

1 **Crowdsourcing: Spatial clustering of low-affinity binding sites amplifies *in***
2 ***vivo* transcription factor occupancy**

3 Justin Malin^{1,3*}, Daphne Ezer^{2,3}, Xiaoyan Ma², Steve Mount¹, Hiren Karathia¹, Seung Gu
4 Park¹, Boris Adryan², Sridhar Hannenhalli^{1*}

5 ¹Center for Bioinformatics and Computational Biology

6 Department of Cell and Molecular Biology

7 University of Maryland, College park, MD

8 ²Department of Genetics

9 University of Cambridge, Cambridge, UK

10 **³These authors contributed equally to the work**

11 ***Co-corresponding authors**

12 Justin Malin

13 3104E Biomolecular Sciences Building (#296)

14 University of Maryland, College Park, MD 20742, USA

15 301 405 8218 (v) 301 314 1341 (f)

16 jmalin@umd.edu

17

18 Sridhar Hannenhalli

19 3104G Biomolecular Sciences Building (#296)

20 University of Maryland, College Park, MD 20742, USA

21 301 405 8219 (v) 301 314 1341 (f)

22 sridhar@umiacs.umd.edu

23 **Key words:** Transcriptional regulation, epigenomics, enhancer, enhancer cluster,

24 chromatin structure, crowdsourcing, information content, binding affinity, cell-

25 specificity, binding site degeneracy, chromatin hub

26

27

28

29

30

31

32

33

34

35

36

37 **Abstract**

38 To predict *in vivo* occupancy of a transcription factor (TF), current models consider only
39 the immediate genomic context of a putative binding site (BS) – impact of the site’s
40 spatial chromatin context is not known. Using clusters of spatially proximal enhancers,
41 or archipelagos, and DNase footprints to quantify TF occupancy, we report for the first
42 time an emergent group-level effect on occupancy, whereby BS within an archipelago
43 experience greater *in vivo* occupancy than comparable BS outside archipelagos, i.e. BS
44 not in spatial proximity with other homotypic BS. This occupancy boost is tissue-specific
45 and scales robustly with the total number of BS, or enhancers, for the TF in the
46 archipelago. Interestingly, enhancers within an archipelago are non-uniformly impacted
47 by the occupancy boost; specifically, archipelago enhancers that are enriched for BS
48 corresponding to degenerate motifs exhibit the greatest occupancy boost, as well as the
49 highest overall accessibility, evolutionary selection, and expression at neighboring gene
50 loci. Strikingly, archipelago-wide activity scales with expression of TFs with degenerate,
51 but not specific, motifs. We explain these results through biophysical modelling, which
52 suggests that spatially proximal homotypic BS facilitate TF diffusion, and induce boosts
53 in local TF concentration and occupancy. Together, we demonstrate for the first time
54 cooperativity among genomically distal homotypic BS that is contingent upon their
55 spatial proximity, consistent with a TF diffusion model. Through leveraging of three-
56 dimensional chromatin structure and TF availability, weak archipelago binding sites
57 *crowdsource* their occupancy as well as context specificity, with coordinated switch-like
58 effect on overall archipelago activity.

59

60

61

62 **Introduction**

63 Eukaryotic transcriptional regulation is critically mediated by the binding of specific
64 transcription factors (TF) to their cognate DNA binding sites in the genome (Spitz and
65 Furlong 2012). A TF's *in vivo* DNA binding varies dramatically over developmental time
66 and across tissues (Yáñez-Cuna et al. 2012; Plank and Dean 2014), and as such, a TF's *in*
67 *vitro* binding preference, or motif, does not accurately predict its *in vivo* binding
68 (Yáñez-Cuna et al. 2012; Zinzen et al. 2009). Thus, a TF's DNA binding motif suffers
69 from being, both insufficiently informative to precisely specify binding in the large
70 genomic substrate and insensitive to the *in vivo* environment, making it essential to
71 characterize additional determinants of *in vivo* TF-DNA binding (Moses et al. 2004;
72 Heinz et al. 2013).

73 Spatio-temporal variation in TF binding has been shown to be, in part, mediated by the
74 local chromatin state of a binding site (BS) (Hesselberth et al. 2009). Recent work
75 highlights three additional features of functional binding: (1) high GC content (White et
76 al. 2013), (2) cooperative binding (Smith et al. 2013; Yáñez-Cuna et al. 2012) and (3)
77 clusters of homotypic clusters of BS for a common TF, or HCTs (Gotea et al. 2010; Ezer
78 et al. 2014a). These three features have been shown to be enriched in gene promoters
79 and distal enhancers and to contribute to functional *in vivo* binding leading to
80 transcriptional activation (Gotea et al. 2010; White et al. 2013; Arvey et al. 2012; Sharon
81 et al. 2012). Still, most BS predicted by current models are not bound *in vivo* (Arvey et
82 al. 2012; Slattery et al. 2014; Moses et al. 2004).

83 To date, research on determinants of functional TF binding, such as those above, have
84 focused on putative BS and immediate flanking sequences. In parallel, the three-
85 dimensional organization of the genome has emerged as an important mediator of
86 transcriptional regulation, where spatial, not genomic, proximity is determinative
87 (Fullwood et al. 2009; Ing-simmons et al. 2014; Babaei et al. 2015; Filippova et al. 2013).
88 Chromatin looping events have been shown to combine at a higher organizational level,
89 where they confer spatial proximity to functionally related genes and their distal
90 regulatory regions (Fullwood et al. 2009; Li et al. 2012; Schwarzer and Spitz 2014;
91 Fraser 2006; Lieberman-aiden et al. 2009). In vertebrates, for example, Hox genes,
92 globin genes, and olfactory receptors, along with their distal enhancers, adopt a spatially
93 clustered conformation, termed as 'regulatory Archipelago' (AP), as a prerequisite for
94 robust transcriptional activation (Schoenfelder et al. 2010a; Vernimmen 2014;

95 Montavon and Duboule 2012; Schwarzer and Spitz 2014; Markenscoff-Papadimitriou et
96 al. 2014). Despite mounting evidence supporting functional criticality of chromatin
97 interactions in context-specific transcriptional regulation, the potential impact of spatial
98 clustering of BS on their individual TF occupancy has not been investigated. Recent
99 findings that spatially clustered enhancers (we borrow the term ‘archipelago’ to refer to
100 such spatially clustered enhancers) often share BS for the same TF, i.e., homotypic sites
101 (Taher et al., 2013; Malin et al., 2013) make such enquiry even more compelling. Given
102 the impact on functional binding of genomic homotypic clusters of BS (He et al. 2012b;
103 Ezer et al. 2014a), we investigated the impact of *spatial* homotypic clusters of BS – that
104 is, spatially clustered but genomically distant BS for a mutual TF.

105 In what follows, it's crucial to distinguish binding affinity of a TF for a BS, which is
106 typically assessed *in vitro*, from TF occupancy at a BS, which is an *in vivo* state and
107 depends on additional factors – most directly, TF concentration (Foat et al. 2006).
108 Importantly, TF concentration and, hence, TF occupancy, may be non-uniformly
109 spatially distributed in the nucleus (Schoenfelder et al. 2010a; Chakalova and Fraser
110 2010). Indeed, as described by facilitated TF diffusion, BS for a common TF in a HCT act
111 together to briefly 'trap' a TF into diffusing back and forth amongst themselves along the
112 chromatin (Brackley et al. 2012; Ezer et al. 2014a, 2014b), resulting in higher-than-
113 expected occupancy in the HCT.

114 Here, based on clusters of spatially proximal enhancers, or APs (Malin et al 2013,
115 Sheffield et al 2013), and using nucleotide-resolution DNase footprints to quantify
116 context-specific *in vivo* TF occupancy data (Neph et al. 2012), we demonstrate a strong
117 group-level effect on TF occupancy whereby individual BS within an AP experience
118 greater *in vivo* occupancy than their counterparts outside APs, i.e., enhancers that are
119 not in spatial proximity with other enhancers. We refer to the differential occupancy as
120 ‘*occupancy boost*’. Strikingly, occupancy boost for a TF in an AP scales robustly with
121 both the number of putative BS per AP enhancer, i.e., homotypicity, as well as the
122 number of accessible enhancers in the AP. TFs with degenerate motifs, which are
123 expected to have abundant putative BS, are consistently among the TFs subject to the
124 greatest occupancy boosts; in large APs, mean boosts for homotypic BS corresponding to
125 highly degenerate motifs are between 2 and 3-fold. Based on these results, we propose
126 that *in vivo* occupancy at particular BS in an AP is amplified by the presence of
127 homotypic BS in spatial proximity, i.e., BS ‘*crowdsource*’ their own occupancy boost
128 along with other homotypic BS in their spatial proximity. We extend the aforementioned

129 biophysical model of facilitated diffusion of TFs within one HCT to show that the
130 observed occupancy boost can result from TFs briefly ‘trapped’, or diffusing, among
131 multiple spatially proximal HCTs. Among AP enhancers, we find that specifically those
132 enriched in degenerate motifs (‘enriched enhancers’) are most affected when compared
133 with AP enhancer depleted for degenerate motifs and comparable non-AP enhancers;
134 they exhibit several-fold greater boost in activity, their neighboring genes exhibit several-
135 fold greater expression and, consistent with functional significance, they exhibit a greater
136 evolutionary conservation. Tellingly, mean AP-wide chromatin accessibility scales with
137 the expression of cognate TFs with degenerate – but, not specific – motifs. Combined
138 with other results, this points to a role for crowdsourcing in: (i) initiating a positive
139 feedback loop whereby greater TF occupancy at enriched enhancer BS increases the
140 overall accessibility at these enhancers, thus facilitating further occupancy; (ii) endowing
141 enriched enhancers with switch-like behavior, activating them in specifically those
142 tissues where chromatin structure and TF availability together result in sufficient
143 occupancy boost

144

145

146 **Results**

147 **Analysis overview**

148 Our analysis is based on previously identified enhancer clusters (Malin et al. 2013)
149 comprising ~1600 enhancers in 40 clusters. Enhancers were clustered based on
150 correlated DNase hypersensitivity (DHS) profiles across 37 cell lines (representing 82
151 cell lines). Enhancers in the same cluster were shown to (i) have functionally related
152 gene neighbors with correlated expression, indicative of coordinated regulation, (ii)
153 share BS for several TFs, and (iii) be spatially proximal to one another. We will refer to
154 such enhancer clusters as 'archipelagos' (APs) borrowing from (Spitz and Furlong 2012).
155 We refined the APs identified in (Malin et al. 2013) to ensure tight spatial proximity
156 among AP enhancers (see Methods). Note that properties (ii) and (iii) above together
157 imply a higher spatial density of homotypic BS within an AP, particularly for TFs with
158 abundant putative BS (Fig. 1A); this generally corresponds to TFs with degenerate
159 motifs.

160 Previous study suggests that genomic clustering of homotypic BS can boost the cognate
161 TF's occupancy at individual BS in the cluster (Brackley et al., 2012; Ezer et al., 2014b).
162 By extending the notion of genomic clusters of homotypic BS to *spatial* clusters of
163 homotypic BS, here we hypothesize that analogously, mediated by conformational
164 changes in chromatin organization, such spatial BS cluster may boost binding occupancy
165 for a cognate TF at individual BS, with potential downstream functional impacts. To test
166 this hypothesis, we **(i)** contrasted *in vivo* TF binding occupancy in AP enhancers to that
167 in a stringently controlled set of '*non-AP*' enhancers in the same tissue, **(ii)** assessed
168 functional impact of occupancy boost at AP enhancers – particularly those enriched for
169 degenerate motif sites – and their putative gene targets, **(iii)** characterized context-
170 specificity of the occupancy boost, and **(iv)** tested context-specific TF availability as an
171 upstream driver of the occupancy boost.

172 In the following analyses, *in vivo* occupancy was estimated for each putative BS using,
173 alternatively, high-resolution curated cell type-specific DNase footprint data (Neph et al.
174 2012) or ChIP-Seq (www.encodeproject.org/ENCODE) (see Methods). For additional
175 validation, key tests were repeated using an alternative set of previously published APs
176 (Sheffield et al. 2013). *In vivo* occupancy of a TF and other functional analyses were
177 primarily performed in each AP's most active cell line out of 34 examined– its so-called
178 'AP-active' tissue (see Methods), which offered 730K BS for analysis and 1.8m more in

179 the alternative AP dataset. Results in the AP-active tissue were then contrasted with
180 those in 'AP-inactive' cell lines. With respect to background for BS-level analyses, for a
181 given TF, each AP enhancer was matched to a non-AP enhancer controlling for
182 chromatin accessibility and the number and type of putative BS (Methods).

183 **Occupancy boost at AP BS increases with homotypic BS density within AP,** 184 **supporting crowdsourcing of *in vivo* TF occupancy**

185 A putative BS for a TF is deemed bound *in vivo* by the TF if the BS overlaps DNase
186 footprint based on stringent criteria (Fig. 1B and Methods). For a collection of homotypic
187 BS, occupancy is quantified as the proportion that are bound *in vivo*. We first found that
188 the *in vivo* occupancy across all BS in all APs in their respective AP-active cell lines
189 (Methods) is significantly higher than the occupancy across matched non-AP enhancers
190 (Wilcoxon test p-value = 2.5E-7), albeit with modest effect size (odds ratio of occupancies
191 = 1.21). However, according to our central hypothesis we expect AP occupancy to be
192 pronounced primarily for the subset of TFs with large numbers of BS in spatial proximity
193 in an AP (Fig. 1A). We therefore examined such TFs explicitly.

194 For a given TF in one AP, we calculated the TF's *coverage* as the total number of cognate
195 BS in the AP, and calculated its *occupancy boost* as the ratio of occupancy at these BS to
196 the occupancy at a matched set of non-AP enhancer BS (Methods). Such TF-AP pairs are
197 the fundamental unit at which the crowdsourcing hypothesis predicts occupancy will be
198 impacted. Because background levels of BS occupancy in the genome are generally low
199 (3-4%), the occupancy is zero in both AP and control non-AP enhancer sets for a
200 majority (65%) of the 25K TF-AP pairs; these pairs were excluded for this analysis. Of
201 the remaining TF-AP pairs, 3.6k have non-zero occupancy in both AP and non-AP,
202 encompassing ~95K enhancer-TF pairs and ~205K BS (we call this the *reciprocal* set),
203 and additional 5k TF-AP have non-zero occupancy in either AP or in matched non-AP BS
204 (*non-reciprocal* set). We analyze the two sets of TF-AP pairs separately.

205 We stratified the reciprocally occupied TF-APs into 8 bins with exponentially increasing
206 coverage cutoffs and calculated the overall occupancy boost for each bin as the mean
207 occupancy boost among member TF-APs. A 95% confidence interval for each bin's mean
208 was estimated using a bootstrap. As shown in Fig. 2A, the occupancy boost robustly
209 increases with the TF coverage in the AP. Specifically, we found a substantial difference
210 in occupancy boost between TF-APs with the highest and lowest 50% coverage (mean of
211 77.7 % versus 2.1 %; Wilcoxon p-value = 1.4e-5). This trend also holds when coverage

212 was alternatively quantified as the number of enhancers in an AP with at least one BS for
213 the TF (Fig. 2B), suggesting that the boost is not due to disproportionate contribution
214 from a few enhancers, but instead relies on widely dispersed BS across the AP's
215 enhancers. Interestingly, the boosts for high coverage TF-APs increase when the digital
216 footprint binding criterion for assessing occupancy is made more stringent
217 (Supplemental Fig. S1A).

218 Abundance of a TF's cognate BS is strongly correlated with its motif degeneracy, as
219 quantified by the motif's relative entropy (RE) (Fig. 2C) – low RE is identified with high
220 degeneracy (Hannenhalli 2008). Given this association, we also directly assessed the
221 relationship between TF motif RE and occupancy and found very similar trends
222 (Supplemental Fig. S1B). Interestingly, the steepest declines in occupancy boost and in
223 the distribution of TF RE echo one another, as shown, in that both occur at RE = 7.
224 When, in an independent set of 11 tissues, we instead used ChIP-Seq to infer occupancy
225 the trend, again, remained very similar (Supplemental Fig. S1C). Taken together, the
226 above analyses strongly suggest that binding sites for high coverage TFs experience a
227 substantial occupancy boost in AP enhancers relative to BS in comparable non-AP
228 enhancers.

229 The overall TF coverage is affected by both the mean number of BS per AP enhancer
230 ('homotypicity') and the number of enhancers per AP ('AP size'). Next, we assessed the
231 relative contributions of these two constituents of coverage on the occupancy boost. As
232 shown in Fig. 2D, for the reciprocal set, AP size and homotypicity independently and
233 robustly impact the magnitude of occupancy boost (p-value = $4.2e-6$). A similar analysis
234 on 5K non-reciprocal TF-AP pairs shows a similar and significant trend (Supplemental
235 Fig. S1D; p-value $8.1E-5$). Our overall conclusion does not change when we used an
236 independent set of 450 AP enhancer clusters reported in (Sheffield et al. 2013) based on
237 1.8 million BS (supplementary Note 1) (Figs. S2A, S2B). Consistent with the expected
238 importance of spatial proximity, as the Hi-C based screen for pairwise distance between
239 fellow AP member enhancers became more stringent, the trend improved
240 (Supplementary Note 1).

241 The observed link between spatial BS abundance and occupancy, we reasoned, may
242 partly be mediated by cooperativity among the bound TFs within an AP (Martinez and
243 Rao 2012). However, as explained in supplementary Note 2, we found that even though
244 heterodimerizing TFs exhibit a somewhat (but significantly) larger boost than other TFs,

245 this difference can largely be explained by their greater BS coverage (Supplemental Fig.
246 S3. Taken together, these results strongly suggest that AP TF's occupancy at a specific BS
247 in an AP is 'crowdsourced' by the collective of its putative BS in the AP.

248 **TF occupancy boost in spatial clusters of BS is consistent with a facilitated-** 249 **diffusion model**

250 Previous studies have shown that a facilitated diffusion model can explain the greater occupancy
251 in homotypic clusters of BS (Brackley et al. 2012). Here we simulated an extended version of the
252 biophysical model for HCTs in isolation to determine whether the crowdsourcing effect is
253 sufficient to explain the observed AP-mediated occupancy boost. All details pertaining to the
254 model, algorithms, and results are provided in Supplemental File 'Facilitated_Diffusion_Model'.
255 First, in a simple scenario, we investigated specifically how occupancy is affected by (i) the
256 spatial distance between enhancers within an AP cluster, (ii) the genomic distance between
257 binding sites within an AP enhancer, and (iii) the number of binding sites within each enhancer.
258 Then, we investigated various 3D organizations of homotypic binding sites that would closely
259 represent the crowding of low-RE binding sites in AP clusters *in vivo*. Our simulation results
260 suggest that both the presence of homotypic clusters and inter-strand jumping between
261 enhancers within an AP can increase the average TF occupancy. For instance, in the case of four
262 enhancers containing pairs of homotypic binding sites, as shown in Fig. 3B, there was a 60% to
263 170% increase in TF occupancy when the enhancers were part of AP clusters (in which
264 enhancers are 100nm to 200nm apart, as compared to being 10000 nm apart, by default), and
265 in the case of eight enhancers containing pairs of homotypic binding sites there was an 118% to
266 277% increase occupancy (Fig. 3D), which suggests that the crowdsourcing effect is a
267 biophysically sound strategy for increasing local TF occupancy in APs at a biologically
268 meaningful scale.

269 **AP enhancers that are enriched for degenerate motifs have greater occupancy** 270 **boost and higher downstream functional impact**

271 Previous studies have shown a trend for AP enhancers to share homotypic BS (Malin et
272 al 2013, Sheffield et al 2013). We reasoned that this could result if AP enhancers were
273 enriched for degenerate, i.e., motifs with low RE, which are expected to be generally
274 abundant, thus increasing the chance of being shared among AP enhancers. We found
275 this to be true (Fig. 4D). Furthermore, our empirical results above and the biophysical
276 model suggest that homotypicity within an enhancer as well as spatial clustering of such
277 enhancers can increase TF occupancy at biologically meaningful scales. Together, these

278 led us to hypothesize that the particular AP enhancers that are enriched for low-RE
279 (degenerate) motif BS may be especially affected by crowdsourcing. In the following
280 analyses, we employ multiple thresholds for RE to define ‘low’ RE motifs.

281 Activity of an enhancer, and ultimately, activation of its target gene, is associated with
282 the amount of TF binding in the enhancer (Smith et al. 2013; Fisher et al. 2012). Hence,
283 given the occupancy boost at low-RE TF sites in AP enhancers, we hypothesized that the
284 enhancers that are enriched for low-RE TF sites will have the largest overall increase in
285 activity, and by proxy, increase in expression of their putative target genes. For each AP
286 enhancer, assuming its closest gene neighbor to be its putative target (Djebali et al.
287 2012), we calculated its ‘expression boost’, similar to occupancy boost above, as the
288 relative change in expression of the gene relative to the target gene of the control non-AP
289 enhancer (Methods). We then compared expression boost for AP enhancers enriched for
290 low-RE motif BS (*enriched* enhancers) to that for AP enhancers depleted in low-RE sites
291 (*depleted* enhancers) (Methods). As shown in Fig. 4, the putative target genes of
292 enriched enhancers have much greater expression in AP-active tissues than their non-AP
293 counterparts, while the depleted enhancers do not. Moreover, as the degree of
294 enrichment increases from top 50% to top 10%, the relative expression boost increases
295 from 62% to 196% (for low-RE cutoff of 5); The difference in AP and non-AP neighbor
296 gene expression for depleted enhancers is nearly as stark – but in the opposite direction:
297 as the degree of enrichment increases from top 50% to top 10%, the drop in AP
298 expression relative to non-AP grows from -32% to -72% (i.e., non-AP expression is 3.5-
299 fold higher). GC content differences cannot explain either trend, as GC content in non-
300 AP enriched (depleted) enhancers is <10% (15%) higher than in corresponding AP
301 enhancers. Also, notably, the highest expression boost (Figure 4A, row4) and the highest
302 occupancy boost (Fig. S1B) occur at similar low-RE cutoffs, thus strengthening the link
303 between the two; at higher low-RE cutoff (being more permissive) the observed effect
304 weakens and eventually disappears.

305 As a complementary test (Supplementary Note 3), we compared TF binding patterns
306 between enhancers within 50Kb of highly expressed genes and enhancers within 50kb of
307 lowly expressed genes. Consistent with above, we found that the low-RE motif usage
308 (defined as the fraction of bound BS that are low-RE) was 1.8 to 3.0 times higher in AP
309 enhancers near highly expressed genes than in those near lowly expressed genes. As a
310 negative control, we did not observe this pattern for matched non-AP enhancers. Taken
311 together, these results suggest that low-RE binding specifically at enriched AP enhancers

312 have a significant impact on downstream gene expression. As an additional
313 ascertainment of the functional importance of enriched AP enhancers, we found such
314 enhancers to be up to 120% (90%) more evolutionarily conserved at RE cutoff of 5 (4)
315 (using 20 species PhastCons score (Siepel et al. 2005)) than matched non-AP enhancers;
316 indeed, the greater their enrichment, the greater the evolutionary constraint we observed
317 (Fig. 4A, row 3). Depleted AP enhancers, by contrast, were at most 40% more conserved
318 than their non-AP counterparts.

319 TF binding and chromatin accessibility are intimately connected; higher accessibility
320 typically leads to higher occupancy, while TF binding can help displace a nucleosome
321 and increase accessibility (Teif and Rippe 2012). Therefore, we assessed whether
322 enriched enhancers exhibit a greater boost (relative to matched non-AP enhancers) in
323 overall accessibility compared with depleted enhancers. For this analysis, we normalized
324 AP enhancer accessibility by that of non-AP enhancers matched one-to-one with AP
325 enhancers for numbers of low and not-low RE BS, as described previously, except DHS
326 was explicitly left uncontrolled. As shown in Fig. 4A (row 2) and Fig. 4B, at RE threshold
327 of 5, the most enriched enhancers exhibit ~10-fold greater DHS boost (with respect to
328 matched non-AP enhancers) than do depleted enhancers; as expected, the trend weakens
329 for higher RE thresholds. To further resolve the effect of BS RE on enhancer
330 accessibility, we tracked changes in accessibility as we increased the number of low-RE
331 (high-RE) sites, while holding relatively constant the number of high-RE (low-RE) sites.
332 As shown in Supplemental Fig. S4, increasing the number of low-RE BS (for relatively
333 constant high-RE BS count) has a substantial positive impact on enhancer's accessibility
334 – especially when the number of high-RE BS is low – while increasing the number of
335 high-RE BS (for a relatively constant low-RE BS count) has an insignificant or negative
336 impact. In addition, we found that histone acetylation levels (H3K27Ac), which are
337 associated with active enhancers, had a ratio between enriched and depleted AP
338 enhancers 3-fold higher than the same ratio in non-AP enhancers (Fig. 4E).

339 Finally, we observed occupancy boosts in enriched enhancers that, unexpectedly, were
340 up to two-fold higher than in depleted enhancers, often in the same AP and for the same
341 levels of coverage (220% vs. 110% at low-RE cutoff of 4 (not shown); 93% vs 39% in
342 alternative APs (Sheffield et al 2013) (Supplementary Note 4; Fig. 4C, Supplemental Fig.
343 S2C). To account for these unexpectedly high boosts, we address the potential for higher-
344 order interactions in enriched enhancers among BS for distinct TFs with degenerate
345 motifs (see Discussion)

346 These results – the relatively higher occupancy boosts, chromatin accessibility,
347 downstream gene expression, and evolutionary constraint in enriched enhancers, along
348 with greater prevalence of enriched enhancers among AP than non-AP enhancers –
349 strongly suggest a hitherto unreported special functional relevance of AP enhancers that
350 are enriched for low-RE binding sites.

351 **Cell type-specificity of AP enhancer occupancy boost and activity**

352 Given the link between occupancy boost and spatial clustering of BS, and given the
353 context-specificity of spatial proximity (Ay et al. 2014), we expect the occupancy boost to
354 exhibit cell type specificity. In addition to identifying the cell type where an AP is deemed
355 active (as employed in analyses thus far), we also identified the cell types where an AP is
356 deemed inactive, namely those where less than 40% of the AP enhancers were DNase I
357 hypersensitive. To offset the paucity of bound sites in inactive tissues, all qualifying
358 inactive tissues for each AP were pooled. We found that for the TF-AP pairs in the
359 highest coverage bin, occupancy boost dropped from ~112% in its AP-active cell type
360 down to 38% in inactive cell types (Fig. 5A). In addition, we estimated tissue specificity
361 of each TF as the cross-tissue dynamic range of its occupancy, defined as the ratio of its
362 occupancy in AP-active tissue(s) to that in AP-inactive tissues, calculated over the
363 identical AP BS. Notably, this provides evidence of the occupancy boost's tight
364 association with coverage without the need for non-AP occupancy as a baseline
365 (Methods). After controlling for DHS across coverage bins, we find that the TF-APs with
366 top 10% coverage display 135% greater occupancy in active relative to inactive tissues,
367 while in the matched non-AP context it is 38% (Fig. 5B). Even larger differentials
368 between AP and non-AP contexts were observed for their respective ratios of non-
369 reciprocal binding in active and inactive tissues (Supplemental Fig. S5A). Interestingly,
370 we found that high coverage TF-AP pairs for heterodimerizing TFs exhibit substantially
371 higher specificity than other TFs (225% vs. 140%) (Supplemental Fig. S5B), particularly
372 TFs in MADS and bZIP domain families, suggesting an augmented level of cooperative
373 binding in APs. This, we suspect, is due to the relatively binary nature of cooperative
374 binding: in response to small increments in TF concentrations, heterodimers exhibit
375 disproportionately large changes in occupancy (Giorgetti et al 2010).

376 **AP activity is correlated with availability of TFs, specifically those with degenerate motifs**

378 Our results thus far suggest that crowdsourcing may be intimately connected to the

379 regulation of AP enhancer-gene complexes, as it provides a way for the cell to prime or
380 induce activity in multiple genomic elements simultaneously, in a specific spatial and
381 tissue context. As shown above, the boost in overall activity (approximated by DHS) of
382 an AP enhancer is in fact far higher in AP enhancers enriched for low-RE (high coverage)
383 BS. However, it is not clear whether the binding of TFs corresponding to the low-RE
384 motifs increases the overall accessibility at the enriched enhancer, or alternatively, an
385 already increased accessibility at enriched enhancer (by some unknown mechanism) is
386 responsible for greater occupancy of particular TFs at those enhancer. In order to resolve
387 this circularity, we tracked tissue-specific gene expression of TFs in 11 cell types and
388 studied its relationship with tissue-specific AP enhancer accessibility (Methods). As
389 shown in Fig. 5C, mean AP enhancer accessibility increases robustly with the expression
390 of TFs comprising high-coverage (red), but not low-coverage, TF-APs (gray), for which
391 there was a slight but significant inverse relationship, also noted in an analysis above
392 (Fig. 2A). No such associations were observed in non-AP enhancer sets controlled for
393 low- and high-RE BS counts (Fig. 5D). Thus, AP enhancer accessibility and activity is
394 highly responsive to the levels of high-coverage AP TFs as they vary across tissues.

395 Together, these results strongly suggest that crowdsourced boosts in TF occupancy,
396 through the context-specific binding of high coverage TFs, may help drive tissue-specific
397 activation of enhancer networks and their target gene complexes.

398 **Discussion**

399 **Summary.** Here, we have shown that a TF's *in vivo* occupancy at a particular cognate BS is
400 much greater when the BS is in spatial clustered with other homotypic BS (i.e., in an AP) than
401 when it is not. Strikingly, the size of the occupancy boost robustly scales with the number of BS,
402 or relatedly, the number of enhancers, in the archipelago (AP), suggesting, for the first time, that
403 the BS in an AP cooperatively crowdsource their own occupancy. To ensure the robustness of
404 our conclusions, we used stringent controls and employed multiple (i) sources for AP enhancers
405 (Malin et al 2013, Sheffield et al 2013), (ii) experimental backgrounds (non-AP enhancers in the
406 AP-active tissue, the same enhancer in AP-inactive tissues), (iii) occupancy scales (per BS, per
407 enhancer), and (iv) types of occupancy data (curated digital footprints, ChIP-Seq. These
408 observations are not adequately explained by current models, however, they closely agree with a
409 standard biophysical model of facilitated TF diffusion that duly accounts for the augmented
410 diffusion of TFs among spatially proximal homotypic BS. Effectively, a collective of homotypic
411 clusters of TF BS (HCTs) cooperatively alter their microenvironment, raising the local

412 concentration of their cognate TF. The accompanying occupancy boost, in turn, contributes to
413 higher overall enhancer activity and putative target gene expression.

414 **Genomic versus Spatial homotypicity.** Our work synthesizes the regulatory roles of HCTs
415 (e.g. Crocker et al 2015), and of stable chromatin structures (e.g. Downen et al., 2014), by
416 showing that it is precisely the interplay of numerous HCTs mediated by chromatin folding that
417 gives rise to the previously undocumented biophysical effect that we have termed
418 *crowdsourcing*. Enhancer-enhancer interactions have been reported in the context of HOX and
419 globin gene regulation as well as in high-throughput ChIA-PET assays, but their functional
420 nature has remained elusive. A notable exception, spatial clustering of enhancers around an
421 olfactory receptor gene have been associated with removal of repressive H3K9me3
422 (Markenscoff-Papadimitriou et al. 2014); it is plausible that crowdsourcing is an upstream
423 trigger of this change – through a general remodeling of the local chromatin state, or through
424 increased binding of a TF that mediates chromatin remodeling.

425 **Higher order impact of crowdsourcing.** Unexpectedly, we observed up to two-fold higher
426 occupancy boost and 10-fold greater normalized chromatin accessibility in AP enhancers
427 enriched for degenerate motifs ('enriched enhancers') than in depleted enhancers. A likely
428 explanation is the emergence of an aggregate occupancy effect among an enriched enhancer's
429 frequent degenerate BS, which serves to remodel the local chromatin state. Under inactive
430 conditions – that is, in AP-inactive tissues or outside of APs – we found that enriched enhancers
431 (which inherently tend toward far lower GC content than depleted enhancers) display
432 substantially higher chromatin accessibility compared to depleted enhancers. This is consistent
433 with previous work suggesting that nucleosomes favor unbound, low GC-content sequence, yet
434 are readily displaced by strongly binding pioneer factors, or, as in the case of crowdsourcing, by
435 an aggregate of distinct TFs (Barozzi et al. 2014; Wasson and Hartemink 2009).

436 In an AP-active tissue, conversely, enriched AP enhancers experience a widespread surge in
437 binding, thereby displacing the nucleosome and boosting occupancy further, in a positive
438 feedback loop (Fig. 5E). Taken together, the markedly divergent accessibility inside versus
439 outside an active AP confer to enriched enhancers switch-like behavior, where their state is
440 determined by their context: harbored in an AP replete with degenerate homotypic BS, their
441 accessibility increases – but only in tissues in which the cognate TFs are available. In light of this
442 highly context-specific activation and the rapid evolutionary gain of BS for degenerate motifs,
443 we suggest that enriched AP enhancers can evolve adaptively relatively free of consequences
444 from spurious binding. This is the first work to highlight the special functional significance of AP

445 enhancers enriched for abundant, low-affinity BS. However, further work is needed to confirm
446 to what extent depleted enhancers, whose neighbor genes are expressed at up to three-fold
447 *lower* levels than those of similar non-AP enhancers, have a unique, perhaps repressive, role.
448 Interestingly, genes controlling cell identity in stabilized chromatin structures were found
449 accompanied by repressed genes that coded for yet other lineage-specifying regulators (Downen
450 et al 2014).

451 Crowdsourcing integrates well with the two prevailing models of coordinated activation of
452 spatially co-localized gene complexes (Fig. 5E), while providing a missing piece of the puzzle.
453 Whether (a) long-range enhancer-gene loops form *de novo* upon (or along with) activation of a
454 gene cluster (Deng et al. 2012), or (b) the loops are pre-formed and are activated by TF
455 availability (Ghavi-Helm et al. 2014), the cell requires TFs to functionally bind and activate
456 elements specifically in a targeted gene cluster. Crowdsourcing of low-affinity BS is well-suited
457 for such targeting, as it can induce specificity through emergent switch-like binding behavior.
458 Interestingly, a recent study showed a strong correlation between pathway-level gene activity
459 and pathway-level spatial proximity across cell types (Karathia, Hannenhalli et al., manuscript
460 under review), suggesting that chromatin structure is intimately connected with gene complex
461 activation, consistent with (a), above. Unlike direct enhancer-gene interaction in the standard
462 model for distal transcriptional regulation, crowdsourcing, interestingly, is not observable at the
463 level of single enhancer-gene interaction, but instead emerges only at higher levels of chromatin
464 organization and co-regulated gene modules.

465 **Tissue specificity and cooperative binding.** We found that crowdsourcing is highly tissue-
466 specific, as high-coverage AP BS exhibit several-fold greater occupancy in AP-active relative to
467 AP-inactive tissues. Such tissue specificity is consistent with the dependence of crowdsourcing
468 on chromatin context and TF availability, where differential TF availability likely acts not only
469 directly but also by influencing higher-order chromatin conformation (Pombo and Dillon 2015).
470 Crowdsourcing endows the cell with a high degree of fine-grained regulatory control, as
471 occupancy boost magnitude is shaped by the collective availability of multiple TFs and
472 conditioned on the chromatin-induced spatial proximity of their cognate sites. Fundamentally,
473 crowdsourcing provides an alternative mechanism of cooperativity to direct cooperative binding
474 of heterodimerizing TFs, an established source of tissue-specificity. Indeed, crowdsourcing acts
475 complementarily to cooperative binding (Supplementary Note 2).

476 **Differential occupancy as a vehicle for specificity.** In contrast to previous work
477 underscoring the functional importance of weak (low occupancy) binding at non-consensus or

478 inherently weak BS that typical ChIP-Seq processing tends to miss due to stringent cutoffs
479 (Tanay 2006; Biggin 2011; Essien et al. 2009), crowdsourcing leverages spatial chromatin
480 context to imbue inherently low-affinity sites with unexpectedly high-occupancy binding.
481 Unidentified crowdsourcing may therefore underlie previous reports linking particular low-
482 affinity sites with context-specific regulation (e.g. Ramos and Barolo 2013), or linking individual
483 HCTs to unusually robust binding, for example in regulatory regions upstream of HOX genes
484 (Crocker et al., 2015). Indeed, occupancy boosts we observed at spatially clustered HCTs were
485 computed with respect to nominally ‘isolated’ HCTs. As shown by Crocker et al (2015),
486 occupancy is more robust where degenerate homotypic sites are located in genomic clusters.
487 HCTs, however, are highly abundant in the genome (Gotea et al. 2010) as well as, by nature,
488 spatiotemporally invariant, which raises a well-known conundrum, viz. how a TF discriminates
489 among a multitude of candidate BS (Stewart and Plotkin 2013; Stewart et al. 2012). In contrast
490 to the static and relatively low specificity of an individual HCT, a large collective of homotypic
491 low-affinity sites can attain high specificity and spatiotemporal responsiveness precisely by their
492 capacity to reconfigure the local TF environment *en masse* – in specific favorable chromatin
493 contexts. Subject to coordinate regulation, a locus is therefore targeted, not through a fixed,
494 individual address on the one-dimensional genome, but as a conditional and collective nexus in
495 the 3-D chromatin topology – a mobile area code to which the motif’s short sequence is
496 appended.

497 **Potential implications for transcription factories, superenhancers.** An archipelago, as
498 described here, represents a group of spatially clustered enhancers and their likely target genes,
499 which are often functionally related (Sheffield et al. 2013; Malin et al. 2013). Meeting this same
500 general description are subnuclear compartments known as transcription factories (Edelman
501 and Fraser 2012). Transcription factories have been shown to concentrate resources such as
502 RNA PolII, core components of transcription, as well as some master TF regulators
503 (Schoenfelder et al. 2010b). However, it is unclear precisely how distinct factories achieve
504 specific and differential concentrations of master regulator TFs (Schoenfelder et al. 2010a).
505 Crowdsourcing offers a possible explanation, and is consistent with a speculated role for
506 resident sequences (Schoenfelder et al. 2010a). While it is generally assumed that high
507 concentrations of TFs are critical in recruiting genes and their distal regulatory regions to the
508 factory, our work suggests alternative causality, as supported by formal biophysical simulations.
509 Although not confirmed, our characterization of archipelagos suggests their operational overlap
510 with factories.

511 Intriguingly, there is also ample overlap between the conditions for crowdsourcing and known

512 features of superenhancers – enhancers 100Kb or longer in length comprising numerous
513 smaller regions, which regulate genes critical to cell identity (Whyte et al. 2013). Recent works
514 have shown that superenhancers are highly tissue-specific and densely occupied by master
515 regulator TFs, which often recognize degenerate motifs (Vahedi et al. 2015; Heinz et al. 2015).
516 Critically, superenhancers also reveal unusually high levels of spatial interaction among their
517 subunits (Heinz et al. 2015). We thus speculate that crowdsourcing may play a role in
518 superenhancer function.

519 **Methods**

520 **Enhancer clusters ('APs').** In previous work, genomically dispersed clusters of
521 enhancers with correlated activity across cell lines showed evidence of spatial proximity,
522 particularly in tissues in which the enhancers were active, where spatial proximity
523 between two genomic segments was inferred from Hi-C (Malin et al. 2013). Starting with
524 previously published 40 enhancer clusters, we iteratively filtered out the enhancers from
525 each cluster whose mean spatial proximity to other enhancers was at least one standard
526 deviation below the original mean across all enhancers in the cluster. This results in 40
527 APs with a total of 1480 enhancers (Supplementary File 1) with ~37 per AP, ranging from
528 6 to 89 enhancers per AP.

529 **Determining *in vivo* occupancy at a BS using digital footprint data.** Putative
530 BS in each enhancer were identified using TRANSFAC vertebrate motifs (Matys et al.
531 2006) and motif scanning tool PWM_SCAN (Levy and Hannenhalli 2002) at 95
532 percentile score cutoff. We identified *in vivo* TF occupancy by overlapping putative BS
533 with the high-confidence genome-wide digital DNase hypersensitivity footprints
534 identified in 38 human cell lines (Neph et al. 2012). Digital footprints are a single-base-
535 pair resolution readout in which the absence of aligned reads in a particular segment of
536 open chromatin has been shown to predict binding of a protein (Neph et al. 2012). For a
537 TF, a particular putative BS was considered bound by the cognate TF if there was specific
538 overlap between the BS and a footprint, with further requirement that (i) the midpoint of
539 a footprint must overlap the BS; (ii) the midpoint of the BS must overlap the footprint;
540 and (iii) $BS\ length + 1 > footprint\ length > BS\ length - 4$. The latter criteria excludes
541 otherwise significant footprints that are either too short or too long to confidently be
542 associated with a given motif instance. When a footprint strongly overlaps sites for
543 multiple TFs, it was included in the analysis for all such TFs. These highly stringent
544 criteria were applied equally to AP and to non-AP data.

545 **AP-active and AP-inactive cell lines.** For each AP, we identified the cell line in
546 which it was most active. Cell lines deemed active for a given AP are those in which at
547 least 80% of the AP's enhancers are in open chromatin regions, based on overlap with
548 DHS narrow peaks. In case of more than one such tissue, except where noted, we
549 selected the tissue with the highest percentage of open enhancers (see Fig. 1A), depicting
550 workflow. Approximately 95 percent of AP enhancers were found to be accessible in an
551 AP's 'most active tissue', which for the 40 APs, span 15 distinct cell types out of 34 tested.

552 **Establishing non-AP control.** To establish a non-AP control, for each combination of
553 TF and AP enhancer we identified a non-AP enhancer (sampled with replacement) with
554 an identical motif profile, *i.e.* the vector containing the number of instances of each motif
555 mapping to the given TF. This is an important control, as the number of homotypic BS in
556 an enhancer that are cognate to a given TF impacts occupancy (He et al. 2012a). We note
557 that AP and non-AP enhancer have very similar distributions of total BS and length.
558 Additionally, for each TF motif and AP, AP enhancers' mean DHS in the AP's most active
559 tissue was matched to within 5% in the corresponding non-AP enhancers' mean DHS in
560 the same tissue. Any TF-AP enhancer pair for which a non-AP could not be found
561 meeting these tight controls was excluded. This procedure yielded 430K AP and non-AP
562 TF-enhancer pairs that harbored 730K BS, of which 31K BS had a DNase footprint
563 suggestive of a binding event.

564 **Determining TF occupancy at enhancer resolution with ChIP-Seq data**

565 We used ENCODE NarrowPeak Chip-Seq data for 11 cell types – AG10803, AoAF, HA-h,
566 HAEpiC, HCM, HEEpiC, HFF, HIPEpiC, HMF, HMVEC-dBl-Ad, HMVEC-dLy-Neo,
567 HVMF, NH-A, NHDF-Ad, NHLF – which gave a total of 135 TFs and 294 TF-cell line
568 pairs. To increase sample size, all AP-active cell types were used for each AP – that is, all
569 cell types for which 90% of the AP's enhancers were DNase hypersensitive (Using a
570 cutoff of, alternatively, 80% or 100% did not change the observed trend). Enhancer
571 occupancy by a given TF was determined based on overlap between a \pm 50bp window
572 surrounding the ChIP-Seq peak and one or more putative motif instances detected
573 within the enhancer. To mitigate concerns over systematic biases stemming from
574 variability in protocols or labs of origin, we note that all ChIP-Seq data had identical *de*
575 *facto* weighting for AP and non-AP enhancer-TF pairs, since these were matched by
576 motif for BS counts.

577

578 **Comparing neighbor gene expression between AP and non-AP enhancers.** As
579 a proxy for an enhancer's target gene, following convention, we used the gene closest to
580 the enhancer. As an extra measure of stringency, in case of non-AP enhancer, we
581 excluded those enhancers that were farther than 50kb from the nearest gene promoter.
582 We paired each AP enhancer with one of the remaining non-AP enhancers while
583 controlling for DHS peak height (within 2%) and numbers of both low-RE and high-RE
584 sites (within 2%), where RE class is based on a (variable) RE threshold. This yielded
585 ~1200 pairs of AP and matched non-AP enhancers; the exact number varied with the RE
586 threshold. For gene expression, five cell types were used for which overall AP activity,
587 calculated as described above, was at or near its maximum as observed in 15 cell types
588 for which we had digital DNase footprint and RNA-Seq data
589 (www.encodeproject.org/ENCODE). These were HSMM, A549, NHLF, Ag04450,
590 and Bj.

591 **Estimating TF's RE.** The relative entropy was calculated for each TF motif (i.e.,
592 position weight matrix) using TRANSFAC (version 2014.3) (Hannenhalli 2008). In
593 cases where there were multiple motifs associated with a particular TF (coming from
594 different publications etc.), the motif with the lowest RE was chosen, because it is
595 expected to numerically dominate the genome-wide BS for the TF, given its higher
596 degeneracy.

597 **Identifying low-RE enriched and low-RE depleted AP enhancers.** A relative
598 entropy (RE) cutoff was used to classify each putative BS in an enhancer as either low-
599 RE or high-RE (complement of low-RE). This cutoff was varied from RE = 4 (classifies
600 ~2% enhancers as low-RE) to RE = 9 (classifies > 50% enhancers as low-RE). For each
601 AP enhancer, after tallying the number of low-RE and high-RE sites, an enrichment p-
602 value was generated by applying a Fisher Exact test comparing the numbers of low/high
603 RE sites in the enhancer to those in the pooled set of control (non-AP) enhancers. Based
604 on this enrichment p-value, enhancers were sorted, and the top and bottom ranked x% of
605 enhancers compared in subsequent analysis (x in {10, 20, 50}).

606 **Calculating a normalized conservation score.** To compare evolutionary
607 conservation of low-RE BS enriched AP enhancers to low-RE depleted AP enhancers, we
608 used PhastCons scores, based on 20 mammalian species (Siepel et al. 2005), which are
609 resolved to the individual base. Mean scores across the two classes of enhancers were
610 normalized with respect to non-AP enhancers matched one-to-one with an AP enhancer

611 having approximately the same number of low-RE and high-RE BS, based on a variable
612 RE cutoff, Additionally, we ensured that non-AP enhancers were within 50Kb of the
613 promoter of a highly expressed gene, i.e. its fpkm > 1.0 – which, depending on the cell
614 type, includes approximately the ten percent most highly expressed genes.

615 **Determining occupancy boost with alternative set of AP enhancers.** We
616 obtained sets of correlated regions generated in (Sheffield et al 2013). Each Sheffield
617 cluster of DNase hypersensitive (HS) regions initially spanned multiple chromosomes.
618 To make them consistent with enhancer clusters from Malin et al (2013), regions from a
619 single Sheffield cluster located on distinct chromosomes were treated as distinct clusters,
620 and we retained at most the two largest such clusters from each Sheffield cluster.

621 Consistent with previous procedures, we derived enhancer clusters from each Sheffield
622 cluster by only retaining the regions that overlapped a putative enhancer represented by
623 a large pooled set of 98,000 P300 ChIP-Seq peaks used previously (Malin et al 2013).

624 To further cull the thousands of resulting enhancer clusters, we excluded those with < 10
625 enhancers or with mean enhancer DHS < 100 in their most active tissue. We further
626 excluded Sheffield clusters in which fewer than 90% of enhancers were DNase
627 hypersensitive in their most active tissue, resulting in 474 clusters – averaging ~16
628 enhancers each, though ranging to over 100. Similar to above (see ‘AP Enhancers’) we
629 used Hi-C data to screen enhancers in each AP that were less spatially proximal, on
630 average, to the remaining members. To prevent excessive removal of additional
631 enhancers, given the already modest mean pre-screen AP size, we implemented the Hi-C
632 screen in a single pass, without recursively updating each enhancer’s mean Hi-C score
633 after removal of a fellow AP member. This resulted in 472 non-empty APs with an
634 average of ~15 enhancers each.

635 For background control, we used the complement of P300 ChIP-Seq peaks overlapping
636 any of the screened set of approximately 2.6M Sheffield et al DNase hypersensitive
637 regions. This resulted in too few putative enhancers, and so to this we added back ChIP-
638 Seq peaks overlapping any cluster (on one chromosome) of hypersensitive regions with
639 fewer than five members and with mean DHS > 50 in its most active cell type; this
640 produced a background pool of ~18K enhancers. Non-AP enhancers from this set were
641 matched with AP enhancers as described above. In order to accommodate the smaller
642 APs in this alternative dataset, we loosened the stringency on DHS control such that
643 mean DHS was matched to within 1% for AP and non-AP sets, overall, with mean DHS

644 for individual TF-AP combinations constrained to less than double or more than half the
645 mean DHS of their non-AP counterpart.

646 **TF expression-AP activity correlation.** This analysis used data from each of 15 cell
647 lines for every AP, encompassing ~2.4 million BS. Each (TF, AP enhancer, cell line)
648 triplet was assigned (i) a DHS value, corresponding to AP enhancer and cell line; (ii) a
649 coverage score, corresponding to AP enhancer and TF; (iii) a normalized RNA-Seq value
650 corresponding to TF and cell line. Analysis was limited to triplets with a coverage score
651 in the top and bottom 20%. In each of these coverage classes, triplets were further sorted
652 based on the TF's expression in the given cell line and screened to include only triplets
653 with top or, alternatively, bottom 20 (or 25 or 50) percent TF expression. For each
654 coverage class, the percentage difference in mean cell-type specific DHS between the low
655 TF expression and high TF expression cohorts was plotted. Confidence intervals for each
656 percentage difference were computed on the basis of 50K bootstrap replicates.

657 **H3K27Ac levels**

658 We downloaded Encode ChIP-Seq peaks from human umbilical vein cells (HUVEC) for
659 histone mark H3K27Ac, known to be associated with active enhancer states (Calo and
660 Wysocka 2013). This cell line was chosen for its combination of available data and a large
661 number of enhancers in APs that are active in the cell line. We compared the ratio in
662 mean ChIP-Seq levels between top 10% enriched and top 10% depleted AP enhancers to
663 the same ratio for non-AP enhancers, matched one-to-one with the AP enhancers as
664 described above. An AP enhancer and its matched non-AP enhancer were included only
665 if the AP enhancer belonged to an AP that was 'active' in HUVEC (>80% of its enhancers
666 was DNase hypersensitive). This resulted in ~40 enriched and ~100 depleted AP
667 enhancers, and the same number of non-AP enhancers.

668 **Supplemental information includes:**

669
670 Facilitated Diffusion Model
671 Extended Experimental Procedures
672 1 Supplemental Data File
673 5 Supplemental Figures and Legends
674

675 **Acknowledgements**

676 JM wishes to thank Ivan Ovcharenko for valuable feedback throughout the process, and Avinash
677 Das, Kun Wang, and Mahfuza Sharmin for technical assistance.
678
679

680 **Author Contributions**

681 Designed analysis: JM with help from SH and SM. Performed analysis JM with help from HK.
682 Biophysical modelling and simulations: DE, XM. Crowdsourcing mechanism and functional
683 implications: JM. Wrote paper: JM, SH, DE. All authors reviewed paper. SGP helped with
684 illustrations.

685
686 **References**

- 687 Arnold CD, Gerlach D, Spies D, Matts J a, Sytnikova Y a, Pagani M, Lau NC, Stark A. 2014.
688 Quantitative genome-wide enhancer activity maps for five *Drosophila* species show
689 functional enhancer conservation and turnover during cis-regulatory evolution. *Nat Genet*
690 **46**: 685–92. <http://www.ncbi.nlm.nih.gov/pubmed/24908250> (Accessed July 9, 2014).
- 691 Arvey A, Agius P, Noble WS, Leslie C. 2012. Sequence and chromatin determinants of cell-type-
692 specific transcription factor binding. *Genome Res* **22**: 1723–34.
693 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3431489&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3431489&tool=pmcentrez&rendertype=abstract)
694 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3431489&tool=pmcentrez&rendertype=abstract) (Accessed July 18, 2014).
- 695 Ay F, Bunnik EM, Varoquaux N, Bol SM, Prudhomme J, Vert JP, Noble WS, Le Roch KG. 2014.
696 Three-dimensional modeling of the *P. falciparum* genome during the erythrocytic cycle
697 reveals a strong connection between genome architecture and gene expression. *Genome*
698 *Res* **24**: 974–988.
- 699 Babaei S, Akhtar W, de Jong J, Reinders M, de Ridder J. 2015. 3D hotspots of recurrent
700 retroviral insertions reveal long-range interactions with cancer genes. *Nat Commun* **6**:
701 6381. <http://www.nature.com/doi/10.1038/ncomms7381>.
- 702 Barozzi I, Simonatto M, Bonifacio S, Yang L, Rohs R, Ghisletti S, Natoli G. 2014. Coregulation of
703 Transcription Factor Binding and Nucleosome Occupancy through DNA Features of
704 Mammalian Enhancers. *Mol Cell* **54**: 844–857.
705 <http://dx.doi.org/10.1016/j.molcel.2014.04.006>.
- 706 Biggin MD. 2011. Animal transcription networks as highly connected, quantitative continua.
707 *Dev Cell* **21**: 611–26. <http://www.ncbi.nlm.nih.gov/pubmed/22014521> (Accessed March 3,
708 2013).
- 709 Brackley C a., Cates ME, Marenduzzo D. 2012. Facilitated diffusion on mobile DNA:
710 Configurational traps and sequence heterogeneity. *Phys Rev Lett* **109**: 1–5.
- 711 Calo E, Wysocka J. 2013. Modification of Enhancer Chromatin: What, How, and Why? *Mol Cell*
712 **49**: 825–837. <http://dx.doi.org/10.1016/j.molcel.2013.01.038>.
- 713 Chakalova L, Fraser P. 2010. Organization of transcription. *Cold Spring Harb Perspect Biol* **2**:
714 a000729.
715 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2926752&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2926752&tool=pmcentrez&rendertype=abstract)
716 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2926752&tool=pmcentrez&rendertype=abstract).
- 717 Crocker J, Abe N, Rinaldi L, McGregor AP, Frankel N, Wang S, Alsawadi A, Valenti P, Plaza S,
718 Payre F, et al. 2015. Low Affinity Binding Site Clusters Confer Hox Specificity and
719 Regulatory Robustness. *Cell* 191–203.

- 720 Deng W, Lee J, Wang H, Miller J, Reik A, Gregory PD, Dean A, Blobel G a. 2012. Controlling
721 long-range genomic interactions at a native locus by targeted tethering of a looping factor.
722 *Cell* **149**: 1233–44.
723 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3372860&tool=pmcentrez&](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3372860&tool=pmcentrez&endertype=abstract)
724 [endertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3372860&tool=pmcentrez&endertype=abstract) (Accessed August 14, 2014).
- 725 Djebali S, Davis C a., Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,
726 Schlesinger F, et al. 2012. Landscape of transcription in human cells. *Nature* **489**: 101–
727 108.
- 728 Downen JM, Fan ZP, Hnisz D, Ren G, Abraham BJ, Zhang LN, Weintraub AS, Schuijers J, Lee TI,
729 Zhao K, et al. 2014. Control of Cell Identity Genes Occurs in Insulated Neighborhoods in
730 Mammalian Chromosomes. *Cell* **159**: 374–387.
731 <http://dx.doi.org/10.1016/j.cell.2014.09.030>.
- 732 Edelman LB, Fraser P. 2012. Transcription factories: genetic programming in three dimensions.
733 *Curr Opin Genet Dev* **22**: 110–4. <http://www.ncbi.nlm.nih.gov/pubmed/22365496>
734 (Accessed August 7, 2013).
- 735 Essien K, Vigneau S, Apreleva S, Singh LN, Bartolomei MS, Hannehalli S. 2009. CTCF binding
736 site classes exhibit distinct evolutionary, genomic, epigenomic and transcriptomic features.
737 *Genome Biol* **10**: R131.
738 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3091324&tool=pmcentrez&](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3091324&tool=pmcentrez&endertype=abstract)
739 [endertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3091324&tool=pmcentrez&endertype=abstract) (Accessed March 14, 2013).
- 740 Ezer D, Zabet NR, Adryan B. 2014a. Homotypic clusters of transcription factor binding sites: A
741 model system for understanding the physical mechanics of gene expression. *Comput Struct*
742 *Biotechnol J* **10**: 63–9.
743 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4204428&tool=pmcentrez&](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4204428&tool=pmcentrez&endertype=abstract)
744 [endertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4204428&tool=pmcentrez&endertype=abstract) (Accessed November 3, 2014).
- 745 Ezer D, Zabet NR, Adryan B. 2014b. Physical constraints determine the logic of bacterial
746 promoter architectures. *Nucleic Acids Res* **42**: 4196–4207.
- 747 Filippova D, Patro R, Duggal G, Kingsford C. 2013. Multiscale identification of topological
748 domains in chromatin. *Lect Notes Comput Sci (including Subser Lect Notes Artif Intell*
749 *Lect Notes Bioinformatics)* **8126 LNBI**: 300–312. Algorithms for Molecular Biology.
- 750 Fisher WW, Li JJ, Hammonds AS, Brown JB, Pfeiffer BD, Weiszmann R, MacArthur S, Thomas
751 S, Stamatoyannopoulos J a, Eisen MB, et al. 2012. DNA regions bound at low occupancy by
752 transcription factors do not drive patterned reporter gene expression in *Drosophila*. *Proc*
753 *Natl Acad Sci* **109**: 21330–5.
754 <http://www.pnas.org/content/early/2012/12/05/1209589110.abstract> (Accessed February
755 28, 2013).
- 756 Foat BC, Morozov A V, Bussemaker HJ. 2006. Statistical mechanical modeling of genome-wide
757 transcription factor occupancy data by MatrixREDUCE. *Bioinformatics* **22**: e141–9.
758 <http://www.ncbi.nlm.nih.gov/pubmed/16873464> (Accessed February 28, 2013).
- 759 Fraser P. 2006. Transcriptional control thrown for a loop. *Curr Opin Genet Dev* **16**: 490–5.
760 <http://www.ncbi.nlm.nih.gov/pubmed/16904310> (Accessed November 3, 2013).
- 761 Fullwood MJ, Liu MH, Pan YF, Liu J, Xu H, Mohamed Y Bin, Orlov YL, Velkov S, Ho A, Mei PH,
762 et al. 2009. An oestrogen-receptor-alpha-bound human chromatin interactome. *Nature*

- 763 **462**: 58–64.
764 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2774924&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2774924&tool=pmcentrez&rendertype=abstract)
765 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2774924&tool=pmcentrez&rendertype=abstract) (Accessed February 28, 2013).
- 766 Ghavi-Helm Y, Klein F a., Pakozdi T, Ciglar L, Noordermeer D, Huber W, Furlong EEM. 2014.
767 Enhancer loops appear stable during development and are associated with paused
768 polymerase. *Nature* **512**: 96–100. <http://www.nature.com/doi/10.1038/nature13417>
769 (Accessed July 9, 2014).
- 770 Gotea V, Visel A, Westlund JM, Nobrega M a, Pennacchio L a, Ovcharenko I. 2010. Homotypic
771 clusters of transcription factor binding sites are a key component of human promoters and
772 enhancers. *Genome Res* **20**: 565–77.
773 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2860159&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2860159&tool=pmcentrez&rendertype=abstract)
774 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2860159&tool=pmcentrez&rendertype=abstract) (Accessed October 24, 2014).
- 775 Hannenhalli S. 2008. Eukaryotic transcription factor binding sites - Modeling and integrative
776 search methods. *Bioinformatics* **24**: 1325–1331.
- 777 He X, Duque TSPC, Sinha S. 2012a. Evolutionary origins of transcription factor binding site
778 clusters. *Mol Biol Evol* **29**: 1059–1070.
- 779 He X, Duque TSPC, Sinha S. 2012b. Evolutionary origins of transcription factor binding site
780 clusters. *Mol Biol Evol* **29**: 1059–70.
781 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3278477&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3278477&tool=pmcentrez&rendertype=abstract)
782 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3278477&tool=pmcentrez&rendertype=abstract) (Accessed June 8, 2013).
- 783 Heinz S, Romanoski CE, Benner C, Allison K a, Kaikkonen MU, Orozco LD, Glass CK. 2013.
784 Effect of natural genetic variation on enhancer selection and function. *Nature* **503**: 487–
785 492. <http://www.ncbi.nlm.nih.gov/pubmed/24121437> (Accessed November 7, 2013).
- 786 Heinz S, Romanoski CE, Benner C, Glass CK. 2015. The selection and function of cell type-
787 specific enhancers. *Nat Rev Mol Cell Biol* **16**: 144–154.
788 <http://dx.doi.org/10.1038/nrm3949>.
- 789 Hesselberth JR, Chen X, Zhang Z, Sabo PJ, Sandstrom R, Reynolds AP, Thurman RE, Neph S,
790 Kuehn MS, Noble WS, et al. 2009. Global mapping of protein-DNA interactions in vivo by
791 digital genomic footprinting. **6**: 283–289.
- 792 Ing-simmons E, Seitan VC, Faure AJ, Flicek P, Dekker J, Fisher AG, Lenhard B, Merckenschlager
793 M. 2014. Spatial enhancer clustering and regulation of enhancer-proximal genes by
794 cohesin.
- 795 Levy S, Hannenhalli S. 2002. Identification of transcription factor binding sites in the human
796 genome sequence. **514**: 510–514.
- 797 Li G, Ruan X, Auerbach RK, Sandhu KS, Zheng M, Wang P, Poh HM, Goh Y, Lim J, Zhang J, et
798 al. 2012. Extensive promoter-centered chromatin interactions provide a topological basis
799 for transcription regulation. *Cell* **148**: 84–98.
800 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3339270&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3339270&tool=pmcentrez&rendertype=abstract)
801 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3339270&tool=pmcentrez&rendertype=abstract) (Accessed February 27, 2013).
- 802 Lieberman-aiden E, Berkum NL Van, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I,
803 Lajoie BR, Sabo PJ, Dorschner MO, et al. 2009. of the Human Genome. **33292**: 289–294.

- 804 Malin J, Aniba MR, Hannenhalli S. 2013. Enhancer networks revealed by correlated DNase
805 hypersensitivity states of enhancers. *Nucleic Acids Res* **41**: 6828–38.
806 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3737527&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3737527&tool=pmcentrez&rendertype=abstract)
807 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3737527&tool=pmcentrez&rendertype=abstract) (Accessed October 26, 2013).
- 808 Markenscoff-Papadimitriou E, Allen WE, Colquitt BM, Goh T, Murphy KK, Monahan K, Mosley
809 CP, Ahituv N, Lomvardas S. 2014. Enhancer Interaction Networks as a Means for Singular
810 Olfactory Receptor Expression. *Cell* **159**: 543–557.
811 <http://linkinghub.elsevier.com/retrieve/pii/S0092867414011829> (Accessed October 24,
812 2014).
- 813 Martinez GJ, Rao A. 2012. Immunology. Cooperative transcription factor complexes in control.
814 *Science* **338**: 891–2.
815 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3621126&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3621126&tool=pmcentrez&rendertype=abstract)
816 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3621126&tool=pmcentrez&rendertype=abstract).
- 817 Matys V, Kel-Margoulis O V, Fricke E, Liebich I, Land S, Barre-Dirrie a, Reuter I, Chekmenev D,
818 Krull M, Hornischer K, et al. 2006. TRANSFAC and its module TRANSCompel:
819 transcriptional gene regulation in eukaryotes. *Nucleic Acids Res* **34**: D108–D110.
- 820 Montavon T, Duboule D. 2012. Landscapes and archipelagos: spatial organization of gene
821 regulation in vertebrates. *Trends Cell Biol* **22**: 347–54.
822 <http://www.ncbi.nlm.nih.gov/pubmed/22560708> (Accessed October 20, 2013).
- 823 Moses AM, Chiang DY, Pollard D a, Iyer VN, Eisen MB. 2004. MONKEY: identifying conserved
824 transcription-factor binding sites in multiple alignments using a binding site-specific
825 evolutionary model. *Genome Biol* **5**: R98.
- 826 Neph S, Vierstra J, Stergachis AB, Reynolds AP, Haugen E, Vernot B, Thurman RE, John S,
827 Sandstrom R, Johnson AK, et al. 2012. An expansive human regulatory lexicon encoded in
828 transcription factor footprints. *Nature* **489**: 83–90.
829 <http://www.ncbi.nlm.nih.gov/pubmed/22955618> (Accessed February 28, 2013).
- 830 Plank JL, Dean A. 2014. Enhancer Function: Mechanistic and Genome-Wide Insights Come
831 Together. *Mol Cell* **55**: 5–14. <http://www.ncbi.nlm.nih.gov/pubmed/24996062> (Accessed
832 July 10, 2014).
- 833 Pombo A, Dillon N. 2015. Three-dimensional genome architecture: players and mechanisms.
834 *Nat Rev Mol Cell Biol* **16**: 245–257. <http://www.nature.com/doifinder/10.1038/nrm3965>.
- 835 Ramos AI, Barolo S. 2013. Low-affinity transcription factor binding sites shape morphogen
836 responses and enhancer evolution. *Philos Trans R Soc Lond B Biol Sci* **368**: 20130018.
837 <http://rstb.royalsocietypublishing.org/content/368/1632/20130018.short?rss=1>.
- 838 Schoenfelder S, Clay I, Fraser P. 2010a. The transcriptional interactome: gene expression in 3D.
839 *Curr Opin Genet Dev* **20**: 127–33. <http://www.ncbi.nlm.nih.gov/pubmed/20211559>
840 (Accessed October 18, 2013).
- 841 Schoenfelder S, Sexton T, Chakalova L, Cope NF, Horton A, Andrews S, Kurukuti S, Mitchell J a,
842 Umlauf D, Dimitrova DS, et al. 2010b. Preferential associations between co-regulated genes
843 reveal a transcriptional interactome in erythroid cells. *Nat Genet* **42**: 53–61.
844 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3237402&tool=pmcentrez&r](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3237402&tool=pmcentrez&rendertype=abstract)
845 [endertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3237402&tool=pmcentrez&rendertype=abstract) (Accessed October 18, 2013).

- 846 Schwarzer W, Spitz F. 2014. The architecture of gene expression: integrating dispersed cis-
847 regulatory modules into coherent regulatory domains. *Curr Opin Genet Dev* **27**: 74–82.
848 <http://www.ncbi.nlm.nih.gov/pubmed/24907448> (Accessed September 22, 2014).
- 849 Sharon E, Kalma Y, Sharp A, Raveh-Sadka T, Levo M, Zeevi D, Keren L, Yakhini Z, Weinberger
850 A, Segal E. 2012. Inferring gene regulatory logic from high-throughput measurements of
851 thousands of systematically designed promoters. *Nat Biotechnol* **30**: 521–30.
852 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3374032&tool=pmcentrez&re
853 ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3374032&tool=pmcentrez&rendertype=abstract) (Accessed December 5, 2014).
- 854 Sheffield NC, Thurman RE, Song L, Safi A, Stamatoyannopoulos J a, Lenhard B, Crawford GE,
855 Furey TS. 2013. Patterns of regulatory activity across diverse human cell types predict
856 tissue identity, transcription factor binding, and long-range interactions. *Genome Res* **23**:
857 777–88.
858 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3638134&tool=pmcentrez&re
859 ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3638134&tool=pmcentrez&rendertype=abstract) (Accessed January 22, 2014).
- 860 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J,
861 Hillier LW, Richards S, et al. 2005. Evolutionarily conserved elements in vertebrate, insect,
862 worm, and yeast genomes. *Genome Res* **15**: 1034–1050.
- 863 Slattery M, Zhou T, Yang L, Dantas Machado AC, Gordân R, Rohs R. 2014. Absence of a simple
864 code: how transcription factors read the genome. *Trends Biochem Sci* **39**: 381–399.
- 865 Smith RP, Taher L, Patwardhan RP, Kim MJ, Inoue F, Shendure J, Ovcharenko I, Ahituv N.
866 2013. Massively parallel decoding of mammalian regulatory sequences supports a flexible
867 organizational model. *Nat Genet* **45**: 1021–8.
868 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3775494&tool=pmcentrez&re
869 ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3775494&tool=pmcentrez&rendertype=abstract) (Accessed July 14, 2014).
- 870 Spitz F, Furlong EEM. 2012. Transcription factors: from enhancer binding to developmental
871 control. *Nat Rev Genet* **13**: 613–26. <http://www.ncbi.nlm.nih.gov/pubmed/22868264>
872 (Accessed September 17, 2013).
- 873 Stewart AJ, Hannonhalli S, Plotkin JB. 2012. Why transcription factor binding sites are ten
874 nucleotides long. *Genetics* **192**: 973–85. <http://www.ncbi.nlm.nih.gov/pubmed/22887818>
875 (Accessed September 19, 2013).
- 876 Stewart AJ, Plotkin JB. 2013. The evolution of complex gene regulation by low-specificity
877 binding sites. *Proc Biol Sci* **280**: 20131313.
878 <http://www.ncbi.nlm.nih.gov/pubmed/23945682>.
- 879 Taher L, Smith RP, Kim MJ, Ahituv N, Ovcharenko I. 2013. Sequence signatures extracted from
880 proximal promoters can be used to predict distal enhancers. *Genome Biol* **14**: R117.
881 <http://www.ncbi.nlm.nih.gov/pubmed/24156763> (Accessed October 30, 2013).
- 882 Tanay A. 2006. Extensive low-affinity transcriptional interactions in the yeast genome. *Genome*
883 *Res* **16**: 962–72.
884 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1524868&tool=pmcentrez&re
885 ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1524868&tool=pmcentrez&rendertype=abstract) (Accessed October 2, 2014).
- 886 Teif VB, Rippe K. 2012. Calculating transcription factor binding maps for chromatin. *Brief*
887 *Bioinform* **13**: 187–201.

- 888 Vahedi G, Kanno Y, Furumoto Y, Jiang K, Parker SCJ, Erdos MR, Davis SR, Roychoudhuri R,
889 Restifo NP, Gadina M, et al. 2015. Super-enhancers delineate disease-associated regulatory
890 nodes in T cells. *Nature*. <http://www.nature.com/doi/10.1038/nature14154>.
- 891 Vernimmen D. 2014. Uncovering Enhancer Functions Using the α -Globin Locus. *PLoS Genet*
892 **10**: e1004668.
893 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4199490&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4199490&tool=pmcentrez&rendertype=abstract)
894 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4199490&tool=pmcentrez&rendertype=abstract) (Accessed October 24, 2014).
- 895 Wasson T, Hartemink AJ. 2009. An ensemble model of competitive multi-factor binding of the
896 genome. *Genome Res* **19**: 2101–2112.
- 897 White MA, Myers CA, Corbo JC, Cohen BA. 2013. Massively parallel in vivo enhancer assay
898 reveals that highly local features determine the cis -regulatory function of ChIP-seq peaks.
- 899 Whyte WA, Orlando DA, Hnisz D, Abraham BJ, Lin CY, Kagey MH, Rahl PB, Lee TI, Young RA.
900 2013. Master Transcription Factors and Mediator Establish Super-Enhancers at Key Cell
901 Identity Genes. *Cell* **153**: 307–319. <http://dx.doi.org/10.1016/j.cell.2013.03.035>.
- 902 Yáñez-Cuna JO, Dinh HQ, Kvon EZ, Shlyueva D, Stark A. 2012. Uncovering cis-regulatory
903 sequence requirements for context-specific transcription factor binding. *Genome Res* **22**:
904 2018–30.
905 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3460196&tool=pmcentrez&re](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3460196&tool=pmcentrez&rendertype=abstract)
906 [ndertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3460196&tool=pmcentrez&rendertype=abstract) (Accessed November 6, 2014).
- 907 Zinzen RP, Girardot C, Gagneur J, Braun M, Furlong EEM. 2009. Combinatorial binding
908 predicts spatio-temporal cis-regulatory activity. *Nature* **462**: 65–70.
909 <http://www.ncbi.nlm.nih.gov/pubmed/19890324> (Accessed March 1, 2013).

910

911

912

913

914

915

916

917

918

919

920

921

922

923 **Figures and Legends**

924

925 **Figure 1. Expectation and testing of differential occupancy**

926 **(A)** The combination of spatial proximity and genomic homotypic clusters of TFBS produce
927 high homotypic TF BS concentration. As illustrated, low-RE (degenerate) motif BS have a higher
928 expected frequency in the genome than high-RE motif BS, including more frequent HCTs. In a
929 spatially proximal chromatin context, effective homotypic BS concentrations are particularly
930 elevated for low-RE motif BS. This effect is further accentuated in archipelagos of enhancers,
931 which have been shown to be enriched for HCTs for shared TFs. High effective homotypic BS
932 concentration is likely a pre-requisite for the crowdsourcing effect. Large ovals denote
933 archipelagos of functionally related enhancers and target genes. Darkness of background color
934 approximates the maximum expected homotypic BS concentration. Not drawn to scale. Green:
935 DNA. Black: BS. BS=binding site; RE=relative entropy; HCT = homotypic (genomic) cluster of
936 TFBS.

937 **(B)** Calculating differential TF occupancy boost based on curated digital DNase footprint data.
938 Shown is the procedure for calculating occupancy boost for each (AP, TF) pair. For each
939 enhancer in an AP, and each TF with one or more putative BS in the enhancer, a non-AP
940 enhancer is chosen (with replacement) after controlling for mean enhancer-wide chromatin
941 accessibility (DHS) in the AP's most active tissue, and for the number of putative BS. For each
942 TF-AP pair, then, occupancy boost is calculated as the percent difference in the number of
943 putatively bound BS, where binding is determined in a binary manner: **1**, if a curated footprint
944 tightly overlaps a given motif instance, **0**, otherwise. If multiple TF motifs tightly overlap a given
945 footprint, conservatively, all are classified as bound. Putative BS are indicated by a '1', or '2',
946 respectively, for example TFs SOX and XBP1. A circle around a BS signifies it is imputed as
947 bound by its cognate TF. Note that the toy calculation of occupancy boost does not correspond
948 to the data displayed. AP = archipelago; TF = transcription factor, BS = binding site. DNase
949 digital footprint scans from Neph et al 2012.

950

951

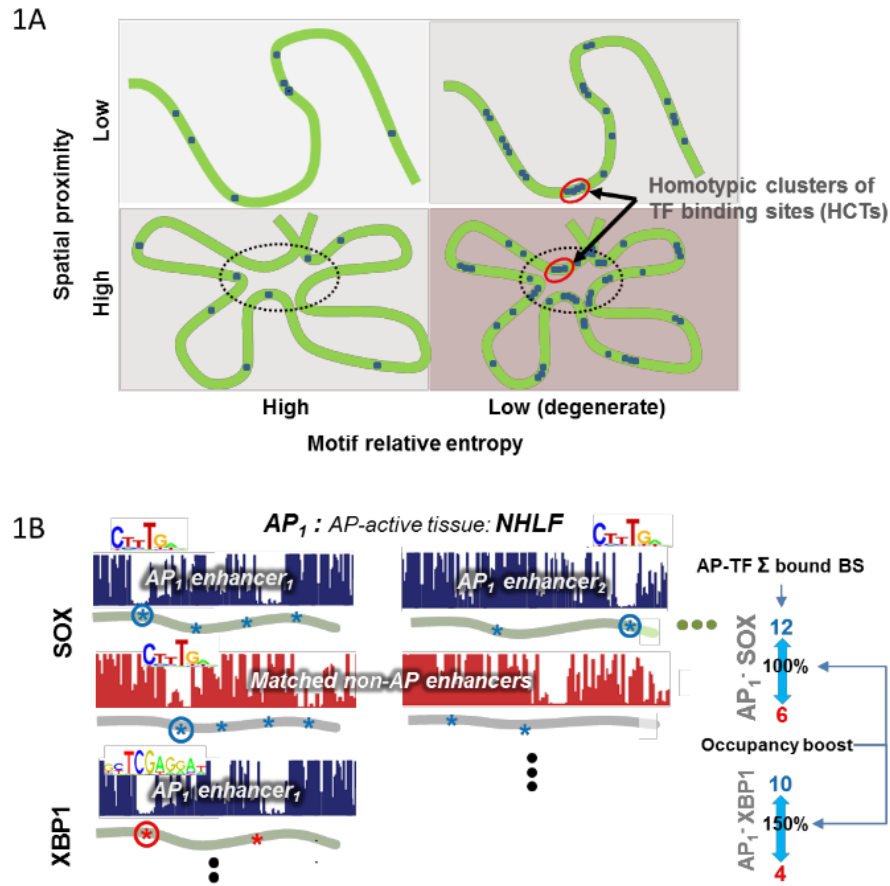
952

953

954

955

956



957

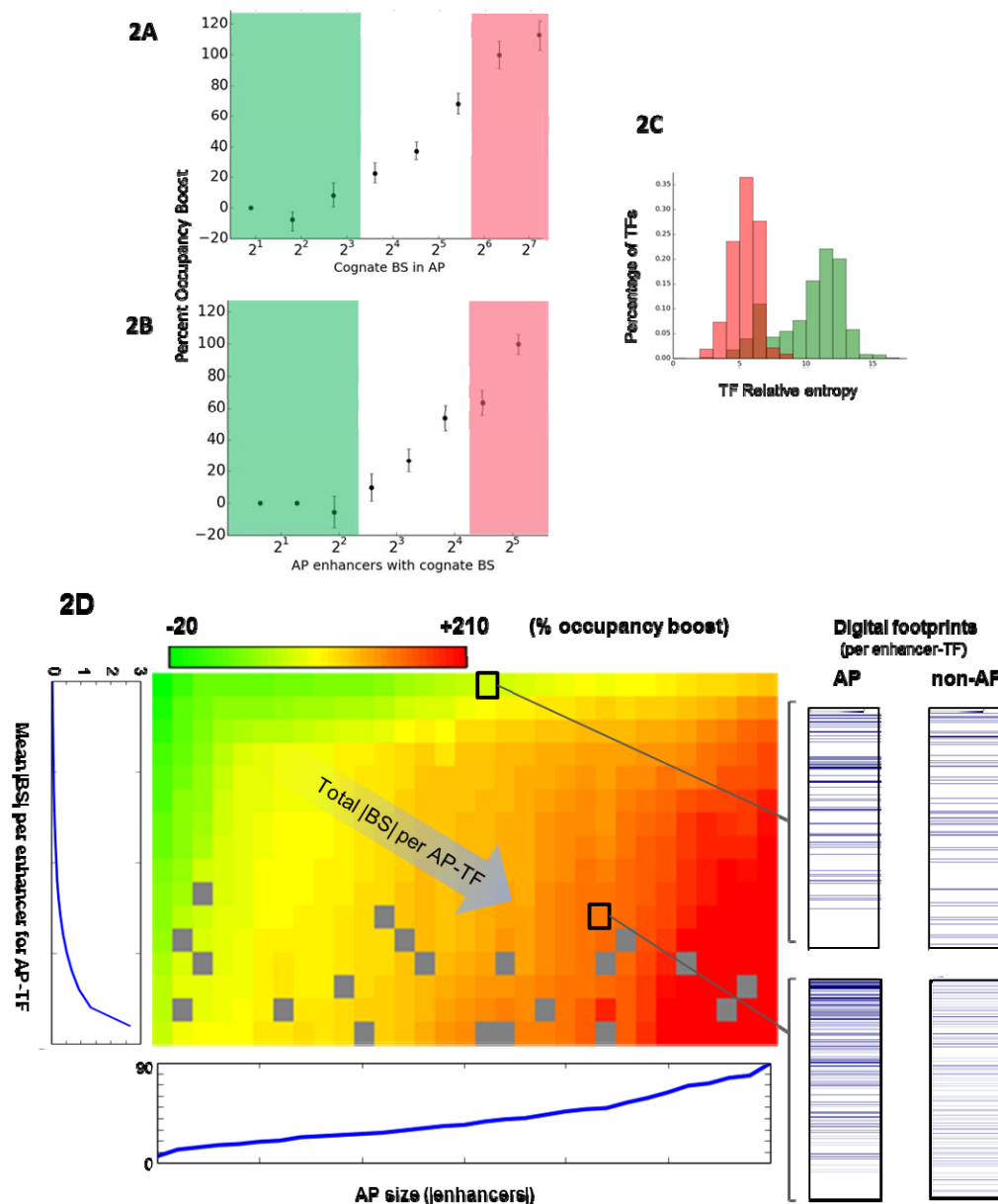
958

959 **Figure 2. Differential AP occupancy ‘boost’ scales with TF coverage in the AP**

960 TF-AP combinations were sorted on the basis of coverage and mean occupancy boost was
 961 determined for each group of TF-APs, where occupancy boost refers to differential occupancy in
 962 AP and non-AP enhancers matched 1-to-1 for the TF’s motif signature (the number and type of
 963 motifs) in a given enhancer, as well as for mean DHS across the AP. Occupancy was calculated
 964 based on the overlap of curated DNase digital footprints (Neph et al 2012) with high-confidence
 965 TRANSFAC motif instances. TF-APs and their non-AP counterparts were included in this
 966 analysis only if they both had non-zero occupancy (See also Figures S1, S2, S3). Coverage was
 967 calculated alternatively as:

968 **(A)** the number of cognate BS for a given TF in a given AP or **(B)** the number of enhancers in
 969 the AP with one or more motif instances cognate to the TF. **(C)** There is a strong inverse
 970 relationship between TF-AP coverage and TF relative entropy, as expected. TF-APs with the top
 971 (salmon) and bottom (green) 20% coverage were mapped to the TFs’ RE values. Colors in 2A
 972 and 2B illustrative and do not coincide with these coverage ranges. AP: archipelago, RE =

973 relative entropy. **(D)** Mean occupancy boost versus coverage that has been decomposed along
 974 two axes. Each TF-AP pair was binned based on the number of enhancers in an AP (column) and
 975 the mean number of BS per AP enhancer (row). Plots to the left of and below the heatplot show
 976 mean boost for each row and column, respectively. Red (green) heatmap cells indicate high
 977 (low) percentage occupancy boost after Lowess smoothing. Grey cells indicate no data. In the
 978 right panel, for all TF-AP pairs in the selected heatmap cell, significant digital DNase
 979 hypersensitivity footprints in member AP and matched non-AP enhancers are shown, where the
 980 numbers of BS for AP and non-AP enhancer-TF pairs are identical; a blue line indicates a
 981 significant footprint overlapping a putative BS. Enhancers are sorted from bottom to top in
 982 order of increasing chromatin accessibility.



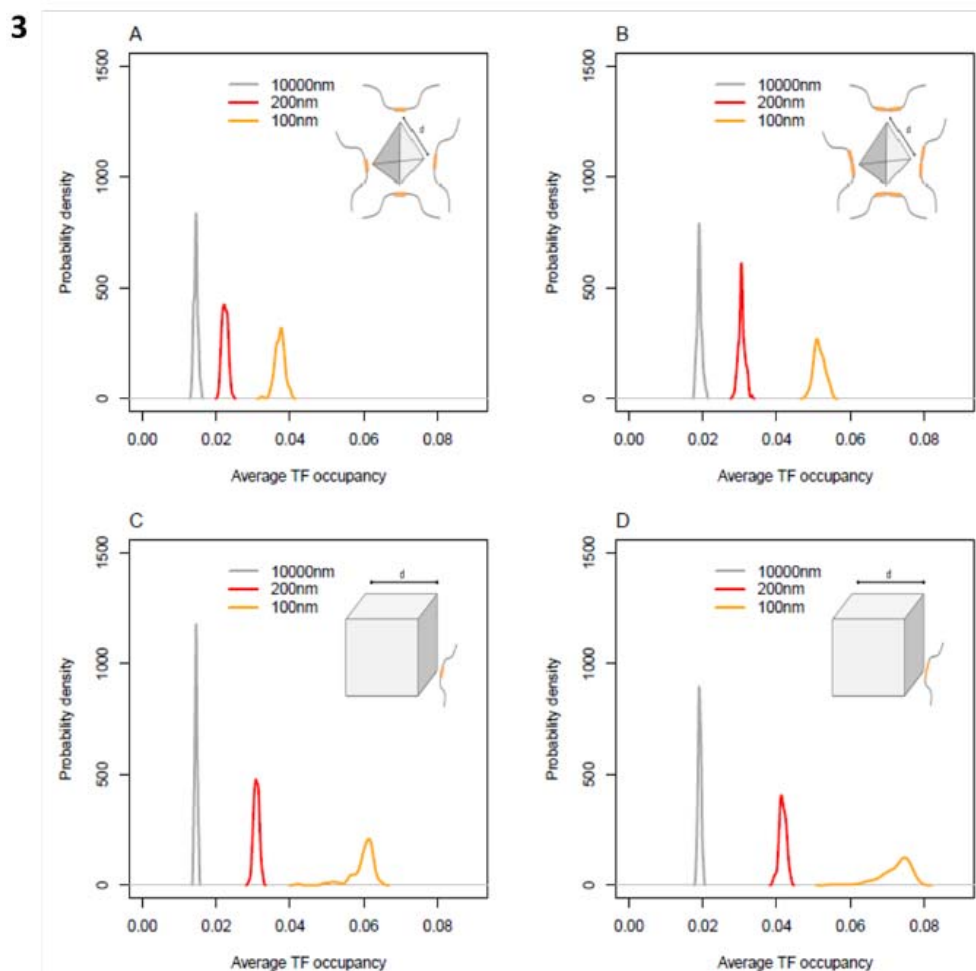
983

984

985 **Figure 3. Biophysically modeling crowdsourcing effect**

986 TF diffusion was simulated for four geometric arrangements of binding sites, and the probability
987 density functions of TF occupancy are shown. The TF occupancy is defined as the average
988 probability that each site is bound. The four simulated scenarios are: a tetrahedron with **(A)** one
989 binding site or **(B)** a pair of binding sites in each corner, which contain 4 or 8 binding sites,
990 respectively; a cube with **(C)** a single binding site or **(D)** a pair of binding sites in each corner,
991 which contain 8 and 16 binding sites respectively. For an additional figure and details on the
992 simulation, see Supplemental Methods.

993



994

995

996

997

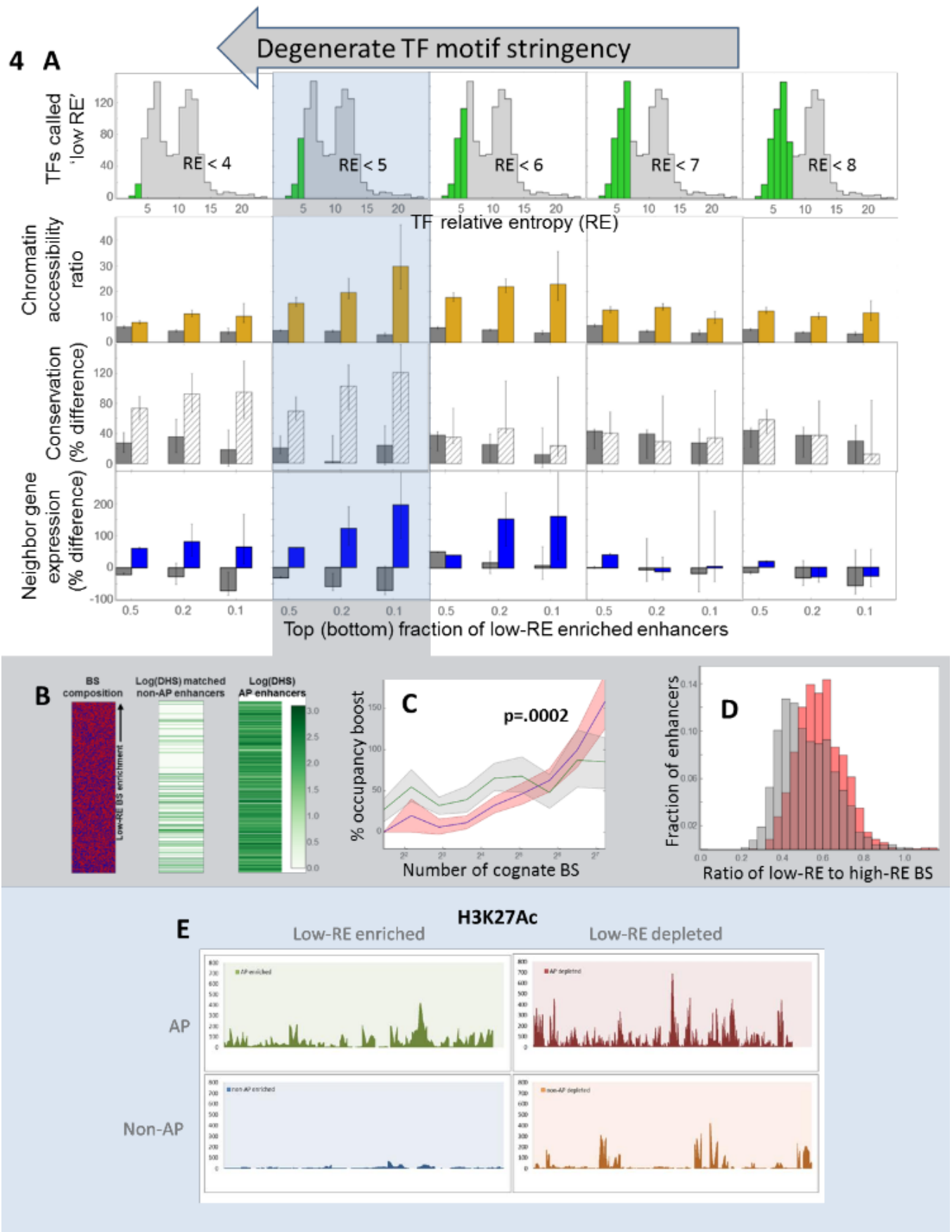
998

999 **Figure 4. Enhancers enriched for degenerate BS are more functional than**
1000 **expected**

1001 **(A)** Enhancers enriched and depleted for low-RE BS were compared in terms of DNase
1002 hypersensitivity (row 2), evolutionary constraint (row 3), and neighbor gene expression (row 4).
1003 Readouts on y-axes indicate values normalized against carefully matched non-AP enhancers for
1004 the given RE cutoff (column). Within each plot, the 10%, 20%, and 50% (x-axis) most enriched
1005 enhancers are indicated in non-grey, while the most depleted enhancers are shown in grey. The
1006 histograms in the top row indicate the fraction (green) of all TFs deemed low-RE for the purpose
1007 of calculating each enhancer's low-RE BS enrichment.

1008 Figures **(B-E)** are for RE cutoff = 5. **(B)** Chromatin accessibility of AP and matched non-AP
1009 enhancers, sorted by low-RE motif enrichment. Aligned heatplots display one enhancer per row
1010 (most enriched at top). (Left) heatplot in which red signifies a low-RE motif and blue a high-RE
1011 motif. Low-RE motif enrichment based on Fisher exact test against a background that included
1012 all non-AP enhancers. For visualization purposes, enhancer lengths and BS lengths
1013 standardized. (left, right). Log of non-AP and AP enhancer DHS, respectively. (See also Figure
1014 S4). **(C)** Percentage occupancy boost is shown as a function of coverage for AP enhancers with
1015 the highest 20% enrichment (blue line) and the highest 20% depletion (green line) for low-RE
1016 BS, along with their 95% confidence intervals. Coverage for a given TF-AP pair was calculated as
1017 the number of cognate BS in in the AP among enriched (depleted) enhancers only. A p-value is
1018 given for a Wilcoxon test comparing boosts among TF-APs with top 20% coverage. **(D)** Ratio of
1019 low-RE to high-RE motifs in AP enhancers vs. non-AP enhancers. AP and non-AP enhancers
1020 were matched one-to-one for DHS in each AP enhancer's most active tissue. Putative BS were
1021 identified based on 95 percentile motif match threshold. The x-axis shows the ratio of low-RE to
1022 high-RE motif sites in each enhancer. Y-axis shows percentage of enhancers analyzed. Red: AP
1023 enhancers, Gray: non-AP enhancers. All confidence intervals were calculated using an
1024 appropriate bootstrap method. AP: archipelago, RE: relative entropy, BS: binding site, TF:
1025 transcription factor. **(E)** Acetylation levels in enriched vs. depleted enhancers. Juxtaposed
1026 views of H3K27Ac ChIP-Seq in HUVEC are shown for 40 (100) AP enhancers in the top row that
1027 are in an AP that is active and in the top 10% for enrichment (depletion). Shown in the bottom
1028 row are views for matched non-AP enhancers.

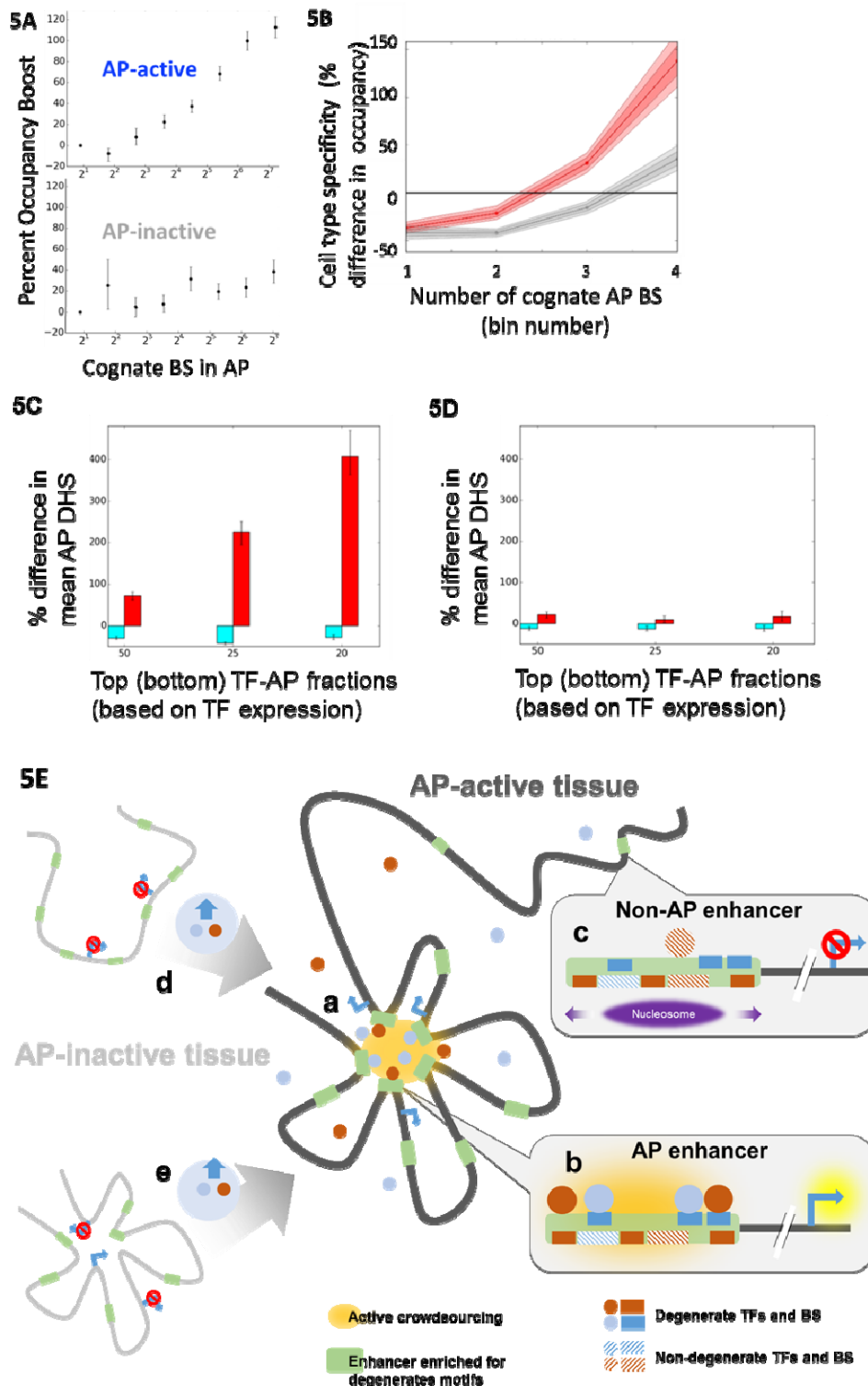
1029
1030
1031
1032



1034 **Figure 5. Mean AP accessibility scales with context-specific availability of TFs with**
1035 **degenerate motifs and not the TFs with specific motifs.**

1036 **(A)** Occupancy boost in cell types with reduced AP activity. Occupancy was computed as a
1037 function of coverage in 'inactive' cell types – those in which fewer than 40% of the AP's
1038 enhancers were DNase hypersensitive (bottom). For comparison, the plot for active cell types
1039 from 2A is reproduced (top) **(B)** Tissue specificity of occupancy for TF-AP and matched non-TF-
1040 AP pairs as a function of TF-AP coverage. Dynamic range (y-axis) for occupancy was calculated
1041 for each TF-AP pair as the percentage difference between mean occupancy in the AP's most
1042 active and inactive cell types. Identical BS in active and inactive cell types were tested. Serving
1043 as a control, dynamic range was also computed for non-TF-APs that were matched to TF-APs. A
1044 TF-AP was required to have non-zero occupancy in both inactive and active cell types. Only TF-
1045 APs shared in AP and non-AP contexts were used for analysis. Shown are results for TF-APs and
1046 for non-TF-APs each sorted into 4 bins with exponentially increasing coverage cutoffs. Red: AP,
1047 Gray: non-AP. 90% and 99% confidence intervals are shown with variable hue. AP-wide DHS as
1048 a function of context-specific TF availability. (See also Supplemental Fig. S5) **(C)** For each TF-
1049 AP, tissue specific DHS was compared across each of 15 tissues for which there was RNA-Seq
1050 data available. (TF, AP, tissue) triplets were segregated into lowest-20%-coverage (cyan) and
1051 highest-20%-coverage (red) classes based on TF-AP, and then further subdivided into low and
1052 high expression based on tissue-specific TF expression. Bar height indicates the percentage
1053 increase in DHS level associated with an increase in TF expression from bottom <x> to top <x>
1054 percentage levels, where <x> is read off the x-axis. **(D)** same as (C) except matched non-AP
1055 triplets were used. **(E)** Model of crowdsourcing effect. **(a)** The yellow highlighted region
1056 represents a regulatory archipelago (AP) consisting of genes and distal enhancers. Within an AP,
1057 spatially proximal binding sites (BS) for a common TF 'crowdsource' an increase in their own
1058 occupancy. Facilitated by increased TF diffusion among large numbers of spatially proximal BS,
1059 a spatial homotypic BS cluster favorably alters TF protein concentration in its
1060 microenvironment. Predictably, TFs with degenerate motifs, and hence pervasive BS, exhibit the
1061 highest occupancy boosts. **(b)** In turn, AP enhancers enriched in degenerate motifs experience
1062 switch-like multi-fold boosts in accessibility and target gene expression. Overall, a context-
1063 specific increase in availability of TFs with degenerate motifs – but not high-specificity motifs –
1064 drives a multi-fold boost in chromatin accessibility, thereby underscoring crowdsourcing's likely
1065 role in AP activation. **(c)** In contrast, a non-AP enhancer does not experience an occupancy
1066 boost and activation. The crowdsourcing mechanism integrates well with the two prevailing
1067 models of context-specific gene module activation: in a targeted tissue, higher expression of TFs

1068 with a degenerate motif may **(d)** induce chromatin loop formation; or alternatively **(e)** facilitate
 1069 release of paused polymerase in pre-formed enhancer-promoter loops. In both cases,
 1070 crowdsourcing ensures a high degree of context-specificity, mitigating spurious occupancy
 1071 outside of or AP-active tissue or AP enhancers enriched for degenerate motifs.



1072