

Quantitative occupancy of myriad transcription factors from one DNase experiment enables efficient comparisons across conditions

Kaixuan Luo,^{1,2,3,11} Jianling Zhong,^{1,2,3} Alexias Safi,^{2,4} Linda K. Hong,^{2,4}
Alok K. Tewari,⁵ Lingyun Song,^{2,4} Timothy E. Reddy,^{1,2,6,7,8} Li Ma,^{1,9}
Gregory E. Crawford,^{1,2,4} Alexander J. Hartemink^{1,2,3,10,*}

¹Computational Biology & Bioinformatics Graduate Program, Duke University,

²Center for Genomic & Computational Biology, Duke University,

³Department of Computer Science, Duke University,
Durham, NC 27708, USA

⁴Department of Pediatrics, Duke University Medical Center,
Durham, NC 27710, USA

⁵Department of Medical Oncology, Dana-Farber Cancer Institute,
Boston, MA 02215, USA

⁶Department of Biostatistics & Bioinformatics, Duke University Medical Center,

⁷Department of Molecular Genetics & Microbiology, Duke University Medical Center,
Durham, NC 27710, USA

⁸Department of Biomedical Engineering, Duke University,

⁹Department of Statistical Science, Duke University,

¹⁰Department of Biology, Duke University,
Durham, NC 27708, USA

¹¹Department of Human Genetics, The University of Chicago,
Chicago, NC 60637, USA

*To whom correspondence should be addressed; E-mail: amink@cs.duke.edu.

Abstract

Over a thousand different transcription factors (TFs) bind with varying occupancy across the human genome. Chromatin immunoprecipitation (ChIP) can assay occupancy genome-wide, but only one TF at a time, limiting our ability to comprehensively observe the TF occupancy landscape, let alone quantify how it changes across conditions. We developed TOP, a Bayesian hierarchical regression framework, to profile genome-wide quantitative occupancy of numerous TFs using data from a single DNase-seq experiment. TOP is supervised, and its hierarchical structure allows it to predict the occupancy of any sequence-specific TF, even those never assayed with ChIP. We used TOP to profile the quantitative occupancy of nearly 1500 human TF motifs, and examined how their occupancy changes genome-wide in multiple contexts: across 178 cell types, over 12 hours of exposure to different hormones, and across the genetic backgrounds of 70 individuals. TOP enables cost-effective exploration of quantitative changes in the landscape of TF binding.

1 **Introduction**

2 Genes are expressed differently in different types of cells and under different conditions. This
3 response of a cell's gene expression to its internal and external context is enacted in large part
4 through the tuned occupancy of transcription factors (TFs) across the genome. To understand
5 how TFs regulate gene expression, it is critical to determine how likely they are to be present at
6 each location in the genome over time, and how that likelihood changes across varying genetic
7 backgrounds, different cell types, and dynamic environmental conditions. We can measure the
8 quantitative occupancy of one TF along the genome using chromatin immunoprecipitation fol-
9 lowed by high-throughput sequencing (ChIP-seq), provided that a selective antibody exists for
10 the TF. While the ENCODE consortium has generated such data for more than 100 human TFs,
11 the data are typically from only a small number of cell types because of a major limitation of
12 ChIP-seq: a separate experiment is required for each TF in each cell type under each condition.
13 Profiling the time-varying genome-wide occupancy of a large set of TFs across a broad range of
14 cell types and conditions is currently impractical since it would require thousands of antibodies
15 and millions of separate ChIP experiments.

16 An alternative strategy for profiling genome-wide TF occupancy is to exploit DNase-seq or
17 ATAC-seq data, which many groups and consortia have generated for a large number of cell
18 types and experimental conditions¹⁻⁴. The primary advantage of this strategy is that a single
19 DNase-seq or ATAC-seq experiment can be used to profile the occupancy of many different
20 TFs at once, and a number of methods employing this strategy have been proposed in recent
21 years⁵⁻¹⁶.

22 Although multiple methods have been developed to predict TF binding (see Supp. Table
23 S1 for an overview of the modeling frameworks used by a number of these methods), many of
24 them require data types beyond DNase¹³⁻¹⁶, making them less efficient at profiling TF occu-

25 pancy across multiple cell types, conditions, or individuals with different genetic backgrounds
26 than methods requiring only one data type. Furthermore, most existing methods model TF oc-
27 cupancy in a binary fashion—each TF is simply considered present or absent at each location in
28 the genome—ignoring the wealth of quantitative information available in the data¹⁷. While this
29 modeling assumption is consistent with the pervasive practice of binary peak-calling in high-
30 throughput sequencing data, it is inconsistent with our knowledge that at different genomic
31 locations, TFs exhibit different levels of occupancy (likelihood of being bound at that location
32 across the cells in a population) in accordance with prevailing thermodynamic and energetic
33 conditions^{2,18–20}. It is also inconsistent with growing evidence that quantitative levels of TF
34 occupancy can play a significant role in regulating gene expression^{21–24}. Therefore, it is impor-
35 tant that statistical models be developed with a quantitative perspective, allowing us to monitor
36 subtle changes in TF occupancy over time across different genetic backgrounds, cell types, and
37 conditions.

38 Here, we describe a novel method called TF Occupancy Profiler (TOP) that integrates
39 DNase-seq data with information about TF binding specificity (in our case, PWM motifs) to
40 predict the quantitative occupancy of multiple TFs genome-wide. In contrast to earlier methods
41 like CENTIPEDE⁵, PIQ⁷, and msCentipede⁹, TOP is supervised, meaning we can use available
42 ChIP-seq data to train it to high accuracy. Importantly, and in contrast to earlier methods like
43 MILLIPEDE⁶ and BinDNase¹¹, TOP employs a Bayesian hierarchical regression framework,
44 which allows it to obtain both TF-specific and TF-generic model parameters by borrowing in-
45 formation across the full spectrum of training TFs and cell types. The hierarchical nature of
46 TOP is significant because it enables us to predict the occupancy of TFs for which we lack
47 training data, including ones that have never before been profiled with ChIP.

48 We used TOP to predict the genome-wide quantitative occupancy of ~1500 TF motifs across
49 178 human cell types, increasing our cell-type-specific view of TF occupancy in human cells

50 over 200-fold relative to the ENCODE ChIP-seq data used to train the model. We have made
51 these predicted TF occupancy profiles freely available for the community. We used them to
52 construct a cell-type specificity map for different TFs, and identified TFs with selective binding
53 and differential occupancy across cell types. To demonstrate TOP's ability to elucidate the
54 dynamics of TF occupancy, we collected DNase-seq data from A549 cells at 12 time points
55 over 12 hours of glucocorticoid exposure²⁵ and used TOP to efficiently screen nearly 1500
56 TF motifs for increased or decreased occupancy throughout the genome following treatment;
57 we did follow-up ChIP experiments for six of those factors to validate our predictions. We
58 show similar results in separate cells stimulated with androgen or estrogen, two other steroid
59 hormones that act through closely related mechanisms. In another application, we predicted
60 quantitative TF occupancy for the same ~1500 TF motifs across 70 Yoruba lymphoblastoid
61 cell lines (LCLs), and mapped thousands of genetic variants associated with quantitative TF
62 occupancy across individuals (which we term 'topQTLs'). These topQTLs suggest specific
63 mechanistic explanations for the functional impact of genetic variants within regulatory regions.
64 In summary, TOP offers a cost-effective strategy for profiling the occupancy of multiple TFs in
65 a single experiment, markedly enhancing our ability to explore subtle quantitative changes in
66 TF occupancy across cell types, conditions, and genetic variants.

67 **Results**

68 **Bayesian hierarchical regression accurately predicts quantitative TF occu-** 69 **pancy from DNase-seq data**

70 Training TOP entailed two basic steps, as illustrated in Fig. 1. First, we used motif matches to
71 enumerate candidate binding sites, and extracted DNase and ChIP data centered on each site
72 for training. Second, we used MCMC to fit our Bayesian hierarchical regression model on
73 spatially-binned DNase data. Owing to its hierarchical nature, once TOP is trained, we can use

74 it to predict occupancy for any TF in any cell type or condition for which we have DNase-seq
75 data, regardless of whether any ChIP-seq data have ever been collected for that TF.

76 We predicted the quantitative number of ChIP-seq reads around candidate TF binding sites
77 using a site-centric approach, as employed by CENTIPEDE⁵ and its successors. Specifically,
78 for each TF, we first identified candidate binding sites by motif scanning with a permissive
79 threshold (using FIMO with P -value $< 10^{-5}$)²⁶. Then, for each cell type, we considered DNase
80 cleavage events occurring within 100 bp of the candidate binding site. Similarly, when training
81 TOP, we counted the number of ChIP-seq reads within 100 bp of the candidate binding site
82 to serve as the target of our regression. Both DNase and ChIP-seq counts were normalized
83 by library size to account for differences in sequencing depth. We simplified the DNase data
84 into five predictive features using bins that aggregate the number of cleavage events occurring
85 within the motif itself, as well as within two non-overlapping flanking regions upstream and
86 downstream; this is the same binning scheme used in the MILLIPEDE model⁶, and markedly
87 reduces the potential impact of DNase digestion bias^{6,9,10,27}.

88 As an alternative, we tried extracting DNase features using wavelet-transformed multi-scale
89 signals from coarse to fine spatial resolution. However, after variable selection using LASSO,
90 we found only the coarsest resolutions yielded significant features for predicting TF occupancy,
91 while fine resolution features were essentially irrelevant (Supp. Fig. S1). Moreover, the simpler
92 MILLIPEDE binning scheme achieved comparable or better prediction accuracy than optimally
93 selected wavelet features (Supp. Fig. S2). As an added benefit, when fitting TOP to a large
94 number of different TFs across many diverse cell types, the five-bin scheme demonstrated su-
95 perior computational efficiency and better generality in capturing common features across TFs
96 and cell types. Thus, the results that follow are all based on DNase data aggregated into five
97 bins.

98 We chose to use a Bayesian hierarchical model because it allows statistical information to

99 be borrowed across TFs and cell types. TOP's hierarchical structure has three levels (Supp.
100 Fig. S3). The bottom level of the hierarchy contains model parameters specific for each TF \times
101 cell-type combination for which ChIP-seq training data is available. In the middle level, one set
102 of TF-specific but cell-type-generic model parameters is shared across all training cell types for
103 each TF. Finally, the top level has one set of TF-generic parameters jointly learned from all TFs.
104 In other words, we obtain more general model parameters as we move to higher levels of the
105 hierarchy. Once TOP's parameters have been trained, we can predict occupancy for any TF in
106 any cell type in which we have collected DNase-seq data by using a model from the appropriate
107 level of the hierarchy.

108 We evaluated TOP's performance in terms of its fit to quantitative TF occupancy as mea-
109 sured experimentally by ChIP (Fig. 2). TOP predicted quantitative occupancy with varying
110 degrees of accuracy across different TFs (Figs. 2A and 3). In light of technical differences and
111 possible batch effects between DNase-seq data generated in different ENCODE labs, we trained
112 two separate hierarchical models for data from Duke and from Washington (UW), achieving
113 comparable performance between them (Figs. 2B and 3). In general, while bottom level mod-
114 els achieved the highest prediction accuracy (median correlation of 0.70 for Duke and 0.75
115 for UW), middle level models performed equally well (0.70 and 0.75), and top level models
116 performed nearly as well (0.68 and 0.74) (Fig. 2B). This indicates that for a TF that has been
117 profiled with ChIP in some cell type, we can use the TF's middle level model to predict its occu-
118 pancy in any other cell type with available DNase data. In addition, for TFs that have never been
119 profiled with ChIP, the top level TF-generic model will still tend to provide good predictions
120 of quantitative occupancy. Our predicted occupancy accurately matched quantitative ChIP-seq
121 occupancy in various cell types, and allowed us to explore TF occupancy in cell types like the
122 embryonic stem cell line H9ES in which no TF ChIP data have been reported (Fig. 2C). The
123 quantitative predictions produce composite landscapes that sensitively reflect cell-type-specific

124 changes in TF occupancy.

125 To compare with alternative existing methods, since our goal is to efficiently and accurately
126 predict quantitative TF occupancy for candidate binding sites using only a single DNase exper-
127 iment, we focused our comparisons on CENTIPEDE⁵ and msCentipede⁹ (see Methods for a
128 discussion of why these were chosen). These predict TF binding in a site-centric framework
129 but generate only predicted TF binding probabilities rather than ChIP-seq read counts. How-
130 ever, since the CENTIPEDE paper showed a substantial correlation between its TF binding
131 predictions (posterior log odds) and ChIP-seq read counts (sqrt transformed), we could use the
132 posterior log odds of TF binding as a proxy for quantitative ChIP-seq predictions. Our results
133 indicate TOP achieves significantly greater accuracy than both CENTIPEDE and msCentipede
134 for both Duke and UW DNase data (Fig. 2B).

135 **TOP reveals a spectrum of predictability across TFs and cell types**

136 Across TFs, we observed a spectrum of predictability of TF occupancy, as indicated by the blue
137 squares in Fig. 3. Predictability was correlated with the degree of DNase depletion at the motif
138 (Supp. Fig. S4). For TFs with higher prediction accuracy, like NRF1 and ATF1, we observed
139 clear profiles of depletion within motif regions and elevation at nearby flanking regions (Supp.
140 Fig. S5), suggesting direct TF–DNA contact. Many of these TFs have previously been classified
141 as pioneer factors⁷. In contrast, TFs with lower prediction accuracy in the ENCODE data, like
142 STATs and SREBPs, showed less marked elevation at nearby flanking regions, and weak or
143 no depletion at motif regions (Supp. Fig. S5). Weaker DNase depletion profiles may result
144 from transient binding with short residence time—known to occur with nuclear receptors and
145 the AP-1 complex^{20,27–29}—or from ChIP data that include many indirect binding events. For
146 some TFs, we observed a high prediction accuracy in most cell types, but a lower prediction
147 accuracy in just one or two cell types. DNase profiles in the latter cases exhibited markedly

148 weaker depletion (Supp. Fig. S6).

149 TOP uses PWM scores to provide information about which sites are more or less likely to
150 be bound in any cell type or condition. However, in the absence of genetic variation, the PWM
151 score of a particular site does not change across cell types or conditions, so TOP's ability to
152 quantify changes in TF occupancy in such situations depends entirely on changes in the DNase
153 data. As expected, when we compared them as single features, the overall level of DNase
154 cleavage was almost always more correlated with ChIP-seq occupancy across cell types than
155 was the PWM score (Supp. Fig. S7).

156 Having established the reliability of TOP's predictions, we applied it to data from differ-
157 ent contexts to illustrate the biological insights that arise from its ability to efficiently compare
158 quantitative occupancy for myriad TFs across conditions; each of the remaining three subsec-
159 tions explores one of these applications: changes in TF occupancy across different cell types,
160 in response to dynamic environmental conditions, and in the context of genetic variation.

161 **TOP maps out the cell-type specificity of TF occupancy**

162 TFs regulate gene expression in a cell-type-specific manner. To assess TF occupancy differ-
163 ences across cell types, we constructed a cell-type differential occupancy map for multiple
164 TFs to reveal distinct patterns in how TFs direct the gene regulation programs of different cell
165 types. For each TF, we calculated the percentage of candidate sites in each cell type show-
166 ing occupancy significantly higher or lower than the mean across cell types (FDR < 10%);
167 we then clustered TFs on the basis of this measure of cell-type specificity (Fig. 4A). Some
168 TFs—including TAL1, GATA1, and NRF1—displayed large differences in occupancy among
169 cell types, whereas the occupancy of other TFs—like the SPs—was quite cell-type-invariant
170 (Fig. 4B). Lending credence to these results, we successfully recovered TFs known to be specif-
171 ically or differentially expressed in certain cell types. For instance, as expected, we saw that

172 POU5F1 (also known as Oct4) occupancy was significantly higher in stem cells, HNFs (hepa-
173 tocyte nuclear factors) were higher in liver cells, GATAs were higher in K562, REST was lower
174 in medulloblastoma, etc.

175 To explore the relationship between a TF's concentration (here approximated by its gene
176 expression level) and its occupancy, we computed the correlation between each TF's average
177 level of occupancy in each cell type with its gene expression level in that same cell type, and
178 observed several categories of TFs with different relationships (Fig. 4C). Many TFs showed
179 significant positive correlations between their gene expression level and average occupancy,
180 most of which are known to be cell-type-specific TFs, such as FOSL2, HNF4A, FOXA1, and
181 POU5F1 (Supp. Fig. S8A). Surprisingly, three TFs (BATF, BHLHE40, and ZEB1, all known
182 repressors) showed significant negative correlations (Supp. Fig. S8B). Since changes in pre-
183 dicted occupancy reflect changes in the DNase-seq data, we suspect that these repressors, upon
184 binding to the DNA, cause the local chromatin state to become inaccessible to other factors.
185 BATF is active in the immune system and known to interact with IRF4. Interestingly, IRF4 and
186 BATF both had high expression in lymphoblastoid cells, yet we predicted high IRF4 occupancy
187 and low BATF occupancy in those cells (Supp. Fig. S8). Thus, certain sets of cofactors may
188 be utilized to up- or down-modulate the occupancy of related TFs in a cell-type- or condition-
189 specific manner.

190 **TOP monitors the dynamics of TF occupancy during hormone response**

191 Nuclear hormone receptors are TFs specifically activated in response to hormone exposure.
192 Once activated, they bind to specific hormone response elements (HREs) where they regulate
193 gene expression, often in conjunction with the binding of cofactors and remodeling of the chro-
194 matin structure. Glucocorticoid receptor (GR), androgen receptor (AR), and estrogen receptor
195 (ER) are type I nuclear receptors, playing critical roles in immune response or reproductive

196 system development, and are heavily involved in many types of cancer. To investigate TF occu-
197 pancy dynamics in response to glucocorticoid, androgen, or estrogen stimulation, we predicted
198 TF occupancy using DNase-seq data collected under each of these treatment conditions. For
199 glucocorticoid (GC) treatment, we conducted DNase-seq experiments in A549 cells (human
200 alveolar adenocarcinoma cell line) over 12 time points from 0 to 12 hours of GC exposure²⁵.
201 For androgen treatment, we collected DNase-seq data in LNCaP cells (human prostate adeno-
202 carcinoma cell line) over 4 time points from 0 to 12 hours following androgen induction²². For
203 estrogen treatment, we used published DNase-seq data before and after estrogen induction in
204 two kinds of cells: Ishikawa (human endometrial adenocarcinoma cell line) and T-47D (human
205 ductal carcinoma cell line)³⁰.

206 We identified sites with significantly differential TF occupancy before and after estrogen
207 induction, as well as over the full time courses for GC and androgen treatment. We then ranked
208 TFs based on the percentage of sites showing significantly increased or decreased occupancy in
209 response to treatment. We grouped TFs with similar motifs together using RSAT clusters³¹ and
210 present results for all significant clusters in Fig. 5 (results for individual TFs in Supp. Fig. S9).

211 We observed different sets of TFs enriched in response to GC, androgen, and estrogen. In the
212 list of most dynamic clusters for GC response (Fig. 5A), GR was ranked at the top—consistent
213 with recent results showing that motif-driven GR binding is the most predictive feature of GC-
214 inducible enhancers^{25,32}—followed closely by C/EBP²⁵. FOX and GATA clusters appeared
215 next, and in both cases, while we identified more sites whose occupancy increased over the
216 time course, we also detected a significant number that decreased.

217 Among TFs whose occupancy was predicted to be significantly responsive to androgen treat-
218 ment (Fig. 5B), AR was at the top of the list, followed by the FOX cluster. Clusters exhibited
219 very few sites with decreasing occupancy along the time course. These observations are consis-
220 tent with our previous findings that androgen induction mainly leads to an increase in chromatin

221 accessibility, and that AR and FOXA1 are key TFs with increased occupancy²². The fact that
222 the occupancy of many AR and FOXA1 sites increased gradually over the duration of the time
223 course (Supp. Fig. S10A) highlights the importance of a quantitative perspective of TF occu-
224 pancy.

225 In the case of estrogen induction (Fig. 5C), ER was ranked at the top, followed closely by
226 the NFY cluster. Intriguingly, the DNase digestion profiles flanking NFYA binding sites showed
227 striking oscillation patterns similar to those observed within nucleosomes³³ (Supp. Fig. S11).
228 This is consistent with previous reports that NFYA has nucleosome-like properties, and plays
229 an important role in maintaining chromatin structure^{7,34}.

230 That different TFs were enriched in these lists may be partly due to cell type differences,
231 but also suggests different utilization of cofactors for GR, AR, and ER binding in response to
232 hormone stimulation. Interestingly, we observed that PWM scores were significantly higher in
233 sites with increased occupancy than sites of unchanged occupancy for GR, AR, and ER, but not
234 for CEBPB, FOXA1, or NFYA (Supp. Fig. S10B), indicating that motif strength for GR, AR,
235 and ER may play a role in prioritizing the selection of binding sites in response to hormone
236 stimulation. This accords with recent results indicating that GR motif strength is predictive of
237 GC-induced enhancer function³².

238 To independently validate our occupancy predictions with data not seen during training, we
239 compared our predictions throughout the GC time course with ChIP-seq data collected in the
240 same experiment²⁵ (Fig. 5D). We computed the correlation between measured and predicted
241 occupancies for CTCF, JunB, FOSL2, cJun, CEBPB, and GR. Across all six TFs and 12 time
242 points, average correlation was 0.70. Over the time course, it was lowest before treatment
243 (0.63) but otherwise consistent (between 0.68 and 0.72). Among TFs, predictions were the
244 most accurate for CTCF (0.91)—not surprising given how predictable we observed it to be
245 (Fig. 3)—and least for GR (0.52). Two reasons for the lower accuracy of GR are that we used

246 a top level model because GR was not profiled as part of ENCODE, and that GR is known to
247 have a weak DNase footprint²⁹. The correlation is particularly low before treatment (time point
248 0), consistent with observations that many GR binding sites occur at regions of the genome that
249 are already open prior to GC exposure³⁵.

250 **TOP identifies genetic variants associated with predicted TF occupancy** 251 **(topQTLs) and provides mechanistic interpretations for dsQTLs**

252 A large majority of genetic variants associated with complex traits are located in non-coding
253 genomic regions³⁶, suggesting roles in transcriptional regulation. To elucidate this, it is im-
254 perative that we continue to identify genetic variants affecting TF occupancy and chromatin
255 dynamics. To examine whether TOP is capable of sensitively distinguishing quantitatively dif-
256 ferential TF occupancy across individuals or genetic variants, we predicted CTCF occupancy in
257 lymphoblastoid cell lines (LCLs) from two trio studies, one from a CEU (CEPH Utah) family
258 and one from a YRI (Yoruba from Ibadan) family^{21,37}. TOP successfully identified differential
259 CTCF occupancy between individuals across CEU and YRI families (Fig. 6A), and was sensi-
260 tive enough to capture quantitative differences in CTCF occupancy between allele genotypes at
261 allele-specific sites within CEU and YRI families (Fig. 6B).

262 Encouraged by this result, we extended our predictions of genome-wide quantitative oc-
263 cupancy to nearly 1500 TF motifs across 70 Yoruba LCLs using TOP applied to previously
264 published genotype and DNase-seq data³⁸. With the resulting TF occupancy profiles across
265 70 individuals, we applied a QTL mapping strategy to identify genetic variants whose geno-
266 types were significantly associated with changes in predicted TF occupancy, which we called
267 ‘topQTLs’. Since genetic variants that change TF motifs often affect TF binding occupancy
268 by changing DNA binding affinity³⁹, we mapped three versions of topQTLs: within 2 kb and
269 200 bp *cis* testing regions around motif matches, as well as SNPs that lie strictly within those

270 motif matches. We sought to compare our topQTLs with previously reported DNase I sensitivity
271 quantitative trait loci (dsQTLs)³⁸. We estimated the heritability and enrichment of heritability
272 for topQTLs and dsQTLs using stratified LD score regression (S-LDSC) on publicly available
273 GWAS summary statistics of multiple human diseases and complex traits. As shown in Fig. 6C
274 and Supp. Fig. S12, the heritability and enrichment estimates for topQTLs are similar among
275 different window sizes, with slightly higher enrichment for topQTLs near motif locations, con-
276 sistent with our understanding of TF binding mechanisms. Heritability estimates are similar
277 between topQTLs and dsQTLs across traits, though dsQTLs tend to show higher enrichment.

278 We then focused our attention on SNPs within TF motif matches because these have the
279 highest potential for causal interpretation. We compared topQTLs within motif matches to a
280 subset of dsQTLs that we call ‘localizable dsQTLs’, dsQTLs that fall inside the 100 bp windows
281 with which they are linked and also lie within TF motif matches. Of the 1230 reported dsQTLs
282 that were localizable, 943 of them were topQTLs (this number increased to 1000 when using
283 $FDR < 20\%$, while 1141 (93%) were associated with a significant change in predicted TF
284 occupancy under a less stringent threshold of $P\text{-value} < 0.05$). Thus, topQTLs provide a direct
285 mechanistic interpretation for nearly all localizable dsQTLs by revealing the identity of TFs
286 likely to drive the observed changes in chromatin accessibility. Moreover, and importantly,
287 we identified more than six thousand additional topQTLs that were not reported as dsQTLs.
288 Among RSAT-clustered motifs, CTCF, STAT, SP, PU, AP-1, POU and NF- κ B, and RREB1
289 motifs had the greatest number (more than 200) of topQTLs (Fig. 6D); most of these factors are
290 known to be active in LCLs and critical for immune cell development^{38,40}. Fig. 6E shows three
291 sample topQTLs, one for NF- κ B that is a non-localizable dsQTL, another for NF- κ B that is a
292 localizable dsQTL, and one for CTCF that is not reported as a dsQTL.

293 That CTCF had the largest number of topQTLs, over 1300, is noteworthy because CTCF
294 plays a key role in chromosomal looping and commonly demarcates the boundaries of topologi-

295 cally associating domains (TADs)⁴¹. A genetic variant that disrupts a CTCF motif not only may
296 have a significant impact on occupancy at loop anchor sites, but also could disrupt TAD bound-
297 aries. Such disruption has been demonstrated experimentally and pathologically to dysregulate
298 the chromatin landscape and the expression of genes within the affected TAD⁴².

299 **Discussion**

300 We introduce TOP to accurately predict quantitative ChIP-seq occupancy using DNase-seq
301 data. TOP effectively learns both TF-specific and TF-generic model parameters among TFs and
302 across cell types using a Bayesian hierarchical regression framework. TOP employs a super-
303 vised learning strategy, trained with existing TF binding specificity, DNase-seq, and ChIP-seq
304 data, yet can accurately predict TF occupancy for new conditions, cell types, or TFs due to
305 its hierarchical structure. In contrast to traditional ways of analyzing ChIP-seq data through
306 peak-calling to label genomic regions as bound or unbound, TOP adopts a quantitative per-
307 spective, allowing us to predict the level of TF occupancy along a continuum. This opens up a
308 new way to investigate quantitative changes in TF occupancy across cell types, treatment con-
309 ditions, and developmental time courses. TOP is general in that it can predict occupancy for
310 any sequence-specific TF of interest with any new DNase-seq data in any cell type or condition
311 without requiring a new ChIP-seq experiment. TOP's ability to use time-course DNase data
312 over 12 hours of GC treatment served as a cost-effective strategy to study the temporal dynam-
313 ics of TF occupancy. By doing one DNase-seq experiment at each time point, we obtained
314 occupancy predictions for 1500 TF motifs, allowing us to screen for TFs showing significant
315 changes in occupancy. For example, TOP results suggest a significant role for FOX and GATA
316 factors in GC-induced transcriptional response (Fig. 5A). Although developed for DNase-seq
317 data, TOP can easily be extended to ATAC-seq data, and though it was trained on human data,
318 it is equally applicable in other organisms. As an example demonstration, we showed that it can

319 successfully predict quantitative Reb1 occupancy across the yeast genome (Supp. Fig. S13).
320 As a resource for the community, we provide genome browser tracks of predicted occupancy
321 for nearly 1500 motifs across 178 cell types, throughout a 12-hour time course of GC exposure,
322 and across 70 LCLs. The occupancy map of TF \times cell-type combinations alone expands the
323 total output of ENCODE TF ChIP-seq efforts over 200-fold.

324 Recently, methods have emerged for the imputation of missing epigenomic data (histone
325 modifications, chromatin accessibility, etc.) using the many other types of available data gen-
326 erated by the ENCODE consortium^{43–45}. Our approach shares a similar goal with these impu-
327 tation methods in trying to predict unmeasured data using models trained on existing datasets
328 across multiple cell types. However, we note some major distinctions between our approach and
329 these recent imputation strategies. First, our approach requires only DNase-seq (or ATAC-seq)
330 and TF ChIP-seq data for training, and requires only DNase-seq (or ATAC-seq) data to make
331 predictions. In contrast, existing imputation methods often require a large variety of existing
332 assays (DNase-seq, RNA-seq, histone modification ChIP-seq, etc.), which may not be readily
333 available, especially in studies to profile new cell types or treatment conditions. Second, our
334 strategy predicts TF occupancy only at candidate binding sites (based on low stringency motif
335 matches), whereas existing imputation approaches attempt to impute a TF's ChIP-seq signal
336 across the entire genome, devoting statistical power and computational effort to genomic lo-
337 cations where a given TF is not likely to bind, which is the vast majority. Third, TOP uses a
338 Bayesian hierarchical regression framework to model DNase digestion features, which allows
339 for easier interpretation than more complex methods, especially those involving deep neural
340 networks. Last but not least, our hierarchical model is able to predict the binding of TFs that
341 have never been measured by ChIP-seq, a significant advantage over imputation methods that
342 require ChIP-seq training data for TFs of interest.

343 The fact that TOP predicts TF occupancy only at candidate binding sites is also a limitation

344 because it has been observed that for many TFs, a large number of their ChIP-seq peaks do not
345 have motif matches³⁹. On the other hand, this does have some benefits, since distinguishing
346 direct from indirect binding can be difficult using ChIP assays. By focusing only on motif
347 matches, our results can be viewed as predictions of TF occupancy that are explainable by
348 direct binding. Indeed, TOP could be used to suggest direct-binding TFs that may be mediating
349 the indirect binding of other TFs in ChIP experiments. As a last observation, since motif quality
350 directly affects which genomic locations we select as candidate binding sites, and since PWM
351 scores also factor into TOP predictions, better TF binding affinity models are important, and
352 should improve TOP's predictions in the future.

353 Our approach can be viewed as complementary to ChIP-based exploration of TF occupancy.
354 Instead of doing one ChIP-seq experiment for every TF in a particular cell type or condition,
355 TOP needs only one DNase-seq experiment to predict the genome-wide occupancy of many
356 TFs. TOP can therefore be used to screen and identify TFs showing significant changes in oc-
357 cupancy, enabling the prioritization of future ChIP experiments for a small number of key TFs.
358 The modeling strategy we present here offers a foundational and cost-effective approach for pro-
359 filing the quantitative occupancy of myriad TFs across diverse cell types, dynamic conditions,
360 and genetic variants.

361 **Methods**

362 **Candidate binding site selection**

363 We defined candidate TF binding sites by PWM scanning across the genome using FIMO²⁶.
364 When training or applying our model, we included as candidate sites all motif matches with
365 P -value $< 10^{-5}$. Similar to CENTIPEDE⁵ and MILLIPEDE⁶, we filtered out candidate sites if
366 more than 10% of the nucleotides in the surrounding window (100 bp flanking each side of the
367 motif) were unmappable.

368 When training the regression model, if the training TF had more than one motif, we manu-
369 ally selected one based on which was the most representative motif for that TF in the Factorbook
370 database⁴⁶. After training, we used the model parameters estimated at various levels of the TOP
371 hierarchy to make occupancy predictions for 1496 motifs, including both JASPAR core motifs
372 (2014 version) and those used by Sherwood and colleagues⁷. The motifs selected for training
373 the model, along with the full list of all motifs used for prediction in this paper, are provided in
374 Supp. Tables S2 and S3.

375 **Normalization and data preprocessing**

376 To account for differences in sequencing depth across experiments in different cell types or
377 conditions, DNase-seq and ChIP-seq data were normalized by library size (scaled to library
378 sizes of 100 million mapped reads for DNase-seq and 10 million mapped reads for ChIP-seq
379 data). This simple library size normalization is flexible for downstream analysis. We considered
380 other types of normalization methods, including quantile normalization, trimmed mean of M -
381 values (TMM), etc. However, these methods usually assume the same distribution of reads
382 across all the peaks (or a subset of common peaks) among all the experiments, which is too
383 strong an assumption in our case—especially, for example, when comparing hormone receptor
384 binding before and after hormone induction—and leads to a high number of false negatives in

385 the GR, AR, and ER analyses.

386 **DNase feature extraction using binning vs. wavelet coefficients**

387 We systematically evaluated different features of the cleavage events ('cuts') arising from DNase
388 digestion, in an attempt to avoid overfitting⁹ and any possible influence of DNase digestion
389 bias. First, we tried to extract multi-resolution features of DNase digestion data using wavelet
390 multi-resolution decomposition. Wavelet methods provide a natural approach to extract the
391 multi-resolution information contained in both DNase cut magnitude and detail profiles. Here,
392 we decomposed DNase-seq data using Haar wavelets with the `wavethresh` package in R.
393 The detail signals were extracted at different resolution levels through the mother wavelet co-
394 efficients while the scales of cuts at different resolution levels were represented by the father
395 wavelet coefficients. We started with windows of size 128 bp around the motif center, but later
396 focused on 64 bp windows around the motif centers, because that was where the majority of the
397 largest mother wavelet coefficients were located. Then we fit regression models with mother
398 wavelet coefficients and log-transformed father wavelet coefficients at multiple resolution lev-
399 els as predictors, together with PWM score, and conducted variable selection with LASSO.
400 Interestingly, variable selection results suggested the scale of DNase cuts (represented by the
401 father wavelet coefficients) was the most significant feature for predicting TF occupancy (Supp.
402 Fig. S1), consistent with previous findings^{6,10,47}. In contrast, very few spiky DNase signals
403 (represented by the mother wavelet coefficients) were selected. Worse, some of the fine details
404 in the DNase signal in the motif region might arise from sequence-specific DNase digestion
405 bias^{6,10,27,47}.

406 Based on these empirical observations of DNase digestion profiles around motifs, we sim-
407 plified the process of DNase feature extraction by using a more flexible binning scheme in place
408 of the rigid dyadic splitting of the wavelet framework. In previous work, we developed a model

409 called MILLIPEDE that divides the motif region and its flanking regions upstream and down-
410 stream into various distinct bins⁶ (Fig. S2A). Following the binning scheme of MILLIPEDE, we
411 compared different binning models from the most complicated M12 model to the simplest M1
412 model, and evaluated their performance in comparison to an optimally-selected wavelet model.
413 Fig. S2B shows the prediction performance of all these models for four TFs in K562 cells using
414 5-fold cross-validation. In summary, different binning models led to roughly similar predic-
415 tion performances and were generally comparable to a model using optimally-selected wavelet
416 features. In agreement with our earlier results⁶, M5 binning—which effectively summarizes
417 the number of DNase cleavage events in the motif region, nearby flanking regions, and distal
418 flanking regions on both sides of the motif—is complex enough to capture the DNase digestion
419 features sufficient to predict TF occupancy at high accuracy, but is also simple enough to fit into
420 the Bayesian hierarchical regression framework and still yield easily interpretable TF-specific
421 and TF-generic signatures. Fig. 2 shows the prediction performance of the TOP model using
422 M5 binning on the training data. We observed very close agreement in prediction performance
423 using training data vs. test results from cross-validation for the TFs listed in Fig. S2B.

424 **Bayesian hierarchical regression model**

425 We designed the hierarchical model to have three levels, with cell types nested within TF
426 branches. (In principle, we could expand the hierarchical model to have an additional branch
427 with parameters for each cell type, i.e., a cell-type-specific but TF-generic model. However,
428 we expect a TF to have similar model parameters in different cell types. Also, in most of the
429 ENCODE tier 3 cell types, very few TFs have been profiled with ChIP-seq, so we would likely
430 have insufficient data to estimate cell-type-specific parameters for most cell types.)

431 ChIP-seq count data are typically fit using a negative binomial distribution, which uses an
432 extra parameter to model the overdispersion in ChIP-seq data better than a Poisson distribution.

433 However, we found a simpler Gaussian linear model on asinh-transformed ChIP-seq data to be
434 a better choice for fitting our Bayesian hierarchical model (asinh transformation is similar to log
435 transformation but handles zero values more gracefully; we used it successfully in our recent
436 NucID model³³). This choice has the added benefit of applying to non-integer data, which
437 arise whenever we average counts over replicate experiments or conduct data normalization.
438 We compared the prediction accuracy of our Gaussian linear model on asinh-transformed data
439 against the alternative of negative binomial regression on integer-rounded data, and observed
440 very close agreement. We ultimately decided to use the Gaussian distribution because it has a
441 nice conjugacy property, allowing posteriors to be estimated through Gibbs sampling, thereby
442 providing a computational advantage over a negative binomial distribution.

The basic regression model for modeling the asinh-transformed ChIP-seq occupancy $y_{t,c,i}$ that is observed when TF t occupies its candidate binding site i in cell type c can be briefly summarized as:

$$y_{t,c,i} \sim \text{Normal}(\mu_{t,c,i}, \nu_{t,c})$$

where

$$\mu_{t,c,i} = \beta_{t,c}^{(0)} + \sum_{j=1}^J \beta_{t,c}^{(j)} \times D_{i,j} + \beta_{t,c}^{J+1} \times \text{PWM}_i$$

443 The $D_{i,j}$ variable represents DNase feature j for site i , while J is the number of DNase features
444 in the model (in our final model, we use M5 binning so $J = 5$).

The Bayesian hierarchical model is specified as follows:

$$\beta_{t,c}^{(j)} \sim \text{Normal}(b_t^{(j)}, 1) \quad \forall j \in \{0, 1, \dots, J + 1\}$$

$$b_t^{(j)} \sim \text{Normal}(B^j, 1) \quad \forall j \in \{0, 1, \dots, J + 1\}$$

$$B^j \sim \text{Normal}(0, 1) \quad \forall j \in \{0, 1, \dots, J + 1\}$$

$$\frac{1}{v_{t,c}} \sim \text{Gamma}(\tau_t^2, \tau_t)$$

$$\tau_t \sim \text{Gamma}(T^2, T)$$

$$T \sim \text{Gamma}(1, 1)$$

445 We used the consensus Monte Carlo algorithm⁴⁸, a parallel technique to reduce the running
446 time of the Gibbs sampler while maintaining predictive performance. Briefly, we split all data
447 randomly into ten equal parts. Gibbs samplers were run on each part separately in parallel for
448 10^6 iterations. Posterior samples from the first 8×10^5 iterations were discarded for burn-in.
449 Each of the remaining 2×10^5 posterior samples from the ten Gibbs samplers were averaged to
450 get the final posterior samples for each model's parameters.

451 **Comparison of prediction accuracy with existing methods**

452 Both CENTIPEDE and msCentipede predict TF binding probabilities using an unsupervised
453 generative framework to model the DNase digestion profiles around candidate sites (motif
454 matches) without ChIP-seq training data. msCentipede improves on CENTIPEDE by using
455 a multi-scale model framework to better model heterogeneity across sites and replicates. We
456 ran CENTIPEDE and msCentipede on DNase cuts data in each TF-cell type combination under
457 default parameter settings. CENTIPEDE was run on DNase data after pooling the replicate
458 samples. msCentipede was run on individual DNase replicates to better capture heterogeneity
459 (its authors demonstrated how beneficial replicates are to msCentipede accuracy). Because the
460 CENTIPEDE paper showed a substantial correlation between its TF binding predictions (post-

461 rior log odds) and ChIP-seq read counts (sqrt transformed), we computed Pearson correlations
462 between measured ChIP-seq counts (library size normalized and sqrt transformed) and posterior
463 log odds from CENTIPEDE and msCentipede, as well as to the quantitative predictions made
464 by TOP models at each level of the hierarchy.

465 We did not include PIQ⁷ in our comparison, as msCentipede has already been shown to
466 significantly outperform PIQ when it has access to DNase replicates⁹. GERV¹² is a statistical
467 method that learns a k-mer based model to predict TF binding using ChIP-seq and DNase-seq
468 data and scores genetic variants by quantifying the changes of predicted ChIP-seq reads between
469 the reference and alternative allele. Like TOP, it tries to predict quantitative TF occupancy, but
470 its main goal is to score genetic variants that affect TF binding, and it treats DNase signals as a
471 binary feature (open vs. closed), which would not be effective in capturing quantitative changes
472 in DNase signals across dynamic conditions. Also, as a k-mer based method, it does not adopt
473 the motif-centric framework that we and the other methods do. For these reasons, we did not
474 include GERV in our comparison.

475 **Differential occupancy comparison across cell types**

476 We used the `edgeR` package⁴⁹ to identify sites with significantly differential occupancy across
477 cell types. For each TF at each candidate binding site, we tested the cell-type effect by con-
478 trasting the predicted occupancy in each cell type (using DNase replicate samples) against the
479 cell-type mean. Sites with predicted occupancy less than 1 read per million were filtered out
480 from the test, and then sites with a significant cell-type effect (FDR < 10%) were selected.

481 When comparing predicted occupancy across cell types, potential influences from copy
482 number variation (CNV) could lead to false positives. However, since our method predicts
483 TF occupancy using DNase data, and since CNV affects both DNase-seq and ChIP-seq counts
484 in a consistent manner (CNV would lead to higher occupancy in both measured and predicted

485 ChIP-seq in a higher copy number region), our predictions should still agree with measured
486 occupancy. To deal with CNV influences while comparing across cell types, instead of directly
487 correcting CNV on both DNase-seq and ChIP-seq data within the regression model, it is easier
488 to do CNV adjustment as a post-processing procedure on the predicted occupancy using input
489 ChIP-seq data. However, since not all these cell types have input ChIP-seq data available, we
490 did not perform CNV corrections in this study (input correction could be performed in those
491 cell types for which input ChIP-seq data are available).

492 **DNase-seq data across hormone treatment conditions**

493 DNase-seq data from LNCaP cells exposed to androgen were collected in our labs. Data from
494 before induction (time point 0) and after 12 hours were already previously published²² and are
495 available from the Gene Expression Omnibus (GEO) under accession GSE34780. DNase-seq
496 data from the 45 minute and 4 hour treatments, along with more samples from before induction,
497 were generated for this study and will be deposited at GEO prior to publication. LNCaP cells
498 were obtained from ATCC. Cells were maintained using the protocol described at [http://
499 genome.ucsc.edu/ENCODE/protocols/cell/human/LNCaP_Crawford_protocol.
500 pdf](http://genome.ucsc.edu/ENCODE/protocols/cell/human/LNCaP_Crawford_protocol.pdf). Prior to stimulation with either androgen (R1881, methyltrienolone) or vehicle (ethanol)
501 for varying time durations, cells were grown in RPMI-1640 medium with 10% charcoal:dextran
502 stripped medium for 60 hours. Androgen was added to culture medium for final concentra-
503 tion of 1 nM in all experiments. Isolation of total DNA, cleavage with DNase I (henceforth,
504 DNase), and subsequent preparation of sequencing libraries were carried out as previously de-
505 scribed (Song and Crawford, Cold Spring Harbor Protocol 2010). Replicates from 12 hours of
506 androgen exposure were previously sequenced on the Illumina GAIIx platform, whereas repli-
507 cates from the 45 minute and 4 hour time points were sequenced for this study on the Illumina
508 HiSeq2000 platform. Sequenced reads were aligned to the genome and further processed as

509 previously described^{22,50,51}.

510 DNase-seq data from A549 cells exposed to the glucocorticoid hormone dexamethasone
511 were collected in our labs. Detailed methods are provided in our paper²⁵.

512 DNase-seq data from Ishikawa and T-47D cells before and after estrogen exposure were
513 collected by others and previously published³⁰; we downloaded their published data.

514 **Differential occupancy comparison across hormone treatment conditions**

515 In the androgen treatment analysis, we combined DNase-seq data from an earlier study²² with
516 three replicates of uninduced samples and two replicates of 12 hour androgen induced samples
517 (using the Illumina GAIIX sequencing platform), and DNase-seq data generated in this study
518 with two replicates of uninduced samples, two replicates of 45 minute induced samples, and two
519 replicates of 4 hour induced samples (using the Illumina HiSeq2000 sequencing platform). AR
520 ChIP-seq data collected in an earlier study with 4 hour androgen induction in LNCaP cells⁵²
521 matched with our DNase-seq data of 4 hour androgen induction were included in the training
522 dataset for AR in the hierarchical model. For each TF, we used edgeR to test linear, quadratic,
523 and cubic trends of TF occupancy changes over the time course of uninduced, 45 minute, 4
524 hour, and 12 hour induced conditions, adjusting for the batch effect from different sequencing
525 platforms (GAIIX vs. HiSeq sequencing). Sites with predicted occupancy less than 10 were
526 filtered out from the test, and then sites with significant linear, quadratic, or cubic trend of TF
527 occupancy over the time course (FDR < 10%) were selected. Very few sites were found to have
528 a significant quadratic or cubic trend, so we focused on sites with a significant linear trend.

529 In the estrogen treatment analysis, we used previously published DNase-seq and ChIP-seq
530 data generated in Ishikawa (endometrial cancer cell line; previously mislabeled as ECC-1) and
531 T-47D (breast cancer cell line) cells before and after estrogen induction³⁰. ER ChIP-seq data
532 from estrogen induced conditions were matched with the corresponding DNase data and in-

533 cluded in the training dataset for ER in the hierarchical model. Occupancy predictions were
534 made for each TF using its middle level parameters in DNase-seq replicate samples in Ishikawa
535 and T-47D, before and after estrogen stimulation. For each TF, we used `edgeR` to test for
536 differential occupancy, where we considered both cell-type effect (Ishikawa vs. T-47D) and
537 treatment effect (estrogen induced vs. uninduced). Sites with predicted occupancy less than 10
538 were filtered out from the test, and then sites with treatment effect significantly higher or lower
539 than zero ($FDR < 10\%$) were selected.

540 In the glucocorticoid (GC) treatment analysis, we used DNase-seq data collected in our labs
541 from A549 cells (human alveolar adenocarcinoma cell line) over 12 time points from 0 to 12
542 hours following exposure to the glucocorticoid hormone dexamethasone²⁵. For each TF, we
543 used `edgeR` to test linear, quadratic, and cubic trends of TF occupancy changes over the 12
544 time points of GC treatment. Sites with predicted occupancy less than 10 were filtered out from
545 the test, and then sites with significant linear, quadratic, or cubic trend of TF occupancy over
546 the time course ($FDR < 10\%$) were selected. Very few sites were found to have a significant
547 quadratic or cubic trend, so we focused on sites with a significant linear trend.

548 After selecting sites with significant differential occupancy, we ranked TFs based on the
549 percentage of sites showing significantly increased or decreased occupancy in response to treat-
550 ment. TFs with similar motifs were grouped together using RSAT clusters³¹ to simplify down-
551 stream interpretation and visualization.

552 **topQTL mapping**

553 We predicted genome wide TF occupancy for 1496 motifs using previously published genotype
554 information and DNase data generated from LCLs from 70 individuals³⁸. For each motif, we
555 focused on those motif matches that had a SNP inside. When making predictions across the 70
556 LCLs using both PWM scores and DNase data, we fixed the PWM scores for candidate sites to

557 be the average of the PWM scores calculated from the two homozygous genotypes for that SNP,
558 in order to avoid using PWM scores twice: in both occupancy predictions and QTL association
559 testing. We mapped topQTLs by testing the associations between genotypes and predicted TF
560 occupancy across the 70 individuals using a linear model (with R package `MatrixEQTL`). For
561 each TF motif, we selected the 10% of candidate sites with the highest predicted occupancy
562 for QTL mapping and downstream analysis (we tested top 10%, 20%, ..., 100% sites, and
563 found the top 10% sites tended to maximize the number of QTLs detected after multiple test-
564 ing correction). To facilitate comparison with dsQTLs, we followed the same data processing
565 procedures as described by the authors³⁸, including z-score standardization, quantile normaliza-
566 tion, and regressing out 4 PCs to remove unidentified confounders. Following Degner et al.³⁸,
567 we corrected the effect of GC content by first partitioning candidate windows (motif matches
568 plus 100 bp flanking windows on both sides) into bins according to their GC content and then
569 normalizing each sample by subtracting each bin median from all the windows belonging to the
570 partition of the corresponding bin. We mapped three versions of *cis* topQTLs by testing SNPs
571 in 200 bp and 2 kb windows around candidate sites, as well as SNPs within motif matches. For
572 each of the TF motifs, genetic variants with significant associations to predicted TF occupancy
573 (FDR < 10%) were identified as topQTLs for that TF motif, and were the basis of all subse-
574 quent analysis in Fig. 6 of the manuscript. TFs with similar motifs were grouped together using
575 RSAT clusters³¹ to simplify downstream interpretation and visualization.

576 **Heritability and enrichment analysis of GWAS summary statistics using** 577 **S-LDSC**

578 We partitioned the heritability of complex traits and estimated heritability enrichment of top-
579 QTLs and dsQTLs using S-LDSC⁵³. S-LDSC partitions the heritability of genomic annotations
580 using GWAS summary statistics and estimates the enrichment as a ratio of the proportion of her-

581 itability explained by an annotation divided by the proportion of SNPs in that annotation. We
582 constructed binary annotations containing lead SNPs of topQTLs (using 200 bp and 2 kb *cis*-
583 testing regions around motif matches, as well as SNPs within motifs) and lead SNPs of dsQTLs
584 downloaded from Degner et al.³⁸ (using 2 kb and 40 kb *cis*-testing regions around the 5% of
585 100 bp DNase windows with the highest DNase I sensitivity). An FDR threshold of 10% was
586 used for both topQTLs and dsQTLs. We applied S-LDSC to our QTL-based annotations using
587 separate models for each QTL annotation. In our S-LDSC analysis, we adjusted for various
588 baseline annotations of SNPs using a baselineLD model⁵⁴, including gene annotations (cod-
589 ing, UTRs, intron, promoter), minor allele frequency, and LD-related annotations. We did not
590 include functional annotations such as enhancer marks in our baseline model, since these anno-
591 tations are likely correlated with the QTL features of interest, and including them may bias our
592 estimates. The GWAS traits and corresponding references are listed in Supp. Table. S4.

593 **Data and code access**

594 TOP is implemented in R, and all code will be made available on GitHub upon publication. Pre-
595 computed genome-wide tracks of quantitative TF occupancy and links to all code resources will
596 also be made available from a single location upon publication at <http://www.cs.duke.edu/~amink/software/>.
597

598 **References**

- 599 [1] Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature*
600 **489**, 75–82 (2012).
- 601 [2] Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor
602 footprints. *Nature* **489**, 83–90 (2012).

- 603 [3] Buenrostro, J. D., Wu, B., Chang, H. Y. & Greenleaf, W. J. ATAC-seq: A Method for
604 Assaying Chromatin Accessibility Genome-Wide. *Curr Protoc Mol Biol* **109**, 21.29.1–
605 21.29.9 (2015).
- 606 [4] Klemm, S. L., Shipony, Z. & Greenleaf, W. J. Chromatin accessibility and the regulatory
607 epigenome. *Nat. Rev. Genet.* **20**, 207–220 (2019).
- 608 [5] Pique-Regi, R. *et al.* Accurate inference of transcription factor binding from DNA se-
609 quence and chromatin accessibility data. *Genome Res.* **21**, 447–455 (2011).
- 610 [6] Luo, K. & Hartemink, A. J. Using DNase digestion data to accurately identify transcription
611 factor binding sites. In *Pac. Symp. Biocomputing*, 80–91 (World Scientific, Hackensack,
612 NJ., 2013).
- 613 [7] Sherwood, R. I. *et al.* Discovery of directional and nondirectional pioneer transcription
614 factors by modeling DNase profile magnitude and shape. *Nat. Biotechnol.* **32**, 171–178
615 (2014).
- 616 [8] He, H. H. *et al.* Differential DNase I hypersensitivity reveals factor-dependent chromatin
617 dynamics. *Genome Res.* **22**, 1015–1025 (2012).
- 618 [9] Raj, A., Shim, H., Gilad, Y., Pritchard, J. K. & Stephens, M. msCentipede: Modeling
619 heterogeneity across genomic sites and replicates improves accuracy in the inference of
620 transcription factor binding. *PLoS ONE* **10**, e0138030 (2015).
- 621 [10] He, H. H. *et al.* Refined DNase-seq protocol and data analysis reveals intrinsic bias in
622 transcription factor footprint identification. *Nat. Methods* **11**, 73–78 (2014).
- 623 [11] Kähärä, J. & Lähdesmäki, H. BinDNase: a discriminatory approach for transcription

- 624 factor binding prediction using DNase I hypersensitivity data. *Bioinformatics* **31**, 2852–
625 2859 (2015).
- 626 [12] Zeng, H., Hashimoto, T., Kang, D. D. & Gifford, D. K. GERV: a statistical method for
627 generative evaluation of regulatory variants for transcription factor binding. *Bioinformat-
628 ics* **32**, 490–496 (2016).
- 629 [13] Li, H., Quang, D. & Guan, Y. Anchor: trans-cell type prediction of transcription factor
630 binding sites. *Genome Res.* (2018).
- 631 [14] Keilwagen, J., Posch, S. & Grau, J. Accurate prediction of cell type-specific transcription
632 factor binding. *Genome Biol.* **20**, 505 (2019).
- 633 [15] Quang, D. & Xie, X. FactorNet: A deep learning framework for predicting cell type
634 specific transcription factor binding from nucleotide-resolution sequential data. *Methods*
635 **166**, 40–47 (2019).
- 636 [16] Schreiber, J., Bilmes, J. & Noble, W. S. Completing the ENCODE3 compendium yields
637 accurate imputations across a variety of assays and human biosamples. *Genome Biol.* **21**,
638 82 (2020).
- 639 [17] Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure
640 through genomic segmentation. *Nat. Methods* **9**, 473–476 (2012).
- 641 [18] Wasson, T. & Hartemink, A. J. An ensemble model of competitive multi-factor binding of
642 the genome. *Genome Res.* **19**, 2101–2112 (2009).
- 643 [19] Li, X.-Y. *et al.* The role of chromatin accessibility in directing the widespread, overlapping
644 patterns of *Drosophila* transcription factor binding. *Genome Biol.* **12**, R34 (2011).

- 645 [20] Lickwar, C. R., Mueller, F., Hanlon, S. E., McNally, J. G. & Lieb, J. D. Genome-wide
646 protein-DNA binding dynamics suggest a molecular clutch for transcription factor func-
647 tion. *Nature* **484**, 251–255 (2012).
- 648 [21] McDaniell, R. *et al.* Heritable individual-specific and allele-specific chromatin signatures
649 in humans. *Science* **328**, 235–239 (2010).
- 650 [22] Tewari, A. K. *et al.* Chromatin accessibility reveals insights into androgen receptor acti-
651 vation and transcriptional specificity. *Genome Biol.* **13**, R88 (2012).
- 652 [23] Gertz, J., Siggia, E. D. & Cohen, B. A. Analysis of combinatorial cis-regulation in syn-
653 thetic and genomic promoters. *Nature* **457**, 215–218 (2009).
- 654 [24] Segal, E. & Widom, J. From DNA sequence to transcriptional behaviour: A quantitative
655 approach. *Nat. Rev. Genet.* **10**, 443–456 (2009).
- 656 [25] McDowell, I. C. *et al.* Glucocorticoid receptor recruits to enhancers and drives activation
657 by motif-directed binding. *Genome Res.* **28**, 1272–1284 (2018).
- 658 [26] Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: Scanning for occurrences of a given
659 motif. *Bioinformatics* **27**, 1017–1018 (2011).
- 660 [27] Sung, M.-H., Guertin, M. J., Baek, S. & Hager, G. L. DNase footprint signatures are
661 dictated by factor dynamics and DNA sequence. *Mol. Cell* **56**, 275–285 (2014).
- 662 [28] Voss, T. C. *et al.* Dynamic exchange at regulatory elements during chromatin remodeling
663 underlies assisted loading mechanism. *Cell* **146**, 544–554 (2011).
- 664 [29] Goldstein, I. *et al.* Transcription factor assisted loading and enhancer dynamics dictate the
665 hepatic fasting response. *Genome Res.* **27**, 427–439 (2017).

- 666 [30] Gertz, J. *et al.* Distinct properties of cell-type-specific and shared transcription factor
667 binding sites. *Mol. Cell* **52**, 25–36 (2013).
- 668 [31] Castro-Mondragon, J., Jaeger, S., Thieffry, D., Thomas-Chollier, M. & van Helden, J.
669 RSAT matrix-clustering: dynamic exploration and redundancy reduction of transcription
670 factor binding motif collections. *bioRxiv* 065565 (2016).
- 671 [32] Vockley, C. M. *et al.* Direct GR Binding Sites Potentiate Clusters of TF Binding across
672 the Human Genome. *Cell* **166**, 1269–1281.e19 (2016).
- 673 [33] Zhong, J. *et al.* Mapping nucleosome positions using DNase-seq. *Genome Res.* **26**, 351–
674 364 (2016).
- 675 [34] Nardini, M. *et al.* Sequence-specific transcription factor NF-Y displays histone-like DNA
676 binding and H2B-like ubiquitination. *Cell* **152**, 132–143 (2013).
- 677 [35] Reddy, T. E., Gertz, J., Crawford, G. E., Garabedian, M. J. & Myers, R. M. The Hy-
678 persensitive Glucocorticoid Response Specifically Regulates Period 1 and Expression of
679 Circadian Genes. *Mol. Cell. Biol.* **32**, 3756–3767 (2012).
- 680 [36] Hindorff, L. A. *et al.* Potential etiologic and functional implications of genome-wide
681 association loci for human diseases and traits. *Proc Natl Acad Sci U S A* **106**, 9362–9367
682 (2009).
- 683 [37] 1000 Genomes Project Consortium *et al.* A map of human genome variation from
684 population-scale sequencing. *Nature* **467**, 1061–1073 (2010).
- 685 [38] Degner, J. F. *et al.* DNaseI sensitivity QTLs are a major determinant of human expression
686 variation. *Nature* **482**, 390–394 (2012).

- 687 [39] Deplancke, B., Alpern, D. & Gardeux, V. The Genetics of Transcription Factor DNA
688 Binding Variation. *Cell* **166**, 538–554 (2016).
- 689 [40] Tehranchi, A. K. *et al.* Pooled ChIP-seq links variation in transcription factor binding to
690 complex disease risk. *Cell* **165**, 730–741 (2016).
- 691 [41] Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles
692 of chromatin looping. *Cell* **159**, 1665–1680 (2014).
- 693 [42] Lupiáñez, D. G., Spielmann, M. & Mundlos, S. Breaking TADs: How Alterations of
694 Chromatin Domains Result in Disease. *Trends in Genetics* **32**, 225–237 (2016).
- 695 [43] Ernst, J. & Kellis, M. Large-scale imputation of epigenomic datasets for systematic anno-
696 tation of diverse human tissues. *Nat. Biotechnol.* **33**, 364–376 (2015).
- 697 [44] Durham, T. J., Libbrecht, M. W., Howbert, J. J., Bilmes, J. & Noble, W. S. PREDICTD
698 PaRallel Epigenomics Data Imputation with Cloud-based Tensor Decomposition. *Nat*
699 *Commun* **9**, 1–15 (2018).
- 700 [45] Schreiber, J., Durham, T., Bilmes, J. & Noble, W. S. Avocado: a multi-scale deep tensor
701 factorization method learns a latent representation of the human epigenome. *Genome Biol.*
702 **21**, 81–18 (2020).
- 703 [46] Wang, J. *et al.* Sequence features and chromatin structure around the genomic regions
704 bound by 119 human transcription factors. *Genome Res.* **22**, 1798–1812 (2012).
- 705 [47] Cuellar-Partida, G. *et al.* Epigenetic priors for identifying active transcription factor bind-
706 ing sites. *Bioinformatics* **28**, 56–62 (2012).
- 707 [48] Scott, S. L. *et al.* Bayes and big data: The consensus Monte Carlo algorithm. *Int. J.*
708 *Manage. Sci. Engin. Manage.* **11**, 78–88 (2016).

- 709 [49] Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: A Bioconductor package for
710 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139–
711 140 (2010).
- 712 [50] Boyle, A. P., Guinney, J., Crawford, G. E. & Furey, T. S. F-Seq: A feature density estima-
713 tor for high-throughput sequence tags. *Bioinformatics* **24**, 2537–2538 (2008).
- 714 [51] Yardımcı, G. G., Frank, C. L., Crawford, G. E. & Ohler, U. Explicit DNase sequence bias
715 modeling enables high-resolution transcription factor footprint detection. *Nucleic Acids*
716 *Res.* **42**, 11865–11878 (2014).
- 717 [52] Massie, C. E. *et al.* The androgen receptor fuels prostate cancer by regulating central
718 metabolism and biosynthesis. *EMBO J.* **30**, 2719–2733 (2011).
- 719 [53] Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide
720 association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 721 [54] Gazal, S. *et al.* Linkage disequilibrium-dependent architecture of human complex traits
722 shows action of negative selection. *Nat. Genet.* **49**, 1421–1427 (2017).
- 723 [55] Sheffield, N. C. *et al.* Patterns of regulatory activity across diverse human cell types predict
724 tissue identity, transcription factor binding, and long-range interactions. *Genome Res.* **23**,
725 777–788 (2013).
- 726 [56] Bar-Joseph, Z., Gifford, D. K. & Jaakkola, T. S. Fast optimal leaf ordering for hierarchical
727 clustering. *Bioinformatics* **17 Suppl 1**, S22–9 (2001).

728 **Acknowledgement**

729 The authors would like to particularly thank David MacAlpine, Raluca Gordân, Galip Gürkan
730 Yardımcı, Jason Belsky, Ian McDowell, Chris Vockley, Tony D’Ippolito, and the rest of the
731 Duke GGR team for helpful comments during the development of this work or in response to
732 drafts of the manuscript. This work was funded in part by NIH grants U01-HG007900 and
733 R01-GM118551.

734 **Author Contributions**

735 K.L., G.E.C., and A.J.H. designed the study. K.L., J.Z., L.M., G.E.C., and A.J.H. conducted and
736 supervised analyses. A.S., L.K.H, A.K.T., L.S., T.E.R., and G.E.C. conducted and supervised
737 experiments. K.L., J.Z., and A.J.H. wrote the paper.

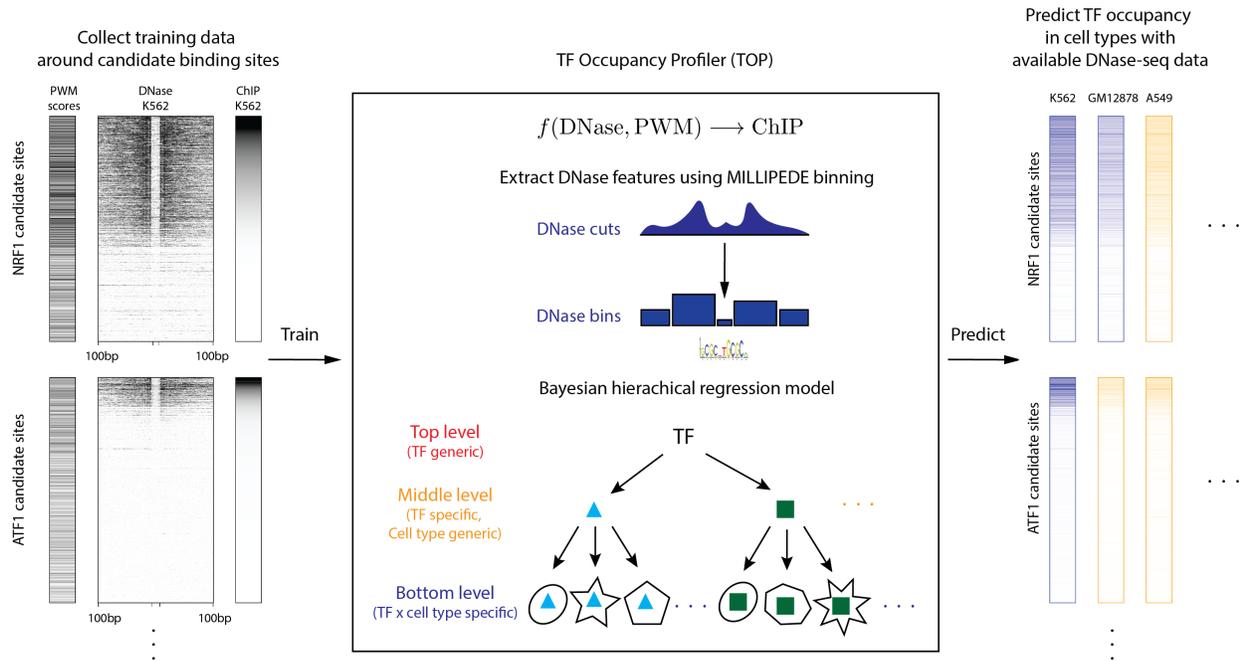


Fig. 1. Schematic outline of the TF Occupancy Profiler (TOP) workflow. (Left) Collect training data. For a sequence-specific TF with a known PWM, compute its candidate binding sites throughout the genome. Then, around each of those sites, collect ChIP-seq and DNase-seq data from the same cell type. (Center) Extract DNase features using MILLIPEDE binning and fit a Bayesian hierarchical regression model to the training data. Bottom level models in the hierarchy make predictions in a TF x cell-type-specific manner, middle level models extend prediction in a TF-specific manner to new cell types, and the top level model extends prediction in a TF-generic manner to new TFs. (Right) Predict occupancy for TFs across cell types. Blue columns indicate a cell type where ChIP-seq measurements are available, allowing us to evaluate the predictive accuracy of our bottom level models. Orange columns indicate a cell type in which we make novel predictions of TF occupancy using middle level parameters of the hierarchical model.

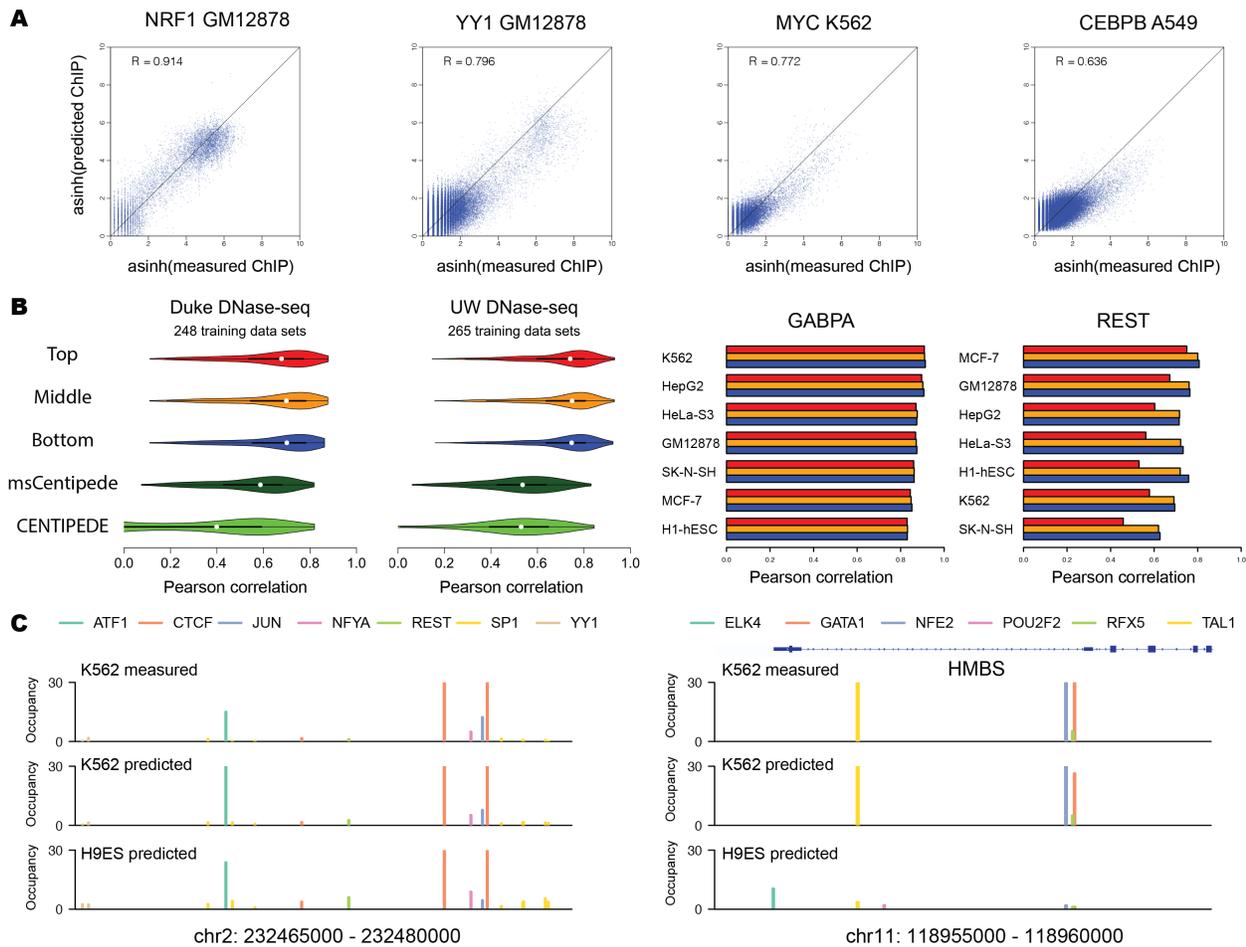


Fig. 2. Evaluation of TOP results. (A) Scatter plots show predicted vs. measured occupancy of a specific TF in a specific cell type, with each dot representing one candidate binding site. The four examples are chosen to represent a range of model performance, from better to worse. (B) Left: Separately for Duke and UW DNase data, violin plots show distribution of Pearson correlations between predicted and measured TF occupancy (sqrt transformed) across TFs and cell types. Predictions were made with TOP models at each level of the hierarchy, as well as CENTIPEDE and msCentipede (using posterior log odds as quantitative measurements of TF occupancy). Right: For many TFs, GABPA being one, all three levels exhibited similar correlation across various cell types. In contrast, for a few TFs, REST being one, the top level model performed markedly worse than bottom and middle level models, suggesting that TF-specific parameters enable more accurate prediction in such cases. (C) Predicted TF occupancy landscapes for two genomic regions in K562 and H9ES cell types. For K562, ChIP-seq data for these TFs are available and are displayed for comparison; for H9ES, no ChIP-seq data are available so TOP provides a novel view of TF occupancy in this embryonic stem cell line. Left: an example genomic region where the occupancy landscape did not change markedly between K562 and H9ES. Right: an example genomic region near the HMBS gene (involved in heme biosynthesis) where GATA1, TAL1, and NFE2 exhibited clear cell-type-specific occupancy.

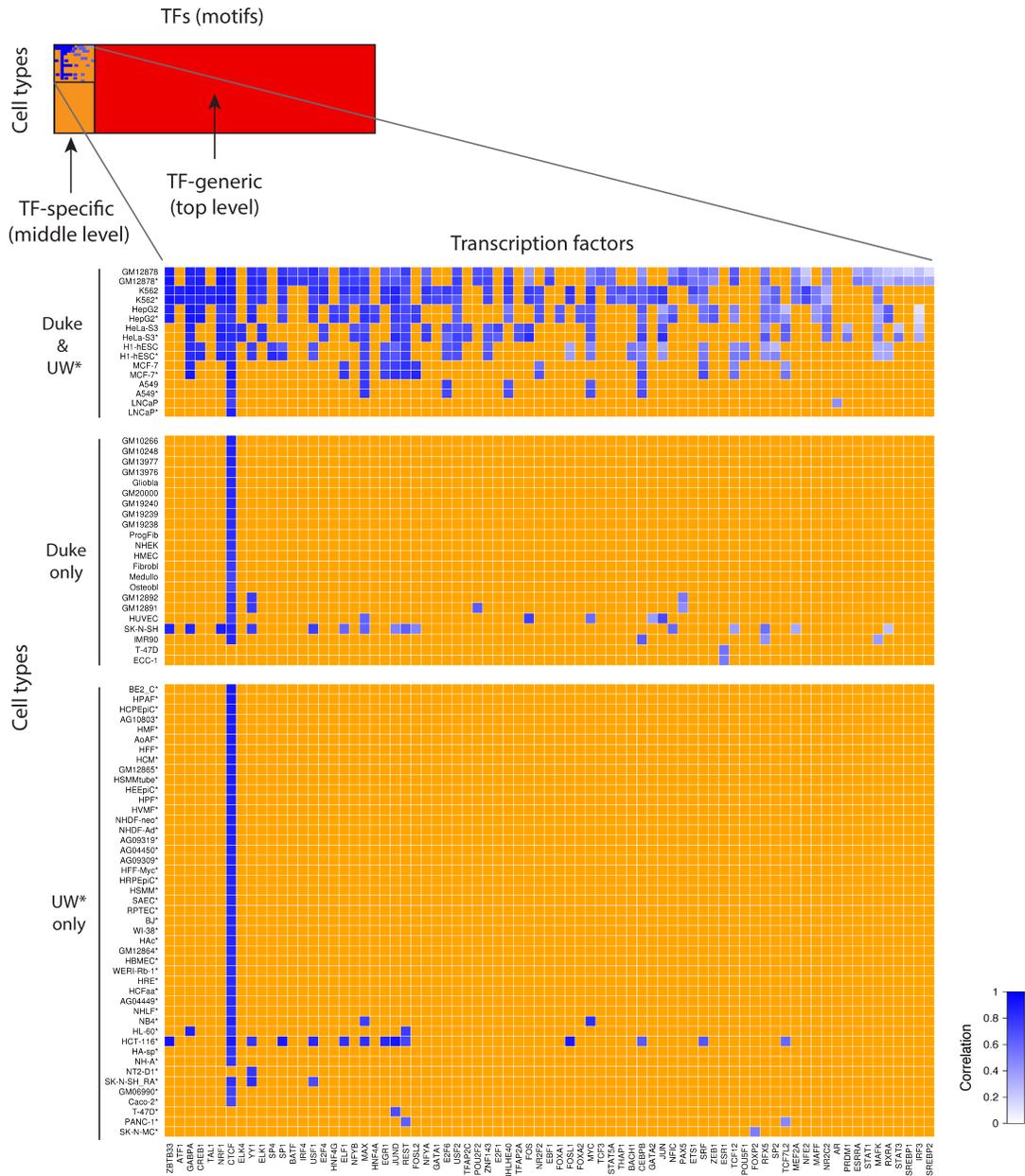


Fig. 3. TOP provides quantitative TF occupancy profiles for nearly 1500 TF motifs in 178 cell types. (Top left) Blue squares represent the TF × cell-type combinations profiled with ChIP-seq as part of the ENCODE project. For each of these TFs, we used a middle level (TF-specific, cell-type–generic) TOP model to generate new occupancy predictions across the rest of the 178 cell types (orange squares). We then used a top level (TF-generic) model to generate new occupancy predictions for the remainder of the 1496 TFs (red squares). (Zoomed inset) In the case of TF × cell-type combinations with ChIP-seq data, we computed the accuracy of TOP predictions; shades of blue indicate the correlation between predicted and measured occupancy. In this submatrix, columns (TFs) and rows within each block (cell types) were sorted by average accuracy, revealing a spectrum of predictability. TFs toward the left were on average more predictable, while TFs to the right were less. Row order is less informative because, except in the top block, it was mainly driven by trivial fluctuations in the predictability of CTCF (in most cell types, CTCF is the only factor whose occupancy was profiled by ENCODE).

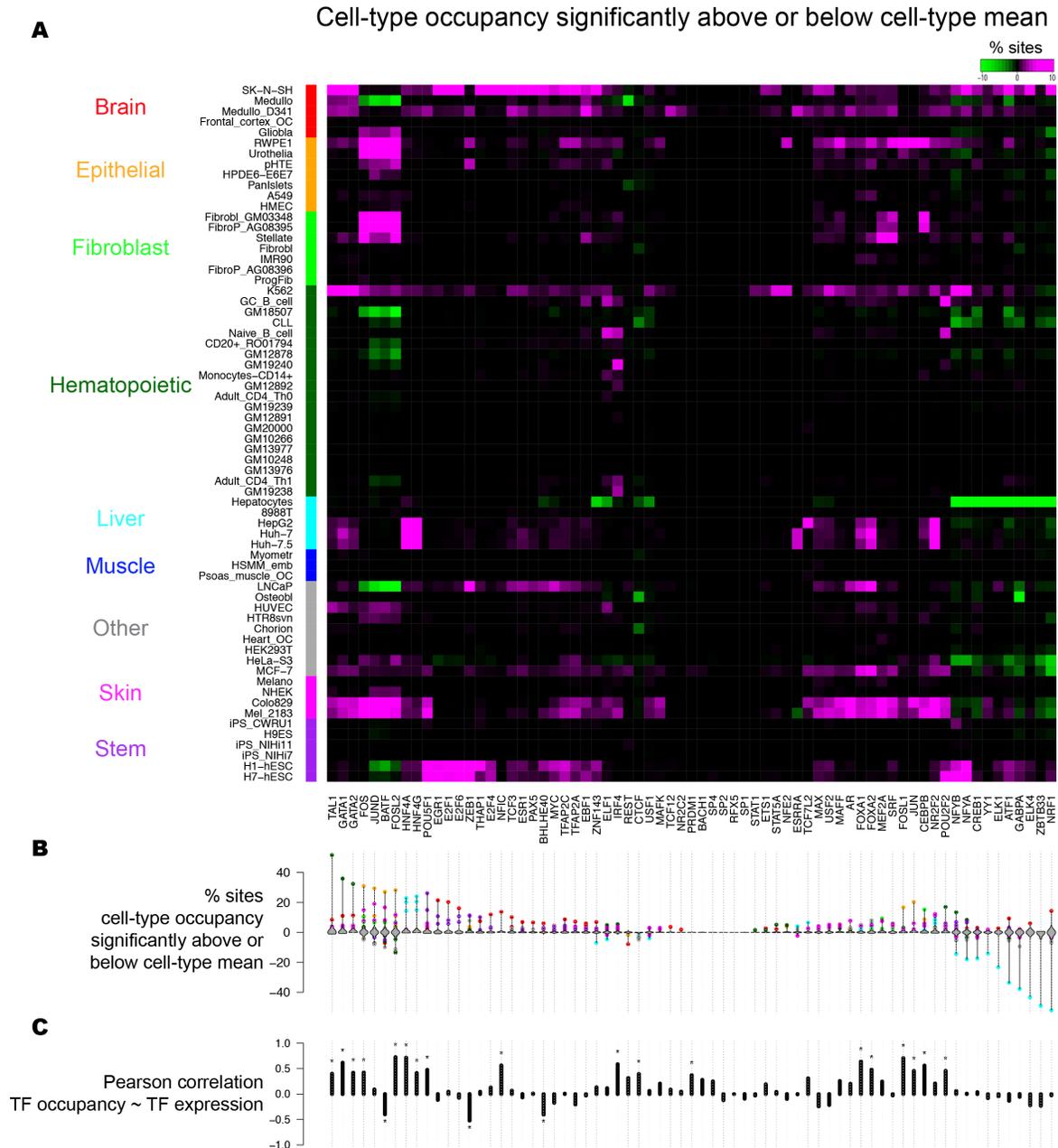


Fig. 4. Cell-type specificity matrix of TFs. (A) Percentage of sites with cell-type occupancy significantly above or below cell-type mean occupancy (FDR < 10%). Cell types were first grouped by lineage (ordered alphabetically)⁵⁵, and within each lineage group were ordered by hierarchical clustering. TFs were ordered by hierarchical clustering (with optimal leaf ordering⁵⁶). (B) Violin plots show for each TF the distribution across cell types of the percentage of sites exhibiting significantly differential occupancy. Colored dots highlight cell types with at least 3% of sites exhibiting significant differential occupancy (color reflects lineage of cell type; for instance, liver and brain exhibit frequent differential occupancy). (C) Pearson correlation across cell types between average predicted occupancy and gene expression of each TF (in this plot, we used only Duke DNase data because corresponding gene expression was measured in each of the cell types).

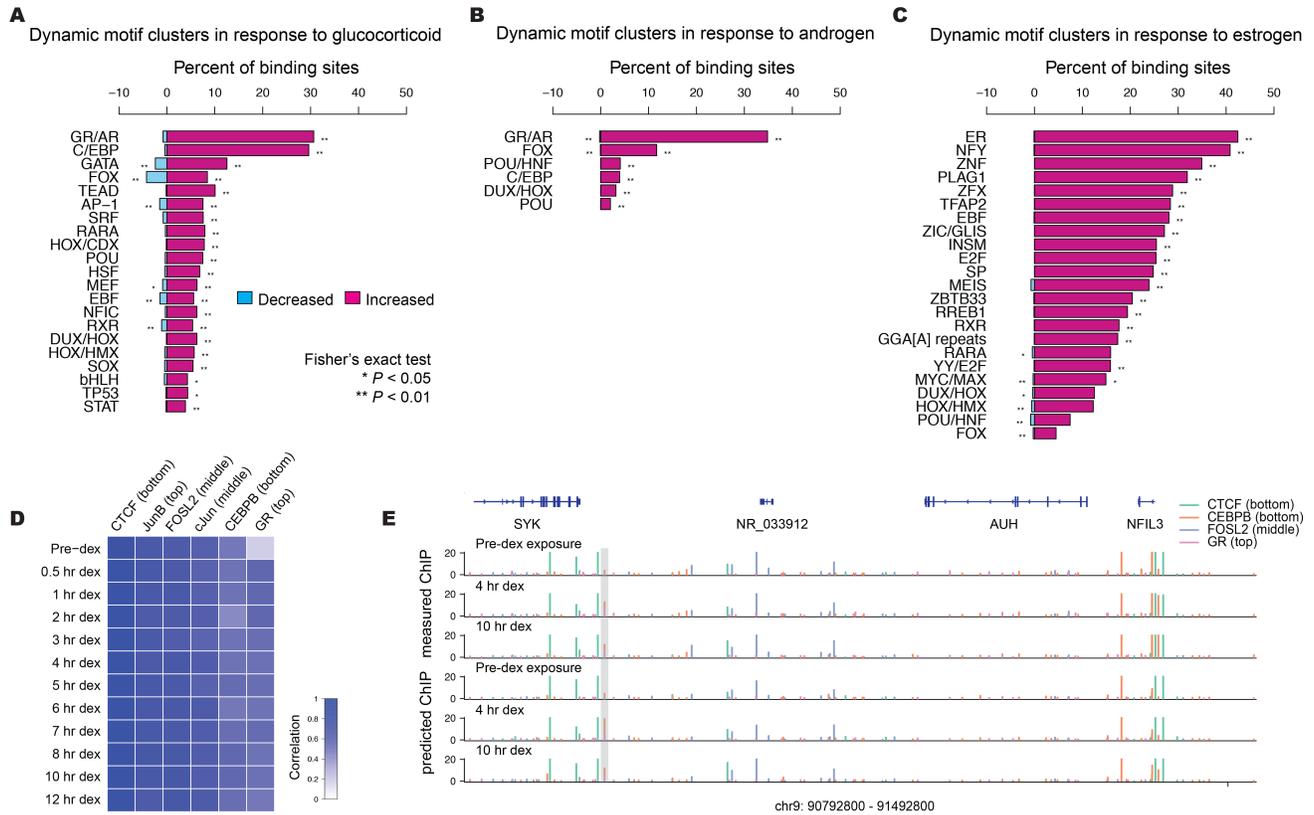


Fig. 5. TF occupancy dynamics in response to hormone stimulation. (A) Motif clusters were ranked by the percentage of candidate sites whose predicted occupancy exhibited either a linear increasing or decreasing trend along the 12 time points of glucocorticoid treatment. Only significant dynamic motif clusters (P -value < 0.05) are listed. (B) Similar to (A), but along the four time points of androgen treatment. (C) Similar to (A) and (B), but before and after estrogen treatment. Because DNase data was collected at 12 time points during treatment with GC, at four time points with androgen, and at only two time points with estrogen, numbers are not necessarily comparable between different experiments in (A), (B), and (C). (D) Prediction accuracy for six TFs was evaluated afterwards using subsequently generated ChIP-seq data²⁵. Shades of blue indicate the correlation between predicted and measured occupancy for each of the six TFs at each time point. Columns (TFs) were sorted by average accuracy across the 12 time points. (E) Measured and predicted TF occupancy landscapes of CTCF, CEBPB, FOXL2, and GR in an example genomic region on human chromosome 9. Predicted occupancy corresponded well with measured occupancy across time, for example revealing in the highlighted region where CEBPB occupancy increased following GC treatment.

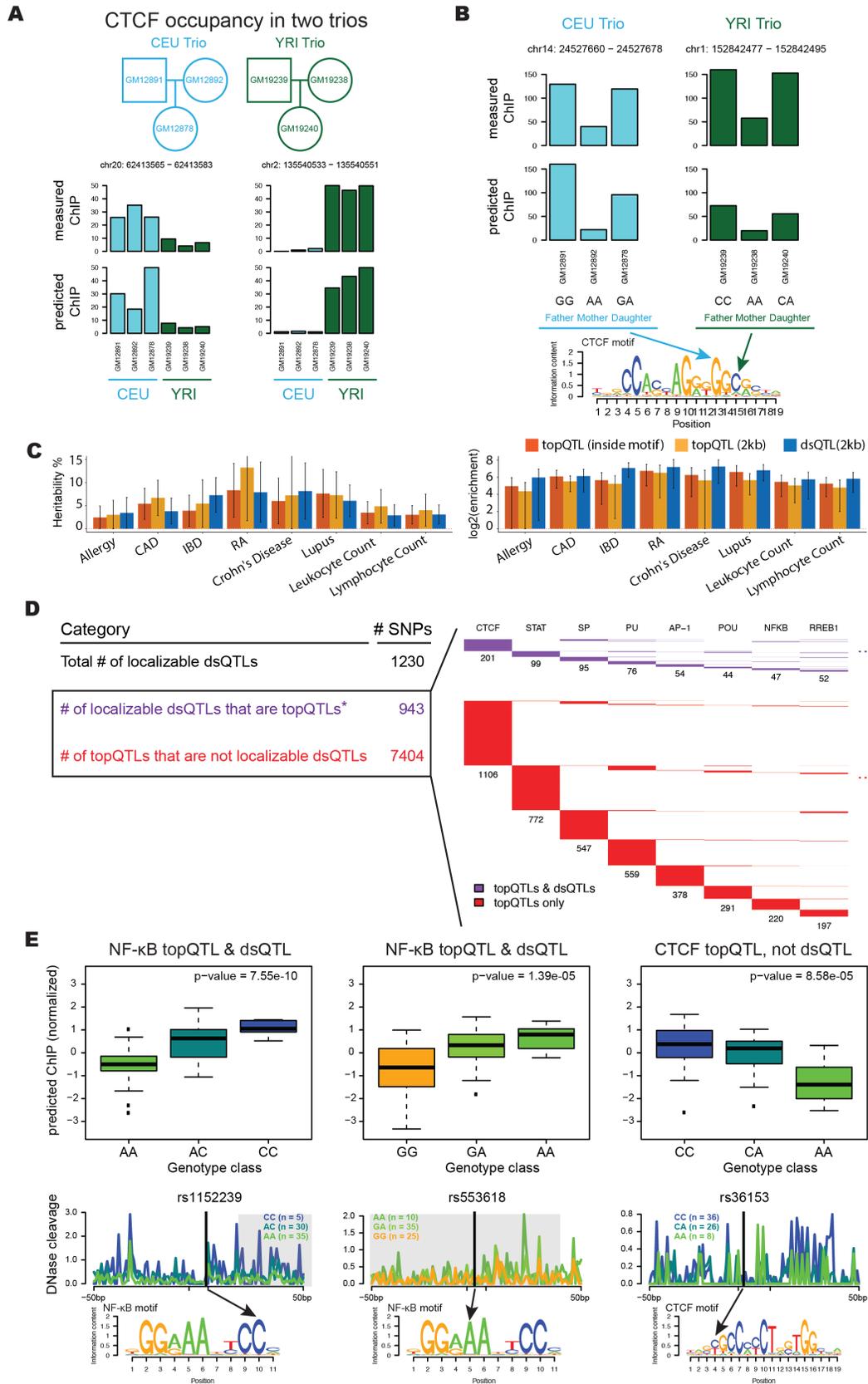


Fig. 6. TF occupancy profile QTLs (topQTLs). (A) Predicting individual-specific occupancy of CTCF. (B) Predicting quantitative allele-specific occupancy of CTCF. (C) Heritability and $\log_2(\text{enrichment})$ estimates for topQTLs and dsQTLs in different diseases and traits using S-LDSC. Lead SNPs for both topQTLs and dsQTLs were used as binary annotations. Bars represent topQTLs from SNPs within motif matches, topQTLs within 2 kb around motif matches, and dsQTLs within 2 kb of the 100 bp DNase window. Error bars represent 95% confidence intervals and were truncated at 15% in the heritability figure. (D) Intersections of topQTLs with localizable dsQTLs (those within their own 100 bp windows and also within motif matches). topQTLs were defined with $\text{FDR} < 10\%$ (* with $\text{FDR} < 20\%$, the number localizable dsQTLs that are also topQTLs is 1000). (Right) Largest motif clusters for topQTLs are displayed in the matrix; each row represents one topQTL which can be explained by one or more motif clusters in the columns. (E) Examples of topQTLs showing normalized allele-specific predicted occupancy; average DNase digestion profiles within 50 bp of the motif for each allele (significant dsQTL windows shaded in gray); and SNP locations within motifs. The CTCF topQTL overlapped a measured CTCF ChIP-seq peak in multiple LCLs, but was not identified as a dsQTL.