Razavi¹, Mihai Albu¹, Hong Zheng¹, The Codebook Consortium, Ivan V. Kulakovskiy^{6,7}, Hamed S. Najafabadi^{2,3**}, and Timothy R. Hughes^{1,4,**}

¹Donnelly Centre, University of Toronto, Toronto, ON M5S 3E1, Canada

²Department of Human Genetics, McGill University, Montréal, QC H3A 0C7, Canada

³Victor P. Dohdaleh Institute of Genomic Medicine, Montréal, QC H3A 0G1, Canada

⁴Department of Molecular Genetics, University of Toronto, Toronto, ON M5S 1A8, Canada

⁵Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY 10065, USA

⁶Vavilov Institute of General Genetics, Russian Academy of Sciences, 119991, Moscow, Russia and Institute of Protein Research, Russian Academy of Sciences, 142290, Pushchino, Russia

⁷Vavilov Institute of General Genetics, Russian Academy of Sciences, 119991, Moscow, Russia

*These authors contributed equally.

**email: hamed.najafabadi@mcgill.ca and t.hughes@utoronto.ca.

40 **The Codebook Consortium**

41

42 **Principal investigators (steering committee)**

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

45

46 **Analysis Centers:**

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

76

77

78

79

80

81 **SUMMARY**

82 **A long-standing challenge in human regulatory genomics is that transcription**
83 **factor (TF) DNA-binding motifs are short and degenerate, while the genome is**
84 **large. Motif scans therefore produce many false-positive binding site predictions.**
85 **By surveying 179 TFs across 25 families using >1,500 cyclic *in vitro* selection**
86 **experiments with fragmented, naked, and unmodified genomic DNA – a method**
87 **we term GHT-SELEX (Genomic HT-SELEX) – we find that many human TFs**
88 **possess much higher sequence specificity than anticipated. Moreover, genomic**
89 **binding regions from GHT-SELEX are often surprisingly similar to those obtained**
90 ***in vivo* (i.e. ChIP-seq peaks). We find that comparable specificity can also be**
91 **obtained from motif scans, but performance is highly dependent on derivation**
92 **and use of the motifs, including accounting for multiple local matches in the**
93 **scans. We also observe alternative engagement of multiple DNA-binding domains**
94 **within the same protein: long C2H2 zinc finger proteins often utilize modular DNA**
95 **recognition, engaging different subsets of their DNA binding domain (DBD) arrays**
96 **to recognize multiple types of distinct target sites, frequently evolving via internal**
97 **duplication and divergence of one or more DBDs. Thus, contrary to conventional**
98 **wisdom, it is common for TFs to possess sufficient intrinsic specificity to**
99 **independently delineate cellular targets.**

100

101 **Keywords:** DNA binding specificity, Transcription factor, TF, Transcription factor binding
102 site, Position weight matrix, PWM, ChIP-Seq, HT-SELEX, GHT-SELEX, SELEX,
103 Modular binding, C2H2 zinc finger, C2H2, RCADEEM, MAGIX, Codebook, Gene
104 regulation

105

106 INTRODUCTION

107 The DNA-binding sequence preference of a Transcription Factor (TF) is typically
108 referred to as a motif, and is most commonly modeled as a position weight matrix
109 (PWM), which describes the relative preference of the TF for each base in the binding
110 site¹. In human, TF binding motifs are generally short and flexible; PWMs are typically 8-
111 14 bases long²⁻⁴, and multiple bases can be tolerated at many positions^{5,6}. Thus, a
112 typical TF PWM scan with default parameters yields over a million potential binding
113 sites in the 3-billion-base human genome, often with multiple high-scoring matches per
114 gene. Very few of the potential target sites are utilized in cells⁷, however, and the actual
115 number of bound sites, as measured by ChIP-seq⁸⁻¹⁰ or other assays¹¹ is typically much
116 lower than the number of motif matches.

117
118 This deficit in specificity has been resolved conceptually by the widespread cooperative
119 binding and synergy among TFs^{5,6,12,13}, and evidence that the chromatin landscape
120 generally dominates TF binding site selection³⁰, such that TF motif matches only
121 determine binding within permissible regions. In the latter model, only a special class of
122 “pioneer” TFs can access target sequences to control the local chromatin. Indeed, some
123 TFs have been shown to have high inherent specificity: for example, CTCF binds the
124 majority of its strongest motif matches in the genome¹⁴, and repositions the surrounding
125 nucleosomes¹⁵. PRDM9, which controls recombination hotspots, has been reported to
126 independently specify roughly half of its binding sites in the genome¹⁶. Another possible
127 explanation for the generally low apparent specificity of TF motifs, however, is that
128 PWMs are inaccurate, or are used inappropriately, or that the PWM model is
129 fundamentally flawed¹⁷. PWMs are often derived from a non-comprehensive set of
130 bound vs. unbound sequences, and there is ongoing controversy regarding the best
131 methods for derivation, underlying representation, and scanning of TF motifs^{1,18}, as well
132 as the impact of DNA shape¹⁹, dependencies among base positions^{17,20}, multimeric
133 binding^{21,22}, and lower-affinity binding sites²³.

134
135 Many human TFs still lack binding motifs, and prominent among them are hundreds of
136 C2H2 zinc finger (C2H2-zf) proteins²⁴. These proteins recognize DNA sequences that
137 approximate a concatenation of the three or four base specificities of their sequential
138 constituent C2H2-zf domains^{25,26}. Different C2H2-zf proteins can bind very different
139 motifs due to both the malleability of the individual C2H2-zf domains and rearrangement
140 of the individual C2H2-zf domains²⁷. An enigmatic feature of the C2H2-zf proteins is
141 their theoretical capacity to recognize very long sequences: the median number of
142 C2H2-zf domains in human TFs is 11, which could contact up to 33 DNA bases, much
143 more than would be needed to specifically recognize even a single target site in the
144 genome, on average. Indeed, C2H2-zf proteins often use only a subset of their DBDs to
145 contact DNA, and whether and how frequently human C2H2-zf proteins utilize different
146 segments of the C2H2-zf domain array to bind different sequences has also been a
147 long-standing question. In a well-studied example, CTCF binding sites appear to reflect
148 a constitutive “core”, bound by fingers 4-7 of the 11 C2H2-zf domain array, flanked by
149 sequences that are bound by alternative usage of upstream and/or downstream C2H2-
150 zf domains^{28,29}. Analysis of the DNA-binding of C2H2-zf proteins to the genome is also

151 complicated by the fact that they often bind repeat elements such as endogenous
152 retroelements³⁰, and thus the target site similarity is derived both from DNA recognition
153 and the shared ancestry of the binding sites. The limited resolution of ChIP-seq
154 (>100bp) presents a related hindrance. These confounding factors, however, can be
155 ameliorated by incorporating information about the bases that are likely preferred at
156 each position of the binding site, as predicted by a C2H2-zf “recognition code” that
157 relates the C2H2-zf amino acid sequences to their binding preferences. These machine
158 learning-based predictions can assist in identifying the most plausible protein-DNA
159 interactions in such cases, as our earlier work demonstrated³¹.

160
161 Here, we describe GHT-SELEX (Genomic DNA HT-SELEX), a novel implementation of
162 the HT-SELEX³² method for identification of the sequence specificity of DNA-binding
163 proteins. HT-SELEX is a high-throughput implementation of SELEX (Systematic
164 Evolution of Ligands by EXponential enrichment)³³, using multi-cycle, automated affinity
165 capture of protein-bound DNA in microwell plates, coupled to multiplexed Illumina
166 sequencing. HT-SELEX utilizes random-sequence DNA, while GHT-SELEX is instead
167 performed with fragmented human genomic DNA, and uses an associated new
168 statistical analysis method, MAGIX (Model-based Analysis of Genomic Intervals with
169 eXponential enrichment). GHT-SELEX is conceptually similar to Affinity-seq¹⁶ and DAP-
170 seq³⁴, but it incorporates multiple selection cycles, and is thus related to earlier genomic
171 SELEX approaches that utilized Sanger sequencing^{35,36}. The use of barcoding,
172 magnetic affinity beads and laboratory automation makes it possible to run GHT-SELEX
173 in parallel with hundreds of samples. We developed GHT-SELEX in the context of the
174 Codebook consortium project³⁷, which was aimed primarily at analysis of 332
175 uncharacterized putative TFs (together with 61 control TFs), and provides comparison
176 data from several other platforms for the same set of TFs (HT-SELEX, ChIP-seq,
177 Protein Binding Microarrays³⁸, and SMiLE-seq³⁹). We successfully applied GHT-SELEX
178 to 179 human TFs, most of which are poorly characterized, thus providing a major
179 expansion in the number of human TF motifs. For dozens of TFs, including some that
180 are considered well-characterized, GHT-SELEX peaks correspond with *in vivo* binding
181 (measured by ChIP-seq) much more accurately than current models would suggest.
182 GHT-SELEX is particularly effective for C2H2-zf proteins, and shows that they often use
183 alternative subsets of their C2H2-zf domains to engage with different genomic target
184 sites. We explore both explanations and ramifications of these observations.

185 186 **RESULTS**

187 188 **Development and testing of GHT-SELEX**

189
190 GHT-SELEX combines the principles of previous genomic DNA selection protocols^{16,34}
191 with HT-SELEX, a method that has been applied successfully to hundreds of human
192 TFs and is compatible with robotics^{32,40}. We developed GHT-SELEX (**Figure 1A**) to run
193 in parallel with HT-SELEX, in the context of the Codebook project. The intended
194 purpose, initially, was to create a DNA library that contains sufficient representation of
195 long repeat sequences that are common in the human genome (e.g. transposons and
196 endogenous retroelements): we reasoned that the difficulty of obtaining long motifs

197 expected for C2H2-zf proteins may be due to the scarcity of long binding sites in a
198 random pool, since representation of any sequence would decrease exponentially with
199 its length. The GHT-SELEX DNA pool used in this study was produced by nonspecific
200 enzymatic fragmentation of HEK293 DNA to fragments with a median length of ~64 bp.
201 HEK293 DNA was chosen for compatibility with ChIP-seq data generated
202 simultaneously (see accompanying manuscript⁴¹), and the length of the DNA was
203 chosen to mimic standard HT-SELEX procedures and provide relatively high resolution.
204

205 We initially tested GHT-SELEX on the Codebook control proteins. Thirty of the controls
206 represented a sampling of well-studied TFs with different classes of DBDs, most of
207 which were previously analyzed using the independent *in vitro* SMiLE-seq platform³⁹. An
208 additional 31 controls were C2H2-zf proteins for which published ChIP-seq data yielded
209 motifs⁴². At the outset, we assumed that GHT-SELEX would yield continuous read
210 coverage across the genome, given conventional estimates of up to a million PWM hits
211 per TF⁷, such that the data could be analyzed directly for enriched motifs among the
212 reads. Indeed, examination of individual mapped reads revealed that they usually
213 accumulate at sites in which all reads overlap with what appears to be a motif match
214 (**Figure 1B**). Remarkably, it also became apparent that GHT-SELEX data typically has a
215 strong resemblance to ChIP-seq data, forming strong peaks found sparsely across the
216 genome. **Figure 1C** shows raw read density for four control TFs, comparing GHT-
217 SELEX to ChIP-seq, and also to target site predictions based on existing and newly-
218 derived (see below) PWM models for the TFs. This observation prompted us to analyze
219 the data as peaks, instead of raw reads.
220

221 Peak calling from the GHT-SELEX data with conventional algorithms is confounded by
222 the fact that different peaks have very different enrichment ratios across the cycles,
223 presumably due to varying affinity of the TF for different sites, the overall increase in
224 motif occurrences in the pool over the successive cycles, and simultaneous reduction in
225 pool complexity, with the strongest binding sites dominating later cycles. As a
226 consequence, enrichment information is distributed across the read cycles, with weaker
227 peaks first appearing and disappearing, and the strongest peaks dominating in the later
228 cycles. To adapt to these issues, we developed an analytical framework that capitalizes
229 on the added information gained from multiple SELEX cycles (**Figure 2A**; see **Methods**
230 for details). The approach relies on a statistical method that explicitly models the
231 exponential growth of TF-bound genomic regions over the SELEX cycles, which leads
232 to a progressively higher proportion of TF-bound fragments and depletion of relative to
233 genomic background. The fragment abundances, in turn, are modeled as latent
234 variables that determine the number of observed reads through a Poisson process. This
235 hierarchical Bayesian model enables the integration of information across different
236 selection cycles, experiments, and batches, to calculate an estimated enrichment
237 coefficient (**Figure 2B**). We refer to this approach as MAGIX (Model-based Analysis of
238 Genomic Intervals with exponential enrichment).
239

240 Among the 61 control proteins, 40 were deemed as successful on GHT-SELEX (see
241 below and accompanying manuscripts^{37,43} for a description of how success was
242 determined). Analysis of the data for the 40 successful controls by MAGIX resulted in

243 between 13 and 137,718 peaks (median 19,400) with enrichment coefficient exceeding
244 5% FDR (see **Methods**). There is a clear enrichment of the motif occurrences for the
245 corresponding TFs within the peaks, with the number of strong PWM hits, on average,
246 declining rapidly at ~50 bp from peak centre, consistent with the DNA fragment size
247 (**Figure 2C**; similar plots for all TFs analyzed are shown in **Document S1**). In addition,
248 higher PWM scores (which would, in theory, predict higher relative affinity) are clearly
249 associated with a higher GHT-SELEX enrichment coefficient (see below), suggesting
250 that the GHT-SELEX/MAGIX is quantitative to some degree.

251

252 **Application of GHT-SELEX to the Codebook TF set**

253

254 We next performed GHT-SELEX and, in parallel, HT-SELEX using fragmented genomic
255 DNA and random 40N ligands (**Table S1**), respectively, to assess DNA binding activity
256 of 331 poorly characterized putative human TFs, as part of the Codebook project. We
257 analyzed individual TFs with up to three types of constructs, and up to three protein
258 expression strategies (two types of *in vitro* transcription–translation reactions, and
259 expression in HEK293 cells, see **Methods**). Several experimental variables were
260 modulated over the course of the experiments, resulting in improvement of success
261 rates, particularly for TFs with long C2H2-zf domain arrays (see **Methods and Table**
262 **S2**). For each TF, the constructs contained the full sequence of a representative
263 isoform, or either all or a subset of its predicted DBDs. In total, we analyzed 1,315
264 constructs encompassing the 61 control TFs and 331 of the 332 putative TFs in the
265 Codebook set of poorly characterized proteins. With these constructs we performed
266 1,534 GHT-SELEX and 1,578 HT-SELEX experiments (see **Methods and Table S3**).

267

268 In separate parts of the Codebook project, this same set of proteins was analyzed using
269 ChIP-seq, Protein Binding Microarrays³⁸, and SMiLE-seq³⁹, as described in the
270 accompanying manuscripts^{37,41,44}. We gauged the success of each TF in each
271 experiment, including the GHT-SELEX experiments, largely based on whether similar
272 DNA-binding motifs (i.e. PWMs) were obtained from different types of experiments, with
273 all data types considered in aggregate by a team of expert curators. This process
274 produced a list of “approved” experiments, as described in an accompanying study⁴³.
275 Selection of a single PWM for each TF for subsequent analyses is described in
276 accompanying study³⁷. The PWM selections incorporated those generated from all data
277 types. PWMs and logos are available in accompanying study³⁷ and online at
278 <https://codebook.ccb.utoronto.ca>, <https://mex.autosome.org>, and
279 <https://cisbp.ccb.utoronto.ca>⁴⁵.

280

281 In total, 139 previously uncharacterized Codebook TFs had at least one “approved”
282 GHT-SELEX experiment (i.e., were successful in GHT-SELEX), of which 131 were also
283 approved in HT-SELEX, 108 in ChIP-seq, and 102 in all three (**Figure 3A and Table**
284 **S3**). The 139 were comprised mainly of C2H2-zf proteins, which are prevalent in the
285 Codebook set (**Figure 3B**). In contrast, 163 of the putative TFs did not yield motifs in
286 any of these assays, suggesting that they either do not bind DNA with sequence
287 specificity, or require post-translational modifications or cofactors. In particular, only two
288 of 49 proteins tested that lacked a known DBD yielded an approved experiment in GHT-

289 SELEX (discussed in greater detail in the accompanying studies³⁷). Including the control
290 TFs, 24 types of DBDs were present among the approved experiments (**Figure 3B**),
291 illustrating that the method can capture motif-containing genomic target site locations of
292 diverse TF types.

293
294 **Unexpectedly high overlap between TF binding to the genome *in vitro* and *in vivo***

295
296 GHT-SELEX analyzed with MAGIX, like ChIP-seq, produces peaks with a continuum of
297 enrichment coefficient values and other associated statistics. Across both Codebook
298 TFs and controls, there are typically a relatively small number of peaks with
299 exceptionally high MAGIX enrichment coefficient values (hundreds to thousands), but
300 we did not observe bimodal distributions that would imply a natural threshold which
301 could be used to discriminate “bound” from “unbound” loci (examples in **Figure 4A**;
302 distributions for all TFs in **Document S1**). We also examined the correspondence
303 between GHT-SELEX/MAGIX peaks, ChIP-seq peaks, and PWM scores, focusing on
304 the 137 TFs for which both ChIP-seq and GHT-SELEX data were available (101
305 Codebook TFs and 36 controls). In most cases, there was a much higher overlap with
306 ChIP-seq peaks and high PWM scores among the highest-scoring GHT-SELEX/MAGIX
307 peaks (examples are shown in **Figure 4A**, and plots for all TFs in **Document S1**). We
308 did not, however, identify a specific peak enrichment coefficient or significance value
309 across all experiments that uniformly corresponds to high enrichment of PWM hits, or
310 the probability of overlap with ChIP-seq peaks.

311
312 Lack of a universal enrichment coefficient threshold across all experiments could be
313 accounted for by TF-specific parameters in both GHT-SELEX and ChIP-seq assays,
314 including different binding kinetics for both sequence-specific and nonspecific DNA
315 binding, the effective concentration of the TFs, and the ability of the TFs to compete or
316 cooperate with nucleosomes and other cofactors *in vivo*. Given that these parameters
317 are unobserved and difficult to estimate from the data available, we implemented a
318 simple scheme to draw thresholds on both peak sets: by sequentially taking equal
319 numbers of highest scoring peaks on a TF-specific basis, we identified the peak number
320 that maximizes the Jaccard statistic of overlap between the GHT-SELEX/MAGIX peaks
321 and ChIP-seq peaks (**Figure 4B**).

322
323 This approach yielded a very striking result, which is that for many TFs, a peak number
324 can be identified with a surprisingly high Jaccard value (Jaccard median 0.1117) (**Figure**
325 **4C,D** and **Table S4**), indicating that the TF intrinsically (i.e. independently) specifies
326 many of the *in vivo* binding sites above the threshold selected. Peak overlap is a
327 demanding statistic, because random expectation (i.e. from choosing genomic regions
328 at random) is near zero, as only a miniscule fraction of the genome is covered by the
329 peaks in either data type, and both experimental variation and noise in generation of
330 peaks will lead to fluctuation of the rank order of peaks, even for replicates. Indeed, this
331 result is not obtained from permuted peak positions, or permuted experiments (i.e.
332 mismatched TFs) (after permutation, Jaccard median 0.0073; Wilcoxon $p=2.6 \times 10^{-38}$)
333 (**Figure 4D**). The peak numbers yielding these high Jaccard values are often relatively

334 low, and correspond to a wide range of ChIP-seq p-value thresholds and MAGIX
335 enrichment coefficient values (**Table S4**).

336
337 Overall, this outcome contrasts with traditional expectation, which is that individual TF
338 would normally not be able to independently specify their DNA targets in the genome⁷.
339 We note that many of the TFs with highest Jaccard maxima are uncharacterized C2H2-
340 zf proteins with long (and intuitively specific) motifs: Among those with Jaccard > 0.1,
341 78% are C2H2-zf proteins (57 out of 73), vs 42% (27 out of 64) for those with Jaccard
342 below 0.1, and overall, the median Jaccard value for C2H2-zf proteins is 0.1582, vs.
343 0.0616 for non-C2H2-zf proteins (Wilcoxon $p=5.14 \times 10^{-8}$). CTCF, a control protein that is
344 known to possess large number of genomic target sites, unusually high specificity and
345 ability to control nucleosome positions^{14,15}, is among those with high Jaccard values,
346 although it is not the highest scoring in this dataset. Counterintuitively, high Jaccard
347 maxima were also obtained for a subset of TFs with relatively short motifs, including
348 NFKB1, GABPA, NACC2, and several CXXC proteins, such as CXXC4 and KDM2A,
349 that mainly bind CG dinucleotides, as expected⁴⁶ (**Figure 5A**).

350 351 **Multiple explanations for high sequence specificity observed in GHT-SELEX**

352
353 We next asked whether PWM predictions across the genome could achieve a level of
354 correspondence to ChIP-seq that we obtained with GHT-SELEX/MAGIX. To do this, we
355 performed a similar maximization of the overlap score (Jaccard) as described above for
356 GHT-SELEX and ChIP-seq, here sweeping through PWM scores (i.e. using PWMs to
357 predict and score “peaks” in the genome; see **Methods** for details). Remarkably, on
358 average, the overlap between PWM predictions and ChIP-seq peaks is similar to that
359 for GHT-SELEX/MAGIX and ChIP-seq peaks (**Figure 5A, Table S4**), and the numbers
360 of peaks at which the maximum Jaccard was obtained is also typically similar (**Figure**
361 **5B**). The slightly higher Jaccard for PWMs in some cases may be due to the simple
362 PWM models smoothing experimental noise in the GHT-SELEX. In some cases,
363 however, this explanation seems implausible; for example, in several instances, very
364 small PWMs (e.g. that of CXXC4, which is mainly a single CG dinucleotide) yielded
365 higher overlap with ChIP-seq peak locations than GHT-SELEX did.

366
367 To our knowledge, such strong ability of PWMs to predict *in vivo* binding sites, over a
368 large set of TFs, is unprecedented. We attribute two main sources. First, the PWMs
369 used in these analyses were selected from a panel of hundreds to thousands of
370 candidate PWMs, specifically choosing those that performed best across numerous test
371 statistics and several data types. The Jaccard statistics against ChIP-seq and GHT-
372 SELEX were among the selection criteria. Thus, lower maximal Jaccard scores – often
373 vastly lower - are obtained from virtually all other PWMs. Hence, in addition to
374 optimizing the thresholds, part of the explanation for the high Jaccard values we
375 obtained lies in the derivation of the PWM itself.

376
377 The second apparent source of performance increase is the PWM scanning and scoring
378 method. For some TFs, scoring a DNA fragment using the sum of predicted affinity
379 scores over a sequence window (i.e. the sum of the PWM probability scores at

380 individual positions, rather than the log-odds that is output by most PWM scanning
381 tools) results in considerably higher maximum Jaccard value than taking the maximum
382 or sum of log-odds PWM scores (which are generally thought to represent binding
383 energy^{47,48}) (**Figure 5C**). Sum-of-affinity scoring presumably reflects the cooperation of
384 multiple adjacent binding sites, traditionally referred to as “avidity”⁴⁹. The effect is most
385 striking for a subset of TFs that bind short or repetitive sequences, including CG
386 dinucleotides and poly-A stretches (**Figure 5D**, but it also appears to underpin the
387 specificity of NACC2 and ZNF48, which have unique, non-repetitive motifs (**Figure 5E**).
388 Points above the diagonal in **Figure 5A**, where PWM prediction shows higher overlap
389 with ChIP than the GHT-SELEX, may therefore represent the impact of TF binding sites
390 over a larger window influencing ChIP-seq but not GHT-SELEX (we performed the
391 PWM scans with a 200 bp window, while the GHT-SELEX fragments are only ~65 bp).
392 For example, scanning 200-base windows with the short CG motif for CXXC4 may be
393 better suited for detection of CpG islands (which dominate the CXXC4 binding sites³⁷, in
394 which the CG dinucleotides will be distributed over a large region (by definition ≥ 200
395 bp).

396
397 These analyses indicate that PWMs can often predict *in vivo* TF binding sites as
398 effectively as actual measurements of binding to the genome made with GHT-SELEX.
399 **Figure 1C** illustrates the increase in correspondence between PWM predictions and
400 ChIP-seq peaks that can be achieved with carefully-selected PWMs and improved
401 scanning procedures. There are, however, many TFs in which no PWM could be
402 derived that rivals GHT-SELEX data in correspondence of ChIP-seq peaks (those below
403 the diagonal in **Figure 5A**). These TFs are almost entirely proteins with a long array of
404 C2H2-zf domains, which we examine more closely in the next section.

405 406 **Alternate usage of C2H2-zf domains within large arrays**

407
408 The expansive collection of GHT-SELEX, HT-SELEX, and ChIP-seq data for C2H2-zf
409 proteins provided an opportunity to examine the long-standing issue of usage of
410 individual C2H2-zf domains within large arrays. Anecdotally, we observed many
411 instances where the motifs detected for C2H2-zf proteins were much shorter than
412 expected based on the number of C2H2-zf domains, as well as examples in which
413 multiple distinct motifs emerged, suggesting that the TFs might use partial subsets of
414 their DBD array to engage DNA at different locations. Proving differential engagement of
415 the specific C2H2-zf domains is challenging, however, due to low statistical power
416 (there are many possible C2H2-zf domain sub-arrays, and a limited number of highly
417 enriched peaks) and the fact that the genome is highly non-random and repeat-rich. To
418 minimize the impact of these issues, we developed a new method that utilizes the
419 C2H2-zf recognition code to assess which sets of C2H2-zf domains are likely to be
420 engaged at any individual binding sites. We call this method RCADEEM (Recognition
421 Code-Assisted Discovery of regulatory Elements by Expectation-Maximization) (see
422 **Methods** for details). **Figure 6A** shows a schematic, and the results of applying
423 RCADEEM to CTCF, illustrating that it produces a “core” motif recognized by fingers 4-7
424 at all sites, and alternative usage of flanking C2H2-zf domains in a subset of sites, very

425 similar to the differential usage of CTCF C2H2-zf domains that has been previously
426 described²⁸.

427
428 We applied RCADEEM to all 120 C2H2-zf proteins for which we had approved data
429 from GHT-SELEX (**Table S4**). We applied RCADEEM on GHT-SELEX data and
430 separately, if available, on HT-SELEX and ChIP-seq; for GHT-SELEX and ChIP-seq, we
431 applied it both with and without repeat sequences (i.e. removing any peaks that overlap
432 with the UCSC Repeatmasker track). In total, we obtained RCADEEM predictions for 86
433 of them (**Table S4**), all of which are available via the web resources accompanying this
434 paper (<https://codebook.cabr.utoronto.ca/>). (For the remaining 34, the algorithm did not
435 converge, suggesting that the sequence preferences of the protein do not closely follow
436 the recognition code, and thus cannot be analyzed in this way). Most of the 86
437 displayed what appears to represent alternative usage of segments of the C2H2-zf
438 domain array on different DNA molecules (e.g. different genomic loci) within the same
439 experiment. We manually classified the apparent C2H2-zf domain usage into the
440 following categories, examples of which are shown in **Figure 6B-F**, while **Figure 6G**
441 provides an overview of the descriptors and other properties of each of the C2H2-zf
442 proteins. 1) *Canonical (30 instances)* follows the baseline assumption that a TF always
443 uses the same set of C2H2-zf domains to recognize sites that can be described with a
444 single PWM. 2) *Core with extensions (24 instances)*, where all sites share a sequence
445 motif bound by a subset of the C2H2-zf domains, which is supplemented by recognition
446 of flanking sequences by adjacent C2H2-zf domains at some binding sites. 3) *Finger*
447 *shift (14 instances)*, where the TF recognizes a range of tiled target sites by binding with
448 variable subsets of adjacent C2H2-zf domains. 4) *Multiple DBDs (32 instances)*, in
449 which subsets of the C2H2-zf domain array appear to function as independent DBDs.
450 The last three binding modes are not mutually exclusive. For example, ZNF471 displays
451 both multiple DBDs and core with extensions with one of the DBDs (**Figure 6F**), while
452 the long finger shift in ZNF665 (**Figure 6D**) leads effectively to multiple DBDs, as the
453 target sites of most N-terminal and C-terminal ends do not overlap with each other.
454 **Table S4** lists the annotations for all 86 proteins.

455 456 **Evolution of C2H2-zf protein DNA-binding specificities via internal duplication**

457
458 In the RCADEEM outputs, different segments of a C2H2-zf domain array (i.e., different
459 DNA binding regions of the protein) are often predicted to bind similar yet distinct sets of
460 sequences. For example, ZNF775 (**Figure 6E**) binds two types of sites that contain a
461 shared GNWGAA consensus, followed by either TTT or GCA trinucleotides. RCADEEM
462 predicts that these two sites are recognized by C2H2-zf domain arrays 1-4 and 5-8,
463 respectively. Indeed, arrays 1-4 and 5-8, as well as 9-11, are homologous, on the basis
464 of sequence identity (visualized at
465 <https://codebook.cabr.utoronto.ca/details.php?TF=ZNF775>), suggesting that they arose
466 from duplications. All three arrays are present in mammals as distant as the Tasmanian
467 devil, indicating that the duplications predate divergence from marsupials, and have
468 since been conserved. The cellular and physiological functions of this protein are
469 unknown, to our knowledge, but this degree of sequence conservation suggests a
470 conserved role across mammals.

471
472 Another example is ZNF721: RCADEEM indicates that it has three DNA-binding modes,
473 with related but distinct motifs (**Figure 7A**), corresponding to homologous C2H2-zf
474 domain arrays containing fingers 6-13, 12-16, and 18-22 (**Figure 7B**). The distinct
475 sequence preferences of the duplicated ZNF721 arrays are supported by experimental
476 data for partial “DBD1” and “DBD2” constructs, corresponding roughly to the first and
477 second half of the full array, which recognize largely distinct subsets of the genomic
478 sites bound by full length TF in GHT-SELEX (**Figure 7A**) and prefer almost entirely
479 distinct 10-mers in HT-SELEX (**Figure 7C**). The function of ZNF721 has not been
480 determined, but sequences recognized by the first (6-13) and third (18-22) duplicated
481 C2H2-zf domain arrays of ZNF721 are found in the highly numerous Alpha repeats,
482 which are fast-evolving elements found at primate centromeres⁵⁰. ZNF721 itself is
483 present only in primates. ZNF721 also binds thousands of unique loci outside known
484 repeat elements, and associates physically with TRIM28/KAP1⁵¹, suggesting a role in
485 gene silencing or heterochromatin formation.

486
487 To survey the prevalence of internal duplication of C2H2-zf domains, we compared all
488 pairs of individual human C2H2-zf domains occurring in the same protein and found that
489 185 human C2H2-zf proteins (~25%) contain at least one pair of C2H2-zf domains that
490 differ by 3 or fewer edits (substitutions, deletions or insertions; **Table S6**), suggesting
491 that they are derived from recent duplications. Furthermore, as in ZNF775 and ZNF721,
492 there are 140 proteins with apparent internal C2H2-zf domain array duplications,
493 defined as two (or more) adjacent C2H2-zf domains (i.e. an array) related to a second
494 such array with two (or more) C2H2-zf domains, with 5 or fewer edits per C2H2-zf
495 domain. Based on recognition code predictions, C2H2-zf domain arrays within internal
496 array duplications have more diverged sequence specificities from each other than
497 individually duplicated C2H2-zf domains (**Figure 7D, Table S6**). The prevalence and
498 diversification of internal C2H2-zf domain array duplications suggest that they are a
499 common modality for evolution of novel functional roles for this large class of proteins.

500
501
502

503 DISCUSSION

504
505 GHT-SELEX assays direct and unassisted binding of a single TFs to the unmodified and
506 unchromatinized genome *in vitro*, revealing surprisingly specific intrinsic sequence
507 preferences for many human TFs. The assay, and the associated MAGIX analysis
508 pipeline, offers several technical advantages over alternatives, including smaller
509 fragment size and compatibility with the same instrumentation used for HT-SELEX.
510 GHT-SELEX data are often more similar to ChIP-seq data than conventional wisdom
511 would suggest it should be^{7,24}, indicating that, for an apparently large subset of TFs,
512 chromatin and cofactors have less critical influence on where binding occurs. This same
513 observation implies that this subset of individual TFs may have greater ability to
514 overcome the chromatin state than is commonly believed. This apparent discrepancy
515 with expectation can be explained partly by technical shortcomings in previous PWM-
516 based genome scans, which are based on PWMs derived from other methods. HT-
517 SELEX and other *in vitro* approaches utilizing random sequence are powerful in that
518 they are unbiased in terms of sequence composition⁵², but they are inherently limited in
519 sequence length and context that can be surveyed. ChIP-seq is invaluable because it
520 can assay binding within cells, but it does not inherently discern direct, indirect, and
521 non-specific binding. Thus, PWMs derived from ChIP-seq and other *in vivo* approaches
522 are influenced by factors other than the TF, in addition to the biased sequence content
523 of the genome. GHT-SELEX provides a powerful intermediate that can resolve
524 ambiguities of both motif discovery and PWM scanning, and thus provides data that
525 complements both ChIP-seq and *in vitro* assays that utilize random sequences.

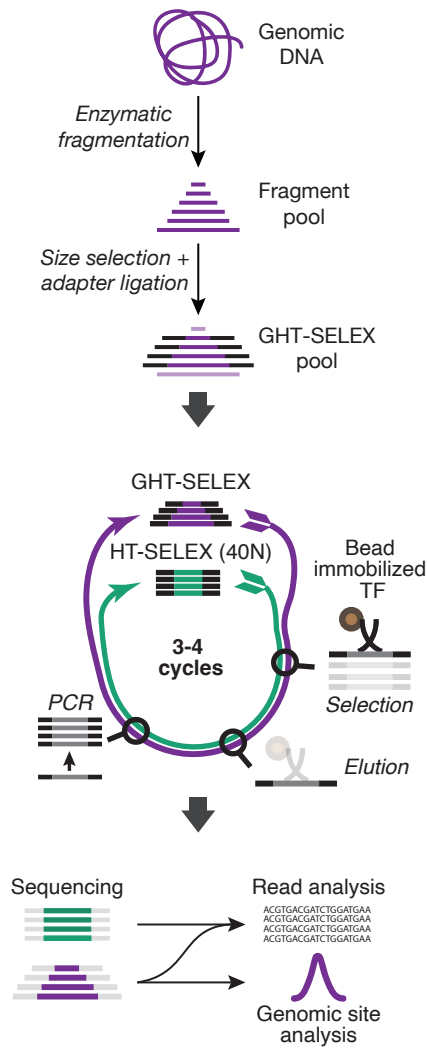
526
527 GHT-SELEX is particularly effective with C2H2-zf proteins, and, together with
528 RCADEEM, has an unprecedented ability to both obtain and dissect *in vitro* the multiple
529 binding modes that are uniquely characteristic of this family, and inherently more difficult
530 to represent as a single PWM. The existence of multiple binding modes also provides a
531 potential explanation for the large number of C2H2-zf domains in each protein. In at
532 least some cases, these large arrays derive from internal duplications of segments of
533 the C2H2-zf domain arrays, possibly facilitating generation of evolutionary novelty via
534 duplication and divergence.

535
536 In contrast to the C2H2-zf family, the most well-studied TFs tend to be in the TF classes
537 such as homeodomain, bHLH, bZIP, nuclear receptor and Sox TFs, because they are
538 the most strongly conserved and often dictate specific biological processes (e.g.
539 morphogenesis, body plan, lineage specification, etc.)²⁴. Our study included some of
540 these TFs (e.g. LEUTX, BATF2, RARA and SRY), and they displayed only limited
541 overlap between GHT-SELEX and ChIP-seq peaks, indicating that many of them cannot
542 independently specify *in vivo* binding locations and hence target genes. It has long been
543 known that TFs controlling chromatin in yeast are largely distinct from those that
544 regulate specific pathways²⁴; we speculate that a similar division may exist in human
545 and other animals.

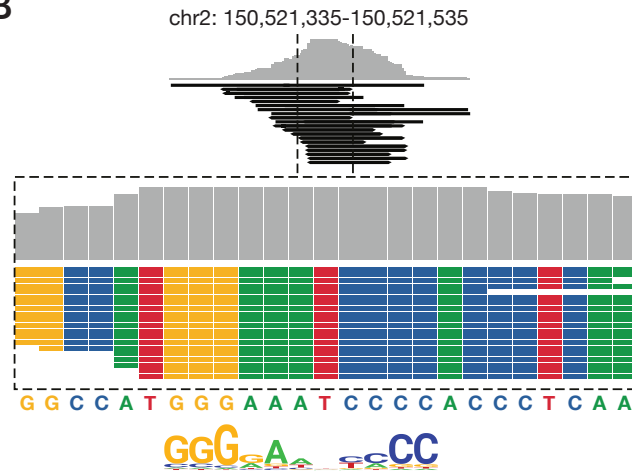
546
547 GHT-SELEX data, together with the larger Codebook dataset, provides an extensive
548 new dataset of TF motifs (i.e. PWMs), encompassing most putative TFs currently

549 lacking them. Accompanying papers provide a thorough analysis of the results of this
550 project, which underscore many challenges and benefits of accurate motif
551 representations. Representation of TF sequence specificity remains an open challenge,
552 more than four decades after the introduction of the standard PWM model⁵³. More
553 accurate representations of large and complex binding sites, in particular for C2H2-zf
554 proteins, could be useful for a variety of purposes, including attributing deep learning
555 filters to individual TFs. Finally, we propose that obtaining data from GHT-SELEX for
556 additional TFs with “known” motifs and genomic binding sites from ChIP-seq will
557 produce a more detailed view of their intrinsic DNA binding abilities, and how this
558 intrinsic ability dictates TF-genome interactions in living cells.

A



B



C

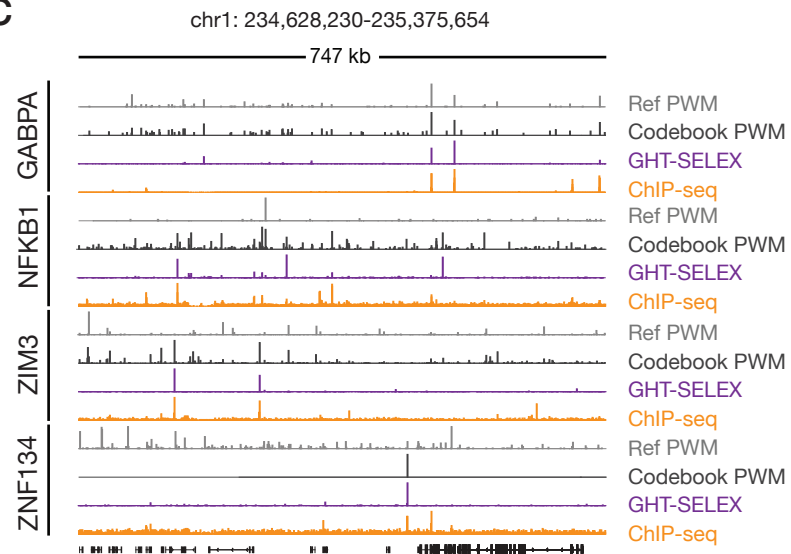
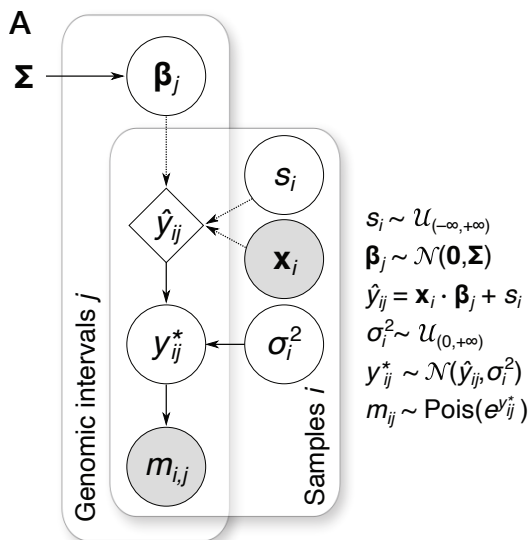
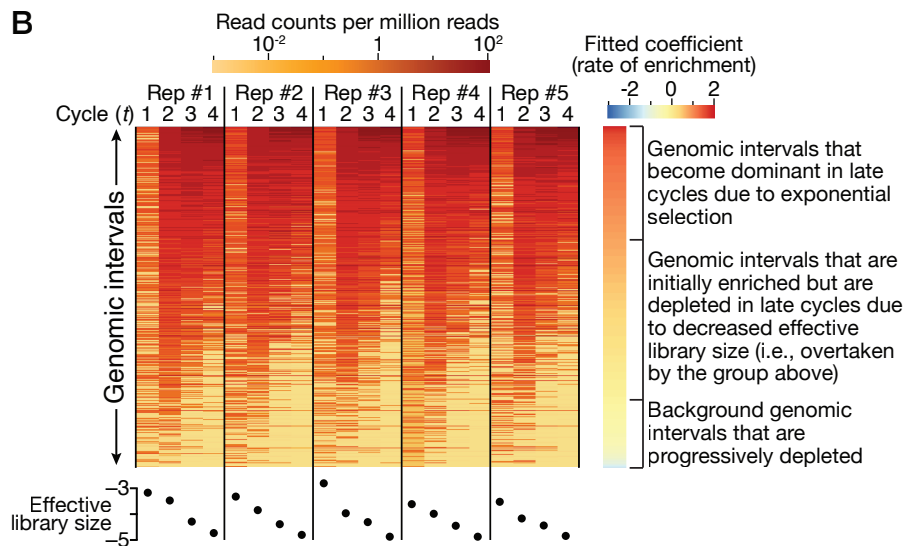


Figure 1. Overview of GHT-SELEX. A. Schematic of GHT-SELEX, showing parallels with HT-SELEX. **B.** Example of read accumulation over a TF motif match for NFKB1. **C.** Genomic binding for four positive control TFs on a genomic region showing (top to bottom) PWM scanning scores (moving average of affinity scores, from MOODS⁶⁰ scan in linear domain, using a window of size 200bp) for literature (Ref) PWMs and Codebook PWMs, followed by read coverage signal observed in GHT-SELEX and CHIP-seq.

A



B



C

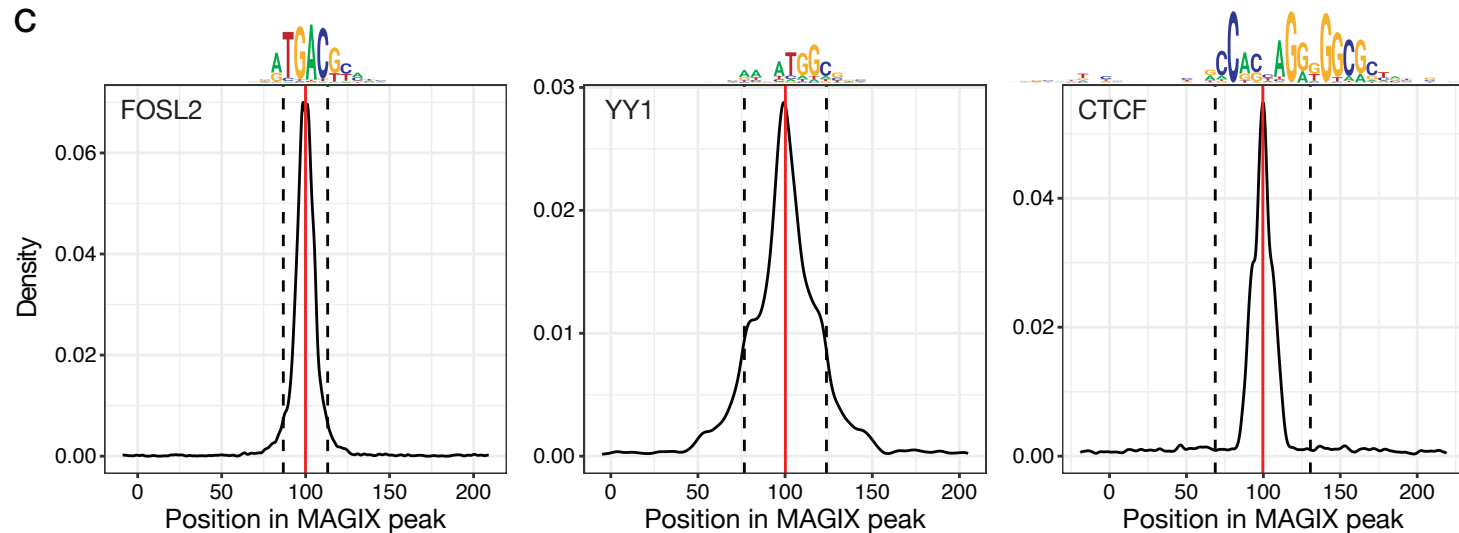


Figure 2. MAGIX method for interpretation of GHT-SELEX data. **A.** A brief overview of the statistical framework of the generative model of MAGIX. Open circles, closed circles, and the diamonds represent latent variables, observed variables, and deterministic computations, respectively. s_i : library size for sample i ; \mathbf{x}_i : vector of sample-level variables for sample i , including an intercept term and a term for the SELEX cycle, in addition to other terms for batch and background effects; β_j : vector of model coefficients for interval j ; m_{ij} : number of observed reads mapping to interval j in sample i . See **Methods** for description of other variables. **B.** Example of actual read count data for CTCF over five replicates of four cycles, illustrating enrichment patterns, fitted coefficients (right), and estimated library sizes (bottom). **C.** Distribution of PWM hits for the top-ranked TF PWM (highest AUROC on GHT-peaks as determined by accompanying study⁴³) within the 5,000 highest scoring MAGIX peaks. PWM hits were identified with MOODS⁶⁰ ($P < 0.0001$). Solid red lines represent the mean PWM hit position within MAGIX peaks and dashed lines represent one standard deviation about the mean.

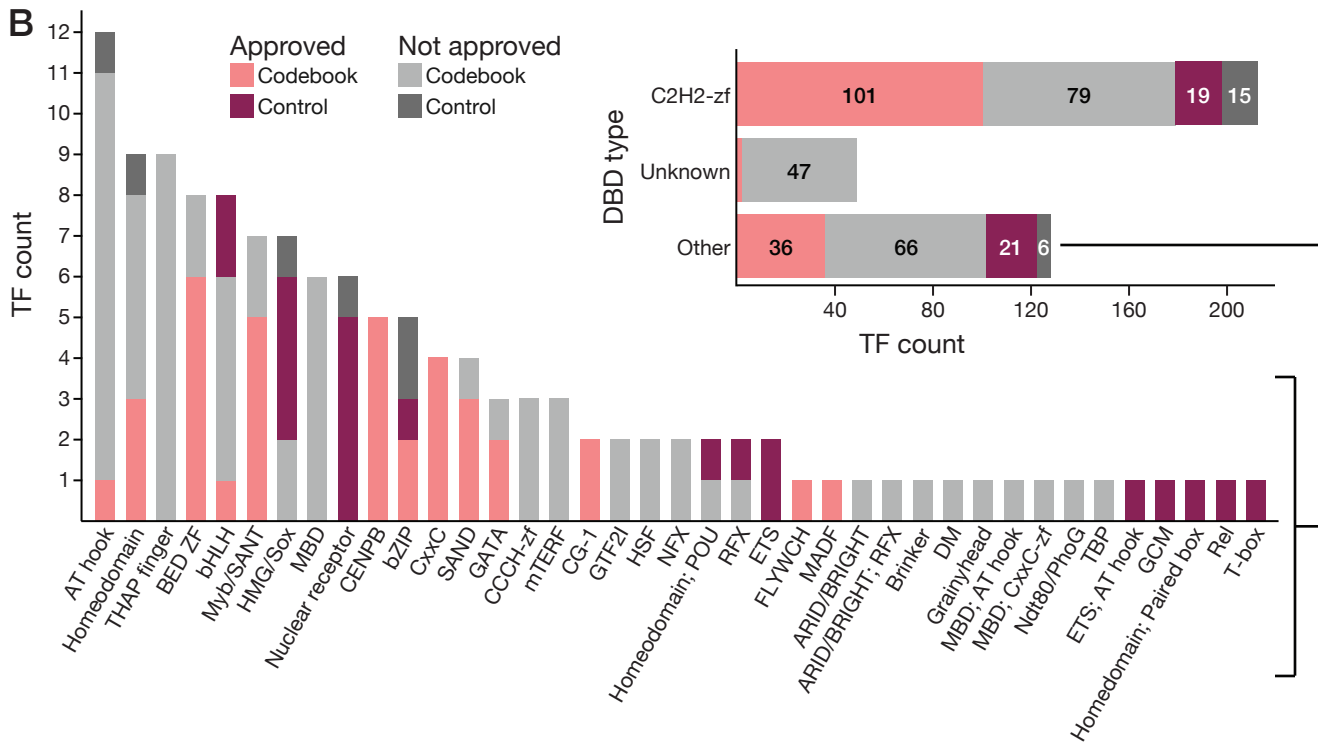
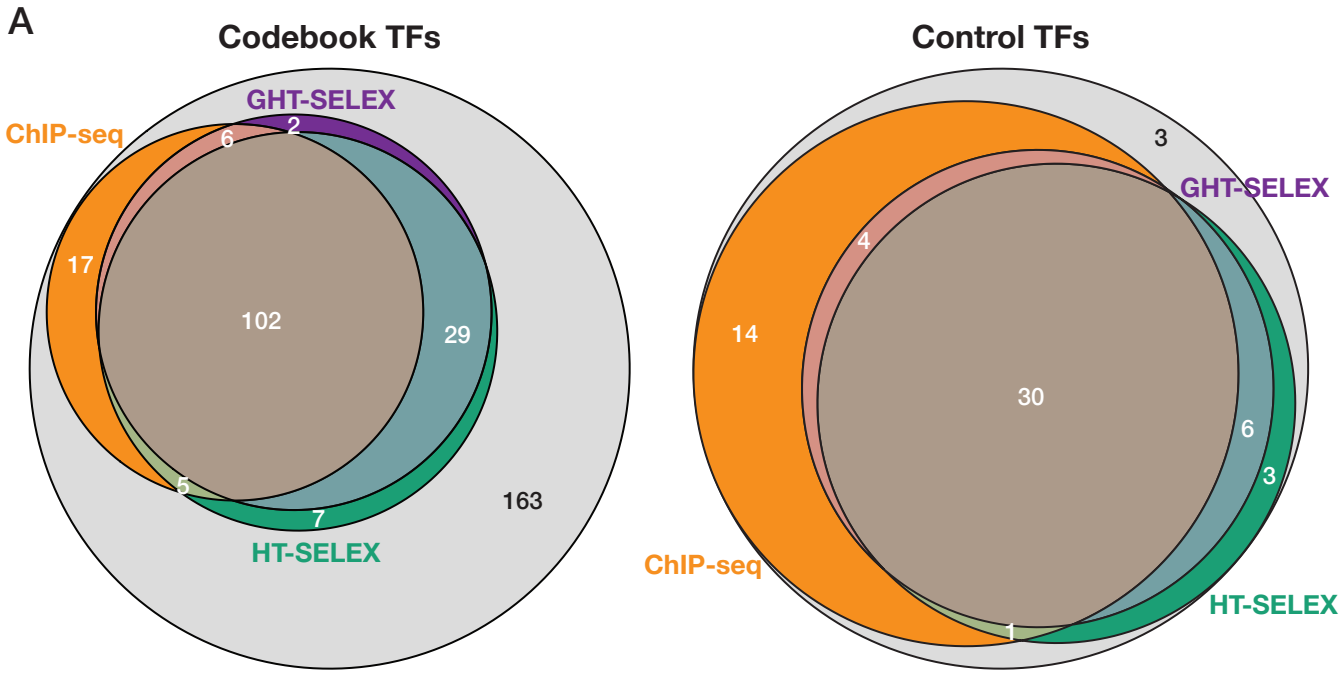


Figure 3. Analysis of 331 Codebook proteins and 61 control TFs using GHT-SELEX. **A.** Venn diagram displays the number of TFs with approved experiments in GHT-SELEX, HT-SELEX, and ChIP-seq for all Codebook TFs (left) and control TFs (right) assayed with GHT- and HT-SELEX. **B.** Bar chart shows the number of TFs with at least one approved GHT-SELEX experiment, categorized based by DBD type. C2H2-zf proteins and those with an unknown DBD (at the beginning of the project) are inset due to large numbers.

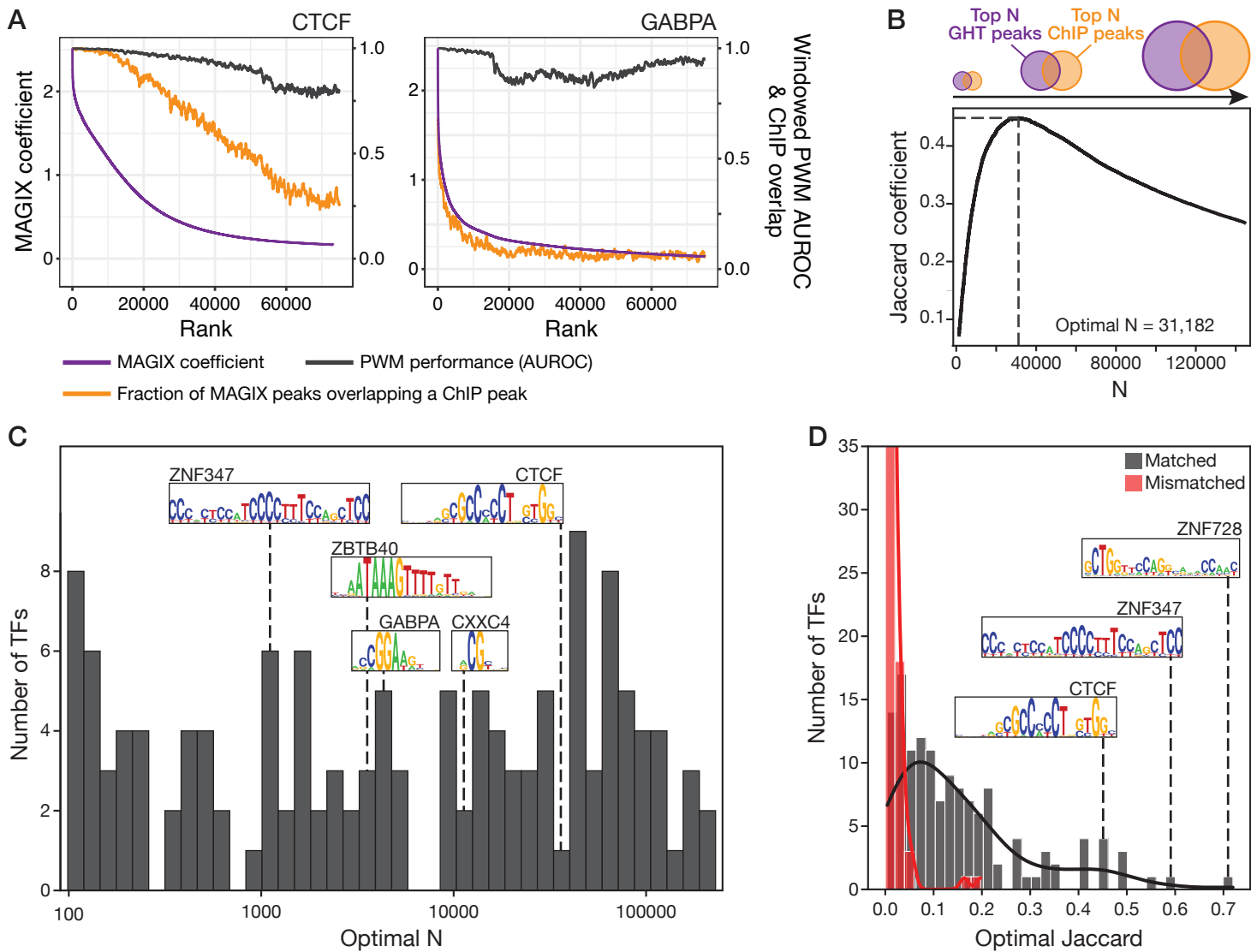


Figure 4. Correspondence between GHT-SELEX and ChIP-seq peaks. A.

Enrichment of ChIP-seq peaks and PWM hits within MAGIX peaks, for two example control TFs. The top 75,000 MAGIX peaks are sorted by their MAGIX enrichment coefficient (purple, left y-axis). Orange line shows the proportion of peaks (in a sliding window of 500 peaks over the ranked peaks, with a step size of 50) that overlap with a ChIP-seq peak (at MACS threshold $P < 0.001$). Black line shows the AUROC for PWM affinity scores (calculated by AffiMx⁵²) of MAGIX peaks in the same window vs. 500 random genomic sites. **B.** Illustration of peak number optimization (for CTCF as an example). **C.** Histogram of the optimal values of N (peak count) for the 137 TFs that have both GHT-SELEX and ChIP-seq peaks. **D.** Histogram of optimal Jaccard values, compared to the maximum Jaccard for mismatched TFs (i.e. between GHT-SELEX for one TF and ChIP-seq for a randomly selected TF).

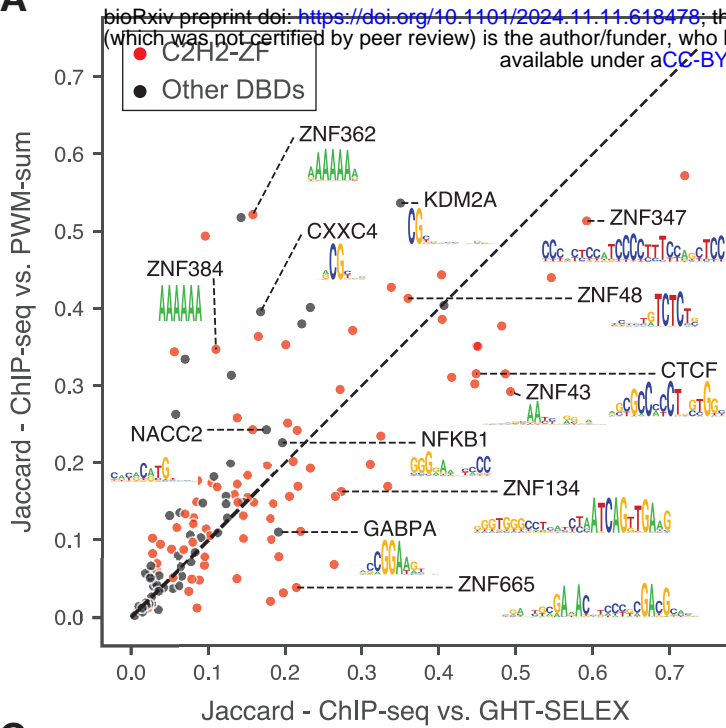
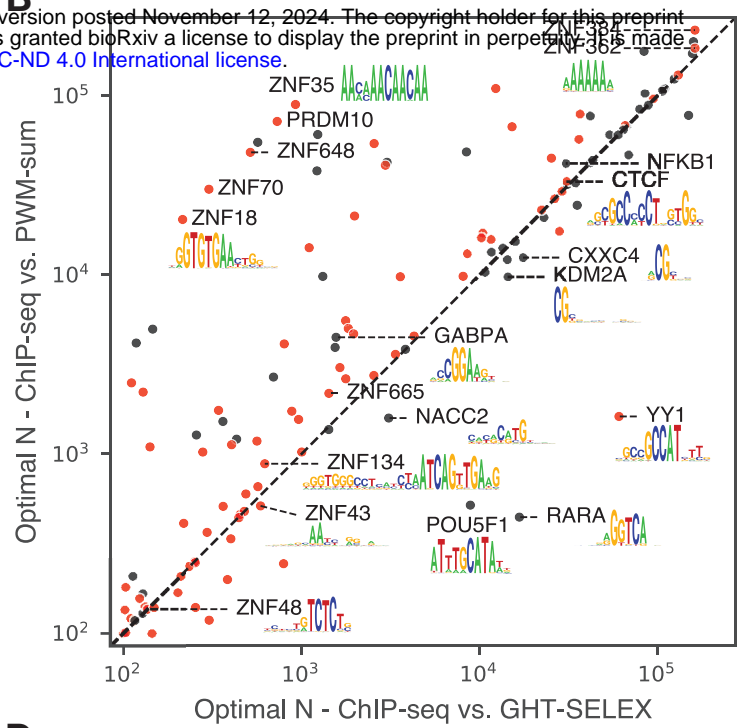
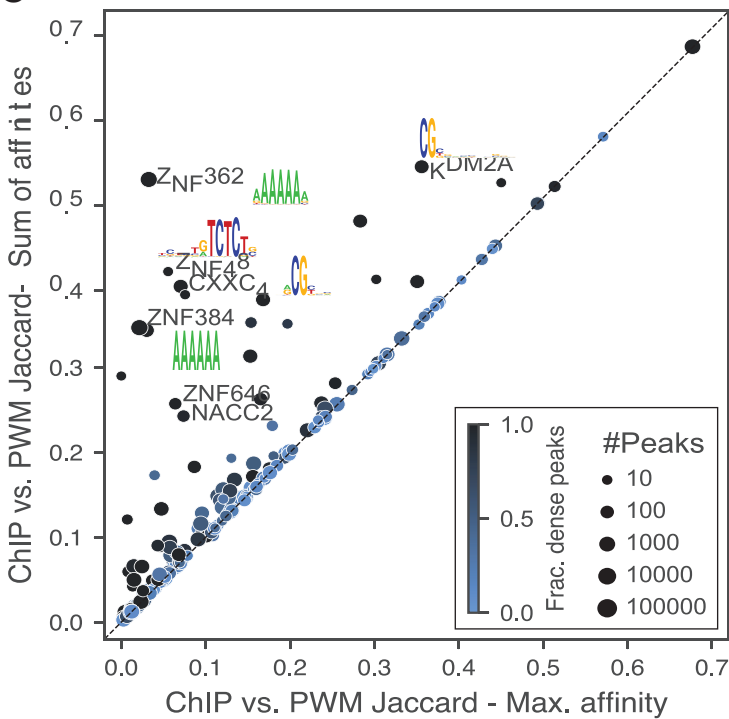
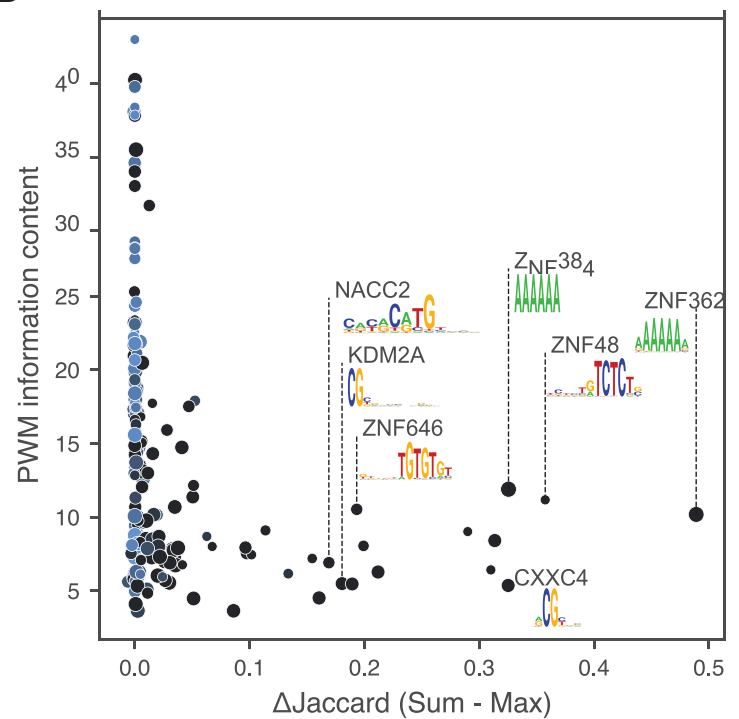
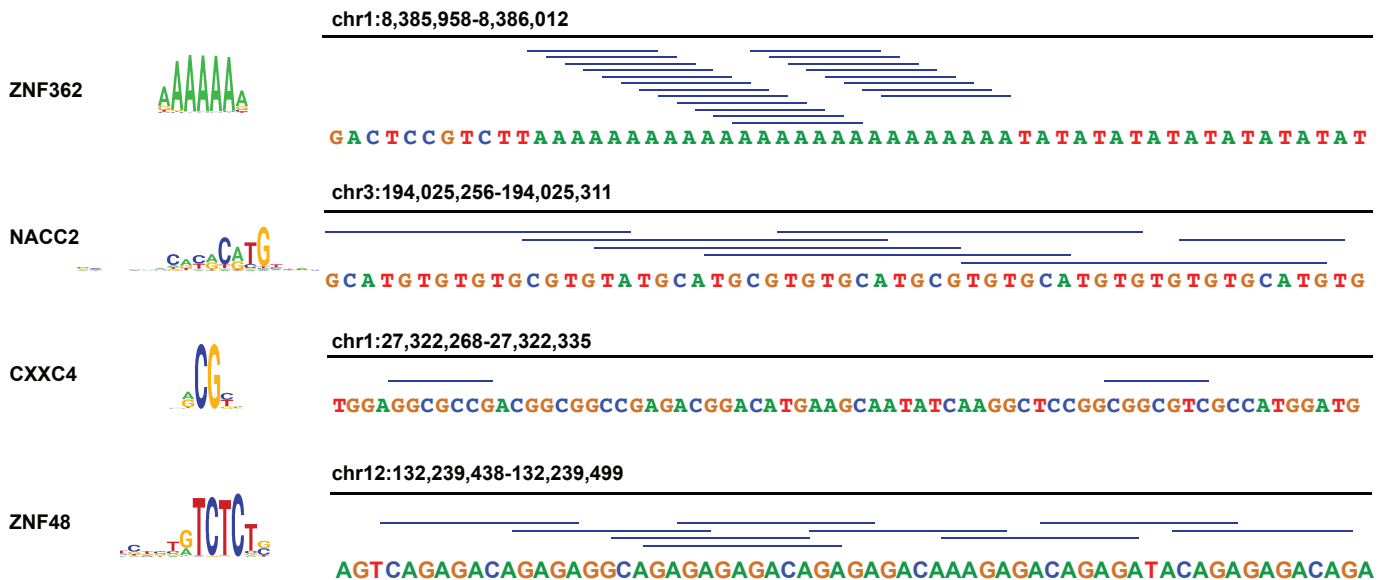
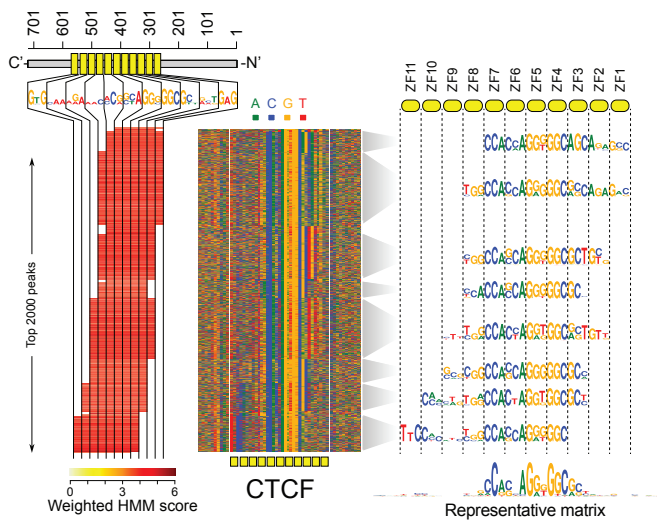
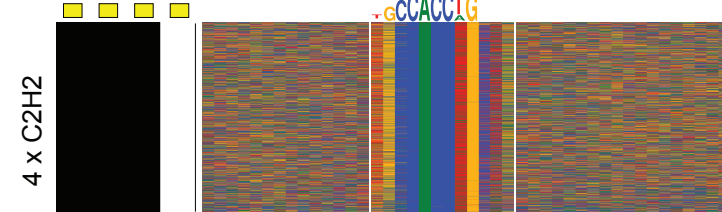
A**B****C****D****E**

Figure 5. High quality PWMs often predict *in vivo* binding sites as effectively as GHT-SELEX peaks. **A.** Scatter plot of optimal Jaccard value between GHT-SELEX peaks and ChIP-seq peaks (x-axis) vs. optimal Jaccard value between PWM-predicted sites and ChIP-seq peaks (y-axis), for all 137 TFs (dots). **B.** Scatter plot of optimal N (peak number) for the same peak set comparisons shown in (A). **C.** Scatter plot showing optimal Jaccard value between PWM-predicted sites and ChIP-seq peaks, for maximum-affinity PWM scoring and sum-of-affinities PWM scoring. Points (TFs) are scaled based on the optimal number of peaks (in the sum scoring), and the color reflects the fraction of binding sites comprised of multiple PWM hits. **D.** Scatter plot of the improvement in the optimal Jaccard value associated with sum-of-affinities PWM scoring vs. information content of the PWM. Points' size and color are the same as panel (C). **E.** Examples of four TFs with multiple motif matches within a single ChIP-seq peak.

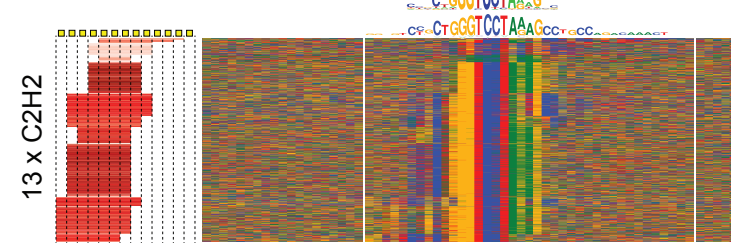
A RCADEEM reveals modular C2H2 DNA binding (CTCF)



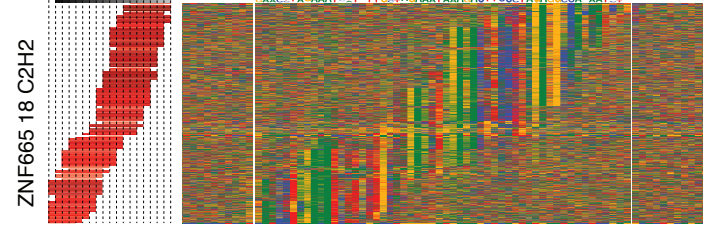
B Canonical (PRDM13)



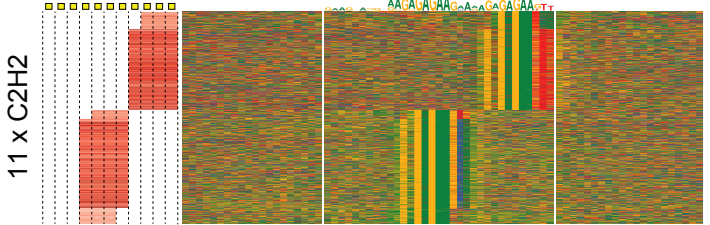
C Core with extensions (ZNF668)



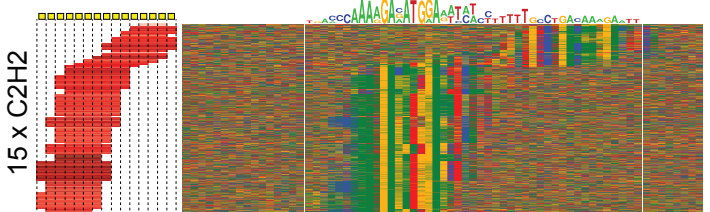
D Finger shift (ZNF665)



E Multiple DBDs (ZNF775)



F Multiple modes (ZNF471)



G

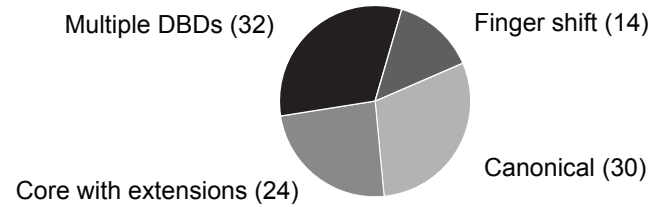
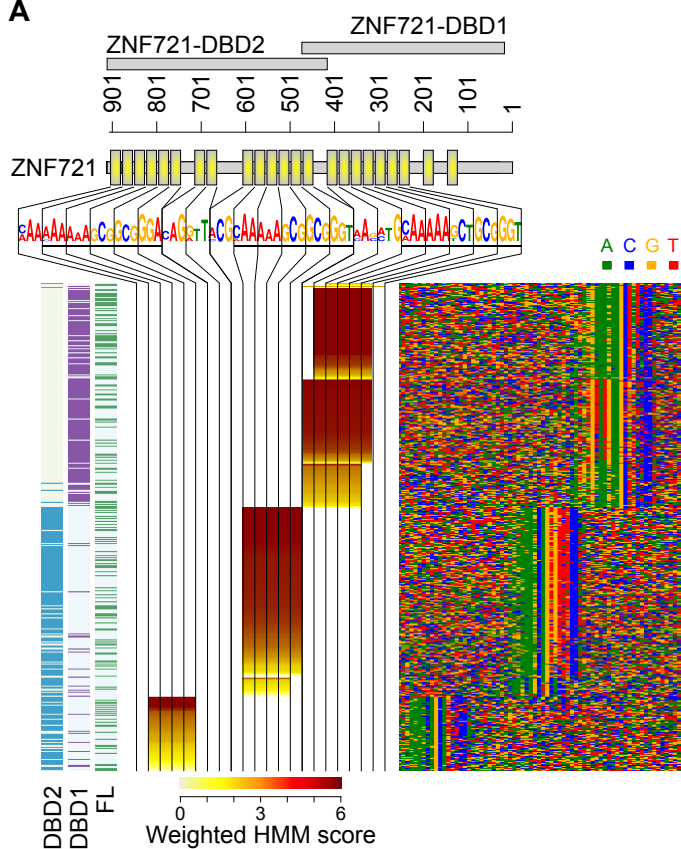
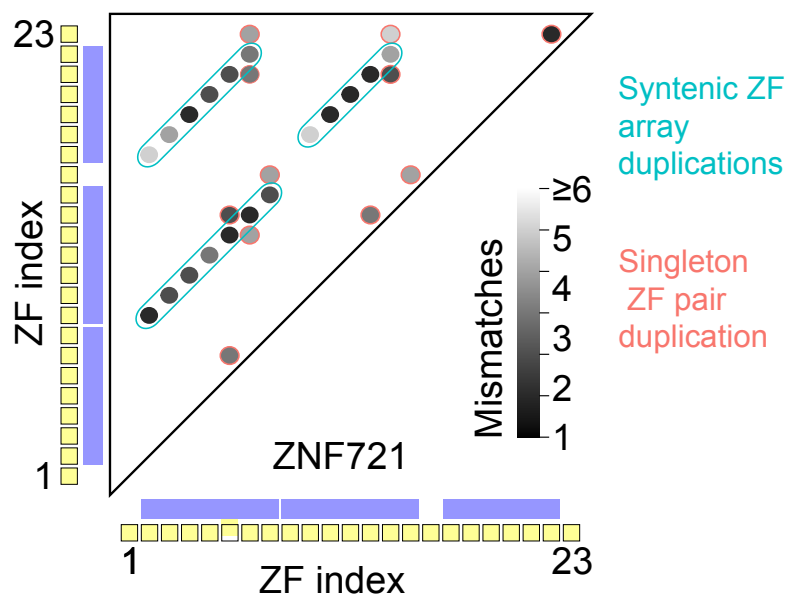


Figure 6. Alternative engagement of individual C2H2-zf domains at genomic binding sites inferred from the recognition code. A. RCADEEM applied to CTCF. *Middle* panel displays the top 2,000 nonrepetitive GHT-SELEX peaks. White vertical bars indicate the region that is expected to contact the DNA based on the assumption that each of the C2H2-zf domains define three contiguous bases. *Left* panel indicates which C2H2-zf domains are inferred to engage each DNA sequence, which is used to determine the row order in the figure. *Right* panel shows motifs for the major sub-sites, derived from base frequencies in the sequence alignment. **B-F,** Top 2,000 non-repeat peak sequences, as in (**A**), for representative TFs with different binding modes, as described in the main text. Above each is shown the sequence logo for the single representative Codebook PWM (*top*) and a motif generated by RCADEEM that represents all the observed sequences (*bottom*). **G.** Number of occurrences of each category among all 86 C2H2-zf proteins for which RCADEEM yielded a significant outcome; note that a TF might appear in multiple or no categories.

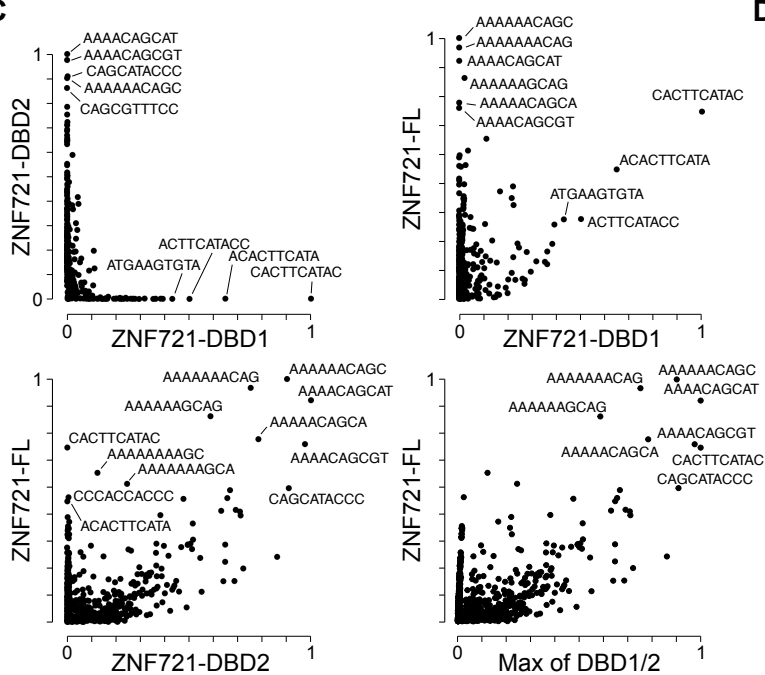
A



B



C



D

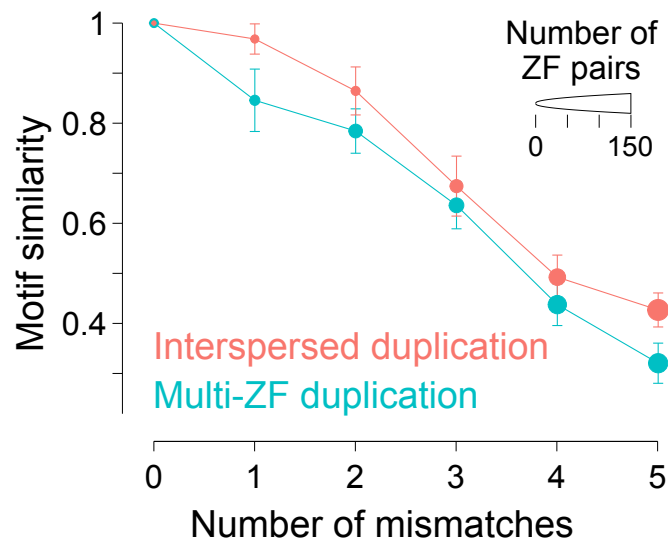


Figure 7. Evolution of C2H2-zf protein DNA-binding specificities through internal duplication of DBDs and DBD arrays.

A. RCADEEM results for the pool of top 500 peaks from full-length ZNF721 and two DBD constructs, after removing peaks that overlap repeats. The construct in which the peak was observed is indicated on the left. **B.** Similarity of C2H2-zf domains of ZNF721 (based on the number of mismatches in the global alignment). Apparent duplicated arrays are encircled by blue border (i.e. syntenic duplications), while single pairs that may also be duplicates are circled in red. **C.** Scatterplots of the HT-SELEX *k*-mer scores⁶¹ (relative counts) across the three ZNF721 constructs. **D.** Comparison of average per-base similarity (correlation of nucleotide frequency) in PWMs predicted by the recognition code, for those present in duplicated arrays vs. those duplicated as individual C2H2-zf domains, with duplicated C2H2-zf domains taken as pairs that are separated from each other by 5 or less edits. DBD pairs have been filtered to contain only combinations where both DBDs are likely to have retained their ability to bind DNA (have DNA binding functionality score⁵⁹ >0.5). Error bars show standard error of the mean.

559
560

561 METHODS

562

563 **TFs and constructs.** Selection of TFs, design of constructs for gene synthesis, and
564 expression vectors are described in accompanying study³⁷. Sequences and other
565 information are available as described below in Data Availability.

566

567 **Protein production.** We used three protein expression systems, which we refer to in
568 **Table S1** and below as *Lysate*, *IVT*, and *eGFP-IVT*, respectively. The *Lysate* system
569 employed recombinant HEK293 cells, created in the accompanying study⁴¹, and in a
570 previous study⁴², which express eGFP-tagged full-length proteins from a Tet-inducible
571 promoter (plasmid backbones pTH13195⁴² and pTH12027). We induced expression by
572 Doxycycline treatment for 24 hours prior to harvest, confirmed via fluorescent
573 microscopy. Whole cell lysates were then harvested from a 10cm plate (~10 million
574 cells) for each line using 1 ml of lysis buffer (50 mM Tris-Cl at pH 7.4 containing 150
575 mM NaCl and 1% Triton X-100), supplemented with protease-inhibitor cocktail (Roche
576 cOmplete mini, 04693159001), as described previously³². Each of the SELEX cycles
577 used 50 ul of lysate. *IVT* used an *in vitro* transcription-translation reaction (PURExpress
578 In Vitro Protein Synthesis Kit, NEB, Cat# E6800L) to express T7-driven, GST-tagged
579 proteins (either full-length or DBDs) (plasmid backbone pTH6838⁴⁵). *eGFP-IVT* employs
580 the TNT SP6 High-Yield Wheat Germ Protein Expression System (Promega, Cat#
581 L3260) to express SP6-driven, eGFP-tagged proteins (either full-length or DBDs)
582 (plasmid backbone pTH16505, an SP6-promoter driven, N-terminal eGFP-tagged
583 bacterial expression vector, modified from pF3A-eGFP³⁹ to contain *Ascl* and *Sbfl*
584 restriction sites after the eGFP. For *IVT* and *eGFP-IVT* production systems, we
585 performed reactions according to kit instructions, but using a smaller volume: 7.5ul of
586 *IVT* or 5ul of *eGFP-IVT* reaction sample was used in each binding reaction of each
587 SELEX cycle.

588

589 **GHT-SELEX and HT-SELEX library preparation.** We fragmented HEK293 genomic
590 DNA (Genscript, USA; Cat. No. M00094) for 45 minutes using NEBNext dsDNA
591 Fragmentase enzyme mix (NEB, M0348S), and then performed a size selection step to
592 reduce the amounts of fragments larger than 200 bp. In the size selection we added
593 0.9X volume of bead suspension (magnetic SPRI beads, supplied with the kit, NEB,
594 E7103S) to the fragmented DNA, mixed the reaction for a minute, and then removed the
595 large DNA fragment bound beads with a magnet, after which we diluted the supernatant
596 5X with water, followed by purification with a PCR purification kit (NEB, T1030S), to
597 recover fragments as small as 25 bp. Next the fragments were converted to an Illumina
598 sequencing compatible library using NEBNext® Ultra™ II DNA Library Prep kit (NEB:
599 E7103S) and NEB E7350 adapters. After adapter ligation, we purified the library with
600 PCR purification kit (NEB, T1030S) and then amplified it for five PCR cycles to convert
601 the partially single stranded adapter flanks to fully double stranded DNA, to increase the
602 amount of the product and reduce the amount of methylated cytosine residues in the
603 initial library. The ninety-six (96) HT-SELEX ligands were prepared as described⁵⁴, with
604 the exception that the reverse primer was replaced with a primer (5'

605 CTGGAGTTCAGACGTGTGCTCTTCCGATCT 3'), that does not contain a T7 promoter
606 sequence, and that HT-SELEX ligands differ from each other by containing a well-
607 specific variable region that flanks the randomized 40 bases indicated in the name of
608 the experiments (e.g. AA40NCCAGTG contains 40-bases flanked by AA and CCAGTG
609 sequences and Illumina adapter sequences). All primers and library preparation
610 schemes are given in **Table S1**.

611

612 **HT-SELEX and GHT-SELEX.** We modified protocols from a previously-described HT-
613 SELEX procedure³². HT-SELEX and the GHT-SELEX ligands contain the same flanking
614 constant regions and thus there were no differences in the selections or sequencing
615 library preparations. We conducted the magnetic bead washing operations below using
616 a Biotek 405TS plate washer fitted with a magnetic carrier. We performed 21 different
617 batches of SELEX, which varied in some technical respects in order to accommodate
618 the three protein production systems and to implement improvements developed during
619 the study (See **Table S2** for description of conditions used in each experimental batch).
620 Protein immobilization was carried out in buffers based either on Lysis buffer (150 mM
621 NaCl and 1% Triton X100 in Tris-Cl, pH 8) or Low stringency binding buffer (LSBB)(140
622 mM KCl, 5 mM NaCl, 1 mM K₂HPO₄, 2 mM MgSO₄, 100 μ M EGTA, 1 mM ZnSO₄ and
623 0.1% Tween20 in 20 mM HEPES-HCl (pH 7). All DNA-protein reactions used LSBB. For
624 GST-tagged proteins, we used glutathione magnetic beads (Sigma-Aldrich G0924-
625 1ML), and for GFP-protein immobilization, we used GFP-Trap Magnetic Agarose"
626 (Chromotek, gtma-100) for initial batches, and Anti-GFP antibody (ab290, Abcam)
627 immobilized to Protein G Mag Sepharose® Xtra (Cytiva, 28-9670-70) for later batches,
628 as the latter showed higher success rate. All selections used 1 μ l of the magnetic bead
629 slurry, a volume that in majority of the cases, according to manufacturers' information,
630 contains excess protein binding capacity but is still visible in microwell plates allowing
631 quality control of the washing steps.

632 **SELEX process:** All of the protocols (described in **Table S2**) followed these general
633 steps: 1) Affinity beads and 96-well plates were blocked with BSA for 15 minutes; 2)
634 Beads and plates were washed to remove unbound BSA; 3) Protein was immobilized
635 into beads for 1h on a shaker; 4) Beads were washed to remove nonspecific proteins
636 and carryover DNA; 5) Protein coated beads were incubated with DNA ligand for 1h to
637 allow the proteins to bind their target sites; 6) Unbound and weakly bound DNA ligands
638 were removed with extensive washing; 7) DNA ligands were eluted by suspending the
639 beads into heat elution buffer (0.4 μ M forward and reverse primers, 1 mM EDTA and
640 1% Tween 20 in 10 mM Tris-Cl, pH 8) transferring the suspension into a conical PCR
641 plate and heat treating it in a PCR machine using a program that cycled between
642 temperatures of 98 and 60°C, in order to denature the proteins and DNA, use
643 convection to drive the DNA into the solution, and to hybridize DNA to the amplification
644 primers; 8) Bead suspension obtained from heat elution was used as template in PCR
645 and qPCR reactions; 9) An additional DNA amplification cycle was performed with 2X
646 more primers and dNTPs to ensure that majority of the ligands are in fully double-
647 stranded state and 10) For batches YWO through YWS, we performed an additional
648 step in which the double-stranded ligands were treated with mung bean nuclease to
649 digest single stranded DNA such as primers or unpaired bases within selection ligands.
650 In each mung bean nuclease reaction, the pH of the solution (PCR reaction) was first

651 lowered by addition of 1:10 volume of 100 mM acetic acid, followed by addition of 1ul
652 (0.75 units) of the enzyme and incubation for one hour at 37°C.

653
654 **Sequencing.** Samples were prepared for sequencing by performing a PCR reaction
655 that indexes each sample and its selection cycle with a unique combination of i7 and i5
656 barcodes, followed by a double stranding reaction with primers that target regions of
657 DNA outside indices (**Table S1**). Following this step, DNA libraries were pooled, purified
658 with a PCR purification kit (NEB, T1030S), and then subjected to Illumina sequencing
659 with 60bp reads at 3M reads per sample (Donnelly Centre sequencing core facility).

660
661 **HT- and GHT-SELEX read processing and mapping.** HT-SELEX reads were filtered
662 by Phred quality score ($Q \geq 30$ in at least 90% of bases). GHT-SELEX reads were
663 parsed with Trimmomatic⁵⁵ to remove the constant regions from genomic fragments that
664 were shorter than the sequencing read length (options:
665 ILLUMINACLIP:CustomAdapters.fa:2:5:5, LEADING:3, TRAILING:3 MINLEN:25). The
666 custom adapters in the fasta file were AGATCGGAAGAGCACACGTCTGAACTCCAG
667 and AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTTA. For GHT-SELEX, we
668 mapped trimmed reads to the human genome build hg38 with bowtie2 (options: --very-
669 sensitive, --no-unal). The mapped reads were further filtered using Samtools (options: -
670 F 1548, version 1.20)⁵⁶.

671
672 **MAGIX statistical framework.** At the core of MAGIX is a generative model that
673 explicitly connects the enrichment of TF-bound genomic intervals to the fragment counts
674 observed across GHT-SELEX cycles. MAGIX, models how TF-bound intervals
675 progressively occupy a higher proportion of selected fragments pool in each cycle
676 relative to genomic background. These fragment proportions, in turn, are treated as
677 latent variables in the model that, together with a sample-specific library size factor,
678 determine the number of observed reads through a Poisson process. Consider the
679 genomic interval $j \in [1, G]$, where G is the total number of unique genomic intervals that
680 we are modeling. Assume that fragments originating from interval j have a starting
681 abundance of a_j in the library. We also assume an exponential enrichment for the
682 fragments, that in each cycle of SELEX, the abundance of these fragments changes by
683 a factor of e^{b_j} , where b_j is the log fold-change in abundance per cycle (referred to as
684 *enrichment coefficient*), conceptually associated with biophysical parameters such as
685 binding energies. Therefore, at cycle t , the abundance of the fragment originating from
686 interval j is given by:

$$687 \quad f_j^t = a_j (e^{b_j})^t = a_j e^{tb_j}$$

688 For convenience, we work with the logarithm of abundance, $y_j = \log f_j$, transforming the
689 exponential equation above to a linear equation as follows:

$$690 \quad E(y_{tj}) = \log(f_j^t) = \log a_j + tb_j,$$

691 which can be seen as the linear multiplication of a feature vector $\mathbf{x}_i = [1 \ t_i]$ for all the
692 samples i (and across different cycles) corresponding to the TF of interest, and the
693 interval-specific parameters $\boldsymbol{\beta}_j = [\log a_j \ b_j]$.

694 We note that to accurately model the enrichment of each fragment per cycle, other
695 factors also need to be taken into consideration, such as background or batch effects,
696 and therefore, the linear equation above needs to be fitted not only to the samples that

697 correspond to the TF of interest, but also samples from other experiments. We embed
 698 these dependencies in a design matrix $\mathbf{X} \in \mathbb{R}^{N \times K}$, where N is the total number of samples
 699 and K is the number of variables to consider, including an intercept term (whose
 700 coefficient will correspond to $\log a_j$ above), a term for the SELEX cycle t (variable for
 701 the samples corresponding to a same TF), and other terms for batch and background
 702 effects. In addition to the variables included in \mathbf{X} , the abundance of each fragment in
 703 each sample depends on a sample-specific scaling factor that is often referred to as the
 704 library size. Assume that this library effect, for each sample $i \in [1, M]$, is the scaling factor
 705 s_i (in logarithmic scale). Therefore:

$$\boxed{E(y_{ij}) = \hat{y}_{ij} = \mathbf{x}_i \cdot \boldsymbol{\beta}_j + s_i}$$

706 Here, \hat{y}_{ij} corresponds to the expected logarithm of the abundance of interval j in sample
 707 i , $\mathbf{x}_i \in \mathbb{R}^K$ is a vector representing the i 'th row of the design matrix \mathbf{X} (i.e., the sample-level
 708 variables for sample i), and $\boldsymbol{\beta}_j \in \mathbb{R}^K$ is an interval-specific vector of coefficients for the K
 709 variables included in the model.
 710

711 We note that the equation above does not have a unique solution. For example, any $\Delta\boldsymbol{\beta}$
 712 can be added to $\boldsymbol{\beta}_j$, followed by subtraction of $\mathbf{X}\Delta\boldsymbol{\beta}$ from \mathbf{s} , without any change in $\hat{\mathbf{y}}$:

$$\boxed{\mathbf{x}_i \cdot \boldsymbol{\beta}_j + \mathbf{s} = \mathbf{x}_i \cdot (\Delta\boldsymbol{\beta} + \boldsymbol{\beta}_j) + (s_i - \mathbf{x}_i \cdot \Delta\boldsymbol{\beta})}$$

714 Therefore, to make the model identifiable, we limit $\boldsymbol{\beta}_j$ so that $\sum_{j \in [1, G]} \boldsymbol{\beta}_j = \mathbf{0}$, where $\mathbf{0}$ is the
 715 zero vector of length K . This constraint is also useful since it means that, across all G
 716 intervals, the mean of each coefficient in $\boldsymbol{\beta}_j$, including the coefficient for the SELEX cycle,
 717 is zero; in other words, the enrichment per cycle for each interval is calculated relative
 718 to the mean of all G intervals.
 719

720 To incorporate the experimental noise in the logarithm of the abundance of interval j in
 721 sample i (i.e. building the real distribution of y_{ij}), we modeled it as a Gaussian random
 722 variable whose mean is given by \hat{y}_{ij} (the linear model above) with a sample-specific
 723 variance σ_i^2 :

$$\boxed{y_{ij} \sim \mathcal{N}(\hat{y}_{ij}, \sigma_i^2)}$$

724 To complete the Bayesian framework, we also assume a multivariate Gaussian prior for
 725 $\boldsymbol{\beta}_j$:

$$\boxed{\boldsymbol{\beta}_j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})}$$

727 Here, $\boldsymbol{\Sigma}$, the covariance matrix of the prior distribution, is shared across all intervals.
 728

729 Altogether, the equations above form the following Bayesian model:
 730

$$\begin{aligned} &\boxed{\boldsymbol{\beta}_j \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})} \\ &\boxed{s_i \sim U(-\infty, +\infty)} \\ &\boxed{\sigma_i^2 \sim U(0, +\infty)} \\ &\boxed{\hat{y}_{ij} = \mathbf{x}_i \cdot \boldsymbol{\beta}_j + s_i} \\ &\boxed{y_{ij} \sim \mathcal{N}(\hat{y}_{ij}, \sigma_i^2)} \end{aligned}$$

736 Assuming that the values for y_{ij} are directly observed, we can obtain the maximum a
 737 posteriori (MAP) estimates of the parameters $\boldsymbol{\beta}_j$ (for $j \in [1, G]$), \mathbf{s} , and σ_i^2 (for $i \in [1, M]$)
 738 through a block coordinate descent algorithm, as previously described⁴⁸.
 739
 740

741 The prior covariance matrix Σ is a hyper-parameter that is obtained using an empirical
742 Bayes approach. More specifically, we first obtain the maximum likelihood estimate
743 (MLE) of the parameters β_j , s_i , and σ_i^2 without assuming any prior on β_j , and then
744 estimate the values of $\sigma_{\beta k}^2$, the variance of each element k in β_j , using the MLE
745 solutions of all β_j coefficients. The covariance matrix Σ is then constructed as
746 $\Sigma = \text{diag}(\sigma_{\beta 1}^2, \dots, \sigma_{\beta k}^2)$.

747

748 We note, however, that the log-abundance values y_{ij} 's are not directly observable in
749 GHT-SELEX data. Instead, we observe m_{ij} , the count of reads mapping to interval j in
750 sample i (see **HT- and GHT-SELEX read processing and mapping**). This parameter
751 adds another step to the framework, leading to the following hierarchical Bayesian
752 model:

753

$$\beta_j \sim \mathcal{N}(\mathbf{0}, \Sigma)$$

754

$$s_i \sim \text{U}(-\infty, +\infty)$$

755

$$\sigma_i^2 \sim \text{U}(0, +\infty)$$

756

$$\hat{y}_{ij} = \mathbf{x}_i \cdot \beta_j + s_i$$

757

$$y_{ij}^* \sim \mathcal{N}(\hat{y}_{ij}, \sigma_i^2)$$

758

$$m_{ij} \sim \text{Pois}(e^{y_{ij}^*})$$

759

760 Here, y_{ij}^* is the logarithm of the true abundance of fragment j in sample i , which is latent.
761 We obtain the MAP estimates of the parameters β_j (for $j \in [1, G]$), s_i , and σ_i^2 (for $i \in [1, M]$)
762 using an expectation maximization (EM) algorithm, in which at each E-step we obtain
763 the expected value of each y_{ij}^* given the observed read count m_{ij} and the current model
764 parameters, followed by re-estimation of the model parameters in the M-step, similar to
765 a previous method established for EM optimization of Poisson-lognormal models⁴⁸.

766

767 **Identifying GHT-SELEX peaks with MAGIX.** The statistical framework described
768 above calculates, for each genomic interval, the rate of enrichment across GHT-SELEX
769 cycles. To identify GHT-SELEX peaks, we used this framework along with the procedure
770 described below to systematically examine all genomic intervals and identify regions
771 with the highest signal (peaks). First, we binned the entire human genome (build hg38)
772 into ~13M non-overlapping intervals of 200 bp and then, for each TF, calculated the
773 read count profiles of these intervals across the GHT-SELEX cycles and replicates.
774 Counts were obtained using bedtools multicov (version 2.30.0)⁵⁷. These read count
775 profiles were used as the input to MAGIX to obtain an enrichment coefficient for each
776 200 bp interval (while controlling for batch effects by including batch-specific pooled
777 controls and variables in the design matrix). Next, the top 200,000 regions with the
778 highest enrichment coefficients were selected as candidate intervals for peak
779 refinement.

780

781 To refine the peak coordinates, we first merged any adjacent candidate intervals, and
782 then calculated the base pair-resolution read count coverage profile across each
783 merged interval (and sum of all GHT-SELEX cycles). The position with the highest read
784 coverage was selected as the candidate peak summit. The read counts overlapping the
785 ± 100 bp around the summits were computed, which were used as input to the MAGIX

786 statistical inference component (above) to recalculate enrichment coefficients, while
787 reusing the library sizes and the empirical hyperparameters estimated from the analysis
788 of all 13M genomic intervals, without re-optimization. For each candidate peak, we also
789 calculated a P-value, representing the statistical significance of the enrichment
790 coefficient (null hypothesis is that the enrichment coefficient is zero). To do so, we
791 obtained maximum likelihood estimate of the model coefficients for each coefficient (i.e.,
792 ignoring the prior distributions), and performed a likelihood ratio test (LRT) against a
793 reduced model in which the enrichment coefficient was restricted to zero.

794
795 To calculate empirical FDRs for the peaks, we first obtained negative peaks by
796 repeating the procedure described above but inverting the cycle labels. In other words,
797 we obtained depleted peaks, relative to the pool, instead of enriched peaks. Then, we
798 calculate the FDR as the fraction of depleted peaks relative to enriched peaks for each
799 coefficient value. The source code for MAGIX is available at
800 <https://github.com/csglab/MAGIX>.

801
802 **Selection of thresholds for peak sets.** We sorted the GHT-SELEX peaks by their
803 MAGIX score (enrichment coefficient, or as named in the peaks BED files, *coefficient.br*,
804 which estimates cycle enrichment). Similarly, we sorted the merged ChIP-seq peaks by
805 P-value. Then, for different values of N (between 100 and the total number of peaks),
806 we took the top N peaks for both peaks sets and calculated the Jaccard index ($= O/(2N -$
807 $O)$, in which O is the intersection of peaks). To eliminate the error in the cases when one
808 peak in a set overlaps with multiple peaks in another set, we used the average of the
809 overlaps for the intersection (i.e. $O=(O1+O2)/2$, in which O1 is the number of peaks in
810 set1 overlapping with any peaks in the set2 and vice versa). The value of N that yielded
811 the maximum Jaccard value was identified, and the threshold for each peak set taken
812 as that which yielded this maximum N. The same process was applied to compare
813 PWM-predicted binding sites and ChIP-seq peaks.

814
815 **Comparing PWM scoring methods.** To create *in silico* predicted binding sites for a TF,
816 we first scanned the genome using the generated PWM (see the Codebook overview
817 manuscript for the details on PWM selection), using MOODS⁵⁸ with a p-value threshold
818 of 0.0001. We then merged the clusters of PWM hits with a distance less than 200bp
819 between neighboring hits, since this is the median length of ChIP-seq fragments, and
820 the task is predicting *in vivo* binding sites; this length is also consistent with the MAGIX
821 bin size. Singleton PWM hits and boundary hits were also expanded to have a width of
822 at least 200bp. The clusters of PWM hits were re-scored using sum-of-affinity (i.e. with
823 PWM log-odds scores at each base converted to linear/probability space, prior to
824 calculation of the sum) and maximum-affinity methods, by either applying a sum or
825 maximum, respectively, over the PWM scores of the cluster members. The resulting
826 sites were sorted by their new score and processed through the same optimization
827 procedure described above for peaks, to maximize their overlap with ChIP-seq peaks.

828
829 **Modeling alternative C2H2-zf binding modes with RCADEEM.** RCADEEM uses a
830 hidden Markov model (HMM) to represent multiple, alternative DNA-binding motifs,
831 each corresponding to the binding preference of a C2H2-zf array. Briefly, the DNA

832 sequences (e.g., GHT-SELEX peaks) are modeled as sequences generated from a
833 discrete Markov process with hidden states that include a background state (S_0) and M
834 motif states S_m ($m \in [1, M]$). The background state, with marginal probability π_0 , emits
835 each nucleotide n with probability $b_0(n)$ ($\sum_n b_0(n) = 1$). The background state can transition
836 to itself (i.e., consecutive DNA nucleotides can be generated from the background state)
837 with probability $a_{0,0}$, or to each motif state S_m ($m \in [1, M]$) with probability $a_{0,m}$
838 ($a_{0,0} + \sum_m a_{0,m} = 1$). Each motif m , with marginal probability π_m , generates a sequence of
839 length l_m , with each nucleotide n at position i emitted with probability $b_{m,i}(n)$ ($\sum_n b_{m,i}(n) = 1$
840 $\forall m \in [1, M], i \in [1, l_m]$). Note that, for each $m \in [1, M]$, the values $b_{m,i}(n)$ form a position-
841 specific frequency matrix (PFM, i.e. the exponential of the classical log-odds PWM) with
842 width l_m , which is fixed to be 3 times the number of zinc finger domains in the array
843 represented by motif m , as each zinc finger domain binds to three nucleotides. Finally,
844 each motif state S_m transitions to the background state with probability $a_{m,0} = 1$.

845
846 We start the model by including the motifs representing all possible consecutive zinc
847 finger domain arrays⁵⁰. We initialize the emission probabilities $b_{m,i}(n)$ for each motif m
848 using the PFM predicted for the associated zinc finger array by a previously created
849 C2H2-zf recognition code⁵¹—this recognition code is a machine learning model that,
850 given the sequence of a zinc finger array, predicts the expected binding preference. The
851 HMM parameters, including all marginal state probabilities, state transition probabilities,
852 and emission probabilities are then optimized via expectation maximization using
853 Baum–Welch algorithm. Then, each of the optimized PFMs are tested for (i) enrichment
854 of the motif in actual sequences compared to dinucleotide-shuffled sequences, and (ii)
855 similarity to the original recognition code-predicted PFM. To achieve (i), for each position
856 x in each DNA sequence k , we calculate $\psi_{k,x}(S_m)$, the probability that it was generated
857 from motif state S_m , using the forward-backward algorithm. The motif score for DNA
858 sequence k is then calculated as $\sum_x \psi_{k,x}(S_m) / l_m$, representing the expected number of
859 times the state S_m is seen in sequence k . For each motif m , these scores are calculated
860 both for actual GHT-SELEX peak sequences and their dinucleotide-shuffled version.
861 Then, the top 100 sequences with the largest scores for each motif are tested to see
862 whether they are enriched in the motif compared to shuffled sequences (Fisher’s exact
863 test, $FDR \leq 0.01$). Motifs that do not pass this cutoff are removed from the model. To
864 achieve (ii), each HMM-optimized PFM is first converted to log-scale (representing a
865 PWM), followed by calculation of Pearson correlation of the PWM entries with those
866 predicted by the recognition code. Pearson correlations are then converted using Fisher
867 transformation in order to calculate a P-value, followed by removal of motifs that do not
868 pass the FDR cut-off ≤ 0.01 . The remaining motifs are then used to reconstruct a smaller
869 HMM, similar to the procedure described above, followed by another round of EM
870 optimization. This procedure is repeated until all motifs pass the cut-offs for enrichment
871 in GHT-SELEX sequences while maintaining significant similarity to the original
872 recognition code-predicted sequences.

873
874 To visualize the binding modes predicted by RCADEEM, the resulting PWMs are used
875 to identify their best match in each of the input sequences using AffiMx⁵². Then, for each
876 sequence, the PWM with the highest weighted HMM score on the best match is kept as
877 the predicted binding mode. To align the sequences, offsets are calculated based on the

878 corresponding C2H2-zf domains (**Figures 6A-F**). C2H2-zf proteins were categorized
879 based on their alternative usage of C2H2-zf domains (i.e., Multiple DBDs, Finger shift,
880 Canonical, and Core with extensions; **Figure 6**) through an expert-curated evaluation
881 (**Table S5**). To make a motif model for each binding mode, we manually selected
882 representative peaks corresponding to each binding mode over the 2000 GHT-SELEX
883 peaks with the highest enrichment coefficient. The sequence (already aligned by
884 RCADEEM) and C2H2-zf domain array coordinates of these peaks were used to create
885 PFMs. The resulting PFMs for those C2H2-zf TFs are available in **Document S2** and
886 online at <https://cisbp.ccb.utoronto.ca>⁴⁵. The logos, coordinates, selected sequences,
887 annotated sequence heatmaps, and associated metadata are available online at
888 <https://codebook.ccb.utoronto.ca>. The source code for RCADEEM is available at
889 <https://github.com/csqlab/RCADEEM>.

890

891 **Comparison of C2H2 DBDs.** C2H2 DBD similarities were compared by pairwise
892 alignment with Needleman-Wunsch algorithm, as implemented in R-package Biostrings
893 and counting substitutions, insertions and unmatched flanking bases as edits. DNA-
894 binding functionality scores and predicted motif similarity for the DBDs were analyzed
895 as described previously⁵⁹.

896

897 **DATA AVAILABILITY**

898

899 The sequencing raw data for the HT-SELEX and GHT-SELEX experiments have been
900 deposited into the SRA database under identifiers PRJEB61115 (HT-SELEX) and
901 PRJEB76622 (GHT-SELEX). Additionally, genomic interval information generated for
902 the GHT-SELEX, has been deposited into GEO under accession GSE278858. The
903 entire Codebook data structure, with many accessory files and browsable results at is
904 available at <https://codebook.ccb.utoronto.ca>. Larger collection of motifs generated for
905 these experiments in an accompanying study⁴³ can be browsed at mex.autosome.org.
906 Source codes for MAGIX and RCADEEM are available from Github
907 (<https://github.com/csqlab/MAGIX> and <https://github.com/csqlab/RCADEEM>).

908

909

910

911 **ACKNOWLEDGEMENTS**

912

913

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

918

919 This work was supported by the following:

- 920 • Canadian Institutes of Health Research (CIHR) grants FDN-148403, PJT-186136,
921 PJT-191768, and PJT-191802, and NIH grant R21HG012258 to T.R.H.
- 922 • CIHR grant PJT-191802 to T.R.H. and H.S.N.
- 923 • Natural Sciences and Engineering Research Council of Canada (NSERC) grant
924 RGPIN-2018-05962 to H.S.N.
- 925 • Russian Science Foundation grant 20-74-10075 to I.V.K.
- 926 • Russian Science Foundation grant 24-14-20031 to F.A.K.
- 927 • Swiss National Science Foundation grant (no. 310030_197082) to B.D.
- 928 • Marie Skłodowska-Curie (no. 895426) and EMBO long-term (1139-2019) fellowships
929 to J.F.K.
- 930 • NIH grants R01HG013328 and U24HG013078 to M.T.W., T.R.H., and Q.M.
- 931 • NIH grants R01AR073228, P30AR070549, and R01AI173314 to M.T.W.
- 932 • NIH grant P30CA008748 partially supported Q.M.
- 933 • Canada Research Chairs funded by CIHR to T.R.H. and H.S.N.
- 934 • Ontario Graduate Scholarships to K.U.L and I.Y.
- 935 • A.J. was supported by Vetenskapsrådet (Swedish Research Council) Postdoctoral
936 Fellowship (2016-00158)
- 937 • The Billes Chair of Medical Research at the University of Toronto to T.R.H.
- 938 • EPFL Center for Imaging
- 939 • Institutional funding from EPFL
- 940 • Resource allocations from the Digital Research Alliance of Canada

941

942 SUPPLEMENTARY TABLES AND DOCUMENTS

943

944 **Table S1. HT- and GHT-SELEX ligand sequences and descriptions.** Table lists the
945 oligonucleotide sequences used in the assay and describes how they anneal with each
946 other on the synthesis and amplification steps.

947

948 **Table S2. Experimental batch specific protocol details.** Table lists the reagents and
949 experimental conditions that varied between different experimental batches.

950

951 **Table S3 GHT-SELEX-experiment metadata.** Table lists all GHT- and HT-SELEX
952 experiments performed in this study indicating: unique experiment identifier; human
953 readable identifier; plasmid identifiers; HNGC symbol; experimental batch; construct
954 type; protein production approach, position in the 96-well; sequencing strategy; number
955 of selection cycles; and whether the experiment was approved or not. Note that GFP
956 control experiments (i.e. empty plasmids) are also included in the table (5 GHT-SELEX
957 and 7 HT-SELEX).

958

959 **Table S4: Genomic region overlap of GHT-SELEX and ChIP-seq peaks and PWM-**
960 **predicted target regions.** Table shows the overlap of optimal ChIP-seq peaks with
961 GHT-SELEX/MAGIX and PWM based predictions for each of the TFs where both
962 datasets were available. Columns show the highest Jaccard coefficient between each
963 pair of datasets and the number of peaks that yielded it.

964

965 **Table S5: C2H2-zf protein DNA-binding mode annotation.** Table lists the 86 C2H2
966 TFs for which RCADEEM result was obtained (out of 120 total C2H2-zf TFs with GHT-
967 SELEX data available) with information of: Total number of C2H2 zinc finger domains;
968 amino acid gaps between these DBDs; number of distinct motifs bound by the TF;
969 modular binding activity annotated for it; whether the protein is likely to contain zinc
970 fingers obtained from internal duplications and whether data was obtained from
971 experiments that expressed different subsets of the TFs C2H2-zf domains.

972

973 **Table S6: Intra-protein C2H2-zf domain duplication dataset.** Table displays all pairs
974 of human C2H2 DBDs that are separated from each other by five or less edits.

975

976 **Document S1: Motif centrality and enrichment in GHT-SELEX/MAGIX peaks and**
977 **its correspondence with ChIP-seq peaks.** Same plots as in **Figure 2C** and **Figure 4A**
978 for all the TFs and DBD constructs in this study with approved GHT-SELEX
979 experiments. *Top*, top-ranked TF PWM (highest AUROC on GHT-peaks as determined
980 by ⁴³). *Middle*, Distribution of PWM hits within the 5,000 highest scoring MAGIX peaks.
981 Solid red lines represent the mean PWM hit position within MAGIX peaks and dashed
982 lines represent one standard deviation about the mean. *Bottom*, Enrichment of ChIP-
983 seq peaks and PWM hits within MAGIX peaks. Orange line shows the proportion of
984 peaks (in a sliding window of 500 peaks over the ranked peaks, with a step size of 50)
985 that overlap with a ChIP-seq peak (at MACS threshold $P < 0.001$). Black line shows the
986 AUROC for PWM affinity scores of MAGIX peaks in the same window vs. 500 random
987 genomic sites.

988 **Document S2: PFMs of C2H2-zf proteins with alternative binding modes.** PFMs
989 representing the different binding modes of C2H2-zf proteins.
990
991

992 REFERENCES

- 993 1. Stormo, G.D. & Zhao, Y. Determining the specificity of protein-DNA interactions. *Nat Rev*
994 *Genet* **11**, 751-60 (2010).
- 995 2. Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor
996 binding profiles and its web framework. *Nucleic Acids Res* **46**, D1284 (2018).
- 997 3. Bernard, B., Thorsson, V., Rovira, H. & Shmulevich, I. Increasing coverage of transcription
998 factor position weight matrices through domain-level homology. *PLoS One* **7**, e42779
999 (2012).
- 1000 4. Lambert, S.A. *et al.* Similarity regression predicts evolution of transcription factor
1001 sequence specificity. *Nat Genet* **51**, 981-989 (2019).
- 1002 5. Wunderlich, Z. & Mirny, L.A. Different gene regulation strategies revealed by analysis of
1003 binding motifs. *Trends Genet* **25**, 434-40 (2009).
- 1004 6. Patel, Z.M. & Hughes, T.R. Global properties of regulatory sequences are predicted by
1005 transcription factor recognition mechanisms. *Genome Biol* **22**, 285 (2021).
- 1006 7. Wasserman, W.W. & Sandelin, A. Applied bioinformatics for the identification of
1007 regulatory elements. *Nat Rev Genet* **5**, 276-87 (2004).
- 1008 8. Valouev, A. *et al.* Genome-wide analysis of transcription factor binding sites based on
1009 ChIP-Seq data. *Nat Methods* **5**, 829-34 (2008).
- 1010 9. Marinov, G.K., Kundaje, A., Park, P.J. & Wold, B.J. Large-Scale Quality Analysis of
1011 Published ChIP-seq Data. *G3 (Bethesda)* (2013).
- 1012 10. Consortium, E.P. *et al.* Expanded encyclopaedias of DNA elements in the human and
1013 mouse genomes. *Nature* **583**, 699-710 (2020).
- 1014 11. Rhee, H.S. & Pugh, B.F. Comprehensive Genome-wide Protein-DNA Interactions
1015 Detected at Single-Nucleotide Resolution. *Cell* **147**, 1408-19 (2011).
- 1016 12. Long, H.K., Prescott, S.L. & Wysocka, J. Ever-Changing Landscapes: Transcriptional
1017 Enhancers in Development and Evolution. *Cell* **167**, 1170-1187 (2016).
- 1018 13. Zinzen, R.P., Girardot, C., Gagneur, J., Braun, M. & Furlong, E.E. Combinatorial binding
1019 predicts spatio-temporal cis-regulatory activity. *Nature* **462**, 65-70 (2009).
- 1020 14. Kim, T.H. *et al.* Analysis of the vertebrate insulator protein CTCF-binding sites in the
1021 human genome. *Cell* **128**, 1231-45 (2007).
- 1022 15. Fu, Y., Sinha, M., Peterson, C.L. & Weng, Z. The insulator binding protein CTCF positions
1023 20 nucleosomes around its binding sites across the human genome. *PLoS Genet* **4**,
1024 e1000138 (2008).
- 1025 16. Walker, M. *et al.* Affinity-seq detects genome-wide PRDM9 binding sites and reveals the
1026 impact of prior chromatin modifications on mammalian recombination hotspot usage.
1027 *Epigenetics Chromatin* **8**, 31 (2015).
- 1028 17. Morgunova, E. *et al.* Two distinct DNA sequences recognized by transcription factors
1029 represent enthalpy and entropy optima. *Elife* **7**(2018).
- 1030 18. Stormo, G.D. DNA binding sites: representation and discovery. *Bioinformatics* **16**, 16-23
1031 (2000).
- 1032 19. Rohs, R. *et al.* The role of DNA shape in protein-DNA recognition. *Nature* **461**, 1248-53
1033 (2009).

- 1034 20. Zhao, Y., Ruan, S., Pandey, M. & Stormo, G.D. Improved models for transcription factor
1035 binding site identification using nonindependent interactions. *Genetics* **191**, 781-90
1036 (2012).
- 1037 21. Jolma, A. *et al.* DNA-dependent formation of transcription factor pairs alters their
1038 binding specificity. *Nature* **527**, 384-8 (2015).
- 1039 22. Horton, C.A. *et al.* Short tandem repeats bind transcription factors to tune eukaryotic
1040 gene expression. *Science* **381**, eadd1250 (2023).
- 1041 23. Kribelbauer, J.F., Rastogi, C., Bussemaker, H.J. & Mann, R.S. Low-Affinity Binding Sites and
1042 the Transcription Factor Specificity Paradox in Eukaryotes. *Annu Rev Cell Dev Biol* **35**,
1043 357-379 (2019).
- 1044 24. Lambert, S.A. *et al.* The Human Transcription Factors. *Cell* **175**, 598-599 (2018).
- 1045 25. Wolfe, S.A., Neklodova, L. & Pabo, C.O. DNA recognition by Cys2His2 zinc finger proteins.
1046 *Annu Rev Biophys Biomol Struct* **29**, 183-212 (2000).
- 1047 26. Klug, A. The discovery of zinc fingers and their applications in gene regulation and
1048 genome manipulation. *Annu Rev Biochem* **79**, 213-31 (2010).
- 1049 27. Stubbs, L., Sun, Y. & Caetano-Anolles, D. Function and Evolution of C2H2 Zinc Finger
1050 Arrays. *Subcell Biochem* **52**, 75-94 (2011).
- 1051 28. Nakahashi, H. *et al.* A genome-wide map of CTCF multivalency redefines the CTCF code.
1052 *Cell Rep* **3**, 1678-1689 (2013).
- 1053 29. Kieffer-Kwon, K.R. *et al.* Interactome maps of mouse gene regulatory domains reveal
1054 basic principles of transcriptional regulation. *Cell* **155**, 1507-20 (2013).
- 1055 30. Najafabadi, H.S. *et al.* C2H2 zinc finger proteins greatly expand the human regulatory
1056 lexicon. *Nat Biotechnol* (2015).
- 1057 31. Najafabadi, H.S., Albu, M. & Hughes, T.R. Identification of C2H2-ZF binding preferences
1058 from ChIP-seq data using RCADE. *Bioinformatics* **31**, 2879-81 (2015).
- 1059 32. Jolma, A. *et al.* Multiplexed massively parallel SELEX for characterization of human
1060 transcription factor binding specificities. *Genome Res* **20**, 861-73 (2010).
- 1061 33. Tuerk, C. & Gold, L. Systematic evolution of ligands by exponential enrichment: RNA
1062 ligands to bacteriophage T4 DNA polymerase. *Science* **249**, 505-10 (1990).
- 1063 34. O'Malley, R.C. *et al.* Cistrome and Epicistrome Features Shape the Regulatory DNA
1064 Landscape. *Cell* **165**, 1280-1292 (2016).
- 1065 35. Singer, B.S., Shtatland, T., Brown, D. & Gold, L. Libraries for genomic SELEX. *Nucleic Acids*
1066 *Res* **25**, 781-6 (1997).
- 1067 36. Zimmermann, B., Bilusic, I., Lorenz, C. & Schroeder, R. Genomic SELEX: a discovery tool
1068 for genomic aptamers. *Methods* **52**, 125-32 (2010).
- 1069 37. Jolma, A. *et al.* Perspectives on Codebook: sequence specificity of uncharacterized
1070 human transcription factors. *bioRxiv*, 2024.11.11.622097 (2024).
- 1071 38. Berger, M.F. *et al.* Compact, universal DNA microarrays to comprehensively determine
1072 transcription-factor binding site specificities. *Nat Biotechnol* **24**, 1429-35 (2006).
- 1073 39. Isakova, A. *et al.* SMiLE-seq identifies binding motifs of single and dimeric transcription
1074 factors. *Nat Methods* **14**, 316-322 (2017).
- 1075 40. Jolma, A. *et al.* DNA-Binding Specificities of Human Transcription Factors. *Cell* **152**, 327-
1076 39 (2013).

- 1077 41. Razavi, R. *et al.* Extensive binding of uncharacterized human transcription factors to
1078 genomic dark matter. *bioRxiv*, 2024.11.11.622123 (2024).
- 1079 42. Schmitges, F.W. *et al.* Multiparameter functional diversity of human C2H2 zinc finger
1080 proteins. *Genome Res* **26**, 1742-1752 (2016).
- 1081 43. Vorontsov, I.E. *et al.* Cross-platform DNA motif discovery and benchmarking to explore
1082 binding specificities of poorly studied human transcription factors. *bioRxiv*,
1083 2024.11.11.619379 (2024).
- 1084 44. Galak, A.J. *et al.* Identification of methylation-sensitive human transcription factors
1085 using meSMiLE-seq. *bioRxiv*, 2024.11.11.619598 (2024).
- 1086 45. Weirauch, M.T. *et al.* Determination and inference of eukaryotic transcription factor
1087 sequence specificity. *Cell* **158**, 1431-43 (2014).
- 1088 46. Birke, M. *et al.* The MT domain of the proto-oncoprotein MLL binds to CpG-containing
1089 DNA and discriminates against methylation. *Nucleic Acids Res* **30**, 958-65 (2002).
- 1090 47. Stormo, G.D. & Fields, D.S. Specificity, free energy and information content in protein-
1091 DNA interactions. *Trends Biochem Sci* **23**, 109-13 (1998).
- 1092 48. Weirauch, M.T. *et al.* Evaluation of methods for modeling transcription factor sequence
1093 specificity. *Nat Biotechnol* **31**, 126-34 (2013).
- 1094 49. Kuznetsov, V.A. Mathematical Modeling of Avidity Distribution and Estimating General
1095 Binding Properties of Transcription Factors from Genome-Wide Binding Profiles.
1096 *Methods Mol Biol* **1613**, 193-276 (2017).
- 1097 50. Alexandrov, I., Kazakov, A., Tumeneva, I., Shepelev, V. & Yurov, Y. Alpha-satellite DNA of
1098 primates: old and new families. *Chromosoma* **110**, 253-66 (2001).
- 1099 51. Huttlin, E.L. *et al.* Architecture of the human interactome defines protein communities
1100 and disease networks. *Nature* **545**, 505-509 (2017).
- 1101 52. de Boer, C.G. & Taipale, J. Hold out the genome: a roadmap to solving the cis-regulatory
1102 code. *Nature* **625**, 41-50 (2024).
- 1103 53. Stormo, G.D., Schneider, T.D., Gold, L. & Ehrenfeucht, A. Use of the 'Perceptron'
1104 algorithm to distinguish translational initiation sites in *E. coli*. *Nucleic Acids Res* **10**, 2997-
1105 3011 (1982).
- 1106 54. Laverty, K.U. *et al.* PRIESSTESS: interpretable, high-performing models of the sequence
1107 and structure preferences of RNA-binding proteins. *Nucleic Acids Res* **50**, e111 (2022).
- 1108 55. Bolger, A.M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
1109 sequence data. *Bioinformatics* **30**, 2114-20 (2014).
- 1110 56. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
1111 2078-9 (2009).
- 1112 57. Quinlan, A.R. & Hall, I.M. BEDTools: a flexible suite of utilities for comparing genomic
1113 features. *Bioinformatics* **26**, 841-2 (2010).
- 1114 58. Korhonen, J., Martinmaki, P., Pizzi, C., Rastas, P. & Ukkonen, E. MOODS: fast search for
1115 position weight matrix matches in DNA sequences. *Bioinformatics* **25**, 3181-2 (2009).
- 1116 59. Najafabadi, H.S. *et al.* Non-base-contacting residues enable kaleidoscopic evolution of
1117 metazoan C2H2 zinc finger DNA binding. *Genome Biol* **18**, 167 (2017).
- 1118 60. Korhonen, J.H., Palin, K., Taipale, J. & Ukkonen, E. Fast motif matching revisited: high-
1119 order PWMs, SNPs and indels. *Bioinformatics* **33**, 514-521 (2017).

1120 61. Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of
1121 occurrences of k-mers. *Bioinformatics* **27**, 764-70 (2011).

1122