# Global reference mapping and dynamics of human transcription factor footprints

Jeff Vierstra[1,*], John Lazar[1,2], Richard Sandstrom[1], Jessica Halow[1], Kristen Lee[1], Daniel Bates[1], Morgan Diegel[1], Douglas Dunn[1], Fidencio Neri[1], Eric Haugen[1], Eric Rynes[1], Alex Reynolds[1], Jemma Nelson[1], Audra Johnson[1], Mark Frerker[1], Michael Buckley[1], Rajinder Kaul[1], Wouter Meuleman[1], John A. Stamatoyannopoulos[1-3,*]

[1] Altius Institute for Biomedical Sciences, Seattle, Washington, USA
[2] Department of Genome Sciences, University of Washington, Seattle, Washington, USA
[3] Division of Oncology, Department of Medicine, University of Washington, Seattle, Washington, USA

**Correspondence to:** jvierstra@altius.org and jstam@altius.org

## Abstract

Combinatorial binding of transcription factors to regulatory DNA underpins gene regulation in all organisms. Genetic variation in regulatory regions has been connected with diseases and diverse phenotypic traits[1], yet it remains challenging to distinguish variants that impact regulatory function[2]. Genomic DNase I footprinting enables quantitative, nucleotide-resolution delineation of sites of transcription factor occupancy within native chromatin[3–5]. However, to date only a small fraction of such sites have been precisely resolved on the human genome sequence[5]. To enable comprehensive mapping of transcription factor footprints, we produced high-density DNase I cleavage maps from 243 human cell and tissue types and states and integrated these data to delineate at nucleotide resolution ~4.5 million compact genomic elements encoding transcription factor occupancy. We map the fine-scale structure of ~1.6 million DHS and show that the overwhelming majority is populated by well-spaced sites of single transcription factor:DNA interaction. Cell context-dependent cis-regulation is chiefly executed by wholesale actuation of accessibility at regulatory DNA versus by differential transcription factor occupancy within accessible elements. We show further that the well-described enrichment of disease- and phenotypic trait-associated genetic variants in regulatory regions[1,6] is almost entirely attributable to variants localizing within footprints, and that functional variants impacting transcription factor occupancy are nearly evenly partitioned between loss- and gain-of-function alleles. Unexpectedly, we find that the global density of human genetic variation is markedly increased within transcription factor footprints, revealing an unappreciated driver of cis-regulatory evolution. Our results provide a new framework for both global and nucleotide-precision analyses of gene regulatory mechanisms and functional genetic variation.

1    **Introduction**

2    Genome-encoded recognition sites for sequence-specific DNA binding proteins are the
3    atomic units of eukaryotic gene regulation.  Binding of regulatory factors to their cognate
4    elements in vivo shields them from nuclease attack, giving rise to protected single nucleotide-
5    resolution DNA 'footprints'.  The advent of DNA footprinting using the non-specific nuclease
6    DNase I[7] marked a major turning point in analysis of gene regulation, and facilitated the
7    identification of the first mammalian sequence-specific DNA binding proteins[8].  Genomic DNase
8    I footprinting[3] enables genome-wide detection of DNA footprints chiefly within regulatory DNA
9    regions, but also over other genomic elements where DNase I cleavage is sufficiently dense.

10   DNase I footprints define sites of direct regulatory factor occupancy on DNA and can be
11   used to discriminate sites of direct vs. indirect occupancy within orthogonal data from chromatin
12   immunoprecipitation and sequencing (ChIP-seq) experiments. Cognate transcription factors
13   (TFs) can be reliably assigned to DNase I footprints based on matching to consensus
14   sequences, enabling TF-focused analysis of gene regulation and regulatory networks[9], and the
15   evolution of regulatory factor binding patterns within regulatory DNA[10].  DNase I is a small
16   enzyme, roughly the size of a typical transcription factor that recognizes the minor groove of
17   DNA and hydrolyzes single-stranded cleavages that, in principle, reflect both the topology and
18   the kinetics of DNA-protein interaction.  Previous efforts to exploit this feature[4] were complicated
19   by slight sequence-driven variation in cleavage preferences; however, these have now been
20   exhaustively determined[11], setting the stage for fully resolved tracing of DNA-protein interactions
21   within regulatory DNA.

22   Currently we lack a comprehensive, nucleotide-resolution annotation of small DNA
23   elements encoding regulatory factor recognition sites that are selectively occupied in different
24   cell types.  Such a reference is essential both for analysis of cell-selective occupancy patterns,
25   and for systematic integration with genetic variation, particularly that associated with diseases
26   and phenotypic traits. Here we combine sampling of >67 billion DNase I cleavages from >240
27   human cell types and states to index, with unprecedented accuracy and resolution, human
28   genomic footprints and thereby the sequence elements that encode transcription factor
29   recognition sites.  We leverage this index to comprehensively assign footprints to transcription
30   factor archetypes, define patterns of cell-selective occupancy, and analyze the distribution and
31   impact of human genetic variation on regulatory factor occupancy and the genetics of common
32   diseases and traits.

33   **Comprehensive mapping of human TF footprints**

34   To create comprehensive maps of TF occupancy, we selected and deeply-sequenced
35   high-quality DNase I libraries from 243 cell and tissues types derived from diverse primary cells
36   and tissues (*n*=151), primary cells in culture (*n*=22), immortalized cell lines (*n*=10) and cancer
37   cell lines and primary samples (*n*=60) (**Extended Data Table 1**). Collectively, we uniquely
38   mapped 67.6 billion DNase I cleavage events (mean 278.2 million cleavages mapped per
39   biosample), greatly eclipsing prior studies[4].  On average, 49.7% DNase I cleavage within each
40   biosample mapped to DNase I hypersensitive sites covering 1-3% of the genome.

41    To identify DNase I footprints genome-wide, we developed a novel computational
42    approach incorporating both chromatin architecture and exhaustively determined empirical
43    DNase I sequence preferences to determine expected per-nucleotide cleavage rates across the
44    genome, and to derive, for each biosample, a statistical model for testing whether observed
45    cleavage rates at individual nucleotides deviated significantly from expectation (**Extended Data
46    Fig. 1**, **Extend Data Fig. 2**, and **Methods**). We note that deriving cleavage variability models for
47    each biosample accounts for additional sources of technical variability beyond DNase I
48    cleavage preference.

49    Using this model, we performed *de novo* footprint discovery independently on each of
50    the 243 biosamples, detecting on average 657,029 high-confidence footprints per biosample
51    (range 220,580-1,664,065) (empirical false discovery rate <1%; **Methods**), and collectively
52    159.6 million footprint events across all biosamples. At the level of individual nucleotides, *de*
53    *novo* footprints were highly concordant between replicates of the same cultured cell type or
54    between same primary cell and tissues types from different individuals (median Pearson's *r* =
55    0.83 and 0.74, respectively) (**Extended Data Fig. 3a-c**). The significance of protected
56    nucleotides tracked closely both the presence of known transcription factor recognition
57    sequences and the level of per-nucleotide evolutionary conservation (**Extended Data Fig. 3d-
58    e**). Within each biosample, genomic footprints encompassed an average of ~7.6 Mb (0.2%) of
59    the genome, with a mean of 4.3 footprints per DHS with sufficient read depth for robust
60    detection (normalized cleavage density within DHS ≥1).

61    **A unified index of human genomic footprints**

62    Comparative footprinting across many cell types has the potential to illuminate both the
63    structure and function of regulatory DNA, yet a systematic approach for joint analysis of
64    genomic footprinting data has been lacking. Given the scale and diversity of the cell types and
65    tissues surveyed, we sought to develop a framework that could integrate hundreds of available
66    genomic footprinting datasets to increase the precision and resolution of footprint detection and,
67    furthermore, serve as a scaffold to build a common reference index of TF-contacted DNA
68    genome-wide.

69    To accomplish this, we implemented an empirical Bayes framework that estimates the
70    posterior probability that a given nucleotide is footprinted by incorporating a prior on the
71    presence of a footprint (determined by footprints independently identified within individual
72    datasets) and a likelihood model of cleavage rates for both occupied and unoccupied sites (**Fig.
73    1a** and **Methods**). **Fig. 1b** depicts per-nucleotide footprint posterior probabilities computed for
74    two DHS within the *RELB* locus across all 243 biosamples exposing the nucleotide-resolved TF
75    occupancy architecture for each element. A notable feature of these data is the remarkable
76    positional stability and discrete nature of footprints within each DHS across the tens to hundreds
77    of biosamples. Indeed, plotting individual nucleotides scaled by their footprint prevalence across
78    all samples precisely demarcates the core recognition sequences of diverse TFs (**Fig. 1b**,
79    bottom).

80  To systematically create a reference set of TF-occupied DNA elements genome-wide,
81  we applied the Bayesian approach to all DHSs detected within one or more of the 243
82  biosamples, and applied the same consensus approach used to establish the consensus DHS
83  index[12] to collate overlapping footprinted regions across individual biosamples into distinct high-
84  resolution consensus footprints (**Methods**). Collectively, we delineated ~4.46 million consensus
85  footprints within ~1.6 million individual DHS (**Fig. 1c**). 82.6% of consensus footprints localize
86  directly within the core of a DHS peak (avg. width 203bp) with virtually all residing within 250bp
87  of a DHS peak summit (**Figure 1d**). As expected, consensus-defined footprints were markedly
88  more reproducible than footprints detected independently within a given biosample (avg.
89  Jaccard distance 0.43 vs 0.29, respectively) (**Extended Data Fig. 4d-e**). Consensus footprints
90  had an average width of 16 bp (middle 95%: 7–44 bp; 90%: 7-36 bp; 50%: 9-21 bp), and
91  collectively annotate 2.1% (72 Mb) of the human genome reference sequence, a compartment
92  slightly larger than protein coding sequences (~1.5%).

## Assigning footprints to TF motifs

94   Our understanding of the recognition motif landscape of human transcription factors has
95  undergone dramatic development during the past decade, and recognition sequences now exist
96  for all major families and subfamilies, and for a large number of individual TF isoforms[13–16]. We
97  thus sought to create a reference mapping between annotated transcription factors and
98  consensus human genomic footprints by (i) compiling and clustering all publicly available motif
99  models[13,17,18]; (ii) creating non-redundant TF archetypes by placing closely-related TF family
100 members on a common sequence axis (**Extended Data Fig. 5**, **Extended Data Table 2** and
101 **Methods**); (iii) aligning these archetypes to the human reference sequence at high stringency
102 ($p<10^{-4}$); and (iv) enumerating all potential TF archetypes compatible with each consensus
103 footprint on the basis of overlap and match stringency (**Methods**). In total, 80.7% of the ~4.46
104 million consensus footprints could be assigned to at least one TF recognition sequence (≥90%
105 overlap; **Methods**), of which 860,780 (19.3%) could be unambiguously assigned to a single
106 factor, and 2,038,220 (45.7%) to a single factor with two lower-ranked alternatives (**Extended
107 Data File 2**).

## Primary architecture of regulatory regions

109  Despite intensive efforts over the past three decades the primary architecture of
110 regulatory regions has remained elusive, with the singular exception of the interferon
111 'enhanceosome'[19]. A prerequisite for understanding the primary architecture of active
112 regulatory DNA is accurate tracing of the TF:DNA interface over an extended interval. Because
113 transcription factor engagement within DNA major or minor grooves creates subtle alterations in
114 DNA shape and protects underlying phosphate bonds from nuclease attack via steric
115 hindrance[5], we asked to what extent fluctuations in corrected DNase I cleavage rates within
116 consensus footprints accurately reflect the topology of the TF:DNA interface. Poly-zinc fingers
117 are the most prevalent class of human transcription factors and have recognition interfaces that
118 potentially cover tens of nucleotides[16]. The DNA recognition domain of the genomic master
119 regulator CTCF comprises 11 zinc fingers, potentially encoding 33bp of sequence (or DNA
120 shape[20]) recognition. We identified 25,852 footprints that coincided precisely with CTCF motifs.

121    Transposing the average corrected per-nucleotide cleavage propensity with an extended co-
122    crystal structure of CTCF[21] accurately traced all features of the protein:DNA interaction
123    interface, including focal hypersensitivity within bent hinge region between zinc fingers 7 and
124    9[22,23] (**Fig. 2a** and **Methods**). A similar result was obtained for widely divergent classes of DNA
125    binding domains such as the paired-box containing TF PAX6[24] (**Fig. 2b**) and other TFs with
126    extant co-crystal structures (not shown). Critically, these topological features are evident at the
127    level of individual TF footprints on the genome (**Fig. 2a-b** and **Extended Data Fig. 6**), indicating
128    that the extended profile of corrected per-nucleotide DNase I cleavage across entire regulatory
129    regions should, in principle, provide a snapshot of the primary structure of active regulatory
130    DNA.

131    **Distinguishing independent vs. cooperative modes of TF occupancy**

132    Transcription factors compete cooperatively with nucleosomes for access to regulatory
133    DNA[25,26]. Although fundamental to eukaryotic gene regulation, it is currently unknown whether
134    nucleosome-enforced TF cooperativity derives primarily from local protein-protein interactions or
135    results from the synergistic effect of independent TF:DNA binding events[26]. We reasoned that
136    the number, relative spacing, and morphology of TF binding events within individual regulatory
137    elements could be used to gain insight into the mechanistic basis of TF cooperativity. We
138    observed that the average footprint width for diverse TFs tightly corresponded to the total width
139    of its recognition sequence (Spearman's $\rho=0.82$, *p*-value=0.001) indicating that DNase I
140    cleavage precisely delineates the boundaries between occupied and unoccupied DNA at
141    nucleotide resolution (**Fig. 3a**).

142    Since the width of genomic footprints tightly tracks the physical structure of individual
143    TFs bound to DNA (**Fig. 2a-b**), and direct TF:TF interactions are dependent on close
144    proximity[19], as such interactions should be reflected in larger footprints that harbor multiple TF
145    recognition sites.  Conversely, independent TF:DNA interaction events should be reflected by
146    compact and widely-spaced footprints harboring single TF recognition sites.  As such, the
147    prevalence of cooperativity mediated by direct TF:TF interactions vs. synergy of independent
148    binding events should be reflected in relative proportion of wide multi-motif footprints vs. well-
149    spaced single footprints. Larger footprints are overwhelmingly associated with two (or more)
150    recognition sequences (**Fig. 3b**), yet such footprints represent only 8% of the global footprint
151    landscape. By contrast, 92% of footprints harbor a single TF recognition site (**Fig. 3c**).

152    Transcription factors can distort DNA upon engagement; as such, the spacing of
153    transcription factors can be critical for establishing an active regulatory structure.  To quantify
154    global footprint spacing patterns, we first binned each DHS by its average accessibility across
155    all samples (because footprint discovery depends on total DNase I cleavage; **Extended Data**
156    **Fig. 1b**), and for each bin we computed the mean number of footprints present per element and
157    their relative edge-to-edge spacing. The density of footprints within the most deeply sampled
158    DHSs genome-wide plateaued at an average of 5.5 per 200 bp (average width of a DHS peak)
159    (**Fig. 3d**), which is in remarkable agreement with a theoretical prediction of the number of
160    human TFs required to destabilize a canonical nucleosome[26] and to encode specificity[27]. Within
161    DHSs, footprints exhibit an average edge-to-edge spacing of ~21bp (middle 50%, 12-35bp)

162  (**Fig. 3e**). Taken together, these results are compatible with the observed lack of evolutionary
163  constraint on the spacing and orientation of TF motifs[28] and provide strong evidence that
164  regulatory DNA marked by DHSs is chiefly instantiated by independent but synergistic TF
165  binding modes (**Fig. 3f**).

**Cell-selective TF occupancy patterns**

167  Analysis of footprint occupancy across all biosamples revealed strong enrichment for the
168  recognition sequences of key regulatory TFs in their cognate lineages (**Extended Fig. 7a** and
169  **Methods**; for example: HNF4A in fetal intestine and liver; GATA factors in erythroid and
170  placental/trophoblast cells and tissues; NEUROG1 in brain; myogenic regulatory factors (e.g,
171  MYF6, MYOD, etc.) in muscle and lung; MEIS1 in developing eye, brain, and muscle tissues;
172  and PAX6 in fetal eye).  In total, we identified 609 motif models matching footprinted sequences
173  (**Methods**); these models encompassed 64 distinct archetypal transcription factor recognition
174  codes (**Extended Data Table 2**), representing virtually all major DNA-binding domain families.
175  For degenerate motifs where the same sequence is recognized by many distinct TFs, we
176  observed highly cell-selective occupancy patterns that could be decomposed into coherent
177  groups corresponding with cell type and function (**Extended Data Fig. 7b-d**).

**Most DHSs encode a single regulatory topology**

179  Given that a significant fraction of DHSs are shared across two or more cell/tissue
180  types[12,29], we next asked whether differential TF occupancy within the same regulatory DNA
181  region (vs. differential actuation of entire DHSs) could be a major driver of cell-selective
182  regulation. Nucleotide-resolution DNase I cleavage provides a topological fingerprint of each
183  DHS, reflecting its unique combination and ordering of occupying TFs. Although detectable on
184  manual inspection[4], systematic analysis of differential TF occupancy has previously not been
185  possible due to dominance of intrinsic cleavage propensities when many data sets are
186  combined. To enable unbiased detection of differential footprint occupancy, we developed a
187  statistical framework to test for differences in relative DNase I cleavage rates at individual
188  nucleotides across many samples, analogous with methods developed for the identification of
189  differentially expressed genes (**Methods**). To estimate the proportion of differentially regulated
190  footprints within DHSs of a given cell/tissue type, we compared footprint occupancy within DHSs
191  broadly accessible in both nervous-system derived samples ($n$=31) with non-nervous-system
192  derived samples ($n$=212). We selected 67,368 DHSs that were highly accessible in at least 10
193  nervous and non-nervous derived samples, and for each DHS, performed a per-nucleotide
194  differential test (**Fig. 4a,b** and **Extended Data Fig. 8**). This analysis identified only a small
195  proportion of DHSs (1,720 DHSs; 2.5%) containing a differentially footprinted element (**Fig. 4c**).
196  Most of these DHSs harbored a single differentially regulated footprint, while a small fraction
197  contained 2-4 differentially occupied elements (**Fig. 4c**). Nonetheless, differentially occupied
198  footprints were significantly enriched recognition sites for known nervous system regulators
199  such as REST, NFIB, ZIC1, and EBF1 (**Fig. 4c**, bottom right and **Extended Data Fig. 9**) and
200  tissue-selective occupancy events paralleled expression of nearby genes (in the case of REST
201  occupancy) (**Extended Data Fig. 10**).

202    Collectively, the above results indicate that the vast majority of regulatory DNA regions
203    marked by DHSs encode a single structural topology reflecting a fixed pattern of footprint
204    occupancy. Nonetheless, at a small minority of elements, DHSs provide a scaffold for cell
205    context-specific TF occupancy that is typically confined to a single or small number of
206    footprinted elements.

**Functional variants in TF footprints**

208    Identifying genetic variants likely to impact regulatory function has remained
209    challenging[2]. Deep sequence coverage at DHSs enables *de novo* genotyping of regulatory
210    variants and simultaneous characterization of their functional impact on local chromatin
211    architecture by quantifying and comparing DNase I cleavage for each allele of a given
212    element[2,4]. The 243 biosamples we analyzed were derived from 147 individuals. *De novo*
213    genotyping of all samples (**Methods**) revealed 3.76 million single nucleotide variants within
214    DHSs, of which 1,656,597 were heterozygous and had sufficient read depth (≥35 overlapping
215    reads) to accurately quantify allelic imbalance.

216    Across individuals, we conservatively identified 117,626 chromatin altering variants
217    (CAVs) that impacted DNA accessibility on individual alleles (median 2.4-fold imbalance) (**Fig.
218    5a**, **Extended Data Fig. 11a** and **Methods**). Within DHSs, CAVs were markedly enriched in
219    core consensus footprints, even after controlling for the increased detection power within this
220    compartment (**Fig. 5b** and **Extended Data Fig. 12**).

221    In protein-coding regions, most functional genetic variation is expected to be deleterious,
222    with rare gain-of-function alleles[30]. Protein-DNA recognition interfaces are likewise presumed to
223    be susceptible to disruption at critical nucleotides, predisposing to loss-of-function alleles.
224    Strikingly, we found CAVs to be nearly evenly partitioned between loss- (disruption of binding)
225    and gain-of-function (increased or *de novo* binding) alleles (**Fig. 5c-d** and **Extended Data Fig.
226    10c**). Homozygosity for either the reference or alternative allele paralleled results from
227    heterozygotes and further revealed that structural changes due to TF occupancy were precisely
228    confined to the DNA sequence recognition interface (**Fig. 5c**, bottom). In many cases, SNVs
229    detected in both heterozygous and homozygous configurations showed strong agreement
230    between allelic ratios and relative footprint strength (**Fig. 5e**; Spearman's $\rho$=0.9, *p*-value < $10^{-5}$).
231    Variants residing within core recognition motifs in footprints were markedly enriched for
232    imbalance vs. non-footprinted motifs; were localized to high-information content positions within
233    the recognition interface (**Fig. 5c**, bottom and **Extended Data Fig. 13**); and paralleled the
234    predicted energetic effect of the variant on the TF binding site (**Fig. 5f** and **Extended Data Fig.
235    14**), thus providing a direct quantitative readout of functional variation impacting TF occupancy.

**DNA elements encoding footprints are hypermutable**

237    We next explored the global distribution of human genetic variation relative to consensus
238    footprints. Transcription factor binding sites appear to be gradually remodeled over evolutionary
239    time via sequential small mutations[31] that could ultimately affect function and phenotype[32].
240    However, patterns of genetic variation within regulatory DNA have not been characterized with

241 high precision. To quantify these, we calculated nucleotide diversity (π) within and around
242 consensus genomic footprints using whole-genome sequencing data compiled from >65,000
243 individual under the TOPMED project[33] (**Methods**).  Canonically, reduced levels of π reflect the
244 elimination of deleterious alleles from the population by natural selection, and hence are
245 typically indicative of functional constraint[34]. Surprisingly, we found a dramatic increase in
246 nucleotide diversity centered precisely within the core of genomic footprints (**Fig. 5a**), and thus
247 that these elements are highly polymorphic in human populations.  This result eclipses prior
248 global analyses indicating that transcription factor occupancy sites are generally not under
249 substantial purifying selection[4] both in the magnitude of the observed effect, and in its
250 nucleotide-precise localization within the core of genomic footprints.

251 The focal increase in genetic diversity within footprints indicated that the nucleotides
252 encoding footprinted elements may have an increased mutational load vs. immediately adjacent
253 sequences. To explore this possibility, we focused on variants with extremely low allele
254 frequencies in human populations (minor allele frequency < $10^{-4}$); such variants are assumed to
255 result from *de novo* germline (ie., non-segregating) mutation and are often used as a surrogate
256 for mutation rate in humans[35]. We found that the distribution of extremely rare variants around
257 and within genomic footprints mirrored that of nucleotide diversity, compatible with context-
258 driven increased mutation rate in the sequences underlying footprints (**Fig. 5b**). Of note, many
259 transcription factors favor recognition of dinucleotide combinations such as CpGs that are
260 intrinsically hypermutable. Conversely, *de novo* mutations have been implicated in the genesis
261 of TF recognition sites[36,37]. Thus, hypermutation within genomic footprints may fill a key
262 evolutionary role by favoring variability in TF occupancy and hence natural variation in gene
263 regulation.

### GWAS variants are enriched in TF footprints

265 Given the above, genetic variation within genomic footprints should, in principle, be a
266 key contributor phenotypic variation; however, to date this has defied accurate quantification.
267 We therefore next resolved the large set of variants strongly associated (nominal *p*-value <
268 $5x10^{-8}$) with diverse diseases and phenotypic traits from the NHGRI/EBI GWAS Catalogue[38] to
269 consensus genomic footprints. To account for the baseline increase in genetic variation present
270 within genomic footprints described above, we performed exhaustive (1,000x) sampling
271 matched variants (by minor allele frequency, linkage-disequilibrium (LD) structure, and distance
272 to the nearest gene) from the 1,000 Genome Project[39] (**Methods**). Additionally, we expanded
273 both GWAS catalogue and matched sampled variants to include variants in perfect LD ($r^2$=1).
274 Within DHSs, aggregated GWAS catalogue SNPs were enriched within footprints but not in non-
275 footprinted subregions, and enrichment within footprints increased monotonically with footprint
276 strength (**Fig. 6c** and **Extended Data Fig. 15**).

277 The GWAS catalogue aggregates hundreds of traits, with corresponding expected
278 diversity in cognate cell/tissue types.  To gain a more accurate view of the enrichment of trait-
279 associated variants in footprints, we compared SNP-based trait heritability of individual
280 traits[40,41].  Using summary statistic data from individual GWAS studies from the UK BioBank, we
281 applied partitioned LD-score regression to compute the relative heritability contribution of

282    variants within all DHSs and footprints collectively vs. DHSs and footprints therein from the
283    expected cognate cell type for the trait (**Fig. 6d-e**). This analysis revealed striking enrichment of
284    variants that account for trait heritability in footprints generally (>5-fold) and most prominently in
285    footprints from the cognate cell type (up to >45-fold) (**Fig. 6d-e**).  Collectively, we conclude that
286    the genetic signals from disease- and trait-associated variants within DHSs emanate from
287    transcription factor footprints, and that variants within footprints are major contributors to trait
288    heritability.

289    **Discussion**

290        Our report details the highest resolution view to date of regulatory factor occupancy
291    patterns on the human genome measured across an expansive range of cell and tissue
292    contexts sampled from >140 genotype backgrounds.  The scale and breadth of the data have
293    enabled the delineation of a reference set of ~4.5 million genomic sequence elements that
294    collectively define nucleotides critical for genome regulation and function and form the building
295    blocks of regulatory DNA. These footprints now provide a ready and extensible nucleotide-
296    precise reference for diverse analyses, particularly those involving genetic variation.

297        The preferential localization of disease- and trait-associated variation within regulatory
298    DNA has heretofore been described in terms of entire regulatory regions demarcated by DHSs
299    or clusters thereof.  Our results now show that genetic association and heritability signals from
300    regulatory DNA overwhelmingly emanate from indexed transcription factor footprints, which
301    should greatly facilitate the connection of disease- and trait-associated genetic variation with
302    genome function.

303        Perhaps most strikingly, we report that human genetic variation is itself concentrated
304    within transcription factor footprints, owing apparently to a combination of mutation propensity
305    and the evolved sequence recognition repertoire of human transcription factors, which favors
306    hypermutable nucleotide combinations (e.g., CG dinucleotides).  Given that human and mouse
307    TFs share the large majority of their recognition landscapes, concentration of variation within TF
308    occupancy sites has likely played a considerable role in shaping human – and indeed all
309    mammalian – gene regulation. It implies, furthermore, that the genome is heavily primed for
310    regulatory evolution, providing a possible mechanism underlying facilitated phenotypic
311    evolution[42]

## References

1. Maurano, M. T. *et al.* Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).
2. Maurano, M. T. *et al.* Large-scale identification of sequence variants influencing human transcription factor occupancy in vivo. *Nat. Genet.* **47**, 1393–1401 (2015).
3. Hesselberth, J. R. *et al.* Global mapping of protein-DNA interactions in vivo by digital genomic footprinting. *Nat. Methods* **6**, 283–289 (2009).
4. Neph, S. *et al.* An expansive human regulatory lexicon encoded in transcription factor footprints. *Nature* **489**, 83–90 (2012).
5. Vierstra, J. & Stamatoyannopoulos, J. A. Genomic footprinting. *Nat. Methods* **13**, 213–221 (2016).
6. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human epigenomes. *Nature* **518**, 317–330 (2015).
7. Galas, D. J. & Schmitz, A. DNAse footprinting: a simple method for the detection of protein-DNA binding specificity. *Nucleic Acids Res.* **5**, 3157–3170 (1978).
8. Dynan, W. S. & Tjian, R. The promoter-specific transcription factor Sp1 binds to upstream sequences in the SV40 early promoter. *Cell* vol. 35 79–87 (1983).
9. Neph, S. *et al.* Circuitry and dynamics of human transcription factor regulatory networks. *Cell* **150**, 1274–1286 (2012).
10. Stergachis, A. B. *et al.* Conservation of trans-acting circuitry during mammalian regulatory evolution. *Nature* **515**, 365–370 (2014).
11. Lazarovici, A. *et al.* Probing DNA shape and methylation state on a genomic scale with DNase I. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 6376–6381 (2013).
12. Meuleman, W., Muratov, A., Rynes, E., Halow, J. & Lee, K. Index and biological spectrum of accessible DNA elements in the human genome. *BioRxiv* (2019).
13. Jolma, A. *et al.* DNA-binding specificities of human transcription factors. *Cell* **152**, 327–339 (2013).
14. Yin, Y. *et al.* Impact of cytosine methylation on DNA binding specificities of human transcription factors. *Science* **356**, (2017).
15. Najafabadi, H. S. *et al.* C2H2 zinc finger proteins greatly expand the human regulatory lexicon. *Nat. Biotechnol.* **33**, 555–562 (2015).
16. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **175**, 598–599 (2018).
17. Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res.* **46**, D1284 (2018).
18. Kulakovskiy, I. V. *et al.* HOCOMOCO: towards a complete collection of transcription factor binding models for human and mouse via large-scale ChIP-Seq analysis. *Nucleic Acids Res.* **46**, D252–D259 (2018).
19. Panne, D., Maniatis, T. & Harrison, S. C. An atomic model of the interferon-beta enhanceosome. *Cell* **129**, 1111–1123 (2007).
20. Rohs, R. *et al.* The role of DNA shape in protein-DNA recognition. *Nature* **461**, 1248–1253 (2009).
21. Yin, M. *et al.* Molecular mechanism of directional CTCF recognition of a diverse range of genomic sites. *Cell Res.* **27**, 1365–1377 (2017).

22. Arnold, R., Burcin, M., Kaiser, B., Muller, M. & Renkawitz, R. DNA bending by the silencer protein NeP1 is modulated by TR and RXR. *Nucleic Acids Res.* **24**, 2640–2647 (1996).

23. MacPherson, M. J. & Sadowski, P. D. The CTCF insulator protein forms an unusual DNA structure. *BMC Mol. Biol.* **11**, 101 (2010).

24. Xu, H. E. *et al.* Crystal structure of the human Pax6 paired domain-DNA complex reveals specific roles for the linker region and carboxy-terminal subdomain in DNA binding. *Genes & Development* vol. 13 1263–1275 (1999).

25. Svaren, J., Klebanow, E., Sealy, L. & Chalkley, R. Analysis of the competition between nucleosome formation and transcription factor binding. *J. Biol. Chem.* **269**, 9335–9344 (1994).

26. Mirny, L. A. Nucleosome-mediated cooperativity between transcription factors. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 22534–22539 (2010).

27. Wunderlich, Z. & Mirny, L. A. Different gene regulation strategies revealed by analysis of binding motifs. *Trends Genet.* **25**, 434–440 (2009).

28. Vierstra, J. *et al.* Mouse regulatory DNA landscapes reveal global principles of cis-regulatory evolution. *Science* **346**, 1007–1012 (2014).

29. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature* **489**, 75–82 (2012).

30. Fu, W. *et al.* Analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants. *Nature* **493**, 216–220 (2013).

31. Payne, J. L. & Wagner, A. The robustness and evolvability of transcription factor binding sites. *Science* **343**, 875–877 (2014).

32. Prud'homme, B. *et al.* Repeated morphological evolution through cis-regulatory changes in a pleiotropic gene. *Nature* **440**, 1050–1053 (2006).

33. Taliun, D., Harris, D. N., Kessler, M. D. & Carlson, J. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *BioRxiv* (2019).

34. Kellis, M. *et al.* Defining functional DNA elements in the human genome. *Proc. Natl. Acad. Sci. U. S. A.* **111**, 6131–6138 (2014).

35. Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation rate heterogeneity in humans. *Nat. Commun.* **9**, 3753 (2018).

36. He, X. *et al.* Methylated Cytosines Mutate to Transcription Factor Binding Sites that Drive Tetrapod Evolution. *Genome Biol. Evol.* **7**, 3155–3169 (2015).

37. Zemojtel, T., Kielbasa, S. M., Arndt, P. F., Chung, H.-R. & Vingron, M. Methylation and deamination of CpGs generate p53-binding sites on a genomic scale. *Trends Genet.* **25**, 63–66 (2009).

38. Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).

39. 1000 Genomes Project Consortium *et al.* A map of human genome variation from population-scale sequencing. *Nature* **467**, 1061–1073 (2010).

40. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).

41. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).

42. Gerhart, J. & Kirschner, M. The theory of facilitated variation. *Proc. Natl. Acad. Sci. U. S. A.* **104 Suppl 1**, 8582–8589 (2007).

**Figure Legends**

**Figure 1: A nucleotide resolution atlas of transcription factor occupancy within the human genome.**

**a,** DNase I cleavage at regulatory DNA elements within the *RELB* locus in CD8+ T cells. Top, windowed DNase I cleavage density. Below, per-nucleotide cleavage counts and genomic footprint posterior probabilities at two DNase I hypersensitive sites in CD8+ T cells. **b,** Heatmap of genomic footprint posterior probabilities computed by integrating 243 datasets within the DHSs. Rows and columns correspond to individual biosamples and nucleotides grouped tissue/organ systems. Black ticks left of heatmap indicate whether region is DHS in each sample. Below, genome sequence scaled by footprint prevalence. Grey boxes define consensus footprints present in one or more cell and/or tissue types (footprint posterior>0.99). **c**, A consensus map of TF occupancy was derived from 243 cell and tissues covering 1.6 million DHS results in a comprehensive annotation of cis-regulatory DNA. **d**, Histogram of footprint location relative to DHS peak summit. Dashed red lines represent the average size of a DHS peak (203 bp).

**Figure 2: DNase I footprints reflect the topological structure of individual TF:DNA interactions**

**a**, Top, Physical structure of CTCF zinc fingers 3-11 bound to its cognate DNA recognition sequence (PDB: 5YEF and 5YEL)[21]. DNA is colored by mean ratio of observed vs. expected cleavages at footprinted CTCF motifs in T regulatory cells. Left, heatmap of a relative cleavage at each of the 25,852 footprinted CTCF motifs. Bottom, aggregate DNase I cleavage summed over all CTCF footprints. Right, DNase I cleavage (observed and expected) at three randomly selected footprints. **b**, Same as **a** for paired-box transcription factor PAX6 (PDB: 6PAX)[24].

**Figure 3: Distinguishing modes of transcription factor occupancy within regulatory DNA**

 **a,** The width of footprints for diverse TFs tightly correlates with the width of their recognition sequence (Spearman's $\rho$=0.82, $p$-value=0.001). **b**, Overlap and spatial enrichment of TF recognition sequences within footprints binned by width. Left, Proportion of footprints uniquely overlapped by 0, 1 or 2 or more recognition sequences. Right, density heatmap of motif occurrences around footprints binned by width. **c,** Percentage of footprints that likely represent the occupancy of a single TF (≤30bp) or multiple TFs (>30bp). **d–e,** Footprint density and footprint spacing (distance edge-to-edge) vs. average DNase I density within DHSs. Grey indicates the middle 50%-ile. **f,** A typical regulatory element (DHS) harbors ~5-6 directly bound TFs spaced roughly 20-bp from each other.

**Figure 4: Comparative footprinting identifies cell-selective TF occupancy at nucleotide resolution**

**a,** Comparative footprinting within the *SCAMP5* promoter identifies 3 footprints differentially occupied in nervous cell and tissue types. Top, DNase I cleavage in two exemplar nervous and non-nervous cell types. Bottom, mean differential per nucleotide cleavage(log$_2$) between nervous-system derived (*n*=26) and non-nervous samples (*n*=12). The color of each bar indicates the statistical significance (–log$_{10}$ *p*) of the per-nucleotide differential test. **b,** Differential footprint testing within thousands of DHS accessible in between nervous and non-nervous related biosamples. **c,** The vast majority of tested DHSs encode a single TF binding topology. Top, percentage of the DHSs tested that containing one or more differentially occupied element. Bottom left, distribution of differentially footprinted elements per DHS. Bottom right, selected TF recognition sequences significantly enriched in differentially occupied footprints (binomial test p<0.01). Indicated in parenthesis is the fold-enrichment vs. expected (based on prevalence of footprinted motif in tested regions).

**Figure 5: Functional genetic variation localizes to TF footprints**

**a**, Allelic imbalance was assessed at all variants overlapping a DNase I footprints (consensus footprint probability < 0.01). **b**, Percentage of variants imbalanced in DNase I footprints and DHS peaks (but not in footprints). **b**, Variant rs10171498 results in the gain of a NFIX footprint. Top, allelically resolved per-nucleotide DNase I cleavage aggregated from 56 heterozygous samples. Middle, DNase I cleavage in two selected samples homozygous for either reference or alternative alleles. Bottom, mean differential per nucleotide cleavage (log$_2$) between homozygous reference (*n*=74) and alternative samples (*n*=12). The color of each bar indicates the statistical significance (–log$_{10}$ *p*) of the per-nucleotide differential test (**Methods**). The variant and differentially footprinted nucleotides precisely colocalize to a NFIX recognition element. **d**, Density histogram of allelic ratios for variants overlapping a footprinted NFIX recognition sequence. Grey, all variants tested for imbalance (*n*=7,110). Blue, all variants significantly (n=1,889) imbalanced. **g**, Scatterplot of allelic imbalance computed from heterozygous individuals (*x*-axis) vs. the relative difference in footprint depth between homozygous individuals at variants overlapping an NFIX footprint. Each point pertains to a SNV within a footprinted NFIX binding site imbalanced in heterozygotes and differentially footprinted in homozygotes. Grey line indicates fit linear model. **e**, Allelic imbalance measurements parallels predicted energetic effects of variants within NFIX footprints. Shown is the mean log-odds motif score (reference vs. alternate allele) of all tested variants within footprinted motifs binned by allelic ratios.

**Figure 6: Human genetic variation is broadly enriched within genomic footprints**

**a,** Distribution of genetic variation with respect to consensus footprints. Plotted is the mean per nucleotide diversity determined from whole genome sequencing of 62,784 individuals (TOPMED project). **b,** Histogram of the distribution of rare variation (minor allele frequency<0.0001) within and surrounding genomic footprints. **c,** Enrichment of GWAS variation within or outside consensus genomic footprints over randomly sampled variants from 1,000 Genome Project. Enrichment was computed after expanding both GWAS and sampled variants with those in perfect LD ($r^2$=1.0, central European population). **d,** Enrichment of SNP-based trait heritability using LD-score regression for UK BioBank GWAS traits lymphocyte counts (**c**) and red blood cell counts (**d**). Asterisk denotes statistically significant enrichments (* indicates p-value<0.01).
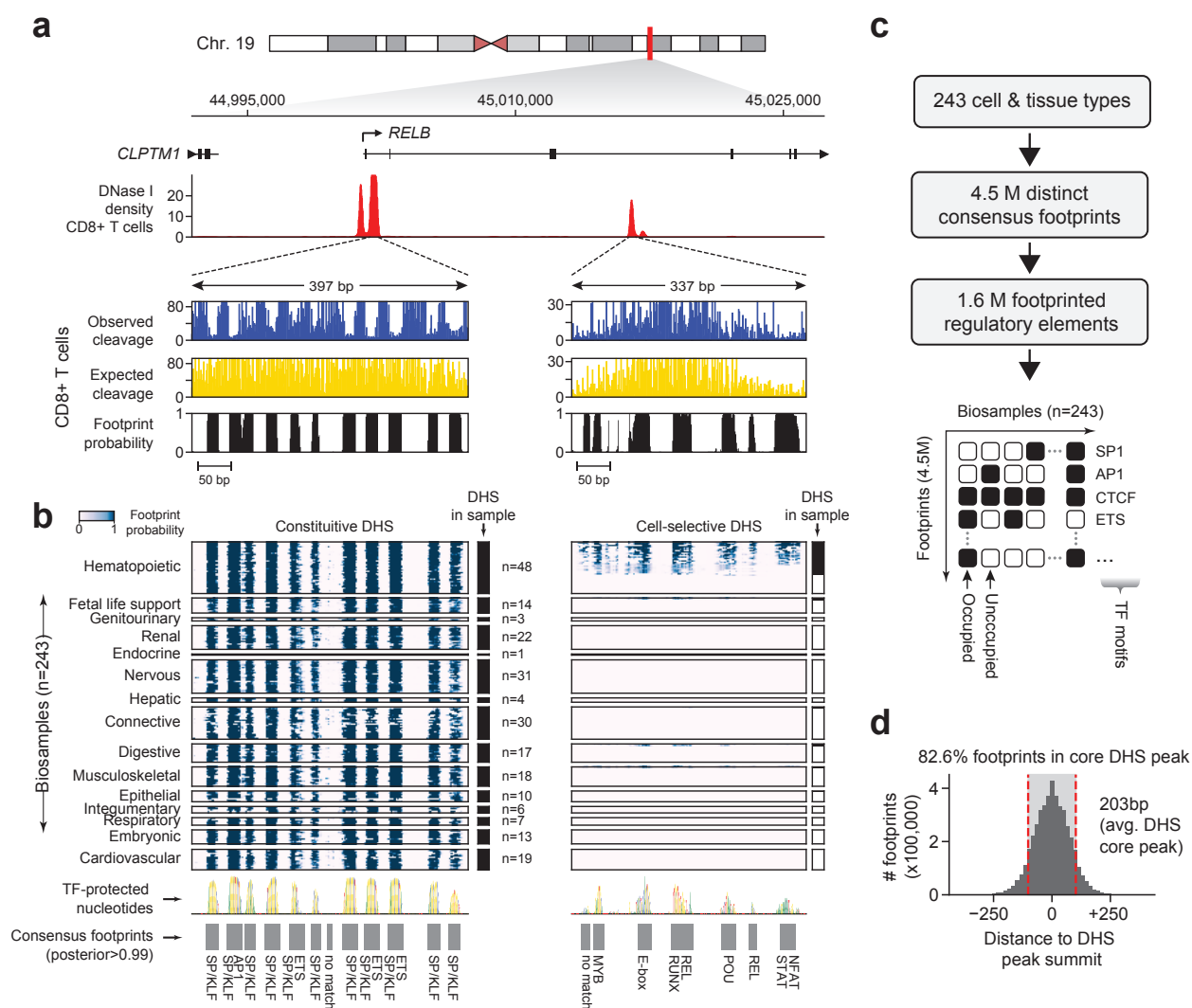
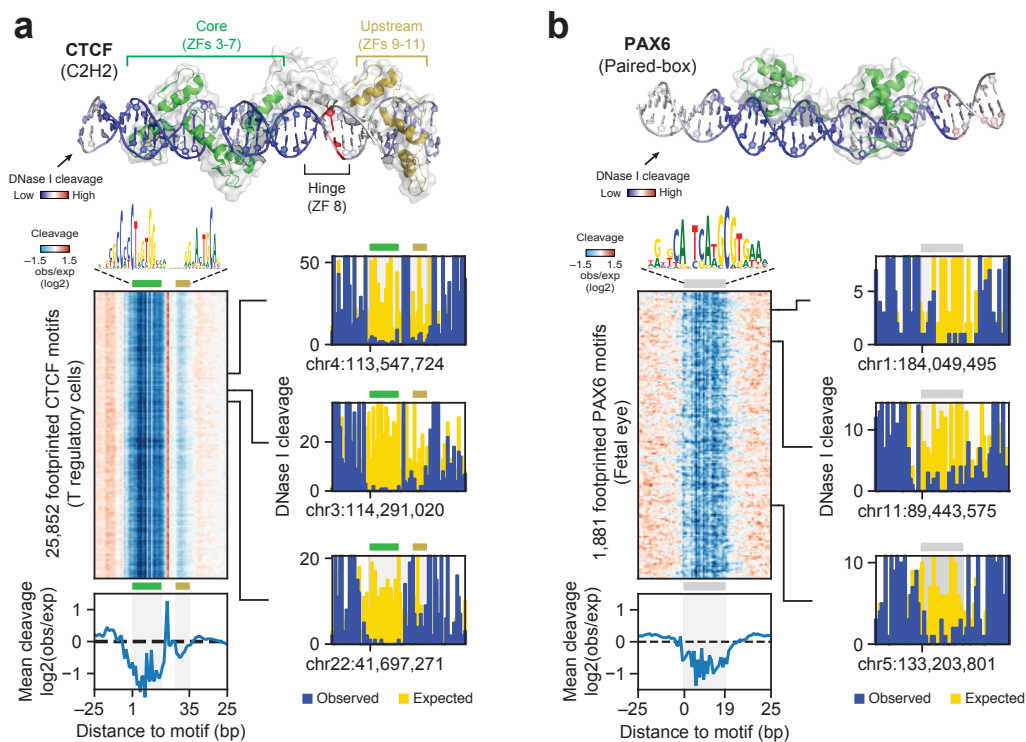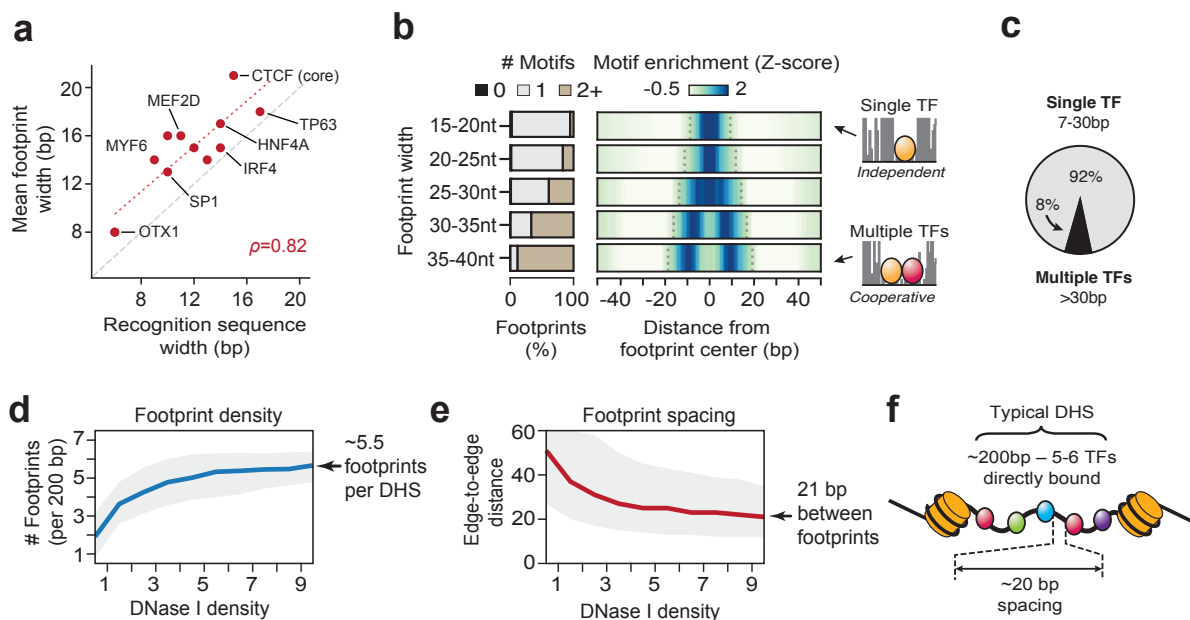# Figure 1

# Figure 2

# Figure 3

# Figure 4

# Figure 5

# Figure 6