

SUPPLEMENTARY MATERIALS

Materials and Methods

Figures S1-S13

Tables S1-S3

References 37-63

MATERIALS AND METHODS

Genomic DNaseI Footprinting

41 of the cell types used in this study were analyzed in Neph et al 2012, and 40 are novel to this study. These 40 novel diverse human cell types were subjected to DNaseI digestion and high-throughput sequencing, following previous methods (37–39) (**Table S1**). To generate genomewide per-nucleotide DNaseI cleavage profiles tags were aligned to the reference genome, build GRCh37/hg19 using Bowtie (40), version 0.12.7 with parameters: *--mm -n 3 -v 3 -k 2*, and *--phred33-quals* for Illumina HiSeq sequencer runs or *--phred64-quals* for Illumina GAI sequencer runs and the 5' ends of the aligned sequencing tags at each position along the genome were summed. Data from additional cell types were utilized from Neph et al. 2012 (13) (**Table S1**). FDR 1% DNase I footprints were identified in each cell type as previously described (13). DNaseI footprints found in any of the 81 cell types were identified using the BEDOPS command *bedops -m* on each of the individual cell-type DNaseI footprint files (41).

Targeted exome footprinting

The two cell types HMF and AG10803 were DNaseI digested and proceeded on to Illumina PE library construction following the previous methods described above for DNaseI footprinting. The DNaseI libraries was amplified by PCR following the Exon Capture SureSelect protocol recommendations with minimal amount of PCR cycles and then purified using Agencourt AMPure XP beads (Beckman Coulter Genomics). Five hundred nanograms of each library was hybridized to Agilent SureSelect Human All Exon Kit (50 Mb) for 24 h at 65 C. The biotinylated probe/target hybrids were captured on DynalMyOne Streptavidin T1 (Invitrogen), washed, eluted, and desalted and purified on a MinElute PCR column (Qiagen) as described in the SureSelect protocol. The eluted captured library was amplified by PCR with minimal amount of PCR cycles. Amplified exon captured libraries were purified using AngencourtAMPure XP beads the samples were then quantified by QubitdsDNA assay (Invitrogen). Samples were diluted to a working concentration of 10 nM. Cluster generation was performed for each sample and loaded on to single lane of an Illumina HiSeq flowcell. Paired end sequencing was performed for 36 cycles according to manufacturer's instructions.

Overlap of DNaseI footprints with coding sequence

Genomic regions annotated as coding were identified using the Consensus CDS (CCDS) (Release 6) (42). Each transcript within this database having two or more exons was utilized to identify first, internal and final coding exons, with first coding exons, by definition, always containing the methionine start codon, and final coding exons, by definition, always containing the stop codon. The density of footprints within exons was computed by first calculating the number of DNaseI footprints within a cell type that overlap coding sequence by at least 50% and dividing that by the total number of mappable bases within the coding sequence. The density of DNaseI footprints surrounding splice acceptor sites, splice donor sites, start codons and stop

codons was calculated by summing the number of DNaseI footprints that overlap each base surrounding these genomic features in each of the 81 cell types.

Generation and analysis of genome-wide methylation patterns

Whole genome methylation data was collected on fetal small intestine, as previously described. Alignment and methylation calling of bisulfite-seq reads was performed using Bismark (43). Methylation calls from Bismark were converted to BED format for subsequent comparisons using the BEDOPS tools (41). The methylation level of each genomic CpG was calculated as the number of distinct fragments representing an unconverted cytosine on either strand, divided by the total number of distinct fragments overlapping the CpG.

Allelic chromatin imbalance measurements

Allelic chromatin imbalance was calculated as previously described (13). Briefly, known autosomal single nucleotide variants (SNVs) were downloaded from the 1000 Genomes Project and SNVs mapping within 36 bases of each other were removed. DNaseI-seq reads overlapping a SNV were counted only if they contained no more than one mismatch, excluding the SNV, and duplicate DNaseI-seq reads were removed from this analysis. To test for allelic chromatin imbalance we first combined DNaseI cleavage reads overlapping a particular SNV in all cell types heterozygous at that SNV. The difference in DNaseI cleavage reads containing each of the two alleles was tested using a two-sided binomial test, with SNVs having $P < 0.01$ considered as showing significant chromatin allelic imbalance. To test whether two sets of SNVs significantly differed in the proportion showing allelic imbalance, a read depth distribution was derived from each set, and the intersection was determined to generate a read-depth-matched random sample as large as possible. At each particular read depth, all sites from the set with fewer instances of that depth were included, and a random sample without replacement was taken from the set with more instances. Finally, we counted sites in each set showing allelic imbalance with two-sided binomial test $P < 0.01$. The difference between these counts was tested for significance with a one-sided Fisher's exact test. The SIFT (44) and PolyPhen-2 (45) scores at each non-synonymous SNV were determined using Ensembl's Variant Effect Predictor(46).

Distribution of common disease- and trait-associated lead SNPs within DNaseI footprints

Using HapMap CEU SNPs for a null model, we found the GWAS SNPs to be significantly enriched in both coding ($P < 10^{-11}$, binomial) and noncoding ($P < 10^{-26}$, binomial) regions of DNase I footprints. For these calculations, we used release 27 of HapMap (47), which contained 4,029,798 SNPs. 5,873 of these lie within coding regions of footprints, while 34 of the 5,386 unique GWAS SNPs lie in these regions. We used the R function `pbinom(x; n, p)`, which returns $\text{Prob}(X > x)$, to compute the P-value; because we want $\text{Prob}(X \geq x)$, we used `pbinom` with $x = 33$, $n = 5,386$, and $p = 5,873/4,029,798$. 373,006 HapMap CEU SNPs lie within noncoding regions of footprints, as do 744 GWAS SNPs. For the corresponding P-value, we used `pbinom` with $x = 743$, $n = 5,386$, and $p = 373,006/4,029,798$.

Transcription factor occupancy at SNPs in linkage with rs495337

Genome wide association study (GWAS) lead SNPs were downloaded from Maurano et al 2012 (48). The synonymous coding variant rs495337, which is associated with Psoriasis (49, 50), was identified as overlapping DNaseI footprints in NB4 cells using BEDOPS (41). To identify other SNPs linked with rs495337 that may also be contributing to the psoriasis

association signal we utilized whole genome sequencing data from 267 individuals of Northern and Western European (CEPH), Finnish (FIN) and English and Scottish (GBR) ancestry (51), corresponding to the geographic regions used in the original GWAS studies. The 98 SNPs with a genotype R-squared greater than 0.8 were used for further analysis. The allele at each of these SNPs in linkage with the rs495337 psoriasis risk allele (G allele) was also determined using the whole genome sequencing data from these 267 individuals. Of these 99 SNPs (including rs495337), 14 overlapped a DNaseI footprint in at least one of the 81 cell types (**Table S3**). In total three of these SNPs were associated with allelic chromatin imbalance at the overlapping DNaseI footprints (rs495337 in NB4 cells; rs492702 in NB4 cells; and rs2281217 in RPMI cells) and two of these were synonymous coding variants (rs495337 and rs492702).

Distribution of DNaseI footprinted TF recognition sequences

Human genome build hg19 was scanned for predicted TRANSFAC (52), JASPAR Core (53) and UniPROBE (54) motif-binding sites using FIMO (55), version 4.6.1, with a maximum p value threshold of 10^{-5} and defaults for other parameters. We marked a putative binding site as being occupied within a cell type if it overlapped a DNaseI footprint within that cell types by at least 3 nt, as previously described in (13). The distribution of DNaseI footprinted TF binding elements within and surrounding coding exons was computed by calculating the distance of the mid-point of the TF binding element to the beginning of the overlapping or neighboring coding exon, relative to the length of that coding exon. TF recognition sequence motif logos were generated using Weblogo 3 and all of the exonic binding sites for that TF (56). Start codon, stop codon, splice acceptor site and splice donor site motifs were generated by aligning all relevant genomic elements using CCDS coding sequence annotations and calculating base enrichments using Weblogo 3 (56).

Protein domain architecture overlapping TF binding sites

The position of NRSF binding sites (JASPAR Core model MA0138.2) relative to the codon frame within the first exon was analyzed using Consensus CDS (CCDS) database gene models (Release 6). Possible alignments included; (1) the first frame of the coding strand; (2) the second frame of the coding strand; (3) the third frame of the coding strand; (4) the first frame of the template strand; (5) the second frame of the template strand; and (6) the third frame of the template strand. Of the 382 footprinted NRSF recognition sequences within the first exon, 253 were found aligned with the third frame of the coding strand. The protein domain architecture overlapping these binding elements were analyzed using; (1) signal peptide domain predictions from the SignalP 4.1 Server (57); (2) transmembrane domain predictions from the TMHMM Server v. 2.0 (58); (3) leucine-rich nuclear export signal predictions from the NetNES 1.1 Server (59); and (4) the Superfamily database of structural and functional protein domain annotations (60). Logos of the frequency of amino acids overlapping TF recognition sequences were generated using Weblogo 3 (56).

For calculation of CTCF splice acceptor site conservation, we calculated the average phyloP score at the 10 bases immediately upstream (relative to reading frame) of all internal exons, or internal exons with splice acceptor sites overlapping footprinted CTCF binding elements. To test if the differences in average phyloP conservation scores for those splice sites overlapping a footprinted CTCF binding element were likely to be observed by sampling error alone, we ran the same calculation on sets of randomly-selected splice sites one million times. On each trial, the same number of splice sites as were observed actually overlapping a

footprinted CTCF binding element were drawn randomly without replacement from the total set of splice sites, and the average phyloP score of bases within those sampled sites was calculated. If the absolute difference between this sample mean and the mean phyloP score for all splice sites exceeded the difference observed on real overlaps, a trial was counted as a hit. The number of such random hits divided by the total number of trials estimates the probability that we could have observed differences of at least that magnitude if there were no relationship between footprinted binding elements and conservation within splice sites. For SREBP1 occupied splice donor sites a similar strategy was used except it focused only on the 10 bases immediately downstream (relative to reading frame) of all first coding exons, as well as first coding exons with splice donor sites overlapping footprinted SREBP1 binding elements.

Evolutionary constraint at footprinted coding sequences

Evolutionary constraint at TF binding sites was calculated using phyloP evolutionary conservation scores (61). 4-fold degenerate bases were identified based on the sequence features of each codon (e.g. the third position of the following codons: CTA, CTT, CTG, CTC, GTA, GTT, GTG, GTC, TCA, TCT, TCG, TCC, CCA, CCT, CCG, CCC, ACA, ACT, ACG, ACC, GCA, GCT, GCG, GCC, CGA, CGT, CGG, CGC, GGA, GGT, GGG, GGC). Non-degenerate bases were identified based on the sequence features of each codon (e.g. the first and second position of every codon except TTA, TTG, CTA, CTT, CTG, CTC, AGT, AGC, TCA, TCT, TCG, TCC, AGA, AGG, CGA, CGT, CGG, CGC. And the second position of TTA, TTG, CTA, CTT, CTG, CTC, AGA, AGG, CGA, CGT, CGG, CGC). To generate the conservation profile of a TF within exons, we calculated the average phyloP at all 4-fold degenerate bases, or non-degenerate bases, overlapping each position within the TF binding element. Only TFs with 20 or more data points contributing to each position within the binding element were used for further analysis. The number of bases contributing to each position within the binding element is shown in **Fig. S7**. To generate the conservation profile of a TF within promoters, a similar process was performed for all bases overlapping TF binding elements within non-coding promoter regions. Pearson correlations were calculated to determine the similarity of the conservation profile of a TF at promoter elements and coding 4-fold degenerate and coding non-degenerate sites.

Mutational age at footprinted coding sequences

Exome sequences were obtained from 6,515 individuals (4,298 of European American ancestry and 2,217 of African American ancestry) (21). Coding variants, specifically synonymous variants and nonsynonymous variants, were classified according to whether they overlapped a DNaseI footprint in any of the 81 tested cell types. Average mutation age for each category was calculated as previously described (21). Briefly, mutation age was estimated based on a derivation of Griffiths and Tavaré(62) by generating a series of coalescent trees under a specified demographic model for European and African American populations (63). Average mutation age across variants for each category was defined as a weighted average of mutation age, where the weights are calculated according to the site-frequency-spectrum (SFS) in this category. Average mutation age in different categories was compared through permutations to identify significant differences (21).

Codon usage biases and TF footprint trinucleotide frequencies

Coding usage biases were obtained using CCDS gene annotations downloaded from the UCSC genome browser, corresponding to human build GRCh37/hg19 or mouse build

NCBI37/mm9. Individual codon locations were parsed into BED format, excluding start codons and any codons overlapping a splice site or that were ambiguous due to overlapping annotations in different reading frames. Coding annotations containing one or more internal stops in the reference sequence were also excluded. Overlaps of codon locations with footprint calls were determined using BEDOPS. Codons that partially overlapped a footprint were excluded. Non-coding trinucleotide frequencies were obtained using the genomic space uniquely mappable by 36-mer sequencing tags. CCDS coding exons, as well as RepeatMasker annotations also downloaded from the UCSC genome browser, were then subtracted from this space using BEDOPS, and the remaining regions divided by overlap with footprint calls. Finally, all reference-strand genomic 3-mers in the footprint or non-footprint space were tabulated separately.

SUPPLEMENTARY FIGURES

Figure S1. Genomic distribution of DNaseI footprints (A) Shown is the average genomic distribution of DNaseI footprints across all 81 cell types. (B) Histogram showing for each coding base along the genome the number of cell types in which that base overlaps a DNaseI footprint. Y-axis is log-10 scale. (C) The proportion of genes containing DNaseI footprints in each of the 81 cell types studied, or in any of the cell types studied. (D) Histogram showing the number of coding DNaseI footprints per gene. (E) The percentage of coding bases occupied by DNaseI footprints from conventional and targeted DNaseI footprinting.

Figure S2. DNaseI footprints identified using additional cell types Total number of coding bases overlapping DNaseI footprints identified using the published datasets.

Figure S3. Sensitivity of coding DNaseI footprints using capture DNaseI-seq. (A) Summary of Capture DNaseI-seq method. (B) Per-nucleotide vertebrate conservation as well as per-nucleotide DNaseI-seq and capture DNaseI-seq cleavage patterns at coding binding elements for NFIC, CTCF, REST, YY1 and NRF1. (C-D) *Capture DNaseI-seq enables extensive DNaseI footprint identification and superior quantification.* (C) The average number of sequenced DNaseI cleavages surrounding DNaseI footprints using DNaseI-seq and Capture DNaseI-seq data. (D) The average depth of DNaseI footprints using DNaseI-seq and Capture DNaseI-seq data. Note that Capture DNaseI-seq enables the more precise quantification of DNaseI footprints.

Figure S4. TFs preferentially occupy coding bases from expressed genes. (A) The average density of DNaseI footprints within coding sequence and outside of coding sequence. (B) *Long genes contain more TF footprints than short genes.* Box-and-whisker plots showing the association of coding gene length with the number of DNaseI footprints within that coding sequence. R-squared and p-values are from a linear regression of coding gene length vs. the number of DNaseI footprints within that coding sequence. (C-D) *Transcription factors preferentially populate highly expressed genes.* (C) Shown is a box-and-whiskers plot of the gene expression in HMVEC_dBNeo cells for genes with 0, 1-4 and 5+ coding DNaseI footprints. (D) Shown is the correlation of exonic footprints count with gene expression in 47 cell types with DNaseI footprint calls and gene expression data.

Figure S5. TF binding elements impart evolutionary constraint on coding sequence. (A) The percentage of 4-fold degenerate bases above a minimum phyloP conservation level that overlap DNaseI footprints. (B) Average phyloP conservation at 4-fold degenerate (left) and non-degenerate (right) exonic bases within DNaseI footprints found any of the 81 cell types, and within exons harboring DNaseI footprints in any of the 81 cell types, yet outside of the actual DNaseI footprint. P-values were calculated using Wilcoxon rank sum two-sided tests. (C) 4,298 sequenced exomes from individuals of European ancestry were utilized to identify SNVs overlapping DNaseI footprints in any of the 81 cell types. (D) 2,217 sequenced exomes from individuals of African American ancestry were utilized to identify SNVs overlapping DNaseI footprints in any of the 81 cell types. (E) The average mutational age at all (grey), synonymous (brown) and nonsynonymous (red) African American coding SNVs identified within and outside of DNaseI footprints. Mutational ages and p-values were calculated as before (21).

Figure S6. Transcription factors influence codon choice. (A-B) Per-nucleotide phyloP vertebrate conservation and DNaseI cleavage plots at (A) a non-coding and (B) a coding NFIC regulatory element. Note that NFIC imparts a stereotyped pattern of evolutionary constraint when bound at non-coding regulatory elements. (C) Average per-nucleotide conservation profile at footprinted binding elements for ZNF219 (left), REST (second), NFKB (third), CTCF (fourth) and MYF (right) overlapping non-coding bases within promoters (blue), 4-fold degenerate coding bases (brown) and non-degenerate coding bases (red). Pearson correlation values (r) between conservation profiles at promoter bases and 4-fold degenerate bases (top) or non-degenerate bases (bottom) are shown in the upper right corner of each plot.

Figure S7. 4-fold degenerate and non-degenerate bases overlapping DNaseI footprinted TF elements. (A-G) Shown are the number of bases overlapping different positions within footprinted (A) KLF4, (B) NFKB, (C) CTCF, (D) MYF, (E) NFIC, (F) ZNF291, and (G) REST binding elements in any of the 81 cell types. Bases overlapping binding elements are broken into; (left/blue) promoter element bases; (middle/brown) 4-fold degenerate bases; and (right/red) non-degenerate bases.

Figure S8. TF sequence preferences and codon usage biases in *M. musculus*. (A) Comparison of global codon usage preferences in *H. sapiens* and *M. musculus*. (B) Comparison of the TF trinucleotide preferences in *H. sapiens* and *M. musculus*. *H. sapiens* trinucleotide preferences are derived from trinucleotides preferentially localized within non-coding DNaseI footprints in *H. sapiens* B-cells. *M. musculus* trinucleotide preferences are derived from trinucleotides preferentially localized within non-coding DNaseI footprints in *M. musculus* B-cells.

Figure S9. TFs are influenced by and can exploit coding features of exons (A) *NFYA*, *AP2* and *SP1* preferentially avoid binding within coding sequence, start codons and splice junctions. The density of (top) *NFYA*, (middle) *AP2* and (bottom) *SP1* DNaseI footprints relative to first, middle and final coding exons. Coding sequence is colored in purple. (B) *NRSF* binding elements preferentially align to the coding strand at the third frame of the codon and exploit start codons. (top) The density of *NRSF* DNaseI footprints relative to first coding exons. (bottom-left) Shown is the *NRSF* motif model as well as a logo of the amino acid sequence at all occupied coding strand *NRSF* binding elements that overlap a start codon. (bottom-right) Shown is the number of *NRSF* binding elements within first coding exons that align to the three different coding positions along either the coding or template strand. (C) Average evolutionary constraint (phyloP) at 4-fold degenerate bases within first coding exon transmembrane (TM) domains that either overlap (top), or do not overlap (bottom) a footprinted *NRSF* binding element (p-value calculated using Wilcoxon rank sum two-sided tests).

Figure S10. TF occupancy at stop codons and splice sites reflects global evolution in TF preferences (A) *CTCF* binding elements exploit splice acceptor sites. (top) The density of *CTCF* DNaseI footprints relative to middle coding exons. (bottom) Shown is the *CTCF* motif model in comparison with the splice acceptor site motif model. (B) *SREBP1* binding elements exploit splice donor sites. (top) The density of *SREBP1* DNaseI footprints relative to first coding exons. (bottom) Shown is the *SREBP1* motif model in comparison with the splice acceptor site motif model. (C) Average evolutionary constraint (phyloP) of the 10 bp non-coding portion of splice

acceptor sites for internal exons, as well as those that overlap footprinted CTCF binding elements (p-value empirically calculated by resampling 1 million times). (D) Average evolutionary constraint (phyloP) of the 10 bp non-coding portion of splice donor sites for internal exons, as well as those that overlap footprinted SREBP1 binding elements (p-value empirically calculated by resampling 1 million times). (E-F) *Transcription factors preferentially avoid occupying splice sites and stop codons.* (E) Shown is the density of DNaseI footprints surrounding start codons, splice donor sites, splice acceptor sites and stop codons. Sequence features of these elements are displayed as motif models and coding sequence is colored in purple. (F) The frequency of the stop codon trinucleotides TAA, TAG and TGA within and outside of non-coding DNaseI footprints.

Figure S11. TF occupancy within coding sequence is modeled by CpG methylation (A) The difference in the percentage of annotated coding and non-coding transcription factor binding elements overlapping a DNaseI footprint in each cell type. Positive values indicate that a greater fraction of coding binding elements are occupied as compared to non-coding elements. (B) The difference in the percentage of CpGs methylated within annotated coding and non-coding transcription factor binding elements in each cell type. Positive values indicate that a greater fraction of CpGs are unmethylated in coding binding elements as compared to CpGs in non-coding elements. (C) Shown is a scatter plot of the preference of 232 TFs for occupying coding vs. non-coding binding elements (y-axis) and being CpG methylated at coding vs. non-coding binding elements (x-axis). Pearson correlation and p-value of a linear regression are shown in the upper right corner.

Figure S12. Coding DNaseI footprints are enriched in variants associated with allele-specific chromatin states Heterozygous coding SNVs associated with allele-specific occupancy are significantly enriched inside DNaseI footprints ($P < 1 \times 10^{-8}$, Fisher's exact test using tag normalized datasets).

Figure S13. Coding variants linked to disease susceptibility can also influence chromatin state (A) Shown is the proportion of coding GWAS variants linked to disease susceptibility that overlap DNaseI footprints in one of the 81 tested cell types. (B-F) *Synonymous coding variants linked to psoriasis susceptibility are associated with chromatin state changes selectively within transformed hematopoietic cells.* (B) SNPs in tight linkage with the psoriasis associated synonymous coding variant (rs495337) were identified using whole genome sequencing data from 267 individuals of Northern and Western European (CEPH), Finnish (FIN) and English and Scottish (GBR) ancestry, corresponding to the geographic regions used in the original GWAS studies. Of the 98 SNPs with a genotype R-squared greater than 0.8 (red points), 14 overlap DNaseI footprints in at least 1 cell type and 3 are associated with allelically imbalanced chromatin state, including the initial lead SNP (rs495337, rs492702 and rs2281217). (C) DNaseI cleavage density profiles surrounding two of the psoriasis linked variants that are associated with allelically imbalanced chromatin state (rs495337 and rs492702) for 16 human cells potentially involved in psoriasis pathogenesis. The genotypes of each cell type at rs495337 and rs492702 are indicated to the right of the plot. (D) DNaseI cleavage pattern surrounding the synonymous SNP rs492702 (left) and the synonymous SNP rs495337 (right) in NB4 cells. Binding elements overlapping DNaseI footprints are indicated below. (E) Shown is the chromatin accessibility associated with either the psoriasis risk or non-risk allele of rs492702 (left) and rs495337 (right)

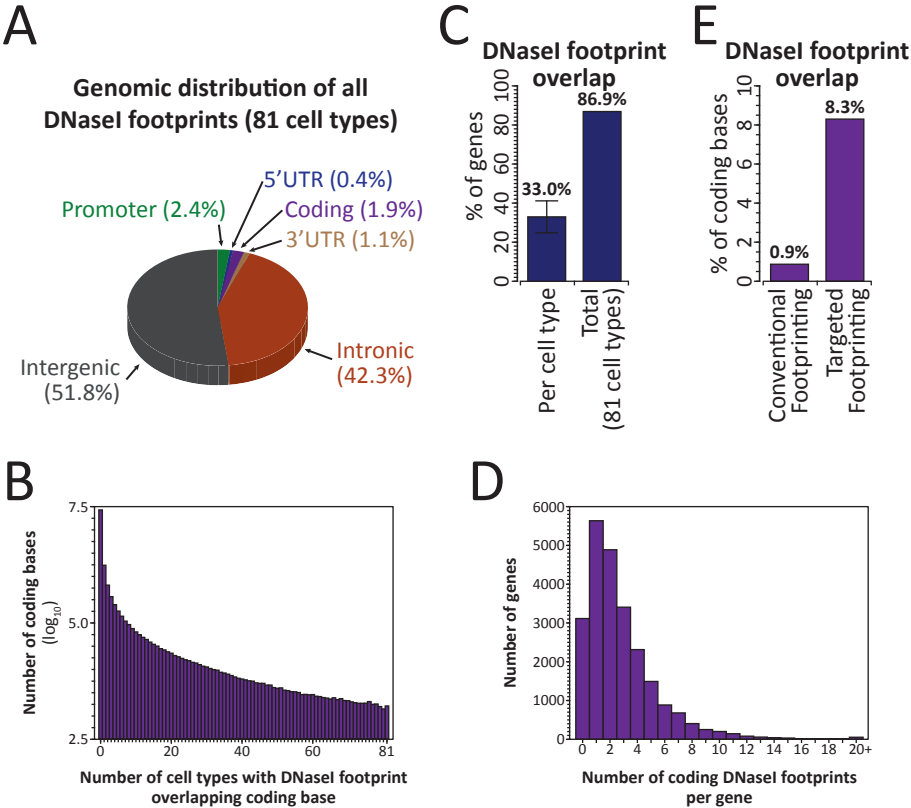
within NB4 cells ($P < 1 \times 10^{-5}$, Fisher's exact test). (F) (top) DNaseI cleavage pattern surrounding the psoriasis-linked non-coding SNP rs2281217 in Melanoma cells (RPMI_7951). (bottom) Shown is the chromatin accessibility associated with either the psoriasis risk or non-risk allele of rs2281217 within Melanoma cells (RPMI_7951) ($P < 1 \times 10^{-2}$, Fisher's exact test).

Table S1. Sample information for the different cell types used in this study

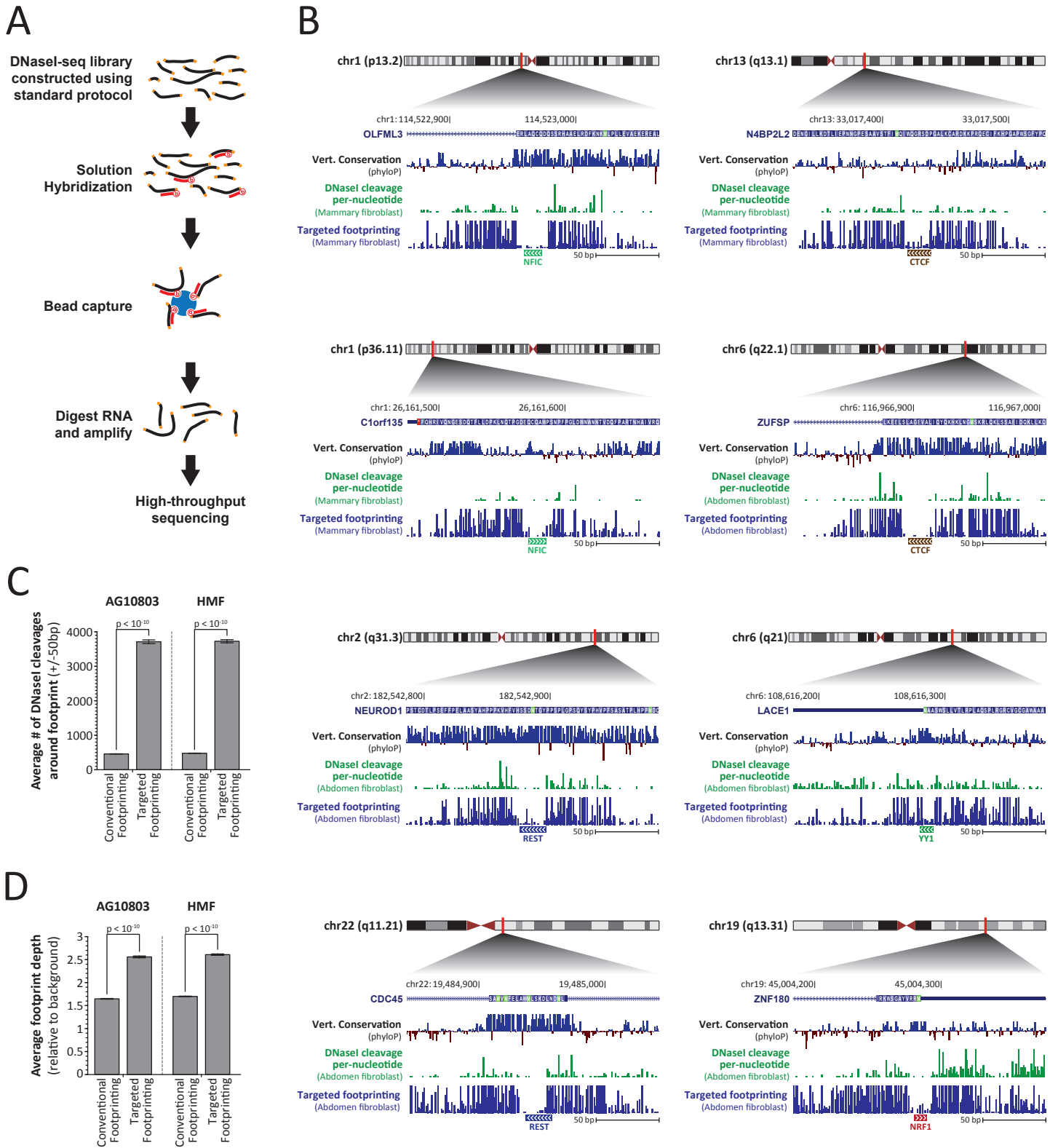
Table S2. Overlap of different codons with DNaseI footprints

Table S3. Regulatory information of SNPs in linkage with rs495337

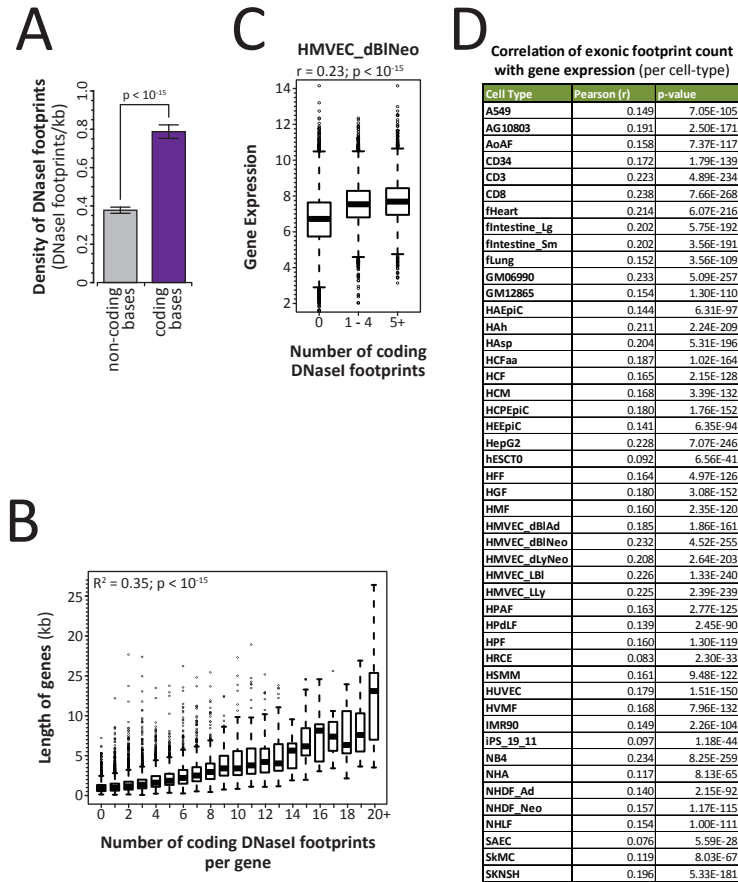
Supplemental Figure 1



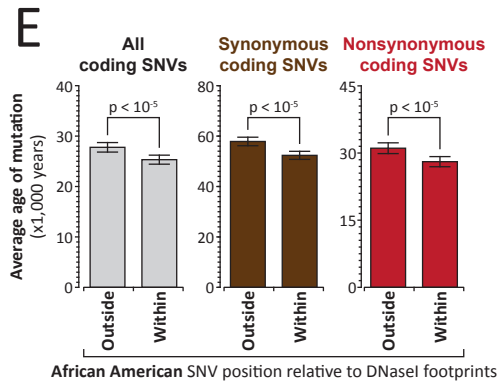
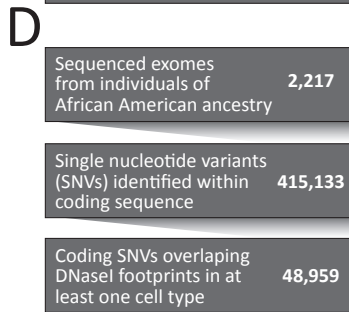
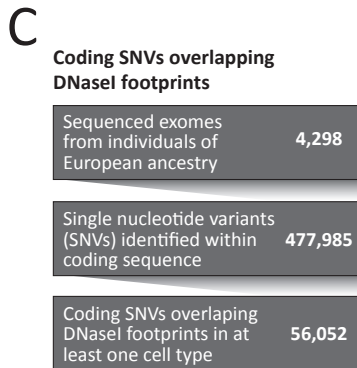
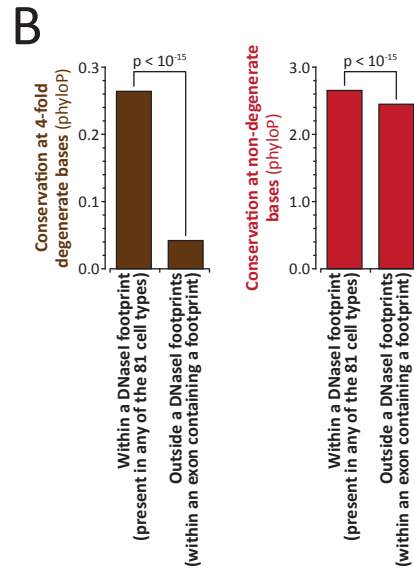
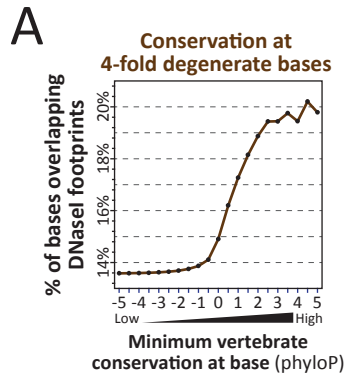
Supplemental Figure 3



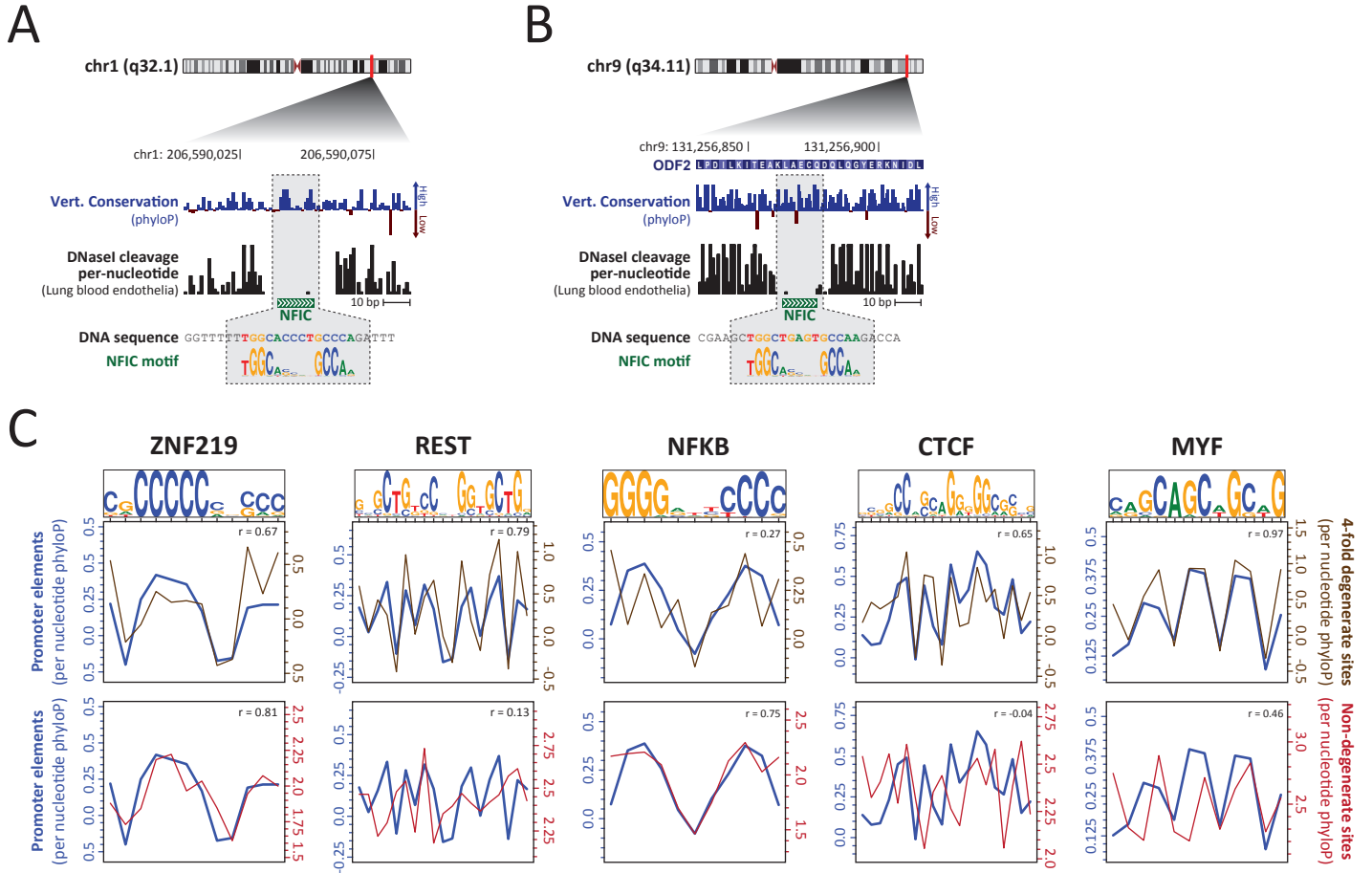
Supplemental Figure 4



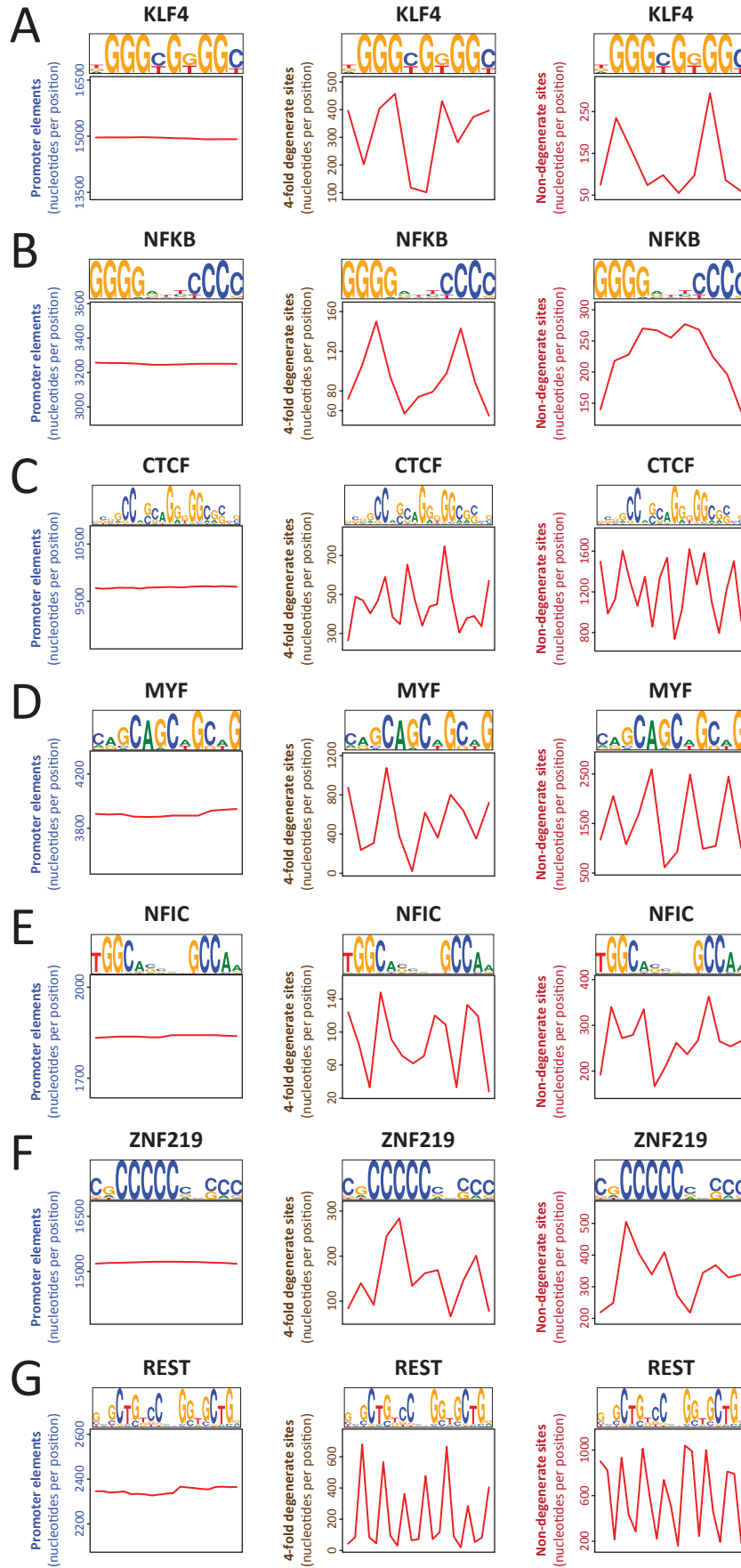
Supplemental Figure 5



Supplemental Figure 6

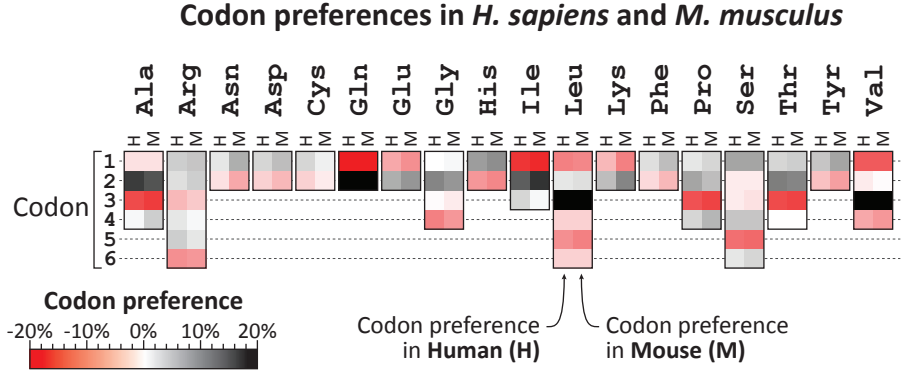


Supplemental Figure 7

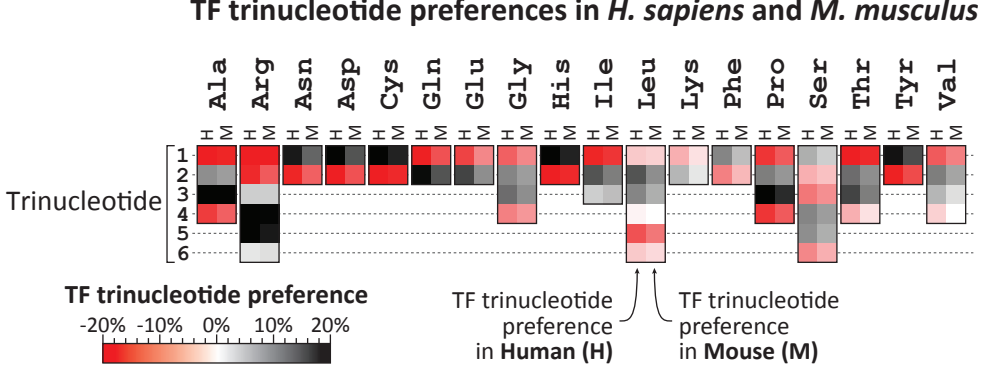


Supplemental Figure 8

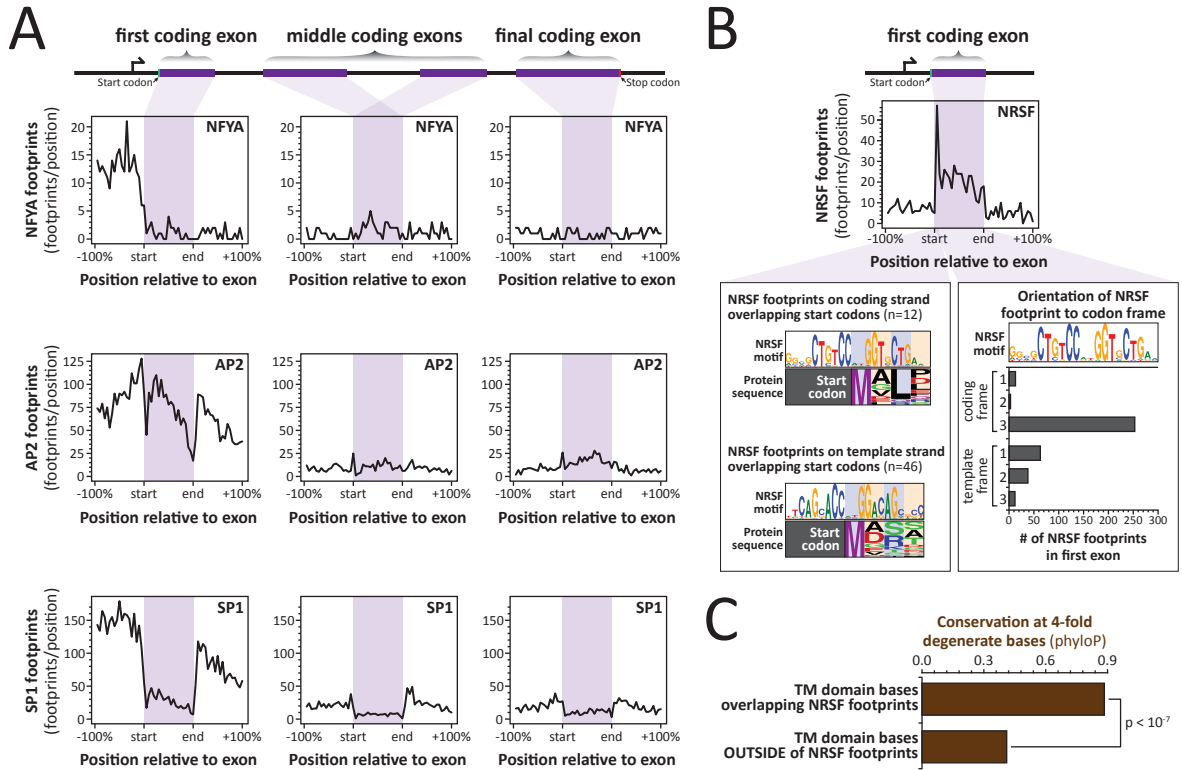
A



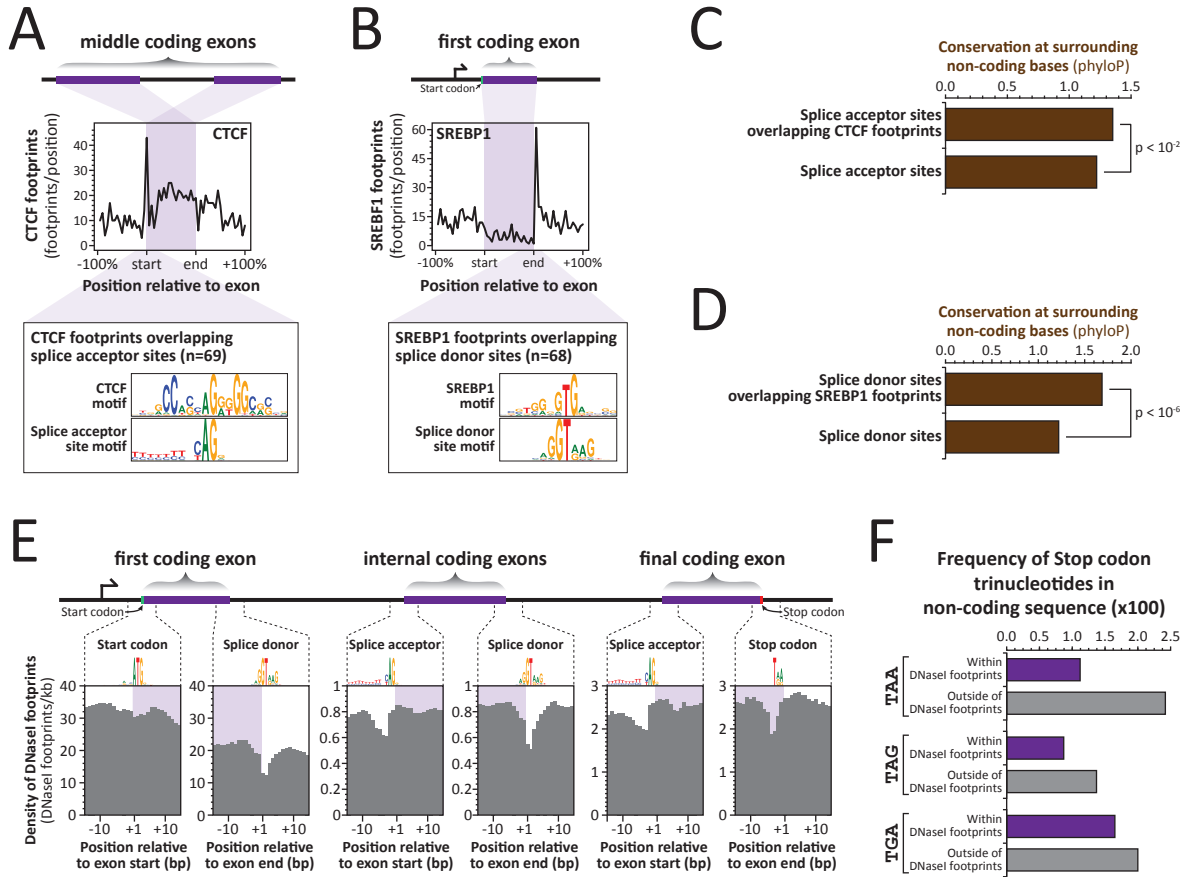
B



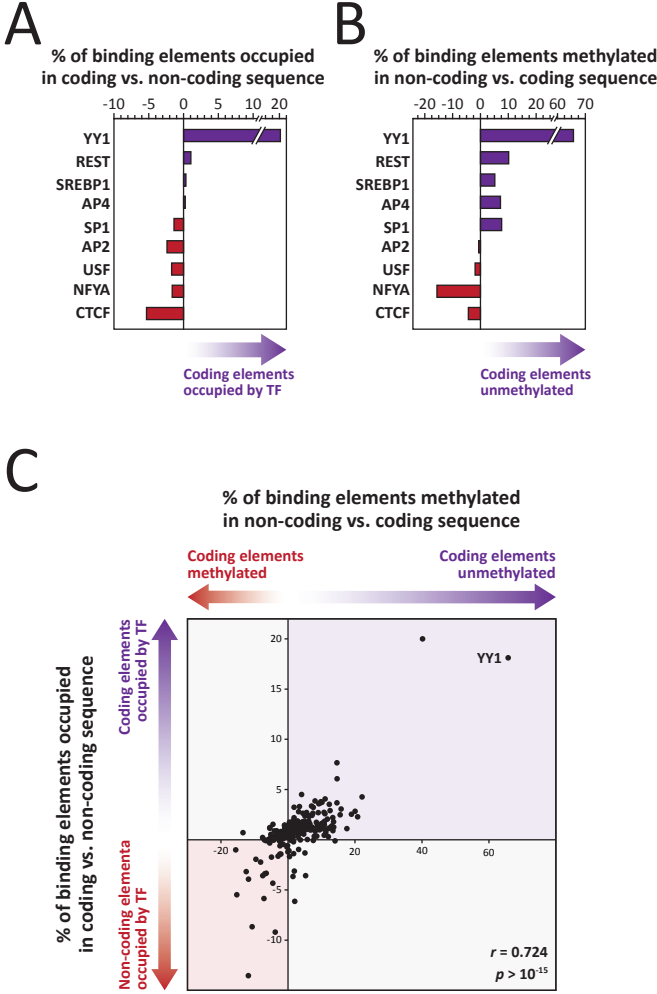
Supplemental Figure 9



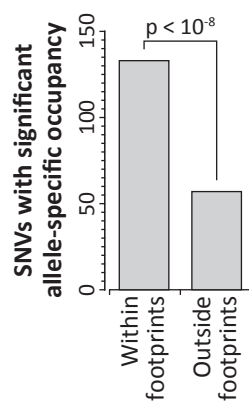
Supplemental Figure 10



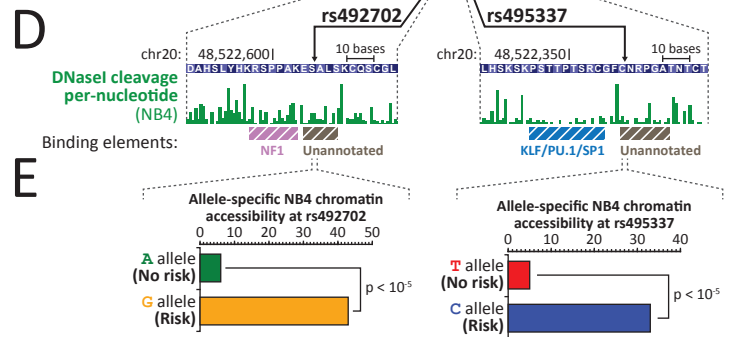
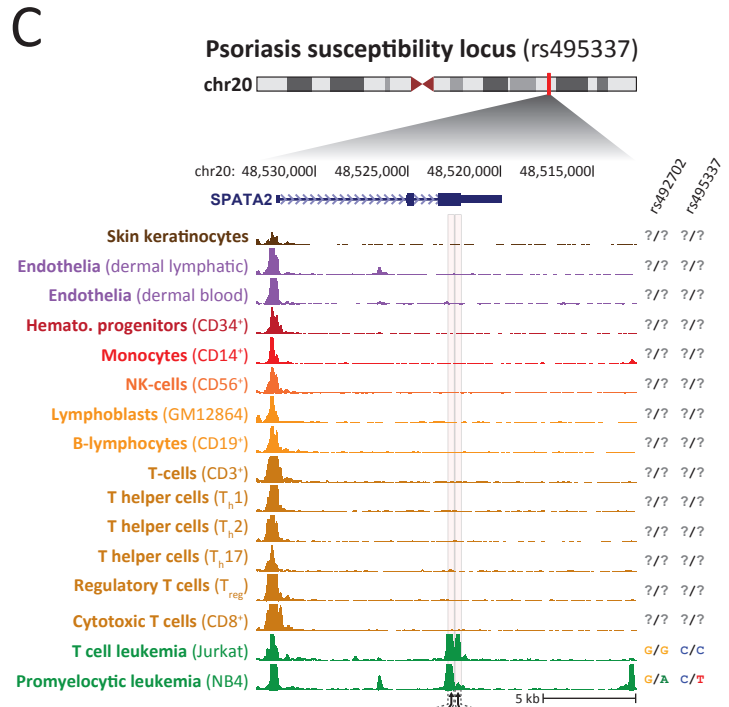
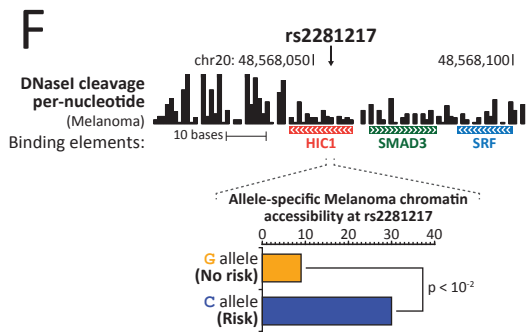
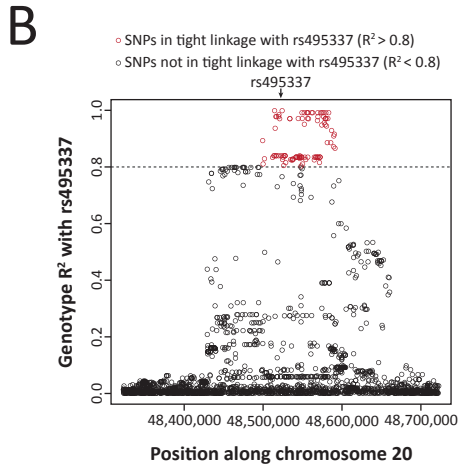
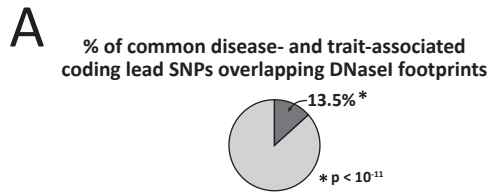
Supplemental Figure 11



Supplemental Figure 12



Supplemental Figure 13



Supplemental Table 1

Cell type	DS number	SPOT	Sequenced reads	Uniquely mapped reads	Total number of DNaseI footprints	Number of coding DNaseI footprints	Neph et al., <i>Nature</i> 2012
A549	DS14289	0.485067	571,473,673	350,629,033	799,061	22,962	No
CD3	DS17198	0.592433	635,018,815	268,414,968	508,821	22,256	No
CD34 (individual 2)	DS12734	0.75638	543,955,153	345,733,458	1,443,618	45,654	No
CD4	DS17881	0.653975	821,450,050	316,205,906	576,410	27,935	No
CD4pos_N	DS14108	0.3923	619,663,019	308,906,610	702,075	19,867	No
CD56	DS17189	0.6606	812,909,023	315,183,940	572,841	24,237	No
CD8	DS17203	0.704	776,670,400	311,987,062	642,374	27,602	No
flintestine_Lg	DS17313	0.43035	533,639,132	392,577,278	1,508,092	51,030	No
flintestine_Sm	DS17317	0.38195	539,120,078	320,341,305	877,043	36,923	No
fMuscle_arm	DS17765	0.500917	252,246,395	204,386,602	1,651,171	66,878	No
fMuscle_back	DS17767	0.488443	599,606,470	282,227,222	991,137	36,532	No
fMuscle_leg	DS20239	0.5809	473,208,353	345,833,051	1,158,459	23,402	No
fMuscle_trunk	DS20242	0.491875	486,203,920	375,241,801	1,180,408	28,667	No
fPlacenta	DS20346	0.442775	1,008,392,565	519,610,800	1,545,783	41,606	No
fSpinal_cord	DS20351	0.4046	494,464,940	367,908,105	836,975	22,823	No
fStomach	DS17878	0.46675	595,242,582	362,167,993	812,854	29,352	No
fThymus	DS20341	0.39645	630,852,138	241,270,763	428,305	15,708	No
HAsp	DS14790	0.525883	376,102,012	292,748,284	1,507,306	32,369	No
HCFaa	DS13480	0.554867	640,379,995	365,088,733	1,543,118	26,867	No
HGF	DS11752	0.4883	521,299,020	336,716,975	1,066,483	17,730	No
HMVEC_LBI	DS13372	0.49885	516,782,159	359,408,173	1,245,886	32,352	No
hTH1	DS18015	0.3854	340,764,122	265,412,022	477,363	18,096	No
hTH17	DS11039	0.2857	382,432,365	207,056,784	533,695	13,274	No
hTH2	DS17597	0.727875	754,042,044	408,764,611	879,882	30,805	No
hTR	DS14702	0.5524	371,137,898	241,377,544	783,257	30,715	No
HUVEC	DS10060	0.386425	1,184,125,568	429,088,276	1,337,725	31,199	No
iPS_19_11	DS15153	0.3743	629,344,480	329,338,603	625,934	20,159	No
iTH2	DS17603	0.4178	513,413,154	383,640,831	731,458	29,241	No
LHCN_M2	DS20548	0.710467	391,281,755	255,134,452	827,462	13,244	No
LHCN_M2_D4	DS20534	0.64325	622,245,602	357,827,356	880,943	19,227	No
M059j	DS20493	0.689833	664,731,412	386,886,267	1,399,284	21,842	No
Mesendoderm	DS19310	0.494517	483,440,991	292,350,792	723,526	18,859	No
MSC	DS21042	0.584033	583,519,364	367,284,952	1,014,611	21,247	No
RPMI_7951	DS20909	0.6849	550,146,566	350,890,993	979,112	15,640	No
Skin_Fibroblasts	DS18224	0.708025	673,387,717	369,970,430	1,565,414	39,398	No
Skin_Keratinocytes	DS18692	0.566975	1,227,133,980	214,248,959	573,146	7,583	No
Skin_Melanocytes	DS18590	0.541433	514,236,139	211,726,949	651,833	11,684	No
T_47D	DS19794	0.58124	632,643,703	350,932,000	694,469	18,047	No
Trophoblast	DS19317	0.298367	486,944,081	289,300,204	365,182	6,231	No
vHMEC	DS18406	0.50275	534,943,942	332,073,525	761,748	14,121	No
AG10803	DS12374	0.71524	406,988,309	284,236,136	1,106,350	15,827	Yes
AoAF	DS13513	0.678227	393,665,970	331,043,131	1,566,105	22,114	Yes
CD20	DS18208	0.571567	338,977,054	240,594,387	603,141	34,313	Yes
CD34 (individual 1)	DS12274	0.704758	287,605,852	221,098,234	902,332	25,716	Yes
fBrain	DS11872	0.71678	247,674,296	202,264,605	1,022,758	15,003	Yes
fHeart	DS12531	0.5493	328,340,544	264,719,957	954,877	22,745	Yes
fLung	DS14724	0.667967	318,544,188	268,295,068	1,181,207	21,843	Yes
GM06990	DS7748	0.6234	238,189,351	137,532,640	434,527	24,825	Yes
GM12865	DS12436	0.56146	416,209,265	263,505,515	811,298	28,230	Yes
HAEPiC	DS12663	0.757515	347,415,753	281,238,046	1,506,375	21,668	Yes
HAh	DS15192	0.5293	491,343,833	367,030,503	966,138	31,141	Yes
HCF	DS12501	0.695743	409,192,151	289,961,235	1,057,693	20,562	Yes
HCM	DS12599	0.73104	380,034,230	274,790,101	1,130,250	23,275	Yes
HCPiC	DS12447	0.758974	305,511,987	253,925,492	1,296,416	24,746	Yes
HEEPiC	DS12763	0.58366	543,479,375	342,975,637	1,263,543	23,587	Yes
HepG2	DS7764	0.5734	248,320,583	168,883,956	448,649	21,562	Yes
hESCT0	DS11909	0.609627	401,363,495	302,050,785	1,279,312	23,013	Yes
HFF	DS15115	0.588125	365,478,871	262,521,646	590,874	13,679	Yes
HIPEPiC	DS12684	0.579325	360,017,644	254,744,863	1,089,872	19,725	Yes
HMF	DS13368	0.74566	384,696,734	311,000,443	1,434,263	20,710	Yes
HMVEC_dBIAd	DS13337	0.719782	292,823,995	239,063,258	1,085,703	22,241	Yes
HMVEC_dBINeo	DS13242	0.56626	395,475,591	293,473,622	1,061,820	31,759	Yes
HMVEC_dLyNeo	DS13150	0.582683	389,329,445	270,345,138	989,591	25,606	Yes
HMVEC_Lly	DS13185	0.62198	453,581,261	313,021,953	872,686	26,578	Yes
HPAF	DS13411	0.700154	320,990,135	255,470,482	1,090,170	19,215	Yes
HPdLF	DS13573	0.674945	371,416,176	304,268,872	1,404,832	18,273	Yes
HPF	DS13390	0.657145	368,365,528	296,713,698	1,175,237	19,056	Yes
HRCE	DS10666	0.599045	284,056,343	236,736,388	1,187,273	11,043	Yes
HSMFM	DS14426	0.65924	467,134,471	367,269,086	1,668,161	24,213	Yes
HVMF	DS13981	0.631767	357,575,517	279,802,866	1,263,772	20,264	Yes
IMR90	DS13219	0.532923	309,171,904	242,507,116	970,254	14,411	Yes
iTH1	DS18018	0.382	462,754,737	355,671,277	499,332	17,891	Yes
K562	DS9767	0.5639	268,452,588	179,970,820	498,615	9,892	Yes
NB4	DS12543	0.56432	439,636,180	323,812,091	1,049,222	35,196	Yes
NHA	DS12800	0.570736	307,812,903	231,589,045	977,890	11,856	Yes
NHDF_Ad	DS12863	0.808153	300,516,213	235,650,107	1,429,337	18,289	Yes
NHDF_Neo	DS11923	0.69876	526,464,624	373,361,757	1,532,787	22,357	Yes
NHLF	DS12829	0.70694	541,170,462	357,163,548	1,567,037	25,608	Yes
SAEC	DS10518	0.58187	296,719,796	243,838,476	1,256,120	14,236	Yes
SkMC	DS11949	0.809203	632,856,867	543,886,965	2,370,625	32,565	Yes
SKNSH	DS8482	0.695	217,691,024	164,615,431	498,908	23,006	Yes

Supplemental Table 1 (continued)

Cell type	Conventional DNaseI footprinting			Targeted DNaseI footprinting		
	Total reads	Coding reads	Coding DNaseI footprints (1% FDR)	Total reads	Coding reads	Coding DNaseI footprints (1% FDR)
AG10803	92,285,815	2,540,236	15,827	44,743,772	9,995,752	190,060
HMF	105,427,896	2,822,710	20,710	34,035,660	8,997,592	156,423

Supplemental Table 2

AA	codon	Total		Outside Footprints			Inside Footprints		
		Counts	Percentage	Counts	Percentage	3 rd base phyloP	Counts	Percentage	3 rd base phyloP
Ala	GCA	158036	22.68%	137349	24.46%	-0.60	17689	14.15%	-0.67
Ala	GCC	280556	40.26%	223397	39.79%	0.30	53651	42.91%	0.34
Ala	GCG	76092	10.92%	44723	7.97%	-1.19	30417	24.33%	-0.22
Ala	GCT	182187	26.14%	155947	27.78%	-0.32	23263	18.61%	-0.37
Arg	AGA	116535	21.04%	104263	23.03%	0.49	10179	10.99%	0.43
Arg	AGG	109274	19.73%	93052	20.55%	0.58	13751	14.84%	0.74
Arg	CGA	61883	11.17%	52934	11.69%	-0.13	8150	8.80%	-0.15
Arg	CGC	105677	19.08%	75776	16.74%	0.15	28655	30.92%	0.38
Arg	CGG	115140	20.79%	88065	19.45%	0.26	25866	27.91%	0.48
Arg	CGT	45341	8.19%	38611	8.53%	-0.54	6061	6.54%	-0.46
Asn	AAC	188446	52.09%	160566	50.79%	1.10	24905	60.41%	1.37
Asn	AAT	173342	47.91%	155596	49.21%	0.32	16324	39.59%	0.32
Asp	GAC	251535	53.35%	209597	51.62%	0.83	38536	64.33%	0.97
Asp	GAT	219908	46.65%	196462	48.38%	0.20	21367	35.67%	0.16
Cys	TGC	122747	53.47%	100642	51.51%	1.13	19138	65.84%	1.30
Cys	TGT	106807	46.53%	94749	48.49%	0.58	9929	34.16%	0.56
Gln	CAA	125287	27.17%	111105	28.24%	0.21	10237	18.58%	0.22
Gln	CAG	335843	72.83%	282347	71.76%	1.26	44870	81.42%	1.44
Glu	GAA	302102	43.32%	268686	44.95%	0.42	28865	33.21%	0.37
Glu	GAG	395223	56.68%	329058	55.05%	0.78	58040	66.79%	1.02
Gly	GGA	160433	25.12%	139913	26.53%	0.02	18405	17.84%	0.01
Gly	GGC	220860	34.57%	172150	32.64%	0.37	45868	44.45%	0.54
Gly	GGG	158362	24.79%	128695	24.40%	-0.08	27713	26.86%	0.20
Gly	GGT	99134	15.52%	86676	16.43%	-0.39	11207	10.86%	-0.32
His	CAC	151888	57.52%	127903	56.13%	0.80	20889	68.25%	1.08
His	CAT	112157	42.48%	99976	43.87%	0.10	9717	31.75%	0.15
Ile	ATA	74740	17.20%	68681	17.87%	0.05	4879	10.77%	0.05
Ile	ATC	199990	46.02%	172015	44.76%	0.87	25617	56.54%	1.06
Ile	ATT	159853	36.78%	143628	37.37%	0.17	14814	32.69%	0.17
Leu	CTA	71581	7.13%	63749	7.50%	-0.32	5972	4.76%	-0.25
Leu	CTC	191801	19.10%	159060	18.71%	0.25	28465	22.70%	0.34
Leu	CTG	395534	39.39%	322523	37.94%	0.83	61818	49.30%	1.01
Leu	CTT	134673	13.41%	118098	13.89%	-0.11	11544	9.21%	-0.16
Leu	TTA	79816	7.95%	72972	8.58%	0.22	4552	3.63%	0.18
Leu	TTG	130746	13.02%	113627	13.37%	1.05	13041	10.40%	1.19
Lys	AAA	250219	44.96%	222456	46.12%	0.61	23534	41.79%	0.61
Lys	AAG	306371	55.04%	259918	53.88%	1.12	32780	58.21%	1.32
Met	ATG	195668	100.00%	171436	100.00%	2.87	19733	100.00%	3.01
Phe	TTC	197471	53.12%	165475	51.71%	1.35	27709	63.57%	1.45
Phe	TTT	174274	46.88%	154511	48.29%	0.67	15882	36.43%	0.59
Pro	CCA	171737	27.37%	151253	28.89%	-0.60	18141	19.15%	-0.67
Pro	CCC	203249	32.40%	166448	31.79%	-0.07	34426	36.35%	0.01
Pro	CCG	72787	11.60%	48657	9.29%	-1.31	23184	24.48%	-0.40
Pro	CCT	179628	28.63%	157156	30.02%	-0.33	18967	20.02%	-0.39
Ser	AGC	198049	23.90%	162287	22.95%	1.03	31855	29.99%	1.17
Ser	AGT	125425	15.13%	110858	15.68%	0.22	12629	11.89%	0.22
Ser	TCA	125015	15.09%	112188	15.86%	-0.59	10352	9.75%	-0.73
Ser	TCC	178300	21.52%	149979	21.21%	0.23	25377	23.89%	0.30
Ser	TCG	45501	5.49%	33087	4.68%	-1.45	11553	10.88%	-0.54
Ser	TCT	156419	18.88%	138806	19.63%	-0.15	14443	13.60%	-0.23
Thr	ACA	147842	28.19%	132832	29.19%	-0.61	12424	20.38%	-0.64
Thr	ACC	183851	35.06%	157462	34.60%	0.19	23809	39.05%	0.30
Thr	ACG	58706	11.20%	46295	10.17%	-1.33	11261	18.47%	-0.70
Thr	ACT	133975	25.55%	118541	26.05%	-0.43	13475	22.10%	-0.45
Trp	TGG	120201	100.00%	101549	100.00%	3.97	16655	100.00%	3.58
Tyr	TAC	147906	54.85%	126426	53.35%	1.21	17736	67.09%	1.50
Tyr	TAT	121736	45.15%	110559	46.65%	0.36	8701	32.91%	0.38
Val	GTA	70420	11.98%	63483	12.59%	-0.30	5702	7.70%	-0.30
Val	GTC	138703	23.59%	117428	23.29%	0.43	19475	26.28%	0.55
Val	GTG	271113	46.11%	227587	45.15%	0.61	38177	51.52%	0.85
Val	GTT	107767	18.33%	95608	18.97%	-0.08	10745	14.50%	-0.11

Supplemental Table 2 (continued)

	Total		Outside Footprints		Inside Footprints		
	tri-nucleotide	Counts	Percentage	Counts	Percentage	Counts	Percentage
Non-Coding Bases	GCA	158036	22.68%	16085315	36.00%	2389977	27.57%
	GCC	280556	40.26%	11500629	25.74%	2684015	30.96%
	GCG	76092	10.92%	1758306	3.94%	918723	10.60%
	GCT	182187	26.14%	15338973	34.33%	2677598	30.88%
	AGA	116535	21.04%	26654892	48.45%	2718974	31.67%
	AGG	109274	19.73%	19671800	35.76%	2868565	33.42%
	CGA	61883	11.17%	2012057	3.66%	499182	5.82%
	CGC	105677	19.08%	1759833	3.20%	910464	10.61%
	CGG	115140	20.79%	2268487	4.12%	998901	11.64%
	CGT	45341	8.19%	2642651	4.80%	588019	6.85%
	AAC	188446	52.09%	17436396	35.28%	2078064	42.62%
	AAT	173342	47.91%	31980316	64.72%	2797683	57.38%
	GAC	251535	53.35%	10814256	41.03%	1671154	48.85%
	GAT	219908	46.65%	15540540	58.97%	1750005	51.15%
	TGC	122747	53.47%	16120832	40.32%	2350619	49.58%
	TGT	106807	46.53%	23865923	59.68%	2390020	50.42%
	CAA	125287	27.17%	22142125	49.86%	2163034	38.82%
	CAG	335843	72.83%	22264852	50.14%	3408480	61.18%
	GAA	302102	43.32%	24962534	58.23%	2900245	51.50%
	GAG	395223	56.68%	17905596	41.77%	2731409	48.50%
	GGA	160433	25.12%	17416918	31.19%	2722195	26.79%
	GGC	220860	34.57%	11496632	20.59%	2655183	26.13%
	GGG	158362	24.79%	14473763	25.92%	3045355	29.97%
	GGT	99134	15.52%	12453105	22.30%	1738846	17.11%
	CAC	151888	57.52%	15847814	41.57%	2278234	50.05%
	CAT	112157	42.48%	22277748	58.43%	2274033	49.95%
	ATA	74740	17.20%	26177302	35.54%	1628736	24.78%
	ATC	199990	46.02%	15486272	21.02%	1853073	28.20%
	ATT	159853	36.78%	31999516	43.44%	3090352	47.02%
	CTA	71581	7.13%	15917018	12.17%	1290231	9.06%
	CTC	191801	19.10%	17875369	13.67%	2796382	19.63%
	CTG	395534	39.39%	22175471	16.96%	3457893	24.28%
	CTT	134673	13.41%	25053558	19.16%	2671903	18.76%
	TTA	79816	7.95%	27833146	21.29%	1756757	12.33%
	TTG	130746	13.02%	21907286	16.75%	2269524	15.93%
	AAA	250219	44.96%	46424991	65.06%	4241843	62.52%
	AAG	306371	55.04%	24932199	34.94%	2542624	37.48%
	ATG	195668	100.00%	22100500	100.00%	2362707	100.00%
	TTC	197471	53.12%	24820252	34.77%	3105419	40.55%
	TTT	174274	46.88%	46567767	65.23%	4551962	59.45%
CCA	171737	27.37%	19310084	34.61%	3088979	30.43%	
CCC	203249	32.40%	14466833	25.93%	3051913	30.07%	
CCG	72787	11.60%	2270381	4.07%	999665	9.85%	
CCT	179628	28.63%	19749821	35.40%	3009067	29.65%	
AGC	198049	23.90%	15295549	14.75%	2545151	18.61%	
AGT	125425	15.13%	19203580	18.52%	2084088	15.24%	
TCA	125015	15.09%	23156768	22.33%	2553818	18.68%	
TCC	178300	21.52%	17348973	16.73%	2860858	20.92%	
TCG	45501	5.49%	2006990	1.94%	509650	3.73%	
TCT	156419	18.88%	26675962	25.73%	3120578	22.82%	
ACA	147842	28.19%	23597149	40.86%	2582610	34.94%	
ACC	183851	35.06%	12396514	21.46%	1803506	24.40%	
ACG	58706	11.20%	2617580	4.53%	602911	8.16%	
ACT	133975	25.55%	19145767	33.15%	2402188	32.50%	
TGG	120201	100.00%	19368891	100.00%	3044710	100.00%	
TAC	147906	54.85%	13939244	34.61%	1078248	41.50%	
TAT	121736	45.15%	26336642	65.39%	1520053	58.50%	
GTA	70420	11.98%	13890032	23.92%	1157080	15.58%	
GTC	138703	23.59%	10835153	18.66%	1720977	23.18%	
GTG	271113	46.11%	15820291	27.24%	2344663	31.58%	
GTT	107767	18.33%	17526547	30.18%	2202690	29.66%	

Supplemental Table 3

SNP	Genomic location	Minor allele frequency	Coding SNP	R ² with rs495337	Allele associated with rs495337 risk allele	Cell types with overlapping DNaseI footprints	Cell types with allelic chromatin imbalance at overlapping footprints
rs495337	chr20:48522329-48522330	T=0.329	Yes	N/A	G	NB4-DS12543.fdr0p01.fps.bed	NB4-DS12543.fdr0p01.fps.bed
rs488255	chr20:48499084-48499085	A=0.426	No	0.894109	C	fMuscle_arm-DS17765.fdr0p01.fps.bed	NONE
rs471054	chr20:48513235-48513236	A=0.482	No	0.840395	C	NB4-DS12543.fdr0p01.fps.bed	NONE
rs492702	chr20:48522584-48522585	T=0.485	Yes	0.840395	G	NB4-DS12543.fdr0p01.fps.bed	NB4-DS12543.fdr0p01.fps.bed
rs4809762	chr20:48533526-48533527	G=0.439	No	0.97146	T	GM06990-DS7748.fdr0p01.fps.bed	NONE
rs6125813	chr20:48541182-48541183	T=0.487	No	0.835503	G	AoAF-DS13513.fdr0p01.fps.bed HVMF-DS13981.fdr0p01.fps.bed	NONE
rs6125819	chr20:48545709-48545710	C=0.487	No	0.835503	T	HAEpiC-DS12663.fdr0p01.fps.bed HVMF-DS13981.fdr0p01.fps.bed M059J-DS20493.fdr0p01.fps.bed	NONE
rs73129264	chr20:48551678-48551679	T=0.327	No	0.992795	G	CD20-DS18208.fdr0p01.fps.bed GM12865-DS12436.fdr0p01.fps.bed hTR-DS14702.fdr0p01.fps.bed	NONE
rs4287819	chr20:48551955-48551956	G=0.446	No	0.951858	T	flntestine_Lg-DS17313.fdr0p01.fps.bed NB4-DS12543.fdr0p01.fps.bed	NONE
rs1056198	chr20:48556228-48556229	T=0.328	No	0.992795	C	hESCT0-DS11909.fdr0p01.fps.bed	NONE
rs2281217	chr20:48568054-48568055	G=0.496	No	0.835503	C	AoAF-DS13513.fdr0p01.fps.bed HGF-DS11752.fdr0p01.fps.bed HMF-DS13368.fdr0p01.fps.bed HPdLF-DS13573.fdr0p01.fps.bed HPF-DS13390.fdr0p01.fps.bed IMR90-DS13219.fdr0p01.fps.bed NHDF_Ad-DS12863.fdr0p01.fps.bed RPMI_7951-DS20909.fdr0p01.fps.bed	RPMI_7951-DS20909.fdr0p01.fps.bed
rs4809769	chr20:48579163-48579164	A=0.330	No	0.992795	C	M059J-DS20493.fdr0p01.fps.bed	NONE
rs6067293	chr20:48581740-48581741	C=0.448	No	0.972231	T	GM12865-DS12436.fdr0p01.fps.bed	NONE
rs6063454	chr20:48590790-48590791	T=0.280	No	0.866789	G	fMuscle_arm-DS17765.fdr0p01.fps.bed fMuscle_trunk-DS20242.fdr0p01.fps.bed HPdLF-DS13573.fdr0p01.fps.bed	NONE