

1

Characterizing the epigenetic signatures of the human regulatory elements: A pilot study

2 Sawyer L. Clement^{1,2} and Hani Z. Girgis*²

¹Department of Biological Science, The University of Tulsa, 800 South Tucker Drive, Tulsa, Oklahoma, USA

²Tandy School of Computer Science, The University of Tulsa, 800 South Tucker Drive, Tulsa, Oklahoma, USA

Email: Sawyer L. Clement - sawyer-clement@utulsa.edu; Hani Z. Girgis* - hani-girgis@utulsa.edu;

*Corresponding author

3 **Abstract**

4 **Background**

5 Chromatin modifications have provided promising clues on how cells that share the same copy
6 of the genome can perform distinct functions. It is believed that enhancers and promoters are
7 marked by a single chromatin mark each, H3K4me1 and H3K4me3, respectively. However, other
8 studies have indicated that enhancers and promoters share multiple chromatin marks, including
9 H3K4me1/2/3 and H3K27ac. Therefore, we asked whether the epigenetic signatures of these
10 regulatory elements consist of a single mark or multiple marks. Repetitive regions, repeats,
11 are usually ignored. However, we found, in public data, that repeats include about 25% of
12 active enhancers. Thus, we asked how the epigenetic signatures of repetitive and non-repetitive
13 enhancers differ. We studied the four marks in IRM90 (human lung fibroblast) and H1 (human
14 embryonic stem cell).

15 **Results**

16 Our results show that enhancers and promoters are enriched significantly with the four marks,
17 which form pyramidal signatures. However, the relative lengths of the marks are different. The
18 promoter signature is directional; H3K4me2/3 and H3K27ac tend to be present downstream of
19 the transcription start site; H3K4me1 tends to be present upstream. H1-specific enhancers have
20 a similar signature to IRM90-specific enhancers; however, it is not the case for active promoters
21 of the two cell types. Interestingly, inactive enhancers show a residual signature that resembles
22 the signature of active enhancers. Finally, the epigenetic signature of enhancers found in repeats
23 is identical to that of enhancers found in non-repetitive regions.

24 **Conclusions**

25 In this study, we characterized the epigenetic signatures of active and inactive enhancers (pyra-
26 midal) as well as active promoters (directional-pyramidal) in two cell types. These signatures
27 consist of four chromatin marks that have been reported to be associated with enhancers and
28 promoters. Interestingly, about one quarter of active enhancers are found in repeats. Active
29 enhancers within repeats and those outside repeats have the same epigenetic signature. These
30 results have great potential to change the way Molecular Biologists think of repeats, and to
31 expand our understanding of gene regulation.

32

33 **Keywords**

34 Epigenetic marks, Chromatin marks, Chromatin modifications, cis-Regulatory elements, En-
35 hancers, Promoters, Repeats, Tissue specificity.

36 **Background**

37 Epigenetic modifications, including DNA methylation and histone modifications, are the
38 primary mechanism for long term regulation of gene expression. As all cells in an organism
39 have the same genome, it falls to differential epigenetic landscapes to determine whether a
40 given cell becomes a neuron, a melanocyte, or a T-cell [1]. While DNA methylation adds a
41 methyl group directly to the DNA strand, histone modifications act on histones, chromatin

42 proteins around which DNA is wrapped [2].

43 These histone/chromatin modifications include methylation, acetylation, ubiquitylation,
44 phosphorylation, and others. Furthermore, these marks have different effects when attached
45 to different histones and at different histone locations; H3K27ac (acetylation of histone 3 at
46 lysine 27) is distinct from H3K18ac (acetylation of histone 3 at lysine 18). Histone mod-
47 ifications alter gene expression by disrupting chromatin organization and by recruiting or
48 blocking the binding of non-histone proteins to DNA. These proteins may include transcrip-
49 tion factors that regulate gene expression or proteins that further modify the chromatin [2].

50 While DNA hyper- and hypomethylation is commonly associated with gene under- and
51 over-expression, the effects of specific chromatin modifications are less clearly understood.
52 However, the scientific community has been making steady progress toward understanding
53 the effects of chromatin marks.

54 Heintzman, et al. [3] reported the association of enhancers and promoters with single
55 chromatin marks: H3K4me1 and H3K4me3, respectively. The pilot phase of the EN-
56 CODE project suggested that enhancers and transcription start sites (TSSs) have “inverse”
57 chromatin patterns. Specifically, enhancers are reported to have high H3K4me1 and low
58 H3K4me3, whereas TSSs are reported to have low H3K4me1 and high H3K4me3. Addi-
59 tionally, five marks, including H3K4me1/2/3, can be used for predicting the TSSs [4]. A
60 third study concluded that H3K4me2 and H3K4me3 are characteristic marks of TSSs. In
61 addition, the same study indicated that “H3K4me1 signal was low but showed some evi-
62 dence of enrichment further downstream from the TSS than H3K4me2 and H3K4me3” [5].
63 In sum, the studies published in 2007 were not in agreement. Some studies suggested that
64 enhancers are marked by the abundance of H3K4me1, whereas promoters are marked by
65 the abundance of H3K4me3 (the single mark hypothesis). However, other studies suggested
66 that H3K4me1/2/3 are present around promoter regions and H3K4me1/3 are present around
67 enhancers (the multiple marks hypothesis).

68 A study by Pekowska, et al. [6] reported the association of H3K4me1/2/3 and the en-
69 hancers specific to the T-cell. Further, the ENCODE project reported that H3K4me1/2
70 are associated with enhancers and H3K4me1/2/3 are associated with promoters; in addi-
71 tion, H3K27ac marks active promoters and active enhancers [7]. At the end of 2012, it was

72 reported that active enhancers and active promoters share H3K4me1/2/3 and H3K27ac.

73 In 2013, a study by Zhu, et al. [8] suggested that H3K27ac and H3K4me3, among other
74 chromatin marks, are indicative of active enhancers. Moreover, another study confirmed
75 that H3K4me1/2/3 and H3K27ac are the best four chromatin marks indicative of active
76 enhancers [9].

77 The multiple marks hypothesis is supported by a larger number of studies than the
78 single mark hypothesis. Nonetheless, the single mark hypothesis is still cited in recent
79 studies [10, 9].

80 Motivated to resolve these discrepancies, we conducted a study on publicly available data.
81 The main purpose of our study is to characterize the epigenetic signatures of enhancers and
82 promoters. Specifically, we studied the distributions of H3K4me1/2/3 and H3K27ac relative
83 to the site and each other with three main questions in mind. First, which of the two
84 hypotheses is more accurate? Second, how does the epigenetic signature of enhancers differ
85 from the signature of promoters? Third, how does the epigenetic signature of enhancers
86 found in repetitive regions (about 25% of active enhancers in the studied cells) differ from
87 that of enhancers outside these regions?

88 **Results and Discussion**

89 This research was conducted to characterize the epigenetic signatures of the main regulatory
90 elements. We focused on the following four chromatin marks: H3K4me1/2/3 and H3K27ac.
91 In this study, we examined the distribution of these marks in the context of promoters and
92 enhancers, with the intention of determining a characteristic epigenetic signature for each.

93 **IMR90-specific enhancers exhibit pyramidal epigenetic signature**

94 We began our investigation by examining the epigenetic signatures around active enhancers
95 (p300 binding sites overlapping DNase I hypersensitive sites (DHSs) provided by another
96 study by Rajagopal, et al. [9]). We chose the IMR90 cell line (human lung fibroblast)
97 for our experiments, as the relevant epigenetic information for this tissue type is publicly
98 available. We began by plotting the distribution of the chosen marks about individual active
99 enhancers (Figures 1a-1f). Eventually, we detected a pattern in the epigenetic distribution,

100 more complex than the simple presence or absence of particular marks. Stretches of each
101 of the four marks appeared in most of the plots, approximately centered about the DHSs.
102 The overlapping stretch of H3K4me1 was usually the longest, followed by H3K4me2 and
103 H3K27ac (roughly equal in length), with the H3K4me3 usually being the shortest. When
104 plotted about a single enhancer, we observed that the four marks form a pyramidal shape.
105 These results confirm that the enhancer epigenetic signature is not defined by the presence
106 or absence of H3K4me1 (the single mark hypothesis), but by the arrangement of the four
107 marks around the enhancer (the multiple marks hypothesis).

108 Next, we sought to determine whether the pyramidal enhancer pattern would persist in
109 a large data set. However, we had no reference for how these signatures compare to the
110 epigenome as a whole. Therefore, we examined the distributions of the four marks around
111 500 segments (each is 500 bp long), spread uniformly throughout the human chromosome 1.
112 We refer to these sequences as control sequences. These random segments have low content
113 of the epigenetic marks, though the actual content varied somewhat between the mark types.
114 H3K4me1 and H3K4me2 appeared in approximately 30% of samples, whereas H3K27ac and
115 H3K4me3 appeared in approximately 15%.

116 The epigenetic signatures of 2000 active enhancers were profiled (Figure 2). We found
117 that each of the four chromatin marks was consistently present at active enhancers (80-95%).
118 Further, the enhancer regions were 3.1-fold more enriched with H3K4me1 than the control
119 sequences (P-value $< 2.2e^{-16}$, Fisher's exact test). Similarly, the other three marks were
120 significantly enriched in the enhancer regions (H3K4me2: 3.4 folds, H3K4me3: 5.2 folds,
121 H3K27ac: 5.2 folds; P-value $< 2.2e^{-16}$, Fisher's exact test). Again, these results show that
122 the four marks are enriched in the enhancer regions.

123 **Epigenetic Marks of Enhancers in repetitive and non-repetitive regions**

124 Interestingly, H3K4me3 was more enriched in the enhancers than was H3K4me1 (5.2 folds
125 vs. 3.1 folds). This observation questions the common assumption that H3K4me1 is the
126 main chromatin mark characterizing enhancers. This assumption is based on the abundance
127 of the mark, not on the enrichment value obtained by comparing the observed abundance in
128 the enhancers to that in the control sequences.

129 The width and the orientation patterns observed in the individual enhancer figures also

130 reappeared. H3K4me1 had the largest average width (5000 bp), followed by H3K4me2 and
131 H3K27ac (3000-4000 bp), with H3K4me3 being the narrowest (1500-2000 bp). Together,
132 these observations support the pyramidal distribution of the studied epigenetic marks about
133 the active IMR90-specific enhancers.

134 **Active enhancers in repetitive regions exhibit the same epigenetic pattern as active** 135 **enhancers outside repetitive regions**

136 A study by Xie et al. [10] conducted on a number of tissues showed (i) transposon subfami-
137 lies have different patterns of hypomethylation across tissue types; (ii) these “differentially-
138 hypomethylated” subfamilies are associated with H3K4me1; (iii) they are associated with
139 the expression of genes in their vicinities; and (iv) the sequences of these subfamilies in-
140 clude binding sites for tissue-specific transcription factors. These four findings suggest that
141 transposon subfamilies have tissue-specific enhancer-like functions. Motivated by these find-
142 ings, we divided active IMR90 enhancers into those overlapping with repetitive regions and
143 those that are not. Out of 25,109 enhancers, 5,925 (23.6%) overlap repetitive regions. We
144 wished to determine whether, if the signature persisted, it would appear on both repetitive
145 and non-repetitive enhancers. Therefore, we profiled the chromatin marks around 1000 non-
146 repetitive enhancers and 1000 repetitive enhancers (Figures 2a and 2b). These figures show
147 that there are no differences between the distribution of these four marks on repetitive and
148 non-repetitive enhancers, indicating that non-repetitive and repetitive enhancers have the
149 same epigenetic signature.

150 **The four epigenetic marks are significantly depleted in inactive enhancers compared** 151 **to active enhancers**

152 We had confirmed that the epigenetic signature of enhancers can be observed in both indi-
153 vidual enhancers and large enhancer sets. However, we had not performed any large-scale
154 analysis of inactive enhancers. As such, we could not be certain whether the pyramidal
155 signature observed previously was indicative only of active enhancers or of enhancers as a
156 whole. To resolve this, we compared the epigenetic signatures of an equal number of inac-
157 tive and highly active enhancers, as determined by eRNA (enhancer RNA) levels (obtained
158 from the Fantom5 project [11]). Figures 3a and 3b show the epigenetic signatures of the

159 active and the inactive enhancers. Active enhancers were significantly more enriched with
160 the four marks than inactive ones (H3K4me1: 1.5 folds, P-value = 0.000259; H2K4me2: 1.6
161 folds, P-value = $3.035e^{-7}$; H3K4me3: 2.8 folds, P-value = $1.677e^{-15}$; H3K27ac: 2.0 folds,
162 P-value = $1.064e^{-8}$; Fisher's exact test). Additionally, the marks were clearly narrower in
163 the inactive enhancers than those found in the active enhancers. Overall, the inactive en-
164 hancers displayed significantly decreased levels of all four epigenetic marks, indicating that
165 the enhancer epigenetic signature is weaker in these enhancers.

166 **Inactive enhancers show residual enrichment of the four epigenetic marks, resembling** 167 **the pyramidal signature**

168 Next, we asked whether the densities of the marks in inactive enhancers are similar to the
169 genome average, i.e. similar to the control sequences. Therefore, we compared these densities
170 in the inactive enhancers and the control sequences (Figures 3b and 3c). Interestingly, inac-
171 tive enhancers were more enriched with the four marks than the control sequences (H3K4me1:
172 1.8 folds, P-value = $9.838e^{-6}$; H2K4me2: 1.9 folds, P-value = $1.289e^{-6}$; H3K4me3: 2.0 folds,
173 P-value = 0.0007712; H3K27ac: 2.5 folds, P-value = $3.752e^{-7}$; fisher's exact test). These
174 results suggest that inactive enhancers exhibit a weaker, yet significant, version of the epi-
175 genetic signature of the tissue-specific active enhancers.

176 **The H1-specific enhancers and the IMR90-specific enhancers have similar epigenetic** 177 **signatures**

178 The study by Rajagopal et al. [9] determined enhancers specific to the H1 cell line experimen-
179 tally (p300 binding sites overlapping DHSs). We asked if the epigenetic signature observed in
180 the enhancers specific to the IMR90 is the same/or similar to that of the enhancers specific to
181 H1. Figures 1g-1i show the four marks around three H1-specific enhancers. The four marks
182 are stacked around the enhancers, suggesting a stacked/pyramidal epigenetic signature.

183 Next, we profiled the epigenetic signature of 2000 H1-specific enhancers (Figure 2). The
184 four studied marks are present around the active enhancers of H1. These profiles differ from
185 those of IMR90 in the width of the densities and in the relative order of the bottom two layers.
186 In general, the densities observed in H1 are narrower than those observed in IMR90. Recall
187 that the signature of the IMR90-specific enhancers consists of these layers: H3K4me1 (the

188 widest), H3K4me2 and H3K27ac (roughly the same width), and H3K4me3 (the narrowest).
189 The signature of the H1-specific enhancers consists of these layers: H3K4me2 (the widest),
190 H3K4me1 and H3K27ac (roughly the same width), and H3K4me3 (the narrowest). The
191 two signatures differ in the relative order of the lower two layers (H3K4me1 and H3K4me2).
192 These results show that the epigenetic signatures of the H1-specific enhancers and the IMR90-
193 specific enhancers are similar, though not identical.

194 **H1-specific enhancers in repetitive regions exhibit the same epigenetic signature as** 195 **those outside repetitive regions**

196 We observed that 23.8% (1,402 out of 5,899) of the H1-specific enhancers overlap repetitive
197 regions. Recall that a similar percentage of the IMR90-specific enhancers overlap repetitive
198 regions as well. Furthermore, the epigenetic profiles of 1000 non-repetitive enhancers and
199 1000 repetitive enhancers of the H1 are almost identical (Figures 2c and 2d). These results
200 confirm the results observed in the non-repetitive enhancers and the repetitive enhancers of
201 IMR90.

202 Given the existence of a pyramidal pattern about active enhancers, our next consideration
203 was whether other regulatory regions, namely promoters, might exhibit a similar signature.

204 **Active promoters exhibit directional-pyramidal epigenetic signature**

205 We examined the epigenetic signature of individual active promoters. First, we defined active
206 promoters as transcription start sites (TSSs) overlapping DHSs. However, a few problems
207 prevented us from receiving immediate results. Initially, we drew random DHS-overlapping
208 promoters from a list of all known TSSs. However, as gene expression in a cell varies
209 greatly over time, not necessarily coinciding with DHS establishment, the resulting figures
210 were inconsistent. Despite this, a few promoters did appear to demonstrate a directional-
211 pyramidal pattern. Therefore, we decided to repeat the experiment with a more precise
212 method for determining active promoters based on gene expression levels.

213 The second time, we chose ten promoters of genes with the highest expression (active
214 promoters) in IMR90 as well as ten promoters of unexpressed genes (inactive promoters).
215 The epigenetic plots made for the inactive promoters showed no apparent pattern, while each
216 plot made for the active promoters demonstrated some version of the directional-pyramidal

217 pattern observed in the first trial (Figures 4a-4c). As with the pattern observed about
218 enhancers, the pattern about promoters involved all four of the studied chromatin marks
219 (H3K4me1/2/3 and H3K27ac). The distinctions were in (i) the order of the layers of the
220 signature and (ii) how these marks were arranged around the TSS, i.e. the directionality.

221 The layers of the directional-pyramidal signature of the active promoters were: H3K4me2
222 (the broadest), H3K4me3, H3K27ac, and H3K4me1 (the narrowest). Recall that the layers of
223 the pyramidal signature of the active enhancers were: H3K4me1 (the broadest), H3K4me2,
224 H3K27ac, and H3K4me3 (the narrowest).

225 Additionally, we found that H3K4me2/3, and H3K27ac all encompassed the TSS-
226 overlapping DHS. Moreover, they were each more present downstream of the promoter than
227 upstream. In contrast, H3K4me1 was mainly present upstream of the DHS, with a small
228 break in the region overlapping the DHS (possibly at the site of the TSS itself). The orienta-
229 tion of the chromatin marks relative to the DHS was inverted on positive strand promoters
230 compared to negative strand promoters.

231 Having observed a distinctive epigenetic pattern in individual promoters, our next step
232 was to determine whether this pattern would appear in a large data set consistently. To
233 this end, we studied the promoters (+/- 250 bp from TSS) of the 100 most expressed (ac-
234 tive promoters) and 100 unexpressed genes (inactive promoters) as determined by RNA-seq
235 values.

236 We compared the four epigenetic marks across the two sets. The active promoters showed
237 clear enrichment of these marks compared to the inactive promoters (H3K4me1: 1.9 folds,
238 P-value = $8.826e^{-07}$; H3K4me2: 2.3 folds, P-value = $5.033e^{-14}$; H3K4me3: 3.6 folds, P-value
239 $< 2.2e^{-16}$; H3K27ac: 3.6 folds, P-value $< 2.2e^{-16}$; Fisher's exact test).

240 As expected, active promoters were also enriched with the four marks compared to the
241 control sequences (H3K4me1: 2.4 folds, P-value = $3.903e^{-16}$; H3K4me2: 3.1 folds, P-value
242 $< 2.2e^{-16}$; H3K4me3: 6.0 folds, P-value $< 2.2e^{-16}$; H3K27ac: 5.6 folds, P-value $< 2.2e^{-16}$;
243 Fisher's exact test).

244 H3K4me1 around active promoters demonstrated a clear drop at the TSS (Figure 5a).
245 This drop mirrored the H3K4me1 TSS breaks observed in the individual promoters.

246 The individual promoter regions show that certain epigenetic marks extend farther on

247 one side of the TSS than the other. The directions of these marks suggest that their dis-
248 tributions are related to the direction of transcription. Therefore, we analyzed the active
249 promoters on the positive and the negative strands separately. Figures 5c and 5d show the
250 epigenetic marks of the 100 active and the 100 inactive promoters. The results matched the
251 epigenetic signature observed in the individual promoters. The directionality was reversed
252 across opposing strands (positive and negative), indicating that the pattern is related to
253 transcription direction. On both strands, H3K4me2/3 and H3K427ac tended to be enriched
254 downstream of the promoter. H3K4me1 tended to be enriched upstream, with the gap at
255 the TSS appearing once again.

256 **Unlike inactive enhancers, the marks around the inactive promoters does not resemble** 257 **the directional-pyramidal signature**

258 Inactive promoters are weakly enriched with three of the four marks (H3K4me2: 1.4 folds,
259 P-value = 0.03304; H3K4me3: 1.7 folds, P-value = 0.01846; H3K27ac: 1.6 folds, P-value =
260 0.04251; Fisher's exact test). However, the densities of these marks do not resemble those
261 of the active promoters, supporting the notion that promoter epigenetic signature is related
262 to gene activation.

263 **Active promoters of the H1 cell line are not marked epigenetically**

264 We studied the active promoters of the H1 cell line using the same procedure used in studying
265 the IMR90 active promoters. However, the directional-pyramidal pattern was not observed
266 in the H1 cell line. Moreover, the densities of the marks around active promoters resemble
267 the genome average. These observations may be due to H1 being a stem cell, as they are
268 known to have a “unique epigenetic signature” [12].

269 **Discussion**

270 In this study, we characterized the epigenetic signatures of enhancers and promoters in the
271 IMR90 and the H1 cell lines. We were motivated by the discrepancies in the literature
272 about how certain chromatin marks are associated with these regulatory elements. Some
273 studies report that H3K4me1 marks enhancers, whereas H3K4me3 marks promoters. Yet
274 other studies report that four chromatin marks (H3K4me1/2/3 and H3K27ac) are present

275 around enhancers and promoters. The first question we considered in this study whether the
276 epigenetic signatures of promoters and enhancers consist of one mark or multiple marks. Our
277 analyses show that these signatures consist of multiple marks, including, but not limited to,
278 all of the four marks studied. We then asked how the epigenetic signature of the promoters
279 differs from the enhancer signature. Additionally, we asked how the signature of active
280 enhancers found in repetitive regions differs from the signature of the ones outside repetitive
281 regions. In sum, this study contributes the following seven findings:

- 282 ● In the IMR90 cell line, H3K4me1/2/3 and H3K27ac form a pyramidal shape around active
283 enhancers. H3K4me1 is the base of the pyramid. H3K4me2 and H3K27ac are the middle
284 layers. H3K4me3 is the top of the pyramid. The regularity with which this pattern appears
285 suggests that it is intrinsically tied to transcription in this cell line.
- 286 ● Active enhancers specific to IMR90 are more enriched with H3K4me3 than with H3K4me1
287 (5.2 folds vs. 3.1 folds), questioning the common assumption that H3K4me1 is the main
288 chromatin mark characterizing active enhancers. This assumption is supported by the
289 abundance, not the enrichment value relative to the genome average, of H3K4me1 around
290 enhancer regions.
- 291 ● Active promoters of IMR90 demonstrate a directional-pyramidal epigenetic signature.
292 H3K4me2 is the base of the pyramid. H3K4me3 and H3K27ac are the middle layers.
293 H3K4me1 is the top of the pyramid. Note that the directional-pyramidal signature of
294 active promoters is roughly the “inverse” of pyramidal signature of active enhancers with
295 regard to chromatin mark length.
- 296 ● Chromatin marks in active IMR90 promoters are unevenly distributed about the TSS.
297 H3K4me2/3 and H3K27ac tend to be more present downstream of the TSS, whereas
298 H3K4me1 tends to be more present upstream, and often has a characteristic gap around
299 the TSS. This transcription-dependent directionality around the TSS suggests that these
300 marks are involved with the initiation of transcription.
- 301 ● The epigenetic signature of the enhancers active in H1 (embryonic stem cell) is similar to
302 the signature of those active in IMR90. However, the promoters of genes active in H1 do
303 not show any recognizable epigenetic pattern consisting of the four studied marks.

- 304 • Inactive enhancers in IMR90 exhibit a residual epigenetic signature that resembles the
305 signature of active enhancers, whereas inactive IMR90 promoters do not exhibit any such
306 signature.

- 307 • The epigenetic signatures of active enhancers in non-repetitive regions and those in repet-
308 itive regions are indistinguishable. As this signature is linked to enhancer activation, this
309 reinforces the notion that repetitive elements have significant regulatory function. As such,
310 we urge the scientific community to stop masking/ignoring repeats and to start studying
311 them.

312 **Materials and Methods**

313 **Data**

314 In this study, we used enhancers experimentally determined by Rajagopal et al. [9]. An
315 enhancer is defined as a DNase I hypersensitive site (DHS) where p300 binds “distal to
316 known UCSC and Gencode” transcription start sites (TSSs). Enhancers studied are specific
317 to the IMR90 cell line (human lung fibroblasts) and the H1 cell line (human embryonic stem
318 cells). Chromatin modification data sets, including those for H3K4me1/2/3 and H3K27ac,
319 were downloaded from the publicly available Human Epigenome Atlas [13]. The DHS data
320 sets were downloaded from the NCBI Gene Expression Omnibus, under the designation
321 Sample GSM468792 [9]. RNA-seq data were downloaded from the ENCODE project [14].
322 Name, TSS, and chromosome location data for all human genes were downloaded from the
323 Ensembl website [15]. As the Ensembl data are in hg38 and the other data sets are in hg19,
324 the Ensembl data were converted into hg19 using the UCSC LiftOver tool. Repeats of the
325 human genome (hg19) detected by RepeatMasker were downloaded from the Institute for
326 Systems Biology website [16]. The eRNA (enhancer RNA) levels for the fetal lung tissue
327 were downloaded from the FANTOM5 project [17, 18, 11].

328 **Individual Enhancers**

329 Individual enhancers (p300 binding sites overlapped by DHSs) were manually selected. For
330 each p300 binding site, the overlapping DHS was located. Then, segments representing
331 the four chromatin marks that overlapped the DHS were detected. These chromatin mark

332 segments were then plotted one on top of another as lines to easily determine how these
333 marks were distributed about the DHS (and by extension the enhancer).

334 **Multi-Enhancers**

335 Active enhancers were taken from a list of experimentally determined enhancers (p300 bind-
336 ing sites overlapping DHSs and distal to known TSSs) [9]. Using a list of repetitive elements
337 in the human genome, these enhancers were separated into two sets: those found in repetitive
338 elements and those found in non-repetitive regions. The first 1000 enhancers from each set
339 were arbitrarily chosen. For each chromatin mark, the length and the distribution about an
340 enhancer were recorded. Then segments representing these marks were summed and plotted
341 together in order to display the distribution of that mark around the enhancers. Separate
342 analyses were done with repetitive and non-repetitive enhancers for purposes of comparison.

343 **Control Sequences**

344 We constructed a set of control sequences by selecting 500 segments distributed uniformly
345 throughout the human chromosome one. Each segment is 500 bp long. The IMR90 chromatin
346 marks overlapping the control sequences were analyzed and summed as done previously.

347 **Active and Inactive Enhancers**

348 The eRNA (enhancer RNA) data was used for selecting a set of enhancers highly active in
349 IMR90 and an equal-sized set of enhancers inactive in IMR90. All enhancers with eRNA
350 expression values 19 or greater were chosen as active enhancers. This cutoff value was
351 chosen because it yielded approximately 100 (102) of the most active enhancers, enough to
352 clearly demonstrate any major epigenetic pattern. The set of inactive enhancers consisted
353 of 102 enhancers with eRNA expression values of 0. The chromatin marks around each
354 enhancer set were analyzed as was done previously. Because these enhancers were not single
355 coordinates (as were the p300 binding sites) but rather representative regions, all chromatin
356 marks overlapping these regions were recorded.

357 **Individual Promoters**

358 Initially, individual promoters were manually selected from a list of all human promoters.
359 These promoters were analyzed the same way as the individual enhancers, plotting the epi-
360 genetic marks around DHS-overlapping promoters. When this provided an unclear pattern,
361 manually chosen promoters were replaced with ten promoters of the most expressed genes
362 in a specific cell line. Additionally, ten promoters of unexpressed genes in the same cell line
363 were selected. The chromatin marks overlapping the promoters were plotted and analyzed
364 as done previously.

365 **Multi-Promoters**

366 We determined the 100 most expressed genes and 100 unexpressed genes in a specific cell line
367 based on the ENCODE gene expression data. As the gene expression data (ENCODE) and
368 the TSS data (Ensembl) were from different sources, some gene names from one did not ap-
369 pear in the other; such genes were removed from consideration. Because the gene expression
370 data were sufficient proof of promoter activation, the requirement for a promoter to overlap
371 a DHS was removed. Instead, all chromatin marks near the promoters (+/- 250 base pair
372 from the TSS) were included in the analysis. For each chromatin mark, the length and the
373 distribution about the promoters of the studied genes were recorded. The chromatin marks
374 around the promoters were analyzed as was done with the enhancers. Separate analyses were
375 done for the 100 most expressed and the 100 unexpressed genes, to contrast the distribution
376 of marks in inactive and active promoters. This procedure was repeated with the additional
377 consideration of gene strand. This time, the 100 most expressed and unexpressed genes on
378 the positive strand and the negative strand were analyzed separately.

379 **List of abbreviations**

380 TSS: transcription start site; DHS: DNase I hypersensitive site; eRNA: enhancer RNA.

381 **Declarations**

382 **Ethics approval and consent to participate**

383 Not applicable.

384 **Consent for publication**

385 Not applicable.

386 **Availability of data and material**

387 All data sets used in this study are publicly available as noted under the Materials and
388 Methods Section.

389 **Competing interests**

390 The authors declare that they have no competing interests.

391 **Funding**

392 This research is supported by an internal grant from the University of Tulsa.

393 **Authors' contributions**

394 SLC designed and implemented the analyses. HZG initiated the study. HZG designed the
395 analyses and calculated the statistical tests. SLC and HZG wrote the manuscript. All
396 authors read and approved the final manuscript.

397 **Acknowledgements**

398 Not applicable.

399 **Authors' information**

400 SLC is a senior student at the University of Tulsa. His major is Biology. HZG is an Assistant
401 Professor of Computer Science at the University of Tulsa. HZG did his postdoctoral work
402 at the National Center for Biotechnology Information (NCBI), the National Institutes of
403 Health (NIH). HZG has been studying regulatory elements and repeats for seven years.
404 HZG majored in Biology and Computer Science while he was undergraduate student at the
405 State University of New York at Buffalo.

References

- 406 1. Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, Sheffield NC, Ster-
407 gachis AB, Wang H, Vernot B, Garg K, John S, Sandstrom R, Bates D, Boatman L, Canfield
408 TK, Diegel M, Dunn D, Ebersol AK, Frum T, Giste E, Johnson AK, Johnson EM, Kutuyavin T,
409 Lajoie B, Lee BK, Lee K, London D, Lotakis D, Neph S, Neri F, Nguyen ED, Qu H, Reynolds
410 AP, Roach V, Safi A, Sanchez ME, Sanyal A, Shafer A, Simon JM, Song L, Vong S, Weaver
411 M, Yan Y, Zhang Z, Zhang Z, Lenhard B, Tewari M, Dorschner MO, Hansen RS, Navas PA,
412 Stamatoyannopoulos G, Iyer VR, Lieb JD, Sunyaev SR, Akey JM, Sabo PJ, Kaul R, Furey TS,
413 Dekker J, Crawford GE, Stamatoyannopoulos JA: **The accessible chromatin landscape of**
414 **the human genome.** *Nature* 2012, **489**(7414):75–82.
- 416 2. Kouzarides T: **Chromatin Modifications and Their Function.** *Cell* 2007, **128**(4):693–705.
- 417 3. Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, Barrera LO, Van Calcar
418 S, Qu C, Ching KA, Wang W, Weng Z, Green RD, Crawford GE, Ren B: **Distinct and**
419 **predictive chromatin signatures of transcriptional promoters and enhancers in the**
420 **human genome.** *Nat Genet* 2007, **39**(3):311–318.
- 421 4. The ENCODE Project Consortium: **Identification and analysis of functional elements in**
422 **1% of the human genome by the ENCODE pilot project.** *Nature* 2007, **447**(7146):799–
423 816.
- 424 5. Koch CM, Andrews RM, Flicek P, Dillon SC, Karaöz U, Clelland GK, Wilcox S, Beare DM,
425 Fowler JC, Couttet P, James KD, Lefebvre GC, Bruce AW, Dovey OM, Ellis PD, Dhimi P,
426 Langford CF, Weng Z, Birney E, Carter NP, Vetrie D, Dunham I: **The landscape of histone**
427 **modifications across 1% of the human genome in five human cell lines.** *Genome Res*
428 2007, **17**(6):691–707.
- 429 6. Pekowska A, Benoukraf T, Zacarias-Cabeza J, Belhocine M, Koch F, Holota H, Imbert J, An-
430 drau JC, Ferrier P, Spicuglia S: **H3K4 tri-methylation provides an epigenetic signature**
431 **of active enhancers.** *EMBO J* 2011, **30**(20):4198–4210.
- 432 7. The ENCODE Project Consortium: **An integrated encyclopedia of DNA elements in**
433 **the human genome.** *Nature* 2012, **489**(7414):57–74.
- 434 8. Zhu Y, Sun L, Chen Z, Whitaker JW, Wang T, Wang W: **Predicting enhancer transcrip-**
435 **tion and activity from chromatin modifications.** *Nucleic Acids Res* 2013, **41**(22):10032–
436 10043.
- 437 9. Rajagopal N, Xie W, Li Y, Wagner U, Wang W, Stamatoyannopoulos J, Ernst J, Kellis M, Ren
438 B: **RFECS: A Random-Forest Based Algorithm for Enhancer Identification from**
439 **Chromatin State.** *PLoS Comput Biol* 2013, **9**(3):1–14.
- 440 10. Xie M, Hong C, Zhang B, Lowdon RF, Xing X, Li D, Zhou X, Lee HJ, Maire CL, Ligon KL,
441 Gascard P, Sigaroudinia M, Tlsty TD, Kadlecik T, Weiss A, O’Geen H, Farnham PJ, Madden
442 PAF, Mungall AJ, Tam A, Kamoh B, Cho S, Moore R, Hirst M, Marra MA, Costello JF, Wang
443 T: **DNA hypomethylation within specific transposable element families associates**
444 **with tissue-specific enhancer landscape.** *Nat Genet* 2013, **45**(7):836–841.

- 445 11. Lizio M, Harshbarger J, Shimoji H, Severin J, Kasukawa T, Sahin S, Abugessaisa I, Fukuda S,
446 Hori F, Ishikawa-Kato S, Mungall CJ, Arner E, Baillie JK, Bertin N, Bono H, de Hoon M, Diehl
447 AD, Dimont E, Freeman TC, Fujieda K, Hide W, Kaliyaperumal R, Katayama T, Lassmann
448 T, Meehan TF, Nishikata K, Ono H, Rehli M, Sandelin A, Schultes EA, 't Hoen PA, Tatum Z,
449 Thompson M, Toyoda T, Wright DW, Daub CO, Itoh M, Carninci P, Hayashizaki Y, Forrest
450 AR, Kawaji H: **Gateways to the FANTOM5 promoter level mammalian expression**
451 **atlas**. *Genome Biol* 2015, **16**:1–14.
- 452 12. Bibikova M, Chudin E, Wu B, Zhou L, Garcia EW, Liu Y, Shin S, Plaia TW, Auerbach JM,
453 Arking DE, Gonzalez R, Crook J, Davidson B, Schulz TC, Robins A, Khanna A, Sartipy
454 P, Hyllner J, Vanguri P, Savant-Bhonsale S, Smith AK, Chakravarti A, Maitra A, Rao M,
455 Barker DL, Loring JF, Fan JB: **Human embryonic stem cells have a unique epigenetic**
456 **signature**. *Genome Res* 2006, **16**(9):1075–1083.
- 457 13. The Human Epigenome Atlas: [<http://www.genboree.org/epigenomeatlas/index.rhtml>]. [Ac-
458 cessed 15 June 2016].
- 459 14. The ENCODE Project: [<https://www.encodeproject.org/data/annotations/>]. [Accessed 15
460 June 2016].
- 461 15. Ensembl: [<http://uswest.ensembl.org/biomart/martview>]. [Accessed 15 June 2016].
- 462 16. The Institute for Systems Biology: [<http://www.repeatmasker.org/species/hg.html>]. [Accessed
463 15 June 2016].
- 464 17. The FANTOM5 Project: [<http://fantom.gsc.riken.jp/5/>]. [Accessed 15 June 2016].
- 465 18. Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X,
466 Schmidl C, Suzuki T, Ntini E, Arner E, Valen E, Li K, Schwarzfischer L, Glatz D, Raithel J,
467 Lilje B, Rapin N, Bagger FO, Jorgensen M, Andersen PR, Bertin N, Rackham O, Burroughs
468 AM, Baillie JK, Ishizu Y, Shimizu Y, Furuhashi E, Maeda S, Negishi Y, Mungall CJ, Meehan
469 TF, Lassmann T, Itoh M, Kawaji H, Kondo N, Kawai J, Lennartsson A, Daub CO, Heutink
470 P, Hume DA, Jensen TH, Suzuki H, Hayashizaki Y, Muller F, Consortium TF, Forrest ARR,
471 Carninci P, Rehli M, Sandelin A: **An atlas of active enhancers across human cell types**
472 **and tissues**. *Nature* 2014, **507**(7493):455–461.

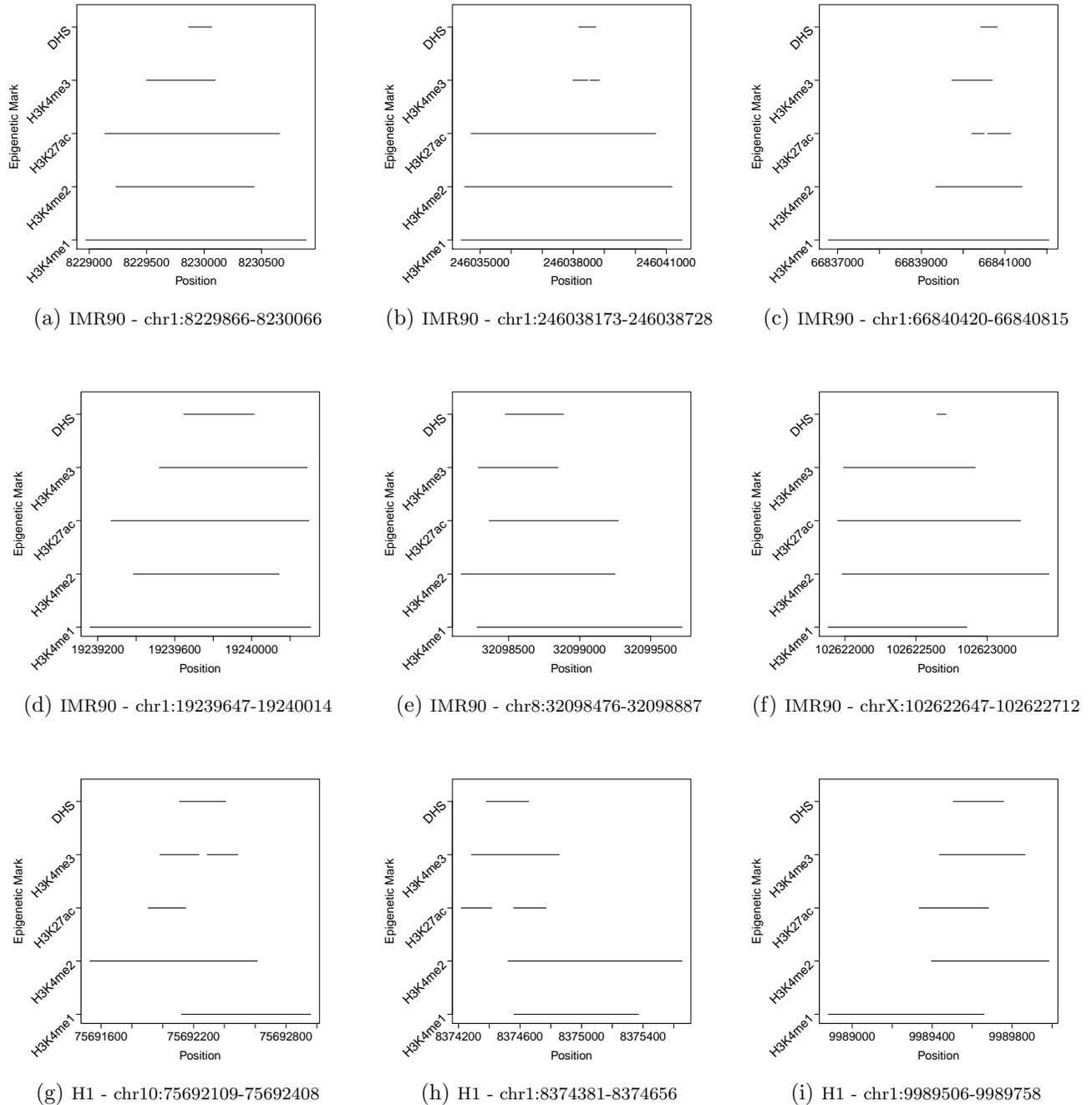


Figure 1: The arrangement of H3K4me1/2/3 and H3K27ac around active enhancers specific to the IMR90 and the H1 cell lines. Active enhancers are defined as p300 binding sites overlapping DNase I hypersensitive sites (DHSs). All coordinates are according to the hg19 assembly. (a-f) The four chromatin marks form pyramidal shape about enhancers specific to IMR90. (g-i) The four marks form stacked/pyramidal shape about enhancers specific to H1. The arrangement of the four marks around the IMR90-specific enhancers is similar, though not identical, to that around the H1-specific enhancers.

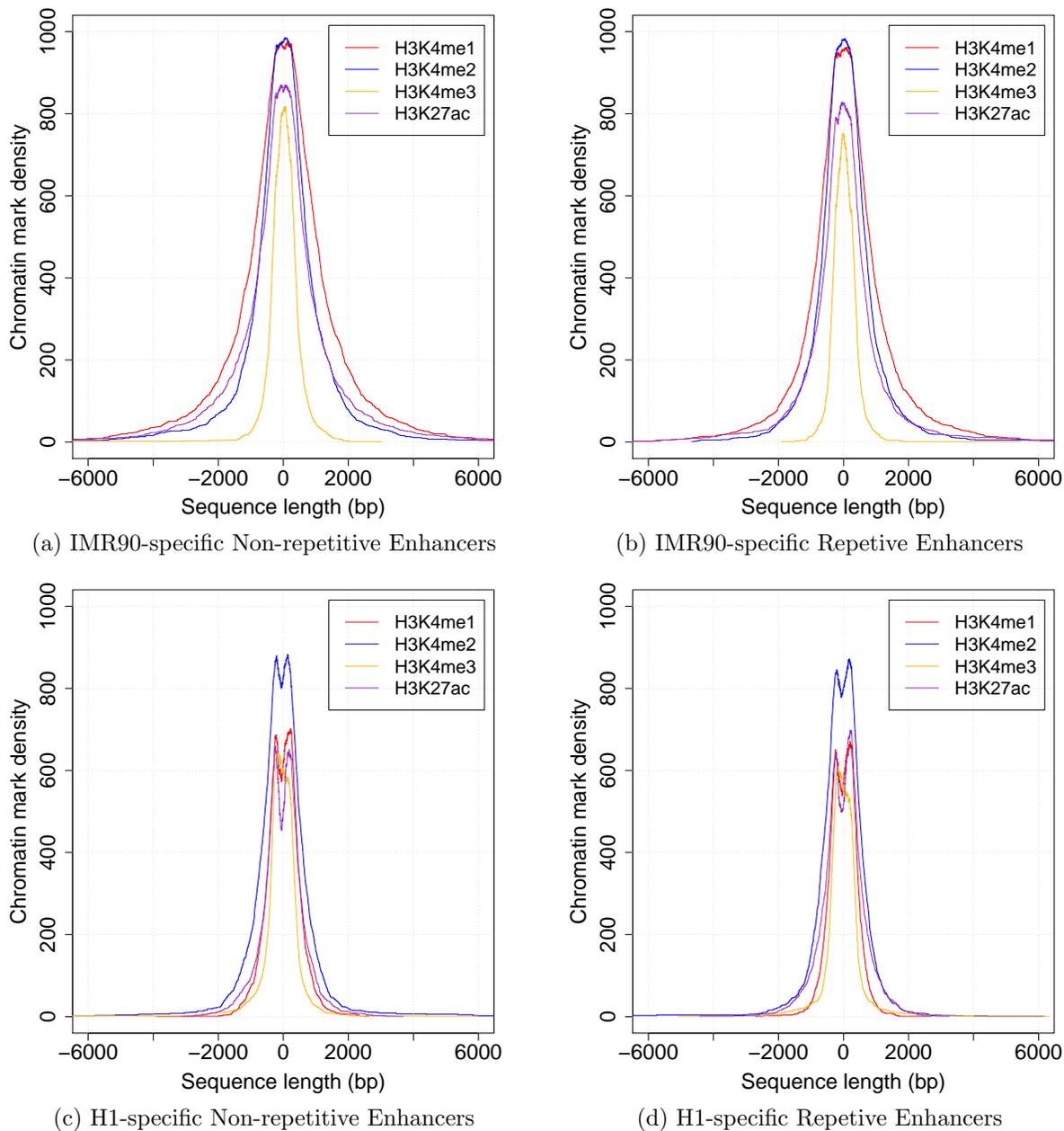
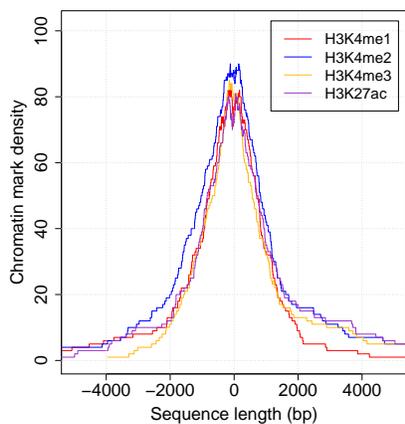
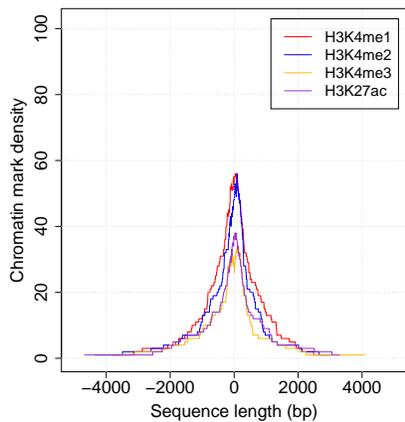


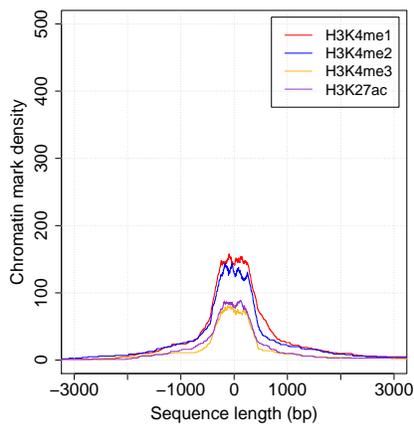
Figure 2: Repetitive (overlapping repetitive elements) and non-repetitive (outside repetitive elements) enhancers exhibit identical epigenetic signatures. Graphs show summed chromatin mark densities of 1000 enhancers centered around the p300 binding sites. (a & b) Comparisons of chromatin mark densities within the repetitive and the non-repetitive IMR90-specific enhancers. The repetitive and the non-repetitive IMR90-specific enhancers conform strongly to the pyramidal epigenetic signature. (c & d) Comparisons of chromatin mark densities within the repetitive and the non-repetitive H1-specific enhancers. The repetitive and the non-repetitive H1-specific enhancers have identical stacked epigenetic signatures.



(a) Active Enhancers



(b) Inactive Enhancers



(c) Control Sequences

Figure 3: Comparisons of chromatin mark densities within (a) active enhancers specific to IMR90, (b) inactive enhancers that are active in other cells, and (c) the control sequences. Chromatin marks around the active enhancers are both broader and more common than those around the inactive enhancers. The inactive enhancers are significantly more enriched with the four marks than the control sequences

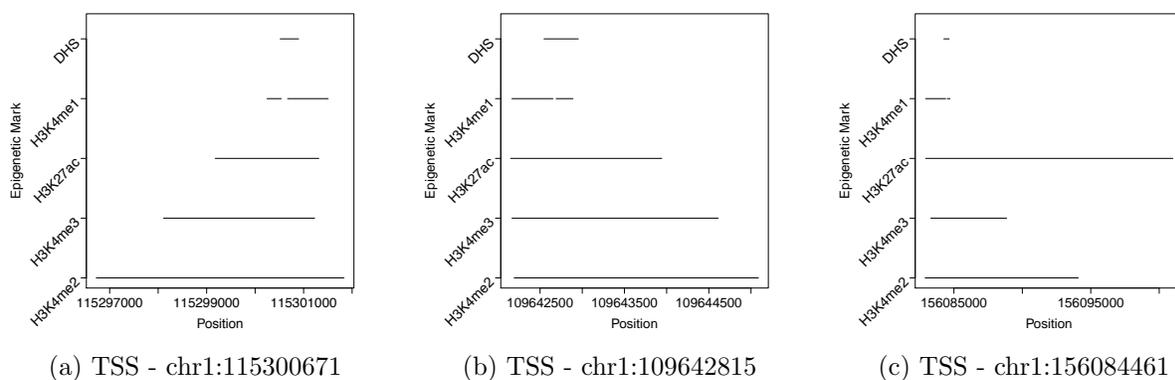


Figure 4: The epigenetic signature of active promoters in IMR90. (a-c) Four chromatin marks (H3K4me1/2/3 and H3K27ac) form directional pyramidal shape about three promoters of highly expressed genes in the IMR90 cell line. All coordinates are according to the hg19 assembly.

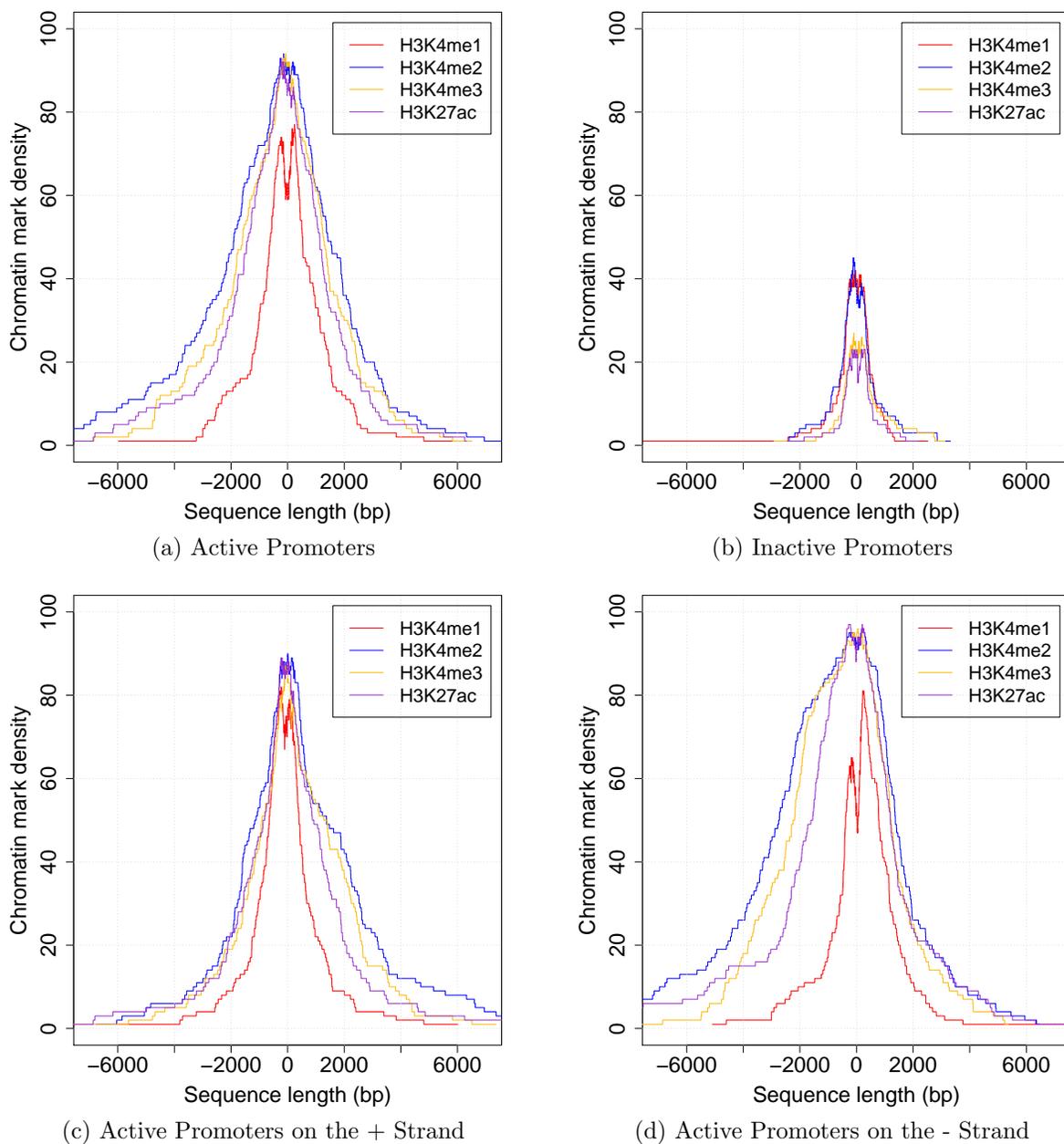


Figure 5: Active promoters in IMR90 exhibit directional-pyramidal epigenetic signature. Each graph shows summed chromatin mark density of 100 promoters centered around the TSS. (a & b) Comparisons of chromatin mark densities between the active (promoters of the 100 most expressed genes) and the inactive promoters (promoters of unexpressed genes in IMR90). Active promoters are significantly enriched with all marks studied. (c & d) Comparisons of chromatin mark densities within the active promoters on the negative and the positive strands. Epigenetic patterns of opposite strands roughly mirror each other, corresponding to opposite directions of transcription. H3K4me2, H3K4me3, and H3K27ac are enriched downstream of the promoter, while H3K4me1 is enriched upstream.