

Human enhancers associated with immune response harbour specific sequence composition, activity, and genome organization

Charles-Henri Lecellier^{1,2,3,*}, Wyeth W. Wasserman³, Remo Rohs⁴,
and Anthony Mathelier^{3,5,6,*}

¹*Institut de Génétique Moléculaire de Montpellier, Université de Montpellier, 1919
Route de Mende - 34293 Montpellier cedex 5, France*

²*Institut de Biologie Computationnelle, 860 rue de St. Priest, 34095 Montpellier
cedex 5, France*

³*Centre for Molecular Medicine and Therapeutics at the Child and Family Research
Institute, Department of Medical Genetics, University of British Columbia, 980 West
28th Avenue, Room 3103, V5Z 4H4, Vancouver, BC, Canada*

⁴*Molecular and Computational Biology Program, Department of Biological Sciences,
University of Southern California, Los Angeles, CA 90089, USA*

⁵*Centre for Molecular Medicine Norway (NCMM), Nordic EMBL partnership,
University of Oslo and Oslo University Hospital, Oslo, Norway*

⁶*Department of Cancer Genetics, Institute for Cancer Research, Oslo University
Hospital Radiumhospitalet, Oslo, Norway*

* Co-corresponding authors: charles.lecellier@igmm.cnrs.fr (CHL) and
anthony.mathelier@ncmm.uio.no (AM)

September 30, 2016

Abstract

The FANTOM5 consortium recently characterized 38,554 robust human enhancers from 808 cell and tissue types using the Cap Analysis of Gene Expression technology. We used the distribution of guanine and cytosine nucleotides at enhancer regions to distinguish two classes of enhancers harboring distinct DNA structural properties. A functional analysis of their predicted gene targets highlighted one class of enhancers as significantly enriched for associations with immune response genes. Moreover, these enhancers were specifically enriched for regulatory motifs recognized by TFs involved in immune response. We observed that immune response enhancers were cell type specific, preferentially activated upon bacterial infection, and with long-lasting response activity. Looking at chromatin capture data, we found that the two classes of enhancers were lying in distinct topologically-associated domains and chromatin loops. Our results suggest that specific DNA sequence patterns encode for classes of enhancers that are functionally distinct and specifically organized in the human genome.

Background

Gene expression is regulated through many layers, one of which being the regulation of the transcription of DNA segments into RNA. Transcription factors (TFs) are key proteins regulating this process through their specific binding to the DNA at regulatory elements, the TF binding sites (TFBSs) [1]. These regulatory elements are located within larger regulatory regions, the promoters and enhancers [2]. While promoters are situated around transcription start sites (TSSs), enhancers are distal to the genes they regulate. The canonical view is that chromatin conformation places enhancers in close 3D proximity to their target gene promoters through DNA looping [3–5]. High-resolution chromatin conformation capture (Hi-C) technology maps genomic regions in spatial proximity within cell nuclei [6]. The Hi-C technology identified specific genomic neighbourhoods of chromatin interactions, the topologically associating domains (TADs), which represent chromatin compartments that are stable between cell types and conserved across species [7, 8].

Studies have shown relationships between the composition of a DNA sequence in guanine (G) and cytosine (C) and chromatin organization, for instance in relation to nucleosome positioning [9, 10]. Furthermore, sequence composition is intrinsically linked to the three-dimensional structure of the DNA. Topological studies have used sequence properties to predict four structural features of DNA: helix twist (HelT), minor groove width (MGW), propeller twist (ProT), and Roll [11, 12]. These topological properties have been shown to inform the analysis of protein-DNA interactions obtained from high-throughput experiments [13–16], emphasizing the importance of DNA sequence composition in transcriptional regulation.

DNA sequence composition and other features of promoter regions have been extensively studied, including such key advances as the discovery of CpG is-

lands. The analysis of promoter regions in the human genome was accelerated by the development of the Cap Analysis of Gene Expression (CAGE) technology [17,18], which identifies active TSSs in a high-throughput manner based on 5' capped RNA isolation. Using CAGE data, a large scale identification of the precise location of TSSs in human [19] led to the classification of promoters into four classes based on G+C content (%GC) [20]. The study highlighted that GC-rich promoters are associated with genes involved in various binding and protein transport activities while GC-poor promoters are associated with genes responsible for environmental defense responses. While promoters overlapping CpG islands are commonly assumed to be ubiquitous drivers of housekeeping genes, comprehensive analysis of CAGE data from > 900 human samples showed that a subset deliver cell type-specific expression [21].

Large-scale computational analyses of enhancer regions have been hampered by a limited set of bona fide enhancers. An advantage of the CAGE technology is its capacity to identify *in vivo*-transcribed enhancers. Specifically, it identifies active enhancer regions in biological samples by capturing bidirectional RNA transcripts at enhancer boundaries [22]. Using this characteristic of CAGE data, the FANTOM5 project identified 38,554 “robust” human enhancers across 808 samples [22]. Sequence property analysis suggested that the enhancers share properties with CpG-poor promoters. The findings shed light on the structure, organization, and function of human enhancers.

As enhancers are distal to the genes they regulate, it is challenging to predict these relationships. Based on cross-tissue correlations between histone modifications at enhancers and CAGE-derived expression at promoters within 1,000 bp, enhancer-promoter links have been shown to be conserved across cell types [23]. As the CAGE technology captures the level of activity for both promoters and enhancers in the same samples, predicting the potential targets of the enhancers was obtained by correlating the activity levels of these regulatory regions over hundreds of human samples from the FANTOM5 consortium [22]. Using the predicted enhancer-gene associations, the authors unveiled that closely spaced enhancers were linked to genes involved in immune and defense responses. These results stress that predictions of enhancer-promoter associations are critical to decipher the functional roles of enhancers.

Here, we used the distribution of G+C nucleotides along the sequences of human CAGE-derived enhancer regions to define two classes of enhancers. The specific sequence features of the two classes encoded for distinct topological DNA shape patterns. The enhancers from the GC-poor class were predicted to be functionally associated with genes involved in the immune response whereas the enhancers from the other class were associated with genes involved in biological processes related to transcription. Accordingly, regulatory motifs associated with immune response TFs like NF- κ B are enriched in the DNA sequence of the immune response-related set of enhancers. Independent functional analysis of histone modification and CAGE data highlighted a cell type specificity of these enhancers along with their activation upon bacterial infection. Moreover, immune system enhancers were observed with a long-lasting response activity pattern following cell stimulation in time-course data sets. Finally, we observed

that the two classes of enhancers tended to be structurally organized in the human chromosomes within distinct TADs and DNA chromatin loops.

Results

Guanine and cytosine nucleotide patterns identified two classes of human enhancers with distinct DNA structural properties

To analyze the sequence properties of human enhancers, we considered the set of 38,554 CAGE-derived enhancers found to be significantly active in at least one primary cell or tissue sample in the FANTOM5 project [21, 22]. We extracted 500 bp DNA sequences 5' and 3' of the mid-point of the enhancers as defined by Andersson *et al.* [22]. We sought to identify distinct classes of enhancers based on the distribution of guanines (Gs) and cytosines (Cs) along the enhancer regions. Specifically, each enhancer was represented by a 1,001 bp-long binary vector with 1s representing G+C and 0s representing adenines (As) and thymines (Ts). We clustered the enhancers by applying the k-means clustering algorithm [24] on the vectors. To select the number of clusters k , we considered silhouette plots, which provide a visual representation of how close each enhancer in one cluster is to enhancers in neighbouring clusters [25]. A visual inspection of cluster silhouettes with $k \in [2, 5]$ revealed that the best clustering was obtained with $k = 2$ (Figure S1). We extracted two classes ($k = 2$) of enhancers with distinct distributions of G+C along the enhancer regions (Figure 1a). The two classes were composed of 14,204 and 24,343 enhancers, hereafter referred to as class 1 and class 2, respectively. While enhancers from class 1 were more GC-rich than enhancers from class 2, separating the enhancers solely based on GC content would have resulted in a different classification (i.e. there is an overlap between the classes in terms of G+C content, as shown in Figure 1b).

As DNA sequence and shape are intrinsically linked, we next considered four DNA shape features computed from DNA sequences with the DNASHape tool [12]: helix twist (HelT), minor groove width (MGW), propeller twist (ProT), and Roll. We applied the k-means clustering algorithm with $k = 2$ to vectors combining DNA shape feature values extracted from the GBshape database [11] at 1,001 bp-long enhancer regions centered around enhancers' mid-points. We obtained two sets containing 15,259 (set 1) and 23,288 (set 2) enhancers, respectively. These sets of enhancers derived from DNA shape features were very similar to classes 1 and 2 that were obtained using G+C patterns at enhancer regions. Indeed, class 1 and set 1 have a Jaccard similarity of 0.85; class 2 and set 2 have a Jaccard similarity of 0.90.

We plotted the distribution of the four DNA shape features along the enhancer regions from the two classes obtained with the G+C pattern-based clustering (Figure 1c-f). Similarly, we plotted DNA shape features for the two sets obtained from the DNA shape-based clustering (Figure S2). We consistently

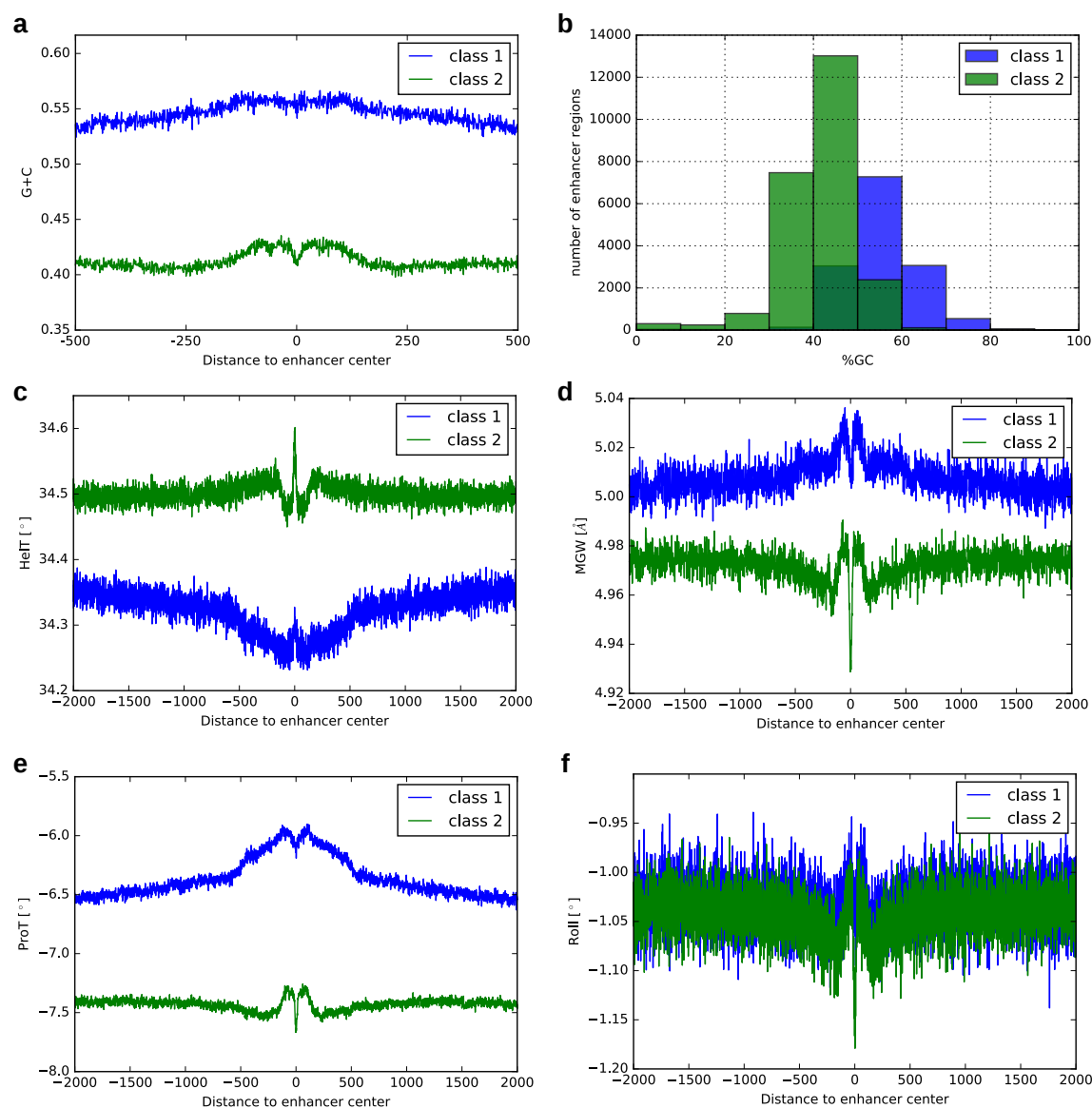


Figure 1: DNA sequence features at enhancers. Features associated with human enhancers from class 1 and class 2 are represented in blue and green, respectively. **a.** G+C values (y-axis) of the k-means cluster centers along DNA regions ± 500 bp centered at enhancer center points (x-axis). **b.** Histogram of the %GC content of the enhancers. **c-f.** Average DNA shape values (y-axis) along the DNA regions $\pm 2,000$ bp centered at enhancer middle-points (x-axis) for DNA shape features HelT (c), MGW (d), ProT (e), and Roll (f).

155 observed that class 1 enhancers harboured lower HelT values at the centre of
 156 the enhancers as well as about 500 bp away from the enhancers' mid-points
 157 (Figure 1c). We observed a symmetrical pattern for MGW with width decrease
 158 at the central positions of the enhancers as well as at the edges (~ 50 -150 bp
 159 away from the mid-points) of the enhancers (Figure 1d). ProT and Roll signals
 160 were also distinct between enhancers from the two classes (Figure 1e-f). The
 161 patterns observed for the DNA shape features were in agreement with the two
 162 distinct patterns of G+C composition computed along the enhancers from the
 163 two classes (Figure 1a).

164 The similarity between G+C- and DNA shape-based clustering stresses that
 165 the G+C pattern is the key discriminant between the two classes of enhancers
 166 while the shape represents a secondary effect of the G+C pattern. We therefore
 167 focused on the two classes of enhancers derived from their G+C pattern in this
 168 report, except otherwise stated. Taken together, these results described two
 169 subsets of human enhancers distinguishable by their distribution of G+C along
 170 their length and reflected in their DNA structural properties.

171 **The two classes of human enhancers associated with specific** 172 **biological processes**

173 Different classes of mammalian promoters, derived from their nucleotide com-
 174 position, were observed to be associated with genes linked to distinct biological
 175 functions [20]. Following the same approach, we sought for a functional inter-
 176 pretation of the classification that we obtained. Based on correlations between
 177 promoter and enhancer activities derived from CAGE data in human samples,
 178 Andersson *et al.* linked enhancers to their potential gene promoter targets [22].
 179 To infer the biological functions of enhancers, we assumed that each enhancer
 180 was associated with the same biological functions as the genes it was predicted
 181 to regulate. Class 1 enhancers were predicted to target 7,713 genes whereas
 182 class 2 enhancers were linked to 7,857 genes (Table S1). In aggregate, the en-
 183 hancers corresponded to a set of 11,271 genes, of which 4,299 were common to
 184 the two classes (representing $\sim 56\%$, $\sim 55\%$ and $\sim 38\%$ of class 1, class 2, and
 185 the combined set of genes, respectively). We submitted the two sets of genes
 186 associated to class 1 and class 2 enhancers to the GOrilla tool [26] to predict
 187 enriched (p-value $< 1 \times 10^{-11}$) gene ontology (GO) biological processes. Note
 188 that the aggregated set of 11,271 genes was used as the background set of genes
 189 for enrichment analyses.

190 Biological processes linked to RNA transcription were found to be enriched
 191 for genes associated with class 1 enhancers (Figures 2a and S3 and Table S2).
 192 Specifically, the directed acyclic graph (Figure S3) of the enriched GO terms
 193 highlighted two leafs corresponding to the terms 'transcription, DNA-templated'
 194 (FDR q-value = 8.7×10^{-13}) and 'regulation of transcription, DNA-templated'
 195 (q = 7.8×10^{-12}). When considering the genes predicted to be regulated by
 196 enhancers from class 2, only two GO biological processes were predicted to be
 197 enriched (Figures 2b and S4 and Table S3): 'immune system process' (q =
 198 6.1×10^{-9}) and 'regulation of immune response' (q = 3.2×10^{-8}).

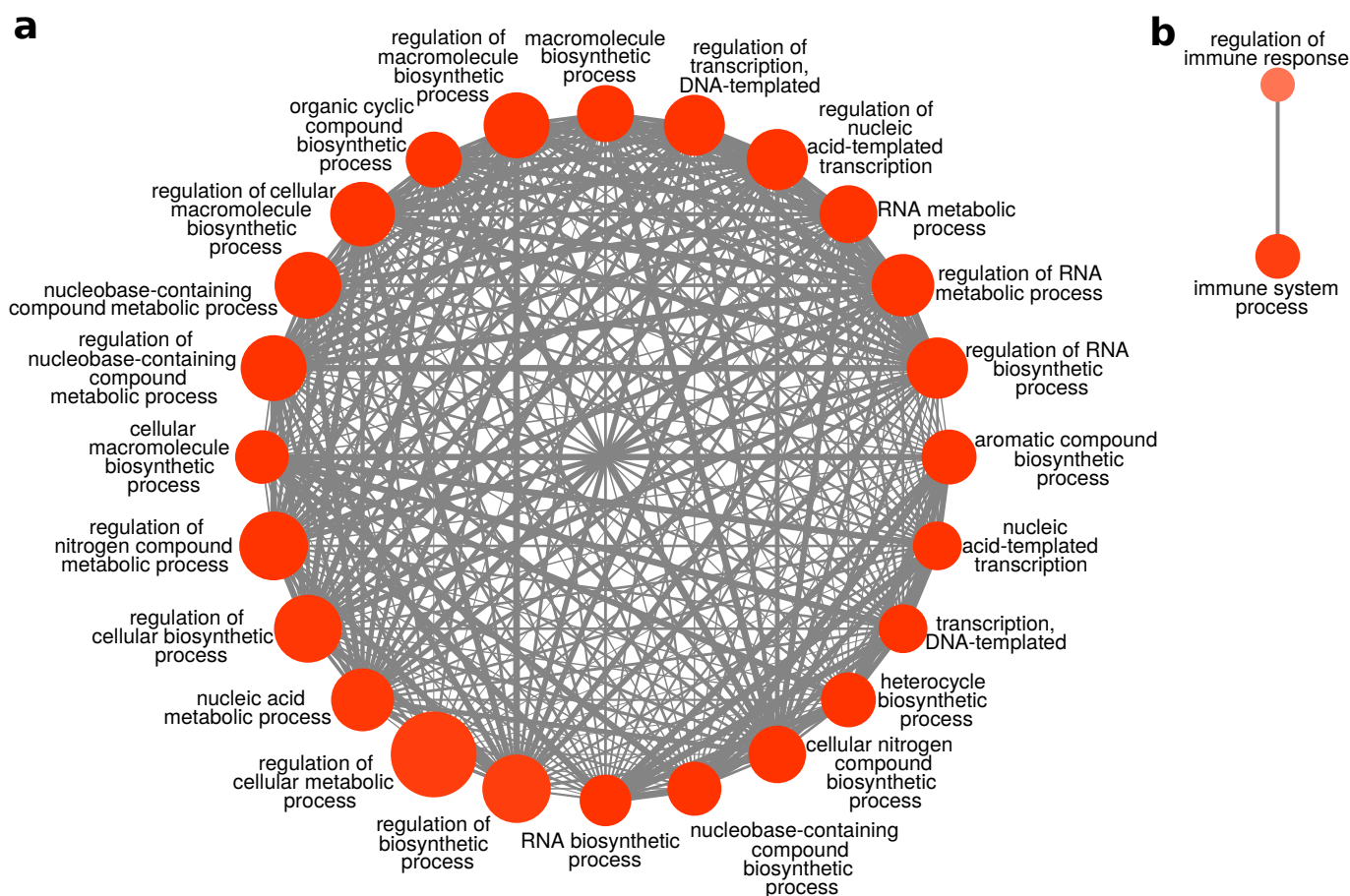


Figure 2: Functional enrichment analysis. Enriched GO biological processes associated with genes predicted to be regulated by enhancers from class 1 (**a**) and class 2 (**b**) were obtained using the GOrilla tool [26]. Nodes in the graphs represent enriched GO biological processes. The color of a node represents the FDR q-value of the corresponding enriched GO biological process, the more red, the lower the q-value (min = 1.37×10^{-13} , max = 3.8×10^{-8}). The size of a node represents the number of genes associated with class 1 (a) and class 2 (b) enhancers in the corresponding GO biological process (min = 761, max = 3,832). Edges between two nodes indicate the number of common genes between corresponding processes. The larger the number of overlapping genes (min = 264, max = 860), the larger the edge between the two corresponding nodes. FDR, false discovery rate; GO, gene ontology.

199 These functional enrichment results were specific to the classification of
200 the enhancers using the pattern of G+C along the enhancer regions. Indeed,
201 we considered a segregation of human enhancer regions solely based on %GC
202 (mean = 46.62, median = 46.15, and standard deviation = 10.56) with GC-
203 poor enhancer regions (%GC < 46.15) assigned to a first set and GC-rich ones
204 (%GC > 46.15) to a second set. We submitted the sets of genes linked to the
205 enhancers from the %GC-based classification to GOrilla and observed that the
206 immune system-related GO terms were found enriched for both sets (Figures S5
207 and S6).

208 Taken together, the functional enrichment results revealed that a classifi-
209 cation based on the distribution of Gs and Cs along human enhancer regions
210 featured two sets of enhancers predicted to be regulating genes with distinct
211 biological functions. While the first class was linked to genes associated with
212 transcription, the second class highlighted enhancers predicted to regulate im-
213 mune system related genes.

214 **Distinct transcription factors predicted to act upon the two** 215 **classes of human enhancers and their predicted promoter** 216 **targets**

217 We sought to identify TF binding motifs enriched within each class of enhancers,
218 to suggest driving TFs for the distinct biological functions. We considered
219 1,001 bp-long DNA sequences centered at the enhancers' mid-points. Posi-
220 tional motif enrichment analyses were performed using the Centrimo tool [27] to
221 predict TF binding motifs over-represented around the enhancers' mid-points.
222 Class 1 enhancer regions were compared to class 2 regions and vice-versa to
223 highlight specific motifs (Figure 3a,c and Data S1). Motifs associated with the
224 Specificity Protein/Krüppel-like Factor (SP/KLF) TFs were enriched in class 1
225 enhancer regions (Figure 3a and Data S1). Members of the SP/KLF family have
226 been associated to a large range of core cellular processes such as cell growth,
227 proliferation, and differentiation [28]. The most enriched motifs in class 2 en-
228 hancer regions were associated with nuclear factor kappa-light-chain-enhancer
229 of activated B cells (NF- κ B)/Rel TFs (Figure 3c and Data S1). As NF- κ B is
230 known to have a central role in immune response [29], the enrichment is consis-
231 tent with an involvement of class 2 enhancers in the immune response biological
232 function (Figure 2b). Other enriched motifs in class 2 enhancers were associ-
233 ated with BACH1/2 TFs, involved in acquired and innate immunity [30], and
234 chromatin remodelling TFs BPTF [31] and SMARCC1 [32].

235 Linked enhancers and promoters were predicted to be driven by similar sets
236 of TF binding motifs when enhancer-promoter links were derived from a dis-
237 tinct collection of CAGE data from time-course studies [33]. We assessed such
238 associations by extending our positional motif enrichment analyses to the pro-
239 moters of genes associated with classes 1 and 2 enhancers, respectively. We used
240 Centrimo to predict motifs locally enriched in 1,001 bp regions centered around
241 corresponding TSSs. Motifs associated with SP1 and NF- κ B TFs were specifi-

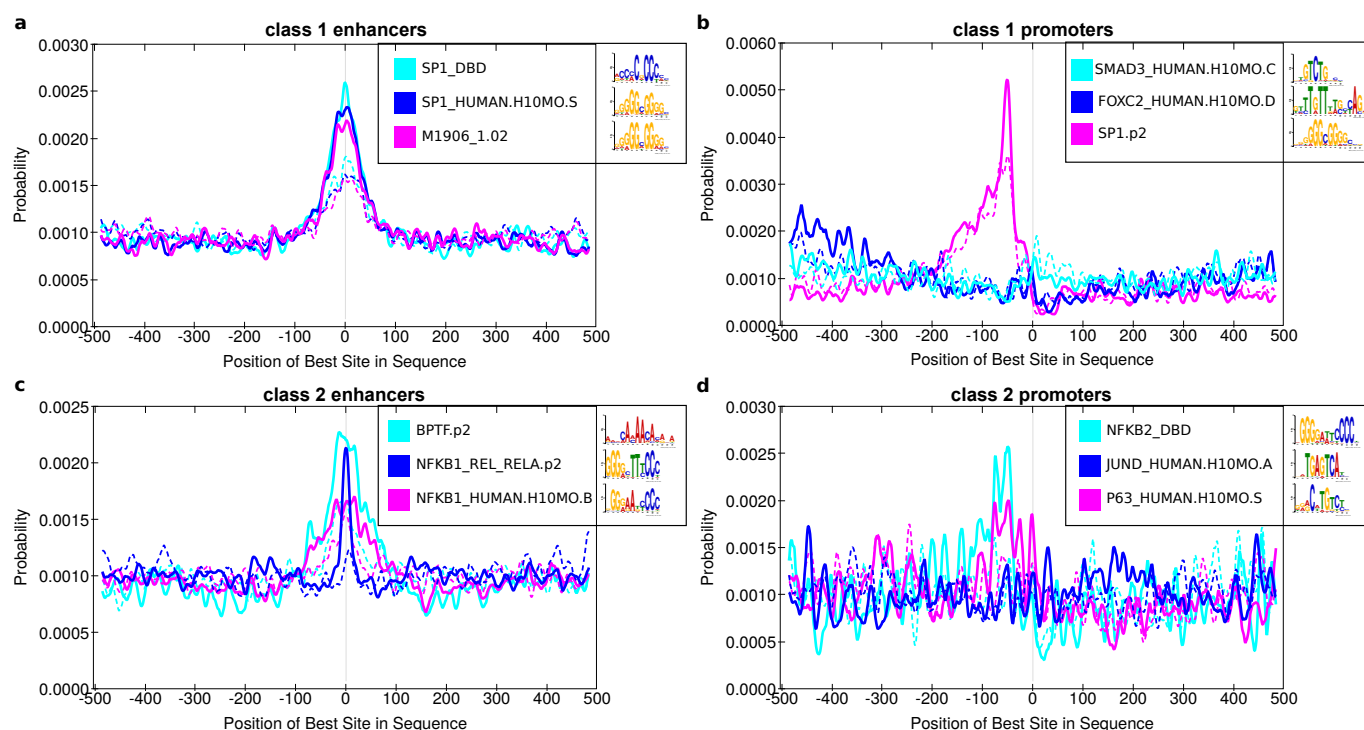


Figure 3: Motif enrichment analysis at enhancer and promoter regions. Regions of ± 500 bp around enhancer mid-points (a, c) and associated genes' TSSs (b, d) were subjected to positional motif enrichment analyses using the Centrimo tool [27]. Enhancers and associated gene targets from class 1 (a, b) and class 2 (c, d) were analyzed separately. The x-axis represents the distance to the enhancer mid-point (a, c) and associated gene TSSs (b, d), respectively. The y-axis represents the probability of predicting TFBSs associated with the motifs given in the legend boxes. Plain lines represent the distribution of predicted TFBSs in the foreground sequences (from class 1 in panels a-b and class 2 in panels c-d). Similarly, dashed lines represent the distribution of predicted TFBSs in the background sequences (from class 2 in panels a-b and class 1 in panels c-d). TSSs, transcription start sites; TFBSs, transcription factor binding sites.

cally enriched about 100 bp and 80 bp upstream of TSSs associated with class 1 and class 2 enhancers, respectively (Figure 3b,d and Data S1). It confirmed that promoters predicted to be targets of enhancers shared the same motifs. Centrimo also predicted SMAD3 and FOXC2 motifs in class 1 promoters and JUN and P63 motifs in class 2 promoters, upstream of TSSs.

We confirmed the motif-based enrichment of NF- κ B binding in class 2 regions by using ChIP-seq data obtained in GM12878 cells for the RELA TF, which is involved in NF- κ B heterodimer formation. By combining data capturing histone modification marks, TF binding, and open chromatin regions from a specific cell type, the ChromHMM [34] and Segway [35] tools segment the genome into regions associated to specific chromatin states. Focusing on predictions from ChromHMM and Segway combined, we found 1,802 ($\sim 12\%$) and 2,813 ($\sim 11\%$) active enhancer regions from classes 1 and 2, respectively. We observed that class 2 enhancers were preferentially bound by RELA. Specifically, 591 active class 1 enhancers and 1,226 active class 2 enhancers overlapped RELA ChIP-seq peaks ($p\text{-value} = 2.3 \times 10^{-13}$). A similar analysis focusing on predicted promoters identified an enrichment for active promoters in class 2 (8,962, $\sim 59\%$, class 1 and 10,752, $\sim 63\%$, class 2 active promoters, $p\text{-value} = 4.5 \times 10^{-9}$). Furthermore, class 2 promoters were preferentially bound by RELA with 1,966 class 1 and 3,179 class 2 active promoters overlapping RELA ChIP-seq peaks ($p\text{-value} < 2.2 \times 10^{-16}$).

Together, these results reinforced the predictions of biological functions specific to class 1 and class 2 enhancers (Figure 2) through the presence of associated TF binding motifs in both enhancers and predicted target promoters.

The two classes of human enhancers exhibited distinct activity patterns

We further investigated the functional differences between the two classes of human enhancers by analyzing their patterns of activity across cell types. In previous studies, enhancer activity has been inferred either from histone modifications or eRNA transcription signatures [5, 34–36]. We considered these two approaches. Namely, we considered histone modification data from 6 cell lines and CAGE data from 71 cell types produced by the ENCODE [37] and FANTOM5 [22] projects, respectively.

We retrieved the segmentation of the human genome obtained using a combination of ChromHMM and Segway in the tiers 1 and 2 cell types from ENCODE [37]. For each cell type, we overlapped enhancers with predicted genome segments to assign an activity state to the enhancer. As an example, Figure 4a presents the proportion of enhancers from classes 1 and 2 that were overlapping with segments associated with active, CTCF, and repressed chromatin states in embryonic stem cells (H1-hESC). We consistently observed that enhancers from class 1 were significantly more active than those from class 2, which were found to be enriched in repressed genomic segments (Figures 4a and S7). Class 1 enhancers were also associated with segments characterized by CTCF binding.

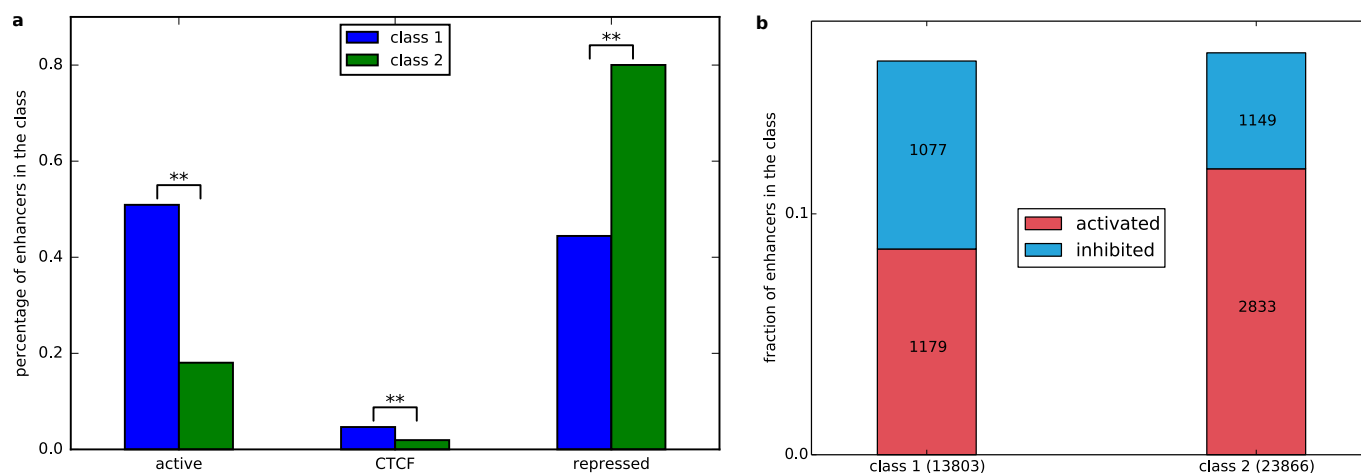


Figure 4: **Human enhancers and genome segmentation.** **a.** Histogram of the proportion of human enhancers (y-axis) in class 1 (blue) and class 2 (green) lying within genome segments (x-axis) as annotated by combined predictions from ChromHMM [34] and Segway [35] on human embryonic stem cells (H1-hESC from the ENCODE project [37]). Statistical significance (Bonferroni-corrected p-value < 0.01) of enrichment for enhancers from a specific class is indicated by '**'. **b.** Stacked histogram of the fraction of human enhancers (y-axis) from class 1 and class 2 predicted to be activated (red) or inhibited (blue). Predictions were obtained using genomic segments predicted by ChromHMM [34] on human dendritic cells before and after infection with *Mycobacterium tuberculosis* [38]. Stacked histogram including unchanged activity is provided in Figure S8.

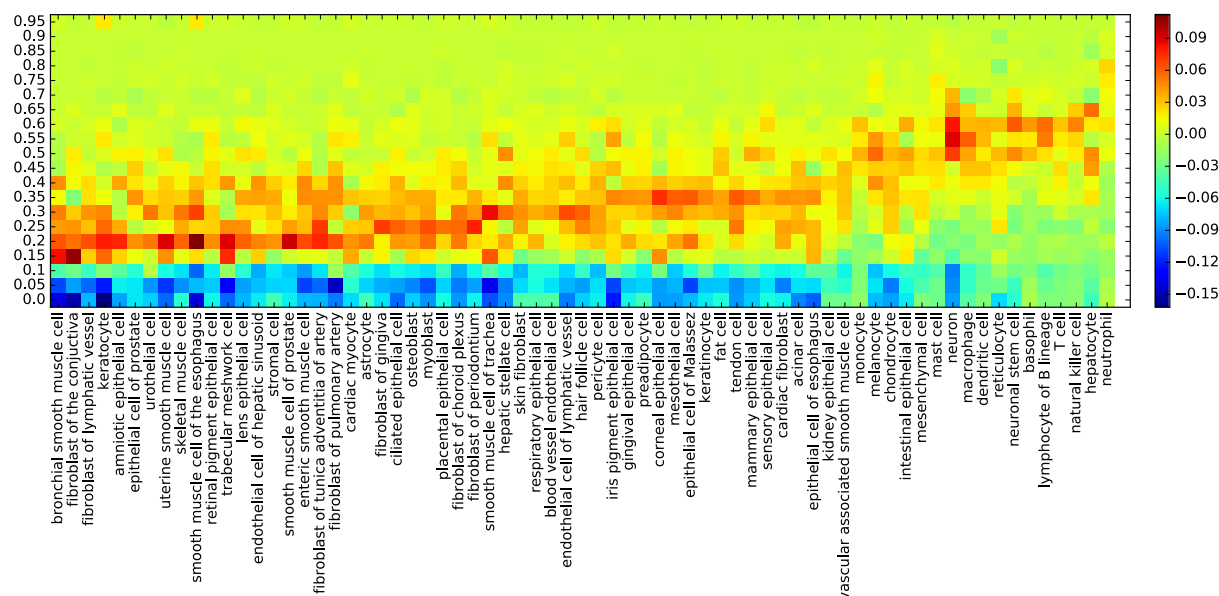


Figure 5: Cell type expression specificities of human enhancers. The difference in cell type expression specificities derived from FANTOM5 CAGE datasets [22] for enhancers in class 1 and class 2 is provided as a heat map. The color (see scale) represents the difference in fraction of expressed enhancers in each cell type (columns) found in each expression specificity range (rows). Positive (respectively negative) values indicate a higher fraction of class 2 (respectively class 1) enhancers. The heat maps corresponding to the enhancers in each class are provided in Figure S9. CAGE, Cap Analysis of Gene Expression.

285 Through CAGE expression analysis in 808 human samples, Andersson *et*
 286 *al.* [22] assigned a cell type-specificity score to human enhancers. Seventy one
 287 cell types were defined by grouping cell and tissue samples [22]. Following
 288 the enhancer expression specificity analysis performed by Andersson *et al.*, we
 289 considered enhancers from class 1 and class 2 separately to highlight potential
 290 activity differences in the 71 cell types (Figure S9). Comparing enhancer activity
 291 specificity over all the cell types between class 1 and class 2, enhancers from
 292 class 2 appeared to be more cell type specific (Figure 5). While immune cells,
 293 neurons, neuronal stem cells, and hepatocytes were previously described to use a
 294 higher fraction of human enhancers [22], the elevated utilization was even more
 295 pronounced for class 2 enhancers (Figures 5 and S9b).

296 Taken together, these results derived from histone marks and transcriptional
 297 data highlighted that enhancers from class 1 were more ubiquitously active over
 298 human cell types than enhancers from class 2, which were more cell type specific.
 299 In our previous functional analyses, we inferred the biological functions of the
 300 two classes of enhancers from the genes they were predicted to regulate. Here, we
 301 further confirmed specific functionalities for the two classes based on enhancer
 302 activity analyses, which corroborated with our functional analysis described
 303 above. Class 1 enhancers were found to be enriched in transcriptional biological
 304 processes, which are required for transcription in all cell types. Enhancers from
 305 class 2 were more cell type specific, with an emphasis in cell types associated
 306 with the immune system, in agreement with the functional enrichment analysis.

307 **Predicted immune system enhancers were activated upon** 308 **cell infection**

309 We sought to further confirm the association of class 2 enhancers with transcrip-
 310 tional control of immune responses. Pacis *et al.* [38] generated genome-wide
 311 DNA methylation, histone marks, and chromatin accessibility data in normal
 312 dendritic cells (DCs) and DCs after infection with *Mycobacterium tuberculosis*
 313 (MTB). The data provided the opportunity to study the chromatin state changes
 314 after infection obtained using the ChromHMM tool [34]. As for the above analy-
 315 sis, we overlapped chromatin state information with the enhancers from classes 1
 316 and 2. To highlight the key epigenetic changes at enhancers, we classified the
 317 transition of activities before and after MTB infection into three groups: acti-
 318 vated (from inactive before MTB infection to active after infection), inhibited
 319 (active to inactive) or unchanged (Figures 4b and S8). We observed that the
 320 enhancers from class 2 were significantly more activated ($p\text{-value} < 2.2 \times 10^{-16}$)
 321 and less inhibited ($p\text{-value} < 2.2 \times 10^{-16}$) when compared to class 1 enhancers
 322 upon MTB infection (Figure 4b). These results reinforced the potential role of
 323 class 2 enhancers in immune response.

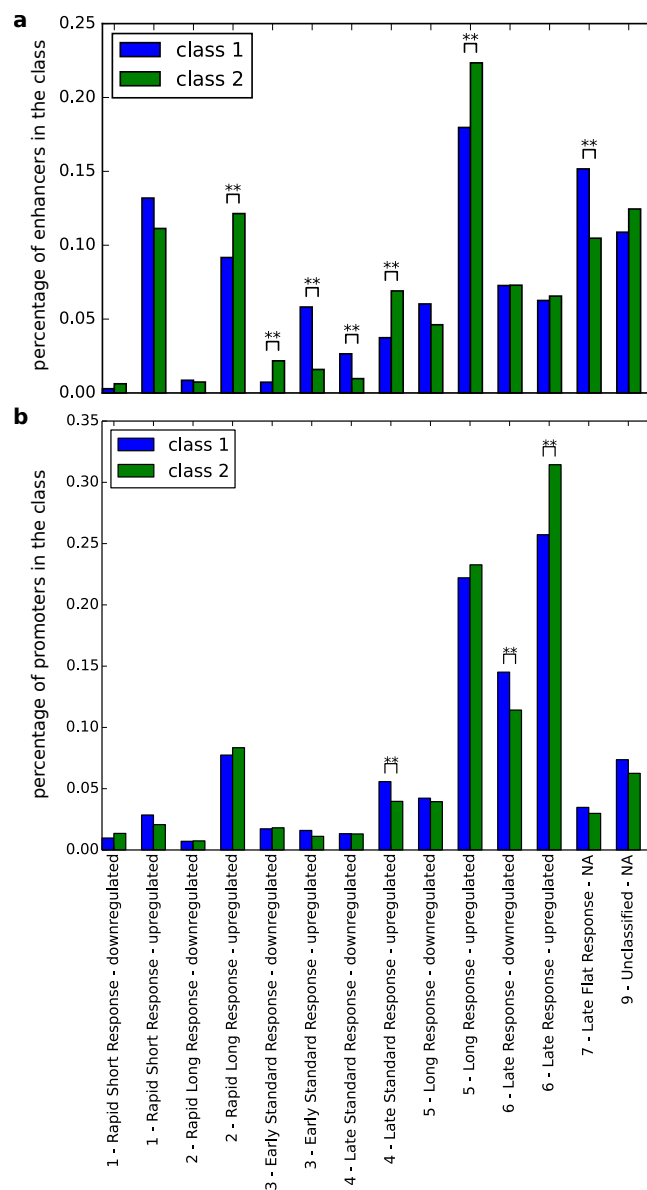


Figure 6: Expression dynamics of human enhancers and associated promoters. Response patterns (x-axis) of human enhancers (a) and promoters (b) in time courses were classified by Arner *et al.* [33]. The percentage (y-axis) of enhancers and promoters from class 1 (blue) and class 2 (green) in each response pattern category are provided as histograms in the two panels. A significant difference (Bonferroni-corrected p-value < 0.01) between class 1 and class 2 enhancers or promoters in a specific category is highlighted by '***'.

324 Predicted immune system enhancers showed long lasting 325 response activity

326 Based on time-courses of differentiation and activation, Arner *et al.* analyzed the
327 transcriptional dynamics of enhancers and promoters [33]. They profiled time-
328 courses with CAGE at a high temporal resolution within a 6 hour time-frame
329 to classify enhancers and promoters into distinct dynamic response patterns
330 of early response activity. We overlaid our classification of human enhancers
331 and their predicted target promoters with the dynamic response pattern data
332 (Figure 6). Within the enhancers associated to any dynamic response pattern
333 ($n = 2,694$; 1,533 and 1,161 from class 1 and class 2, respectively), class 1
334 enhancers were enriched (p-value $< 2.2 \times 10^{-16}$).

335 We focused on the set of 2,694 enhancers classified in the dynamic response
336 patterns. Looking at the peaks of activity specific to early time points ('rapid
337 short response' and 'early standard response'), class 1 enhancers were found
338 to be upregulated while class 2 enhancers were downregulated (Figure 6a).
339 Enhancers from class 2 showed significant activity dynamics corresponding to
340 long lasting and later responses (up-regulated in 'rapid long response', 'late
341 standard response', and 'long response') when compared to class 1 enhancers
342 (down-regulated in 'late standard response') (Figure 6a). The promoters asso-
343 ciated with class 1 were upregulated in the 'late standard response' dynamic.
344 Class 2 promoters exhibited significant up-regulation in the 'late response' dy-
345 namic while class 1 promoters were downregulated in the same dynamic.

346 Taken together, these results identified different dynamics between class 1
347 and class 2 enhancers. Class 1 enhancers were more dynamic than class 2 en-
348 hancers in the FANTOM5 time-course, activated early and for a short period of
349 time while class 2 enhancers harboured long-lasting rapid and late activities. As
350 previously observed [33], the activity of the enhancers were followed by peaks of
351 activity for the associated promoters at later stages (enrichment in late response
352 categories).

353 Enhancers from the same class co-localized within chro- 354 matin domains

355 The organization of the chromatin in cell nuclei is a key feature in gene ex-
356 pression regulation by forming regulatory region interactions within TADs [8].
357 Genes within the same TAD tend to be coordinately expressed across cell types
358 and tissues, and clusters of functionally related genes requiring co-regulation
359 tend to lie within the same TADs [8,40]. Similar to these studies analyzing gene
360 organization observed in chromatin domains, we focused on how the two classes
361 of enhancers were organized with respect to TADs. We compared the distribu-
362 tion of enhancers from the two classes within a set of TADs [7]. Specifically, we
363 assessed whether individual TADs were biased for containing more enhancers
364 associated with a specific class than expected by chance using the Binomial test.
365 The distribution of the corresponding p-values was compared to those obtained
366 by randomly assigning classes 1 and 2 labels to the enhancers. The results high-

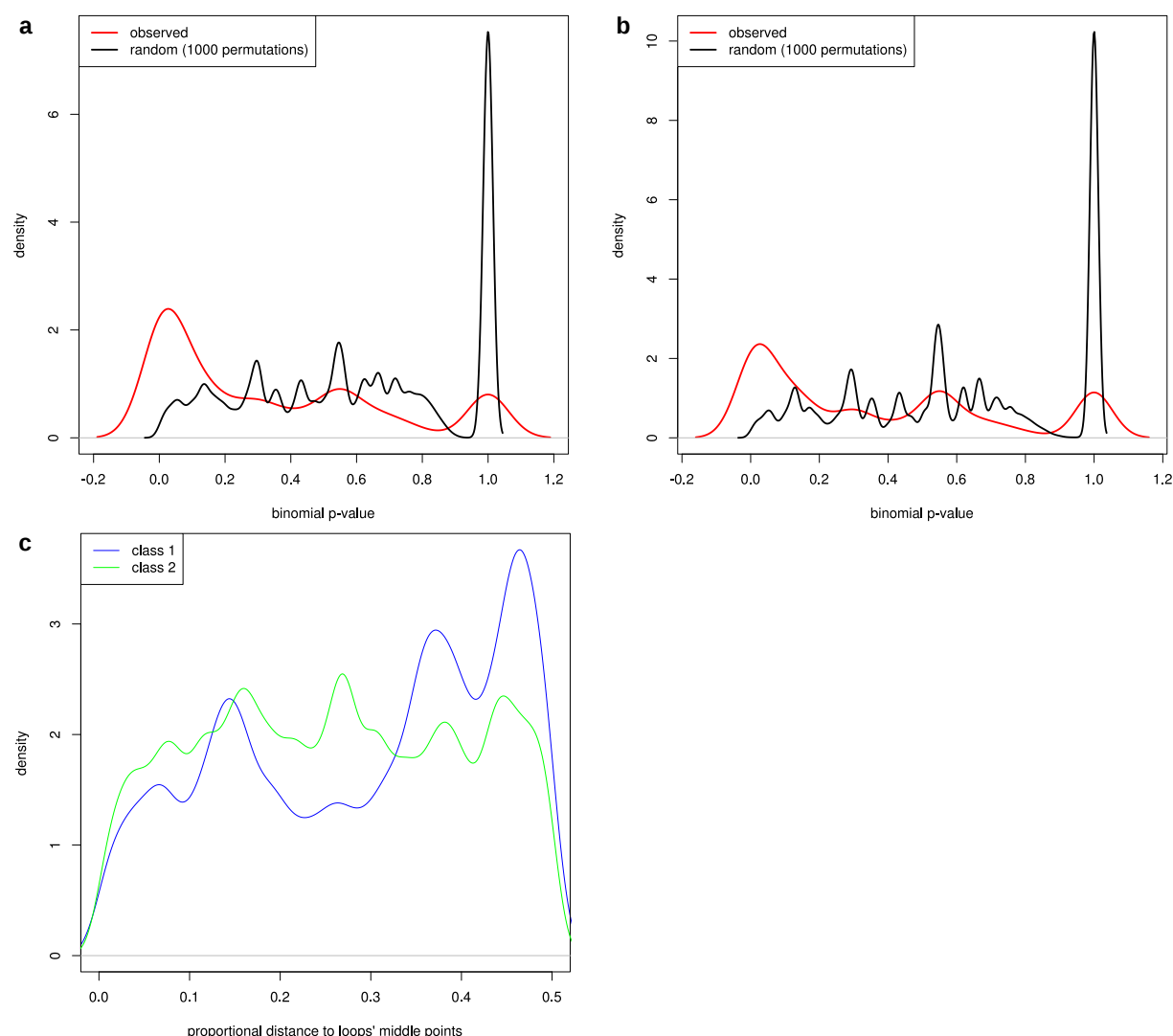


Figure 7: Chromosomal organization of class 1 and class 2 enhancers. **a.** For each TAD [7], we computed the p-value of the Binomial test to assess the enrichment for enhancers from a specific class. The plot compares the density (y-axis) of p-values for Binomial tests (x-axis) applied to classes 1 and 2 enhancers (red) and 1,000 random assignments of class labels to the enhancers (black). **b.** The same analysis as in panel a. was performed using chromatin loops predicted in lymphoblastoid GM12878 cells [39]. **c.** Density (y-axis) of distances (x-axis) between enhancers and chromatin loop centers defined using Hi-C data in GM12878 cells [39]. The distances were normalized by the length of the loops. Enhancers at the center of the loops were found at distance 0.0 while enhancers at chromatin loops boundaries were found at distance 0.5. Results associated with class 1 and class 2 enhancers are depicted in blue and green, respectively.

367 lighted that TADs were enriched for enhancers from a specific class (Figure 7a),
368 showing a genomic organization of human enhancers with respect to chromatin
369 domains.

370 TADs represent interactions within megabase-sized domains of chromatin,
371 which can be subdivided into kilobase-sized chromatin loops of chromatin inter-
372 actions [39]. We refined our analyses of class-based enhancer co-localization by
373 focusing on chromosomal loops derived from 8 cell lines [39]. Similar to what
374 we observed at the TAD level, we found that chromatin loops tended to contain
375 enhancers from a specific class (Figures 7b and S10). Furthermore, class 2 en-
376 hancers were evenly distributed within the chromatin loops whereas enhancers
377 from class 1 were consistently observed to be situated close to the loop bound-
378 aries (Figure 7c). This observation is in agreement with the enrichment for
379 class 1 enhancers in CTCF chromatin segments (Figure 4a) as chromatin loop
380 boundaries are known to be enriched for CTCF binding [39].

381 Discussion

382 We have analyzed the sequence properties of FANTOM5 human enhancers de-
383 rived from CAGE experiments to reveal that a subset with low G+C content is
384 associated with immune response genes. This set of enhancers harbours a G+C
385 pattern that corresponds to characteristic HelT, MGW, ProT and Roll confor-
386 mation of the DNA. The predicted immune system enhancers tend to co-localize
387 within chromatin domains, exhibit cell type specificity, are activated upon in-
388 fection, and are observed with long lasting response activity. In summary, our
389 study of sequence composition patterns along enhancer regions culminates with
390 the identification of human enhancers associated with immune response that
391 harbour specific sequence composition, activity, and genome organization.

392 The analyses of sequence properties in regulatory regions, most prominently
393 CpG islands at promoters, have been key to understanding gene expression
394 regulation [9, 10, 20]. The predicted immune response enhancers exhibit a cell-
395 type specific expression pattern and have low %GC. Nevertheless, it remains
396 unclear how and why immune response enhancers have emerged with these
397 sequence properties.

398 While enhancers predicted to be associated with immune response are GC-
399 poor, a dichotomy of enhancers solely based on GC-content did not highlight a
400 specific set of enhancers associated with immune response genes. This observa-
401 tion might reflect the importance of the DNA sequence, with the DNA shape
402 conformation at enhancer regions as a secondary effect. DNA structural proper-
403 ties were shown to be linked with DNA flexibility, nucleosome positioning, and
404 gene expression regulation [41–45]. The more negative ProT in class 2 (Fig-
405 ure 1e) corresponds to lower GC content [46]. This could relate to 3 hydrogen
406 bonds in G/C pairs versus 2 in A/T pairs, which determines the ability of a base
407 pair to form a ProT angle. Moreover, the differences in DNA shape features be-
408 tween the two classes of enhancers might relate to differences in conformational
409 flexibility. Indeed, we observed class 1 enhancers with less negative ProT, lower

HelT, wider MGW, larger Roll compared to class 2 enhancers (Figure 1c-f). These characteristics all relate to increased flexibility of the DNA [47], which could provide a topological explanation for the differences observed between the two classes.

The classification of human enhancers was performed from basic feature vectors summarizing G+C patterns along enhancer regions. Our observations are reminiscent of the "enhancer-core-promoter specificity" observed in *Drosophila* [48]. Zabidi *et al.* uncovered two classes of enhancers regulating "housekeeping" versus "developmental" genes, which differ in genomic distribution and in the presence of distinct regulatory elements [48]. Each enhancer appears to have acquired specific DNA features to most effectively regulate the particular promoters it has to regulate [49]. This hypothesis is in agreement with our identification of two classes of enhancers that were defined based on DNA sequence composition and predicted to target promoters associated with genes of distinct biological functions. Similar to Zabidi *et al.* [48], we found distinct genomic organizations between two classes of enhancers, one more ubiquitously active while the other was more cell type-specific.

A recent perspective on immunological memory suggested that transient modifications to chromatin along with inducible noncoding RNAs could mediate a "short-term memory" in immune cells [50]. The principle of this short-term memory is that some histone modifications and regulatory molecules like microRNAs and TFs would be persistent after stimulation, even though limited in time. Mediators of response to stimuli were categorized into (i) labile mediators of activation and (ii) long-lasting mediators of short term memory [50]. It represents a way for both the adaptive and innate immune cells to be more effective in responding to secondary stimulations. Our results suggest that enhancers could also be mediators of the immune response to stimuli lying in the two categories presented by Monticelli *et al.* [50]. Indeed, we observed that class 2 enhancers, associated with immune response genes, were showing a long-term response activity to stimuli, as opposed to short-term patterns of activity for class 1 enhancers (Figure 6). We hypothesize that (i) class 1 enhancers are used for a rapid but short response to stimuli, representing labile mediators of activation and (ii) class 2 enhancers correspond to robust long-lasting mediator of short-term memory to memorize which genes need to be activated after stimulation. Dedicated experiments will be necessary to assess this hypothesis.

Materials and Methods

Human enhancers clusterization

We retrieved the hg19 positions of the 38,554 FANTOM5 robust human enhancers in BED12 format from http://enhancer.binf.ku.dk/presets/robust_enhancers.bed [22]. We extracted DNA sequences for regions of 1,001 bp centered at the enhancer mid-points (columns 7-8 of the BED12 file) using the BEDTools [51]. We created binary vectors representing the enhancer sequences

with 1s and 0s corresponding to G or C and A or T, respectively. Note that 7 enhancers were not considered as the 1,001 bp regions contained undefined nucleotides (Ns). The vectors were clustered into $k = 2$ classes using the *k*-means algorithm implemented in the *KMeans* function of the *scikit* Python module [52]. The silhouette plots (Figure S1) were constructed for $k \in [2, 5]$ using the *silhouette_samples* function of the *scikit* Python module. Formally, the silhouette plots display the silhouette coefficient for each enhancer as $\frac{(b-a)}{\max(a,b)}$ where a is the mean intra-cluster euclidian distance and b the mean nearest-cluster euclidian distance.

Similarly, we created vectors representing the enhancers by combining the values of HelT, MGW, ProT, and Roll at the enhancer sequences. The DNA shape values for the hg19 version of the human genome were retrieved as bigWig files from the GBshape database [11]. DNA shape feature values at enhancer regions were obtained by using *bwtool* [53]. It resulted in vectors of 4,004 values each that were submitted to the *Kmeans* function of the *scikit* Python module.

Enhancer gene targets

The enhancer-RefSeq promoter associations were retrieved from http://enhancer.binf.ku.dk/presets/enhancer_tss_associations.bed [22]. The corresponding official gene symbols were considered for the functional enrichment analyses.

DNA shape feature plots

The values of DNA structural features HelT, MGW, ProT, and Roll computed using the DNashape tool [12] were obtained from the GBshape browser [11] as bigwig files at <ftp://rohslab.usc.edu/hg19/>. We retrieved the averaged DNA shape values at the enhancer regions from class 1 and class 2 using the *agg* subcommand of the *bwtool* tool [53].

Gene ontology functional enrichment

Official symbols corresponding to the RefSeq promoters associated with enhancers from class 1 and class 2 were submitted to GOrilla [26] at <http://cbl-gorilla.cs.technion.ac.il/> using the March 5th 2016 update. We used the two unranked list option with genes associated with class 1 or class 2 enhancers as targets and the aggregated set of 11,271 genes associated with the full set of enhancers as background. We searched for enriched GO biological processes with the most stringent p-value threshold ($< 10^{-11}$). The DAG representation of the results in Figures S3-S4 were downloaded from the output page of GOrilla. The visual representation of the results in Figure 2 was constructed manually using Cytoscape 3.4.0 [54]. The same procedure has been applied to genes associated with enhancers classified with respect to their GC-content (Figures S2-S3).

490 Motif enrichment

491 We applied Centrimo [27] from the MEME suite version 4.11.1 with default
 492 parameters to DNA sequences of regions ± 500 bp around the mid-points of
 493 enhancers from class 1 and class 2. Class 1 enhancer regions were used as
 494 foreground and class 2 enhancer regions as background and vice-versa. The
 495 MEME databases of motifs considered for enrichment were derived from [55]
 496 (jolma2013.meme), JASPAR [56] (JASPAR_CORE.2016.vertebrates.meme), Cis-
 497 BP [57] (Homo_sapiens.meme), Swiss Regulon [58] (Swiss_Regulon_human_and_mouse.meme),
 498 and HOCOMOCO [59] (HOCOMOCOv10_HUMAN_mono_meme_format.meme).
 499 The same procedure was applied to promoter regions (± 500 bp around TSSs)
 500 associated with class 1 and class 2 enhancers.

501 Figure 3 has been obtained from the html output of Centrimo by selecting
 502 the 3 most enriched motifs (ranked using the Fisher E-value). We did not
 503 consider inferred motifs in Cis-BP [57].

504 Genome segmentation

505 ENCODE genome segmentation

506 The genome segmentation using the combination of results from ChromHMM [34]
 507 and Segway [35] for ENCODE tier 1 and tier 2 cell types GM12878, H1hesc,
 508 HeLaS3, HepG2, HUVEC, and K562 were retrieved at [http://hgdownload.
 509 cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgSegmentation/](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgSegmentation/).

510 Genome segmentation in dendritic cells

511 The genome segmentation of DCs before and after MTB infection [38] was
 512 computed using ChromHMM [34] and retrieved at [http://132.219.138.157:
 513 8080/DC_NI_7_segments_modID.bed.gz](http://132.219.138.157:8080/DC_NI_7_segments_modID.bed.gz) and [http://132.219.138.157:8080/
 514 DC_MTB_7_segments_modID.bed.gz](http://132.219.138.157:8080/DC_MTB_7_segments_modID.bed.gz).

515 Genome segmentation overlap with enhancers

516 The overlap between enhancers and genome segments were obtained using the
 517 *intersect* subcommand of the BEDTools requiring a minimum overlap of 50%
 518 of the enhancer lengths. We considered enhancers as in active states if they
 519 overlapped the TSS, promoter flank, enhancer, weak enhancer, and transcribed
 520 segments.

521 RELA ChIP-seq data analyses

522 The ENCODE RELA ChIP-seq data in GM12878 cells was retrieved at [http://
 523 hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgTfbsUniform/
 524 wgEncodeAwgTfbsSydhGm12878NfkbTnfaIggrabUniPk.narrowPeak.gz](http://hgdownload.cse.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeAwgTfbsUniform/wgEncodeAwgTfbsSydhGm12878NfkbTnfaIggrabUniPk.narrowPeak.gz). To iden-
 525 tify active FANTOM5 enhancers in GM12878, we considered the overlap be-
 526 tween 1,001 bp-long regions around enhancer's mid-points and genome segments

527 predicted by ChromHMM and Segway combined as enhancer or weak enhancer.
 528 Similarly, active FANTOM5 promoters were obtained by overlapping 1,001 bp-
 529 long regions around TSSs and genome segments predicted as TSS or promoter
 530 flanks by ChromHMM and Segway combined. The identified 1,001 bp-long ac-
 531 tive enhancer and promoter regions were further overlapped with RELA ChIP-
 532 seq peaks. All overlaps were computed with the *intersect* subcommand of the
 533 BEDTools.

534 **Enhancer expression specificity**

The cell-type expression specificity of enhancers was computed as

$$\frac{\text{entropy}(\text{enhancer expression})}{\log_2(\text{number of cell types})}$$

535 in [22]. The binary matrix of enhancer usage across FANTOM5 samples was ob-
 536 tained at [http://enhancer.binf.ku.dk/presets/hg19_permissive_enhancer_](http://enhancer.binf.ku.dk/presets/hg19_permissive_enhancer_usage.csv.gz)
 537 [usage.csv.gz](http://enhancer.binf.ku.dk/presets/hg19_permissive_enhancer_usage.csv.gz). The association between FANTOM5 samples and cell types was
 538 obtained from Tables S10-S11 in [22]. Heat maps in Figure 5 were computed
 539 using the *colormesh* function of the *matplotlib.pyplot* Python module [60].

540 **Enhancer dynamics**

541 FANTOM5 classification in the 14 dynamics displayed in Figure 6 was obtained
 542 from Auxiliary data table S3 in [33]. The classification provided response class
 543 assignments to 1,533 and 1,161 class 1 and class 2 enhancers, respectively. Re-
 544 sponse classes were assigned to 2,311 and 2,407 promoters associated with class 1
 545 and class 2 enhancers, respectively. Note that enhancers and promoters can be
 546 assigned to multiple response classes.

547 Corresponding plots (Figures 6) and enrichment analyses were performed
 548 using *pandas* Python data structure [61] and the *scipy* Python library [62] in
 549 the *IPython* environment [63].

550 **Chromatin conformation data**

551 The enrichment for enhancers associated to a specific class in each TAD or
 552 chromatin domain (see below) was computed using Binomial test p-values as
 553 implemented by the *binom.test* function in the *R* environment [64]. As a con-
 554 trol, we randomly assigned the labels class 1 and class 2 to the enhancers and
 555 computed the corresponding Binomial test p-values; this procedure was applied
 556 to 1,000 random trials.

557 **Topologically associating domains**

558 As TADs have been shown to be conserved between cell types and species,
 559 we retrieved the TADs defined in the first study describing them [7]. The
 560 TADs were predicted in mouse embryonic stem cells and we used the *liftOver*

561 tool from the UCSC genome browser at <https://genome.ucsc.edu/cgi-bin/hgLiftOver> to map them to hg19 coordinates.

563 Chromatin loops

564 The positions of the chromatin loops computed with the HICCUPS tools [39]
565 from Hi-C data on the GM12878, HMEC, HUVEC, HeLa, IMR90, K562, KBM7,
566 and NHEK human cell lines were retrieved from GEO at <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE63525>.

568 Enrichment p-values

569 P-values throughout the manuscript were computed using the Fisher exact test
570 except otherwise stated.

571 Acknowledgements

572 As research parasites [65], we are indebted to the researchers around the globe
573 who generated experimental data and made them freely available. We thank
574 Miroslav Hatas and Georgios Magklaras for systems support, Dora Pak for man-
575 agement support, Chih-Yu Chen for providing the code for enrichment analyses
576 in chromatin conformation data, Tsu-Pei Chiu for his help with DNA shape data
577 and tools and for insightful discussions and Robin Andersson for his help with
578 FANTOM5 enhancer data. CHL was supported by funding from CNRS, *Plan*
579 *d'Investissement d'Avenir* and *Institut de Biologie Computationnelle* (Young
580 Investigator grant). AM and WWW were supported by the Genome Canada
581 Large Scale Applied Research Grant 174CDE. Funding was provided by the
582 Child and Family Research Institute and the British Columbia Children's Hos-
583 pital Foundation, Vancouver, to AM and WWW. AM was also supported by
584 funding from the Norwegian Research Council, Helse Sør-Øst, and the Univer-
585 sity of Oslo through the Centre for Molecular Medicine Norway (NCMM) and
586 the Oslo University Hospital Radiumhospitalet.

587 Author contributions

588 CHL and AM conceived and designed the project. CHL and AM implemented
589 and performed experiments. CHL, WWW, and AM analyzed and interpreted
590 the results. RR was involved in the interpretation of the results relative to the
591 DNA shape features. CHL, WWW, and AM wrote the manuscript.

592 References

- 593 [1] W. W. Wasserman and A. Sandelin, "Applied bioinformatics for the identi-
594 fication of regulatory elements," *Nature Reviews Genetics*, vol. 5, pp. 276–
595 287, Apr. 2004.

- 596 [2] A. Mathelier, W. Shi, and W. W. Wasserman, "Identification of altered cis-
597 regulatory elements in human disease," *Trends in Genetics*, vol. 31, no. 2,
598 pp. 67–76, 2015.
- 599 [3] A. Visel, M. J. Blow, Z. Li, *et al.*, "ChIP-seq accurately predicts tissue-
600 specific activity of enhancers," *Nature*, vol. 457, pp. 854–858, Feb. 2009.
601 00964.
- 602 [4] D. Babu and M. J. Fullwood, "3D genome organization in health and dis-
603 ease: emerging opportunities in cancer translational medicine," *Nucleus*
604 (*Austin, Tex.*), vol. 6, pp. 382–393, Sept. 2015. 00002.
- 605 [5] R. Andersson, A. Sandelin, and C. G. Danko, "A unified architecture of
606 transcriptional regulatory elements," *Trends in genetics: TIG*, vol. 31,
607 pp. 426–433, Aug. 2015. 00013.
- 608 [6] E. Lieberman-Aiden, N. L. v. Berkum, L. Williams, *et al.*, "Comprehen-
609 sive Mapping of Long-Range Interactions Reveals Folding Principles of the
610 Human Genome," *Science*, vol. 326, pp. 289–293, Oct. 2009.
- 611 [7] J. R. Dixon, S. Selvaraj, F. Yue, *et al.*, "Topological domains in mammalian
612 genomes identified by analysis of chromatin interactions," *Nature*, vol. 485,
613 pp. 376–380, May 2012. 01078.
- 614 [8] J. R. Dixon, D. U. Gorkin, and B. Ren, "Chromatin Domains: The Unit
615 of Chromosome Organization," *Molecular Cell*, vol. 62, pp. 668–680, June
616 2016.
- 617 [9] D. Tillo and T. R. Hughes, "G+C content dominates intrinsic nucleosome
618 occupancy," *BMC Bioinformatics*, vol. 10, p. 442, Dec. 2009.
- 619 [10] A. Hughes and O. J. Rando, "Chromatin 'programming' by sequence - is
620 there more to the nucleosome code than %GC?," *Journal of Biology*, vol. 8,
621 p. 96, 2009.
- 622 [11] T.-P. Chiu, L. Yang, T. Zhou, *et al.*, "GBshape: a genome browser database
623 for DNA shape annotations," *Nucleic Acids Research*, vol. 43, no. D1,
624 pp. D103–D109, 2015.
- 625 [12] T. Zhou, L. Yang, Y. Lu, *et al.*, "DNASHape: a method for the high-
626 throughput prediction of DNA structural features on a genomic scale,"
627 *Nucleic Acids Research*, vol. 41, pp. W56–W62, May 2013.
- 628 [13] R. Rohs, S. M. West, A. Sosinsky, *et al.*, "The role of DNA shape in
629 proteinDNA recognition," *Nature*, vol. 461, pp. 1248–1253, Oct. 2009.
- 630 [14] T. Zhou, N. Shen, L. Yang, *et al.*, "Quantitative modeling of transcription
631 factor binding specificities using DNA shape," *Proceedings of the National*
632 *Academy of Sciences*, vol. 112, pp. 4654–4659, Mar. 2015.

- 633 [15] M. Slattery, T. Zhou, L. Yang, *et al.*, “Absence of a simple code: how
634 transcription factors read the genome,” *Trends in Biochemical Sciences*,
635 vol. 39, pp. 381–399, Sept. 2014.
- 636 [16] A. Mathelier, B. Xin, T.-P. Chiu, *et al.*, “DNA Shape Features Improve
637 Transcription Factor Binding Site Predictions In Vivo,” *Cell Systems*, no. 3,
638 pp. 1–9, 2016.
- 639 [17] R. Kodzius, M. Kojima, H. Nishiyori, *et al.*, “CAGE: cap analysis of gene
640 expression,” *Nature methods*, vol. 3, pp. 211–222, Mar. 2006.
- 641 [18] T. Shiraki, S. Kondo, S. Katayama, *et al.*, “Cap analysis gene expression
642 for high-throughput analysis of transcriptional starting point and identifica-
643 tion of promoter usage,” *Proceedings of the National Academy of Sciences*,
644 vol. 100, no. 26, pp. 15776–15781, 2003.
- 645 [19] P. Carninci, T. Kasukawa, S. Katayama, *et al.*, “The Transcriptional Land-
646 scape of the Mammalian Genome,” *Science*, vol. 309, pp. 1559–1563, Sept.
647 2005.
- 648 [20] V. B. Bajic, S. L. Tan, A. Christoffels, *et al.*, “Mice and Men: Their Pro-
649 moter Properties,” *PLOS Genet*, vol. 2, p. e54, Apr. 2006.
- 650 [21] The FANTOM Consortium and the RIKEN PMI and CLST (dgt), “A
651 promoter-level mammalian expression atlas,” *Nature*, vol. 507, pp. 462–
652 470, Mar. 2014.
- 653 [22] R. Andersson, C. Gebhard, I. Miguel-Escalada, *et al.*, “An atlas of active
654 enhancers across human cell types and tissues,” *Nature*, vol. 507, pp. 455–
655 461, Mar. 2014.
- 656 [23] T. R. O’Connor and T. L. Bailey, “Creating and validating cis-regulatory
657 maps of tissue-specific gene expression regulation,” *Nucleic Acids Research*,
658 vol. 42, pp. 11000–11010, Jan. 2015.
- 659 [24] J. MacQueen, “Some methods for classification and analysis of multivariate
660 observations,” in *Proceedings of the Fifth Berkeley Symposium on Mathe-*
661 *matical Statistics and Probability, Volume 1: Statistics*, (Berkeley, Calif.),
662 pp. 281–297, University of California Press, 1967.
- 663 [25] P. J. Rousseeuw, “Silhouettes: A graphical aid to the interpretation and
664 validation of cluster analysis,” *Journal of Computational and Applied Math-*
665 *ematics*, vol. 20, pp. 53–65, Nov. 1987.
- 666 [26] E. Eden, R. Navon, I. Steinfeld, *et al.*, “GORilla: a tool for discovery and
667 visualization of enriched GO terms in ranked gene lists,” *BMC Bioinfor-*
668 *matics*, vol. 10, p. 48, Feb. 2009.
- 669 [27] T. L. Bailey and P. Machanick, “Inferring direct DNA binding from ChIP-
670 seq,” *Nucleic Acids Research*, vol. 40, pp. e128–e128, May 2012.

- [28] J. S. Presnell, C. E. Schnitzler, and W. E. Browne, “KLF/SP Transcription Factor Family Evolution: Expansion, Diversification, and Innovation in Eukaryotes,” *Genome Biology and Evolution*, vol. 7, pp. 2289–2309, Aug. 2015.
- [29] H.-C. Liou, ed., *NF-[kappa]B/Rel transcription factor family*. Molecular biology intelligence unit, Georgetown, Tex. : New York, N.Y: Landes Bioscience/Eurekah.com ; Springer Science+Business Media, 2006. OCLC: ocm68133074.
- [30] A. Itoh-Nakadai, R. Hikota, A. Muto, *et al.*, “The transcription repressors Bach2 and Bach1 promote B cell development by repressing the myeloid program,” *Nature Immunology*, vol. 15, pp. 1171–1180, Dec. 2014.
- [31] M. H. Jones, N. Hamana, and M. Shimane, “Identification and Characterization of BPTF, a Novel Bromodomain Transcription Factor,” *Genomics*, vol. 63, pp. 35–39, Jan. 2000.
- [32] W. Wang, Y. Xue, S. Zhou, *et al.*, “Diversity and specialization of mammalian SWI/SNF complexes,” *Genes & Development*, vol. 10, pp. 2117–2130, Sept. 1996.
- [33] E. Arner, C. O. Daub, K. Vitting-Seerup, *et al.*, “Transcribed enhancers lead waves of coordinated transcription in transitioning mammalian cells,” *Science*, vol. 347, pp. 1010–1014, Feb. 2015.
- [34] J. Ernst and M. Kellis, “ChromHMM: automating chromatin-state discovery and characterization,” *Nature Methods*, vol. 9, pp. 215–216, Mar. 2012.
- [35] M. M. Hoffman, O. J. Buske, J. Wang, *et al.*, “Unsupervised pattern discovery in human chromatin structure through genomic segmentation,” *Nature Methods*, vol. 9, pp. 473–476, May 2012.
- [36] G. Natoli and J.-C. Andrau, “Noncoding Transcription at Enhancers: General Principles and Functional Models,” *Annual Review of Genetics*, vol. 46, no. 1, pp. 1–19, 2012.
- [37] I. Dunham, A. Kundaje, S. F. Aldred, *et al.*, “An integrated encyclopedia of DNA elements in the human genome,” *Nature*, vol. 489, pp. 57–74, Sept. 2012.
- [38] A. Pacis, L. Tailleux, A. M. Morin, *et al.*, “Bacterial infection remodels the DNA methylation landscape of human dendritic cells,” *Genome Research*, vol. 25, pp. 1801–1811, Dec. 2015.
- [39] S. S. P. Rao, M. H. Huntley, N. C. Durand, *et al.*, “A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping,” *Cell*, vol. 159, pp. 1665–1680, Dec. 2014.

- 708 [40] J. H. Gibcus and J. Dekker, "The Hierarchy of the 3D Genome," *Molecular*
709 *Cell*, vol. 49, pp. 773–782, Mar. 2013.
- 710 [41] S. C. J. Parker, L. Hansen, H. O. Abaan, *et al.*, "Local DNA Topography
711 Correlates with Functional Noncoding Regions of the Human Genome,"
712 *Science*, vol. 324, pp. 389–392, Apr. 2009.
- 713 [42] M. Bansal, A. Kumar, and V. R. Yella, "Role of DNA sequence based struc-
714 tural features of promoters in transcription initiation and gene expression,"
715 *Current Opinion in Structural Biology*, vol. 25, pp. 77–85, Apr. 2014.
- 716 [43] T. Raveh-Sadka, M. Levo, U. Shabi, *et al.*, "Manipulating nucleosome dis-
717 favoring sequences allows fine-tune regulation of gene expression in yeast,"
718 *Nature Genetics*, vol. 44, pp. 743–750, July 2012.
- 719 [44] I. Tirosh, J. Berman, and N. Barkai, "The pattern and evolution of yeast
720 promoter bendability," *Trends in Genetics*, vol. 23, pp. 318–321, July 2007.
- 721 [45] K. Struhl and E. Segal, "Determinants of nucleosome positioning," *Nature*
722 *Structural & Molecular Biology*, vol. 20, pp. 267–273, Mar. 2013.
- 723 [46] I. Dror, T. Golan, C. Levy, *et al.*, "A widespread role of the motif envi-
724 ronment on transcription factor binding across diverse protein families,"
725 *Genome Research*, p. gr.184671.114, July 2015.
- 726 [47] F. Comoglio, T. Schlumpf, V. Schmid, *et al.*, "High-Resolution Profiling of
727 Drosophila Replication Start Sites Reveals a DNA Shape and Chromatin
728 Signature of Metazoan Origins," *Cell Reports*, vol. 11, no. 5, pp. 821–834,
729 2015.
- 730 [48] M. A. Zabidi, C. D. Arnold, K. Schernhuber, *et al.*, "Enhancer-core-
731 promoter specificity separates developmental and housekeeping gene regu-
732 lation," *Nature*, vol. 518, pp. 556–559, Feb. 2015. 00049.
- 733 [49] D. S. Lorberbaum and S. Barolo, "Enhancers: holding out for the right
734 promoter," *Current biology: CB*, vol. 25, pp. R290–293, Mar. 2015. 00001.
- 735 [50] S. Monticelli and G. Natoli, "Short-term memory of danger signals and en-
736 vironmental stimuli in immune cells," *Nature Immunology*, vol. 14, pp. 777–
737 784, Aug. 2013.
- 738 [51] A. R. Quinlan and I. M. Hall, "BEDTools: a flexible suite of utilities for
739 comparing genomic features," *Bioinformatics (Oxford, England)*, vol. 26,
740 pp. 841–842, Mar. 2010.
- 741 [52] F. Pedregosa, G. Varoquaux, A. Gramfort, *et al.*, "Scikit-learn: Machine
742 learning in Python," *The Journal of Machine Learning Research*, vol. 12,
743 pp. 2825–2830, 2011.
- 744 [53] A. Pohl and M. Beato, "bwtool: a tool for bigWig files," *Bioinformatics*,
745 vol. 30, pp. 1618–1619, June 2014.

- 746 [54] P. Shannon, A. Markiel, O. Ozier, *et al.*, “Cytoscape: a software environ-
747 ment for integrated models of biomolecular interaction networks,” *Genome*
748 *Research*, vol. 13, pp. 2498–2504, Nov. 2003.
- 749 [55] A. Jolma, J. Yan, T. Whittington, *et al.*, “DNA-Binding Specificities of
750 Human Transcription Factors,” *Cell*, vol. 152, pp. 327–339, Jan. 2013.
- 751 [56] A. Mathelier, O. Fornes, D. J. Arenillas, *et al.*, “JASPAR 2016: a major
752 expansion and update of the open-access database of transcription factor
753 binding profiles,” *Nucleic Acids Research*, vol. 44, pp. D110–D115, Nov.
754 2015.
- 755 [57] M. T. Weirauch, A. Yang, M. Albu, *et al.*, “Determination and Inference
756 of Eukaryotic Transcription Factor Sequence Specificity,” *Cell*, vol. 158,
757 pp. 1431–1443, Sept. 2014.
- 758 [58] M. Pachkov, I. Erb, N. Molina, and E. v. Nimwegen, “SwissRegulon: a
759 database of genome-wide annotations of regulatory sites,” *Nucleic Acids*
760 *Research*, vol. 35, pp. D127–D131, Jan. 2007.
- 761 [59] I. V. Kulakovskiy, Y. A. Medvedeva, U. Schaefer, *et al.*, “HOCOMOCO: a
762 comprehensive collection of human transcription factor binding sites mod-
763 els,” *Nucleic Acids Research*, vol. 41, pp. D195–D202, Nov. 2012.
- 764 [60] J. D. Hunter, “Matplotlib: A 2d Graphics Environment,” *Computing in*
765 *Science Engineering*, vol. 9, pp. 90–95, May 2007.
- 766 [61] W. McKinney, “Data Structures for Statistical Computing in Python,”
767 pp. 51–56, 2010.
- 768 [62] E. Jones, T. Oliphant, P. Peterson, *et al.*, “SciPy: Open source scientific
769 tools for Python,” 2001–.
- 770 [63] F. Perez and B. E. Granger, “IPython: A System for Interactive Scientific
771 Computing,” *Computing in Science and Engineering*, vol. 9, no. 3, pp. 21–
772 29, 2007.
- 773 [64] R Core Team, *R: A Language and Environment for Statistical Computing*.
774 R Foundation for Statistical Computing, Vienna, Austria, 2016.
- 775 [65] D. L. Longo and J. M. Drazen, “Data Sharing,” *New England Journal of*
776 *Medicine*, vol. 374, pp. 276–277, Jan. 2016.

777 Supplementary Figure and Table legends

Figure S1: **Silhouette plots of k-means clusters.** Silhouette plots for clusters obtained using the k-means clusterization algorithm for $k = 2$ (a), $k = 3$ (b), $k = 4$ (c), and $k = 5$ (d). After clusterization of the enhancers using the k-means algorithm, the silhouette score was computed for each enhancer vector. For each k-means clusterization corresponding to each panel, clusters are represented with different colors. The scores range from -1 to 1 with -1 indicating a possible assignment of the enhancer vector to the wrong cluster, 0 indicating that the vector is close to the boundary between two clusters, and 1 indicating that the vector is far away from the boundary between two clusters. The silhouette score is calculated as $(b - a) / \max(a, b)$ where a is the mean intra-cluster distance and b the mean nearest-cluster distance for each sample as implemented in the scikit learn *silhouette_score* function. The red dashed lines represent the average silhouette score over all the enhancer vectors (0.019 for $k = 2$, 0.009 for $k = 3$, 0.005 for $k = 4$, and 0.002 for $k = 5$).

Figure S2: **DNA shape features at enhancers from DNA shape-based clusterization.** DNA shape feature values are provided for human enhancers from set 1 and set 2 in blue and green, respectively. Average DNA shape values (y-axis) along the DNA regions $\pm 2,000$ bp centered at enhancer mid-points (x-axis) for DNA shape features HelT (a), MGW (b), ProT (c), and Roll (d).

Figure S3: **Functional enrichment for genes associated to class 1 enhancers.** Directed acyclic graph of the enriched GO biological processes obtained using GOrilla on the set of genes predicted to be regulated by enhancers from class 1.

Figure S4: **Functional enrichment for genes associated to class 2 enhancers.** Directed acyclic graph of the enriched GO biological processes obtained using GOrilla on the set of genes predicted to be regulated by enhancers from class 2.

Figure S5: **Functional enrichment for genes associated to GC-poor enhancers.** Genes predicted to be regulated by %GC poor (%GC < 46.15) were submitted to GOrilla. The figure represents the directed acyclic graph of the enriched GO biological processes obtained.

Figure S6: **Functional enrichment for genes associated to GC-rich enhancers.** Genes predicted to be regulated by %GC rich (%GC > 46.15) were submitted to GOrilla. The figure represents the directed acyclic graph of the enriched GO biological processes obtained.

Figure S7: **Human enhancers and genome segmentation.** Histograms of the proportion of human enhancers (y-axis) in class 1 (blue) and class 2 (green) lying within genome segments (x-axis) as annotated by combined results of ChromHMM [34] and Segway [35] on human lymphoblastoid (GM12878; **a**), cervical cancer (HeLa-S3; **b**), liver carcinoma (HepG2; **c**), umbilical vein endothelial (HUVEC; **d**), and chronic myelogenous leukemia (K562; **e**) cell lines from the ENCODE project [37]. Statistical significance (Bonferroni-corrected p-value < 0.01) of enrichment for enhancers from a specific class is indicated by

***.

Figure S8: **Enhancer activation upon MTB infection.** Stacked histogram of the fraction of human enhancers (y-axis) from class 1 and class 2 predicted to be activated (red), inhibited (blue), or with unchanged activity (grey). Predictions were obtained using genomic segments predicted by ChromHMM [34] on human dendritic cells before and after infection with Mycobacterium tuberculosis [38].

Figure S9: **Cell type expression specificities of human enhancers.** The cell type expression specificities derived from FANTOM5 CAGE datasets [22] is provided as a heat map for human enhancers in class 1 (**a**) and class 2 (**b**). The color (see scale) represents the fraction of expressed enhancers in each cell type (columns) found in each expression specificity range (rows). CAGE, Cap Analysis of Gene Expression.

Figure S10: **Enrichment of enhancers from a single class within chromatin loops.** For each chromatin loop predicted in HeLa (**a**), HMEC (**b**), HUVEC (**c**), IMR90 (**d**), K562 (**e**), KBM7 (**f**), and NHEK (**g**) cell lines, we computed the p-value of the Binomial test to assess enrichment for enhancers from a single class. The plots compare the density (y-axis) of p-values for Binomial tests (x-axis) applied to classes 1 and 2 enhancers (red) and 1,000 random assignments of class labels to the enhancers (black).

Figure S11: **Organization of enhancers within chromatin loops.** Density (y-axis) of distances (x-axis) between enhancers and chromatin loop centers/anchors. The distances were normalized by the lengths of the loops. Enhancers at the center of the loops were found at distance 0.0 while enhancers at chromatin loops boundaries/anchors were found at distance 0.5. Results associated with class 1 and class 2 enhancers are depicted in blue and green, respectively. HeLa (**a**), HMEC (**b**), HUVEC (**c**), IMR90 (**d**), K562 (**e**), KBM7 (**f**), and NHEK (**g**) cell lines were considered.

Table S1: List of genes associated with class 1 (first column) and class 2 (second column) enhancers derived from Andersson et al. [22].

Table S2: Enriched GO biological processes associated with genes predicted to be regulated by class 1 enhancers.

Table S3: Enriched GO biological processes associated with genes predicted to be regulated by class 2 enhancers.

Data S1: **Centrimo motif enrichment analyses.** Centrimo was applied to regions of 1,001 bp centered around enhancers' mid-points from class 1 (centrimo_class1vs2_enhancers.html) and class 2 (centrimo_class2vs1_enhancers.html). Promoter regions of 1,001 bp centered around TSSs associated with class 1 (centrimo_class1vs2_promoters.html) and class 2 (centrimo_class2vs1_promoters.html) were also subjected to Centrimo. Position enrichment of motifs were focussed around enhancer mid-points for enhancers and all regions for promoters.