

# 1 **Polymer physics and machine learning reveal a combinatorial code linking** 2 **chromatin 3D architecture to 1D epigenetics**

3 Andrea Esposito<sup>1+</sup>, Simona Bianco<sup>1+</sup>, Andrea M. Chiariello<sup>1</sup>, Alex Abraham<sup>1</sup>, Luca Fiorillo<sup>1</sup>, Mattia  
4 Conte<sup>1</sup>, Raffaele Campanile<sup>1</sup>, and Mario Nicodemi<sup>1,2,3 \*</sup>

5 <sup>1</sup>Dipartimento di Fisica, Università di Napoli *Federico II*, and INFN Napoli, Complesso Universitario di Monte  
6 Sant'Angelo, 80126 Naples, Italy.

7 <sup>2</sup>Berlin Institute for Medical Systems Biology, Max-Delbrück Centre (MDC) for Molecular Medicine, Berlin, Germany.

8 <sup>3</sup>Berlin Institute of Health (BIH), MDC-Berlin, Germany.

9 <sup>+</sup>These authors contributed equally

10 \*Contact author's email address: [mario.nicodemi@na.infn.it](mailto:mario.nicodemi@na.infn.it)

11

## 12 **ABSTRACT**

13 The mammalian genome has a complex 3D organization, serving vital functional purposes, yet it  
14 remains largely unknown how the multitude of specific DNA contacts, e.g., between transcribed  
15 and regulatory regions, is orchestrated by chromatin organizers, such as Transcription Factors.  
16 Here, we implement a method combining machine learning and polymer physics to infer from only  
17 Hi-C data the genomic 1D arrangement of the minimal set of binding sites sufficient to recapitulate,  
18 through only physics, 3D contact patterns genome-wide in human and mouse cells. The inferred  
19 binding sites are validated by their predictions on how chromatin refolds in a set of duplications at  
20 the *Sox9* locus against available independent cHi-C data, showing that their different phenotypes  
21 originate from distinct enhancer hijackings in their 3D structure. Albeit derived from only Hi-C, our  
22 binding sites fall in epigenetic classes that well match chromatin states from epigenetic  
23 segmentation studies, such as active, poised and repressed states. However, the inferred binding  
24 domains have an overlapping, combinatorial organization along chromosomes, missing in  
25 epigenetic segmentations, which is required to explain Hi-C contact specificity with high accuracy.  
26 In a reverse approach, the epigenetic profile of binding domains provides a code to derive from  
27 only epigenetic marks the DNA binding sites and, hence, the 3D architecture, as validated by  
28 successful predictions of Hi-C matrices in an independent set of chromosomes. Overall, our results  
29 shed light on how complex 3D architectural information is encrypted in 1D epigenetics via the  
30 related, combinatorial arrangement of specific binding sites along the genome.

## 31 INTRODUCTION

32 The genome of higher organisms has a complex spatial organization within the cell nucleus<sup>1-6</sup> as  
33 revealed by recent technologies<sup>7-13</sup>. Chromosomes are folded in a sequence of 0.5-1.0Mb wide  
34 domains, named TADs<sup>14,15</sup>, in sub-TADs and loops, and in larger structures such as A/B  
35 compartments<sup>7</sup> and meta-TADs<sup>16</sup>. Importantly, such an organization serves vital functional  
36 purposes, as for instance distal enhancers control their target genes by establishing physical  
37 contacts with them, disruptions being linked to human diseases<sup>17-19</sup>. However, how chromatin  
38 architecture is shaped and orchestrated remains mostly unknown.

39 To rationalize the complexity of Hi-C data, polymer models from statistical physics<sup>20-32</sup> and a variety  
40 of computational methods<sup>33-36</sup> have been developed. A class of models, such as the *Strings and*  
41 *Binders* (SBS) model<sup>21</sup>, has focused on the classical scenario where loops and contacts between  
42 distal DNA sites are established by diffusing molecules such as Transcription Factors (TFs), or some  
43 effective interaction potential, bridging cognate binding sites by thermodynamics mechanisms of  
44 phase separation<sup>21,25-32,37-39</sup>. Another interesting classical scenario has been considered by off-  
45 equilibrium polymer models where loops are formed by extrusion, e.g., by molecules that bind to  
46 DNA and extrude a loop<sup>22-24</sup>, based on prior knowledge of the involved molecular factors, such as  
47 CTCF binding sites<sup>29,40</sup>.

48 Here, we use a machine learning approach (PRISMR<sup>41</sup>) to infer from only Hi-C data the genomic  
49 location of the minimal set of binding sites best explaining contact patterns across chromosomes by  
50 only polymer physics via the molecular mechanisms envisaged by the SBS model. While PRISMR  
51 was previously applied to Mb wide genomic regions, we optimize its performance to extend the  
52 approach to the entire genome, improving the statistical power of our method of three orders of  
53 magnitude to understand how complex 3D architectural information is encrypted in 1D epigenetic  
54 signals. Without prior knowledge of binding factors, our approach can infer genome-wide the  
55 specific location of the distinct binding sites whereby DNA contacts are established as captured by  
56 Hi-C data, returning a picture of the key elements underlying chromosomes folding. We show that  
57 the SBS polymer model informed with the inferred binding sites recapitulates Hi-C data across  
58 chromosomes in human<sup>42</sup> and murine<sup>14</sup> cells with high accuracy, illustrating that its minimal  
59 ingredients are sufficient to make sense of a substantial fraction of contact patterns genome-wide.  
60 For sake of simplicity, we focus on the SBS model, but the method can be extended to  
61 accommodate additional mechanisms, such as loop extrusion<sup>22-24</sup>.

62 To test the inferred binding sites of the model and its envisaged folding mechanisms, we compared  
63 its predictions about the impact of mutations on chromosome conformations against independent  
64 experimental data. As a case study, we considered the *Sox9* locus, where cHi-C data are available  
65 for a set of duplications<sup>43</sup>. We implemented in the wild-type chromosome model those duplications  
66 and derived *de novo* the corresponding contact maps that are successfully compared to cHi-C data,  
67 with no fitting parameters available. Our analysis also shows that different genomic variants  
68 produce different neo-TADs around *Sox9* marked by specific enhancer hijackings, hence resulting in  
69 different phenotypes.

70 Importantly, our inference procedure does not exploit previous knowledge on binding factors or  
71 epigenetics marks. Hence, the inferred binding domains can be used to bring together  
72 independently derived information on architecture and epigenetics, e.g., by crossing their genomic  
73 position with ENCODE databases. We find that the different binding domains fall in similarity  
74 classes based on epigenetics, well matching functional chromatin states derived in linear epigenetic  
75 segmentation studies such as active, poised and repressed states<sup>10,44–47</sup>. However, we discover that  
76 they have an overlapping, combinatorial genomic distribution at the current resolution of Hi-C  
77 experiments, lacking in linear segmentation studies, which is shown to be required to explain Hi-C  
78 contacts with high accuracy genome-wide.

79  
80 Finally, we validated the discovered association between machine learned binding sites and  
81 epigenetic features by reversing the approach. In the considered cell types, we used the epigenetic  
82 profiles of the different binding domains of even chromosomes as a code to derive from only  
83 histone marks of odd chromosomes the location and type of their binding sites. Next, those binding  
84 sites were used to inform the polymer models of odd chromosomes, which predicted the  
85 corresponding Hi-C matrices with an accuracy comparable to those directly inferred from Hi-C data.

86  
87 Overall, our results provide insights on how the 1D combinatorial arrangement of a comparatively  
88 small number of binding site types, barcoded by distinctive epigenetic signatures, encodes the  
89 architectural information guiding chromatin organizing factors to establish, through physics, the  
90 multitude of specific 3D contacts across chromosomal scales.

91

92

## 93 RESULTS

### 94 Inferred binding domains explain Hi-C data genome-wide

95 To dissect the molecular mechanisms that contribute to chromatin folding, we used the PRISMR  
96 machine learning procedure<sup>41</sup> to infer the minimal SBS polymer model best explaining *in situ* Hi-C  
97 contact maps in human GM12878 B-lymphoblastoid cells at 5kb resolution across chromosomes<sup>42</sup>  
98 (**Fig. 1B**, **Fig. S1**, Materials and Methods). In the *String and Binders* (SBS) model<sup>21</sup>, a chromatin  
99 filament is modeled as a self-avoiding string of beads, including specific binding sites for diffusing  
100 molecules; such binders can bridge distal cognate sites along the sequence, producing loops and  
101 physical contacts (**Fig. 1C**). In particular, in the SBS model, contact domains of homologous sites are  
102 spontaneously established by their cognate binders via a thermodynamic mechanism known as  
103 micro-phase separation<sup>26,30,39</sup>. Specifically, PRISMR finds the minimal combination of the binding  
104 sites of the SBS model (**Fig. 1D**) that reproduces, within a given accuracy threshold, the  
105 experimental Hi-C contact matrix based only on polymer thermodynamics (Materials and Methods).  
106 The different groups of homologous binding sites are named the *binding domains* of the model.  
107 Importantly, PRISMR uses just Hi-C data as input, with no prior knowledge of binding factors.

108 To check whether the model can explain *in situ* Hi-C data genome-wide, we compared the PRISMR  
109 derived SBS contact matrices to the original data (**Fig. 1E**, **Fig. S1**). In particular, we computed their  
110 Pearson correlation coefficient,  $r$ , their distance corrected Pearson correlation coefficient,  $r'$ , and  
111 their HiCRep stratum adjusted correlation coefficient, SCC<sup>48</sup>. The last two measures were  
112 considered to account for genomic proximity effects (see Materials and Methods). Model and  
113 experimental data were found to be comparatively similar across chromosomes, as  $r$ ,  $r'$  and SCC  
114 range around  $r=0.94$ ,  $r'=0.74$  and  $SCC=0.86$ , respectively. Notably, from the SBS model the  
115 thermodynamics ensemble of chromosomal 3D conformations can also be derived; a snapshot, e.g.,  
116 of chromosome 20 is pictured in **Fig. 1F**.

117 Additionally, to prove the general validity of the method, we tested its performance on a mouse  
118 embryonic stem cells Hi-C dataset at 40kb resolution<sup>14</sup>, finding that the PRISMR inferred and  
119 experimental contact matrices have high correlation values, comparable to those reported above  
120 for the 5kb human data (**Fig. S2**).

121 The model binding domains, i.e., the sets of homologous binding sites along the chromosomes, are

122 the output of PRISMR. The algorithm returns 30 different binding site types per chromosome in  
123 GM12878 (**Fig. 1D**, see Materials and Methods). Interestingly, a single binding domain covers on  
124 average a genomic length comparable to the size of a TAD (3.1+/-1.9Mb), yet the range of their  
125 interactions,  $r_{Int}$ , extends more than one order of magnitude longer, up to tens of Mbs (**Fig. S3A**,  
126 Materials and Methods). The distribution of  $r_{Int}$ ,  $P(r_{Int})$ , is significantly different from a random  
127 control model obtained by bootstrapping the location of binding sites (Materials and Methods) and  
128 is asymptotically consistent with a power-law scaling,  $P \sim 1/r_{Int}$ , typical of hierarchical structures  
129 made of domains within domains, as in Cantor sets (**Fig. S3A**, Materials and Methods). The broad  
130 range of values of  $r_{Int}$  shows that chromatin interactions extend above the size of single TADs, with  
131 higher-order 3D structures formed at scales below and above the A/B compartment level<sup>16</sup>. The  
132 derived 3D structure of chromosomes (**Fig. 1F**) shows indeed that, rather than being a linear chain  
133 of TADs, they tend to fold on themselves in complex structures, such as meta-TADs<sup>16</sup> (see also <sup>31</sup>).

134 Taken together, the high correlations found between the SBS model and Hi-C contact data show  
135 that the 1D binding domains inferred genome-wide by PRISMR contain key information sufficient to  
136 recapitulate 3D contact patterns genome-wide in human and mouse cells. That sheds light on the  
137 molecular mechanisms shaping chromosome architecture, supporting the view that the  
138 combinatorial action of a comparatively small number of TFs, mediating the interactions between  
139 cognate binding sites, can spontaneously fold chromatin in its 3D structure through just the laws of  
140 physics.

#### 141 **Validation of the inferred binding domains against duplications in the *Sox9* locus**

142 To validate the binding domains inferred by our approach, i.e., the determinants of folding and  
143 their envisaged mode of action, we compared our model predictions against previous  
144 independently produced cHi-C data in E12.5 limb bud cells from mice carrying homozygous  
145 structural variants in the *Sox9* locus<sup>43</sup>. We considered three mutations (**Fig. 2**, **Fig. S4**): a 0.4Mb  
146 duplication (*DupS*) in the non-coding DNA region within the *Sox9* gene TAD (intra-TAD duplication),  
147 associated with female to male sex reversal in humans; a 1.6Mb duplication (*DupL*) that  
148 encompasses the neighboring TAD boundary, with no phenotypic effect; and a slightly longer  
149 1.8Mb duplication (*DupC*), associated to limb malformation, which also includes *Kcnj2*, the next  
150 flanking gene. Specifically, we implemented those mutations in the SBS polymer model of the wild-  
151 type region in mESC inferred by PRISMR<sup>30</sup> and derived the novel contact matrices from only

152 polymer physics, with no fitting parameters whatsoever. The Pearson and distance-corrected  
153 Pearson correlation coefficients between the model predicted and cHi-C contact maps across the  
154 three mutations are as high as  $r=0.88$  and  $r'=0.48$ , reflecting their good degree of similarity (**Fig. 2**,  
155 **Fig. S4**).

156 Those results provide a stringent validation to our approach and demonstrate that predictions on  
157 the 3D structure of chromatin based on the inferred binding domains can be accurate to the point  
158 to anticipate ectopic contacts produced by disease-associated mutations.

### 159 **Enhancer hijackings in neo-TADs at the mutated *Sox9* locus link to phenotype**

160 To understand the origin of the ectopic contacts in the mutated systems, within our model we  
161 dissected the interactions of the duplicated from the original DNA sequences and the  
162 corresponding 3D structures, pieces of information inaccessible by Hi-C data (**Fig. 2**, **Fig. S4**, **S5**).

163 *DupS* is fully included within the TAD encompassing *Sox9* (**Fig. S4A**). Within our model, a TAD and its  
164 corresponding enrichment of interactions derive from the presence of a prevailing type of binding  
165 sites in that DNA region (see, e.g., TAD *A*, *B*, *C* in **Fig. 4D-F**). Hence, the duplicated and the original  
166 sequence in *DupS* (region *B2'* and *B2* in **Fig. S5A**) share many homologous binding sites, which  
167 produce the contacts between such regions visible in the interaction matrix mapped along the full,  
168 duplicated genome (**Fig. S4A**, **S5A**). When those contacts are mapped back onto the wild-type  
169 sequence, an excess of interactions appears localized around the mutated region within the  
170 corresponding TAD, but no major changes to the overall contact pattern, as experimentally found in  
171 cHi-C data<sup>43</sup>. The model derived 3D structure of the mutated locus shows, indeed, that the  
172 duplicated region remains well embedded into the rest of the locus (**Fig. S5C**).

173 Conversely, in *DupL* the duplicated region spans two TADs (**Fig. 2**). In our model, those TADs are  
174 produced by different prevailing types of binding sites. Accordingly, the portion of the duplication  
175 within the *Sox9* TAD (region *B1'* in **Fig. 2C**) has enriched contacts with itself and its corresponding  
176 original sequence (*B1*), but less with the portion of the duplication within the flanking TAD and its  
177 original sequence (resp. region *A2'* and *A2*). Since *B1'* is enriched of self-contacts but has  
178 comparatively less interaction with its neighboring genomic regions *A2'* and *A2*, it forms a separate  
179 chromatin domain (termed a 'neo-TAD'<sup>43</sup>) remaining partially isolated from the rest, as seen in a  
180 snapshot of the 3D structure of the locus (**Fig. 2E**, red region). As the isolated neo-TAD does not

181 include main genes, *DupL* has no phenotype<sup>43</sup>.

182 Finally, *DupC* produces a neo-TAD as much as *DupL*, however, it now includes a copy of the next  
183 flanking gene, *Kcnj2* (**Fig. S4**). As seen in the contact matrix of the full genome, within the neo-TAD  
184 the duplicated *Kcnj2* establishes ectopic contacts with the duplicated part of the regulatory region  
185 of *Sox9*. So, *Kcnj2* is mis-expressed, leading to the associated phenotype<sup>43</sup>.

186 In brief, our findings clarify how mutations impact chromatin architecture and the mode of action  
187 of the 3D structure in regulating gene activity. In particular, they explain how the considered  
188 structural genomic variations at the *Sox9* locus differently alter 3D conformation and gene  
189 regulation by specific enhancer hijackings, resulting in distinct phenotypes.

### 190 **Epigenetic profile of the binding domains**

191 To shed light on the nature of the model inferred binding sites (**Fig. 1D**), we correlated their  
192 genomic locations with histone mark tracks available in the ENCODE database<sup>49</sup> for the GM12878  
193 cell line. In particular, we used the binding domains derived from even-numbered chromosomes to  
194 compute such correlations, in order to use the derived barcode linking binding site types and  
195 epigenetics to later independently predict the architecture of odd-numbered chromosomes. In our  
196 analysis, we retained only statistically significant correlation values, i.e., those above a random  
197 control model with sites having bootstrapped genomic positions (Materials and Methods). As the  
198 different binding domains tend to fall in groups with similar epigenetic profiles, we clustered them  
199 to identify genome-wide significantly distinct epigenetic classes (Materials and Methods). The  
200 Akaike Information Criterion<sup>50</sup> (AIC) returns a set of 9 statistically different groups (**Fig. 3A**), a result  
201 also supported by basic hierarchical clustering (**Fig. S6B**).

202 Three classes of binding domains strongly correlate with active chromatin marks (**Fig. 3A**), but they  
203 are distinct from an epigenetic point of view. While class 1 is enriched for only active marks, classes  
204 2 and 3 are both enriched also in H3K9me3. Also, class 3 shows a stronger correlation with  
205 H3K4me1 as compared with class 2, a histone mark associated especially with active enhancer  
206 regions<sup>10,44–47</sup>. Interestingly, the genomic positions of the sites of the first three classes (**Fig. 3B**) are  
207 partially correlated (**Fig. S7C**, Materials and Methods). Their histone signatures are also consistent  
208 with DNA accessibility, early replication time and RNAseq transcription data (**Fig. S6C**). That  
209 supports the view that the binding sites in class 1, 2 and 3 are responsible, genome-wide, especially

210 for specific contacts between transcribed and regulatory regions, mediated by factors such as  
211 active Pol-II, as experimentally demonstrated at a number of loci<sup>40</sup>. Class 4 has the typical signature  
212 of bivalent chromatin, with H3K27me3 combined with active marks. Its binding sites could be  
213 responsible for interactions between regions including, for instance, poised genes and their  
214 regulators, as seen in FISH co-localization experiments<sup>40</sup>. Classes 5 and 6 are significantly correlated  
215 with H3K27me3 and could be responsible for the experimentally observed self-interacting domains  
216 of PRC repressed chromatin<sup>51</sup>. Interestingly, the first six classes are the only ones to correlate with  
217 CTCF binding sites (**Fig. S6C**). That confirms the significance of CTCF in regulating chromatin  
218 architecture and gene activity (see, e.g. <sup>52</sup>), highlighting that its role can be modulated by different  
219 sets of histone marks and molecular factors.

220 Classes 7 and 8 display a lack of active marks, but while class 8 does not correlate with any of the  
221 used histone marks, class 7 shows a correlation with H3K9me3, a mark usually associated with  
222 constitutive heterochromatin and lack of transcription factor binding. Finally, class 9 (named 'low  
223 signal') has a very low correlation with available histone marks. However, consistently with  
224 previous studies<sup>10,44-47</sup>, it covers almost 15% of the genome, while the other classes range from  
225 around 2% to 10% in genomic coverage (**Fig. S7A**). Interestingly, the different classes are  
226 significantly differently enriched over the different chromosomes and not consistent with a uniform  
227 random genomic distribution (**Fig. S7B**, p-values<0.05, Materials and Methods).

228 To understand the relative importance of the different types of binding domains in shaping  
229 chromatin architecture, we conducted a set of *in-silico* experiments with mutant models where  
230 each class, one at the time, is erased. Specifically, from the wild-type chromosome models we  
231 removed the binding domains of a given class. Next, we computed the contact maps of the mutated  
232 model and measured across chromosomes the variation of the Pearson,  $r$ , and distance-corrected  
233 Pearson correlation coefficient,  $r'$ , between the mutated model and wild-type Hi-C contact map  
234 (Materials and Methods). The variation is found to be proportional to the genomic coverage of the  
235 different classes in both cases (**Fig. S7D,E**). That implicates that no binding class has a special role in  
236 holding the architecture of the genome in place. The linear relation whereby the removal of, say,  
237 10% of binding sites genome-wide roughly results in a 10% reduction of  $r$  highlights the structural  
238 stability of the system: the removal of a small fraction of binding sites proportionally alter the  
239 structure but does not produce a sudden collapse of the architecture, as reported by recent  
240 experiments<sup>53-57</sup>.



241 Finally, as a control of the robustness of the association between binding site types and epigenetics,  
242 we applied the same approach to the mentioned mouse ES cells<sup>14</sup>, using the corresponding set of  
243 ENCODE histone modifications in mouse, and found an overall analogous classification (**Fig. S8**).

244 Summarizing, the inferred binding site types have each a specific epigenetic barcode falling in  
245 classes that match well those found by previous epigenetic genome segmentation studies<sup>10,44–47</sup>.  
246 However, our binding domains are inferred from only Hi-C data without prior knowledge of  
247 epigenetics. Hence, they bring together independent information on architecture and epigenetics.  
248 A crucial feature of the model binding domains to explain contact data is that the different types do  
249 overlap with each other along the genome at the resolution of the considered Hi-C data. Therefore,  
250 they naturally provide each DNA window with a distinctive set of binding site types. This is an  
251 important difference with 1D epigenetic segmentation classes: by definition, those have no  
252 genomic overlap, thus each DNA window is associated to only one of such classes. Epigenetic  
253 segmentations have been shown, though, to correlate with Hi-C contacts<sup>28,32,46</sup>.

#### 254 **Epigenetic linear segmentation only partially captures chromatin folding**

255 To deepen our comprehension of the interplay of chromosome epigenetics and folding, we  
256 investigated the architectural information content retained in 1D epigenetic segmentations of the  
257 genome and compared it with the more complex DNA barcoding given by the classes of our binding  
258 domains. As done in previous studies<sup>10,44–47</sup>, we segmented chromosomes in 9 epigenetic classes  
259 based only on ENCODE histone marks (**Fig. 4A,B**). For simplicity, we opted for 9 classes to match the  
260 number of different types of binding domains found above. Such a number of classes is comparable  
261 to those in previous segmentation studies, and our results are not affected by more complex  
262 choices of segmentation (until the scale of the single binding domain is reached). Next, we derived  
263 *in-silico* the contact maps predicted by a polymer model based only on such a 1D epigenetic  
264 segmentation. Specifically, we considered a polymer model where chromatin physical interactions  
265 only occur between homologous 1D-segmented epigenetic regions<sup>28</sup>. Interestingly, while the  
266 overall contact patterns from such a model visually resemble Hi-C patterns (for example,  $r=0.78$  for  
267 chromosome 20), their distance-corrected Pearson correlation,  $r'$ , with Hi-C data is very low (for  
268 chromosome 20  $r'=0.02$ , **Fig. 4A,C** and **Fig. S9**, Materials and Methods). Hence, the patterns derived  
269 from a polymer model constructed from 1D epigenetic segmentation is only partially better than  
270 one where Hi-C pair-wise interactions are replaced by the average value corresponding to that

271 genomic separation. Conversely, an SBS model with 9 types of binding domains, based on  
272 epigenetics classes, genomically overlapping as discussed before, has  $r=0.89$  and  $r'=0.43$  for  
273 chromosome 20; and, as stated, the model with the full set of inferred binding domains has  $r=0.97$   
274 and  $r'=0.84$ .

275 To understand the partial failure of 1D epigenetic segmentation in explaining contact data (**Fig.**  
276 **4B,C**), for each pair of genomic sites we identified the binding domain that mostly contributes to  
277 their pair-wise interaction within the full SBS model (**Fig. 4D,E,F**, Materials and Methods). For  
278 clarity, we focus on a case-study 20Mb-wide region on chromosome 20. Plaid-patterns are visible in  
279 its Hi-C contact map, as expected from A/B compartments (**Fig. 4A**); they are also visible in the  
280 matrix of the most contributing binding domains (**Fig. 4F**), where rich and fine substructures appear  
281 as well. Consider, for instance, the TAD associated to region C in **Fig. 4**. The interactions within that  
282 TAD are mainly related to binding domains in class 7 (magenta, **Fig. 4F**), which is indeed the most  
283 abundant within the genomic region where C is located (**Fig. 4E**). Its interactions with the upstream  
284 region A can be simply traced back to homotypic interactions within class 7 itself, which is also the  
285 most abundant in A. However, the flanking region B, in which class 6 (dark blue) is the 1<sup>st</sup> most  
286 abundant, also interacts with C (**Fig. 4F**). That occurs because class 7 is the 2<sup>nd</sup> most abundant in B  
287 (**Fig. 4E**) and because in C class 6 is, in turn, the 2<sup>nd</sup> most abundant. Such an example illustrates that  
288 a linear epigenetic segmentation model with homotypic interactions fails to account for the  
289 complexity of the observed contact pattern because a homotypic interaction between B and C  
290 would only occur if the two regions belong to the same class. Analogously, the contacts between  
291 regions A and B originate from different overlapping binding domains included in those regions. A  
292 similar reasoning can be extended to the plaid-pattern of A/B compartments (which is a specific  
293 example of a two classes genome 1D segmentation) capturing the overall interactions between  
294 homologous active and repressed regions respectively<sup>7,42</sup>. Yet, a much more complex and finer  
295 structure of contacts exists (including interactions across A and B compartments). Indeed, it has  
296 been shown that polymer models based on a linear epigenetic classification of domains are forced  
297 to include heterotypic interactions to accurately explain Hi-C data<sup>32</sup>.

298 Overall, homotypic interactions between the domains of a coarse-grained linear epigenetic  
299 segmentation of the genome, such as compartment A/B, are not enough to explain the specificity  
300 of Hi-C patterns with high accuracy since a complexity of relevant heterotypic contacts exists  
301 between those regions. The origin of those heterotypic interactions is understood within our

302 analysis showing that multiple binding domains are present in a genomic segment. Their genomic  
303 1D combinatorial overlaps associate a distinctive interaction profile to each DNA segment,  
304 containing the information required to produce through physics the complex details of the system  
305 3D conformations (**Fig. 4**). In turn, the specific set of histone marks barcoding each binding domain  
306 provides a code linking epigenetic to architecture.

### 307 **The epigenetic barcode of binding domains predicts *de novo* chromatin architecture**

308 To validate the identified association between linear epigenetic features and chromosome  
309 conformations, we considered a reverse approach whereby starting from only epigenetics data,  
310 through the mentioned barcode we identify the key binding sites of a set of independent  
311 chromosomes and, next, predict their contact matrices via polymer physics (**Fig. 5A**). Specifically,  
312 we exploited the epigenetic barcoding provided by the classification of the binding domains of  
313 even-numbered chromosomes, as previously described, to identify *de novo* the binding sites of  
314 odd-numbered chromosomes. To determine the locations and types of the binding sites, we  
315 partitioned each 5kb genomic window (5kb is the resolution of Hi-C) of odd-numbered  
316 chromosomes in equal-sized, 0.5kb sub-windows, which we epigenetically profiled by measuring  
317 the abundance of the mentioned key set of histone marks (Materials and Methods). We then  
318 computed the correlations between the epigenetic profile of each sub-window and the centroids of  
319 the epigenetic classes of the binding domains of even-numbered chromosomes (**Fig. 3A**). We focus  
320 on those epigenetics classes because they recapitulate the main functional groups found in  
321 segmentation studies; additionally, considering 9 types of sites is more stringent than considering  
322 all the binding domains found on even chromosomes. Exploiting such a larger set of domains  
323 would, of course, improve our results. Finally, each sub-windows of odd-numbered chromosomes  
324 was assigned with a binding site type corresponding to the epigenetic class having the highest  
325 correlation (**Fig. 5A**).

326 Once obtained the genomic locations of the binding sites along odd-numbered chromosomes, we  
327 computed their contact matrices via the SBS polymer model and compared them with the  
328 corresponding *in situ* Hi-C maps (**Fig. S10A,B**). **Fig. 5B,C** shows, for example, the contact data of  
329 chromosomes 19 and 21 predicted by use of the above defined code that links binding sites, i.e.,  
330 architecture, to epigenetic marks. In all the considered cases, the predicted matrices well capture  
331 the patterns of interactions seen in Hi-C data even at large genomic distances, albeit for simplicity

332 we considered only 9 types of binding domains. Accordingly, the correlation and distance-corrected  
333 correlation coefficients ( $r=0.91$   $r'=0.47$  and  $r=0.91$   $r'=0.63$  for respectively chromosome 19 and 20)  
334 are much higher than those found by 1D epigenetic segmentation, as seen above.

335 Taken together, our results show that the barcode linking epigenetics marks to the binding domains  
336 inferred by PRISMR from Hi-C data, albeit still incomplete, can predict the genome's 3D architecture  
337 to a good level of accuracy. A crucial difference between ours and epigenetic segmentation  
338 strategies to predict chromatin contacts<sup>58</sup> is the intrinsically overlapping nature of binding domains,  
339 lacking in segmentations, which is necessary to recapitulate accurately the complex pattern of  
340 chromatin interactions.

341

## 342 **DISCUSSION**

343 To infer from Hi-C data the different types of DNA binding sites determining chromosome  
344 architecture and their genomic position, we employed an approach based on machine learning<sup>41</sup>  
345 and the physics of the SBS polymer model of chromatin. The SBS model quantifies the scenario  
346 where TFs mediate the interactions between distal cognate binding sites, establishing DNA contacts  
347 and loops<sup>21</sup>. We found that the 3D structures derived by the model informed with the inferred  
348 putative binding domains, and folded through only polymer physics, explain Hi-C data genome-wide  
349 with high accuracy in human GM12878 B-lymphoblastoid<sup>42</sup> and mES<sup>14</sup> cells. That shows that the  
350 basic molecular ingredients considered by the model are sufficient to explain contact patterns  
351 across genomic scales. Thus, the binding domains encode key molecular information required to  
352 fold chromatin and provide an *architectural code* whereby 3D conformations can be established  
353 based on the 1D sequence (**Fig. 6**). To explain folding with high accuracy, they have a combinatorial  
354 organization along chromosomes, which is needed to control the intricate multitude of genomic  
355 interactions captured in Hi-C maps and their functional specificity, via a comparatively smaller  
356 number of molecular factors. Additionally, the non-trivial arrangement of binding domains provides  
357 structural stability to the 3D conformation of the genome, as experimentally reported<sup>53-57</sup>. We  
358 found that binding domains produce chromatin interactions extending across chromosomal scales,  
359 above the size of single TADs and A/B compartments, in a hierarchy of higher-order 3D structures,  
360 as in the meta-TADs<sup>16</sup> picture.

361 Next, we associated each of the Hi-C inferred binding domains to an epigenetic profile based on the  
362 genomic correlation with a few main ENCODE histone marks. The model binding domains turn out  
363 to belong to main epigenetic classes, similar in human and mouse cell types, which well match  
364 known chromatin states (e.g., active, poised, repressed) derived by linear segmentation studies<sup>10,44–</sup>  
365 <sup>47</sup>. However, as stated, the identified binding domains have broad overlaps along the genome, a  
366 feature that is missing in linear segmentations but is required to explain Hi-C accurately. The few  
367 coarse-grained epigenetic classes here discussed constitute only a first, simplified description of the  
368 epigenetic features of the binding domains that shape chromatin architecture inferred by PRISMR.  
369 More generally, their barcode is expected to be associated with a broader set of (still partially  
370 unknown) molecular factors, including histone marks, CTCF<sup>42</sup>, Active/Poised Pol-II<sup>40</sup> and additional  
371 factors, such as PRC1<sup>51</sup>, PRC2<sup>40</sup> and MLL3/4<sup>59</sup>. Furthermore, molecular mechanisms beyond those  
372 envisaged by the SBS model, such as DNA loop extrusion<sup>22–24</sup>, appear to play a role in chromosome  
373 folding and the code can be extended to accommodate them.

374 The inferred binding domains and the associated architectural interaction code were tested by  
375 making predictions on the changes of the 3D structure caused by a set of structural variants at the  
376 *Sox9* locus linked to human diseases. Notably, the predicted contact maps were confirmed by  
377 independent cHi-C data in cells carrying such mutations<sup>43</sup>. This is a stringent validation because  
378 there are no available fitting parameters. The model also helps understanding how the mutations  
379 differently affect the 3D structure of the locus (e.g., forming neo-TADs) and how that differently  
380 impacts gene regulation and, hence, phenotype by enhancer hijackings.

381 Finally, in a reverse approach, based on the discovered link between epigenetics and binding  
382 domains, we identified the binding sites of an independent set of chromosomes from only their  
383 epigenetic marks. Those binding sites were sufficient to predict *de novo*, via the physics of the SBS  
384 model, the contact matrices of those chromosomes with good accuracy, validating our approach.

385 Overall, the agreement between our results and independent experimental Hi-C data strengthens  
386 the scenario where chromatin 3D architectural information is encoded in a 1D combinatorial  
387 arrangement of epigenetically barcoded sites, which can be inferred across chromosomes and cell  
388 types by our computational approach. By integration of different genomic data, it provides a  
389 quantitative picture of the deep cause-effect relationship between epigenetics, architecture and  
390 function, which we tested to predict the phenotypic effect of mutations linked to congenital

391 disorders. That can help the development of computational tools in biomedicine to infer the link  
392 between genotype and phenotype from the features of the genomic landscape.

## 393 MATERIALS AND METHODS

### 394 The String & Binders Switch model of chromatin

395 To investigate the 3D structure of the genome, we employed the String & Binders Switch (SBS)  
396 model<sup>21,26,30</sup>. According to the SBS, a chromatin filament (from small loci to entire chromosomes) is  
397 modeled as a self-avoiding walk polymer chain of beads, a fraction of which, named binding sites,  
398 interacts with diffusing molecular binders. The interaction between binding sites and binders allows  
399 for the formations of loops along the polymer and, therefore, permits its spontaneous folding (**Fig.**  
400 **1C**). Each bead can be bound only by its specific, cognate type of binders and, to fully describe the  
401 complexity of the system, different types of interactions are allowed together with inert sites along  
402 the chain that do not interact with any binder (apart from steric effects). We represent these  
403 different interactions as different “colors” of the system, “gray” beads being the non-interacting  
404 particles (**Fig. 1C**). Key parameters of the model are the concentration,  $c$ , and the binding energy,  
405  $E_{int}$ , of each different type of binder. As a function of  $c$  and  $E_{int}$ , the system of corresponding,  
406 cognate binding sites exhibits a coil-globule phase transition from an open conformation (at low  
407 concentration or energy) to a globule, compact phase (at high concentration or energy) as  
408 extensively discussed in previous studies<sup>26,30,39</sup>. The presence of different sets of binding sites (here  
409 named “binding domains” and represented with different colors) interacting with different, cognate  
410 molecular factors allows the formation of complex 3D structures by microphase separation.

### 411 The PRISMR method

412 To determine the distribution of the different binding sites along the SBS polymer chain, here we  
413 used PRISMR, a previously illustrated machine learning procedure<sup>41</sup>. The PRISMR algorithm is a  
414 polymer physics-based method that, starting from an experimental contact matrix (e.g. Hi-C or  
415 GAM), finds the minimal polymer model that, at equilibrium, best describes the input. Although we  
416 focus on the SBS polymer model to describe a chromatin filament, the PRISMR algorithm can be  
417 easily generalized to different models.

418 A detailed description of the PRISMR method can be found in ref <sup>41</sup>. Here we just summarize the  
419 key points of the algorithm. An SBS polymer model of a genomic region is composed of  $L$  beads,  
420 depending on the resolution of the input contact matrix of the region. For instance, a 10Mb locus at  
421 10kb resolution is partitioned in  $L=1000$  bins. Furthermore, we split each of the  $L$  bins into  $r$

422 different sub-units, considering that a single DNA bin could include many binding sites and interact  
423 with different factors. The SBS polymer is then completely characterized by the arrangement of the  
424 binding sites along the chain. Given the number  $n$  of different types of binding sites, PRISMR finds  
425 the color arrangement along the polymer chain by the minimization, via an iterative Simulated  
426 Annealing (SA) Monte Carlo optimization procedure<sup>60,61</sup>, of a specific cost function made of two  
427 terms. The first term representing the distance between the experimental and the model predicted  
428 contact matrices; the second one is a Bayesian term proportional to the total number of colored  
429 sites of the polymer through a parameter  $\lambda$  and penalizes the addition of new colored beads. In this  
430 way we account for the necessity to fit well the input data and, at the same time, we attempt to  
431 avoid overfitting. After initializing the SBS polymer in a random configuration, by assigning a  
432 random color to each bead, a standard iterative SA procedure is performed, as available in public  
433 software repositories (see e.g.<sup>62</sup>), to optimize the model<sup>60,61</sup>. Schematically, each SA step consists in  
434 randomly changing the color of a polymer bead, compute the average contact matrix of the new  
435 polymer, evaluate the new cost function, compare it with the cost function in the previous step  
436 and, based on it, accept or reject the color change. SA steps are iteratively repeated until  
437 convergence<sup>41</sup>. The entire procedure is repeated many times by varying the polymer initial  
438 configurations and the model parameters  $n$ ,  $r$ , and  $\lambda$ , to find their optimal values.

### 439 **Details on the application of PRISMR genome-wide**

440 In this study, we present the first genome-wide application of the algorithm. Precisely, here we  
441 applied PRISMR over the somatic chromosomes of the human genome, obtaining, for each  
442 chromosome independently, the SBS polymer that best describes its corresponding Hi-C matrix. We  
443 employed published in situ Hi-C data<sup>42</sup> relative to the human GM12878 cell line at 5kb of resolution  
444 and normalized according to the method described in ref<sup>63</sup>. To reduce the local noise in the input  
445 Hi-C data, we applied a gaussian filter with a standard deviation equal to 1 along both  $x$  and  $y$   
446 directions. The optimal value of the parameters of the algorithm has been estimated as already  
447 described in ref<sup>41</sup>, that is, we repeated the SA procedure many times starting from different initial  
448 conditions and different values of  $n$ ,  $r$ , and  $\lambda$  to set these parameters at the values that explain the  
449 input data within a given accuracy. As input data for the optimal parameter evaluation, we used the  
450 contact matrix of chromosome 12, a medium-sized chromosome, obtaining  $n=30$  different types of  
451 binding sites,  $r=30$  and  $\lambda=3 \times 10^{-5}$ . The same values for the parameters  $n$ ,  $r$ , and  $\lambda$  have been used to  
452 obtain the best SBS polymer for all the other chromosomes. **Fig. S1A** shows the comparison



453 between the contact matrices inferred by PRISMR (lower triangular maps) and the in situ Hi-C  
454 matrices (upper triangular maps). The global pattern obtained by PRISMR is highly correlated with  
455 the experimental one as also quantified by the comparatively high values of the Pearson's ( $r$ ),  
456 distance-corrected Pearson's ( $r'$ )<sup>41</sup> and stratum-adjusted (SCC)<sup>48</sup> correlation coefficients (**Fig. S1B**,  
457 see below). In the calculation of  $r$  and  $r'$ , to correct for outliers, we did not consider genomic  
458 distances below 25kb. The PRISMR method is highly generalizable across different experiments and  
459 data resolution. To test that, we also applied our method to genome-wide Hi-C data in mouse  
460 embryonic stem (mES) cells<sup>14</sup> at 40kb resolution (**Fig. S2A**). The correlations between experimental  
461 and model matrices obtained in mouse are as high as the values obtained in human, as shown in  
462 **Fig. S2B**.

### 463 **Structural variants at the Sox9 locus and validation of PRISMR**

464 As a validation of the PRISMR inference method and the SBS model, we implemented in-silico a set  
465 of three previously studied structural variants in E12.5 limb bud cells from mice<sup>43</sup>. Specifically, we  
466 started from a SBS polymer model<sup>3</sup> of the region chr11:109000000-115000000 (mm9, mESC cells)  
467 including the *Sox9* gene and implemented on it, independently, the following duplications: *DupS*, an  
468 intra-TAD duplication of the region chr11:111760000-112160000; *DupL*, an inter-TAD duplication of  
469 the region chr11:110960000-112520000; *DupC*, another inter-TAD duplication of the region  
470 chr11:110760000-112520000. We then computed the PRISMR predicted contact maps for each  
471 duplication, under no adjustable parameters, obtaining the following values of correlations  $r$  and  $r'$ ,  
472 between model and experimental matrices (excluding the effect of strong outliers <5th and >95th  
473 percentile):  $r=0.88$  and  $r'=0.48$  in *DupS*;  $r=0.82$  and  $r'=0.41$  in *DupL*;  $r=0.82$  and  $r'=0.47$  in *DupC* (**Fig.**  
474 **S4**).

### 475 **Matrix similarity evaluation**

476 The agreement between experiment and model matrices has been quantified using Pearson's  
477 correlation coefficient,  $r$ . We also used two additional measures: 1) the distance corrected Pearson  
478 correlation coefficient, denoted by  $r'$ , that is the Pearson's correlation coefficient between the two  
479 matrices where we subtracted from each diagonal (corresponding to a given genomic distance)  
480 their average contact frequency; 2) the stratum-adjusted correlation coefficient, denoted by SCC,  
481 from the HiCRep<sup>48</sup> method with a smoothing parameter  $h=10$  and an upper bound of interaction

482 distance equal to 5Mb. These two measures have been used to put aside the expected decreasing  
483 trend of the pairwise contact frequency with genomic distance, which tends to dominate in the  
484 simple Pearson correlation value.

### 485 **Molecular Dynamics simulations**

486 To obtain 3D conformations of the PRISMR derived SBS models, shown in **Fig. 1F**, **Fig. 2D,E** and **Fig.**  
487 **S5C,D**, we performed Molecular Dynamics (MD) simulations. To this aim, we proceeded as  
488 described in ref <sup>30</sup>. Briefly, the polymer chain and the binders move in the system according to the  
489 Langevin equation, integrated with the LAMMPS software<sup>64</sup>, using standard dimensionless  
490 parameters<sup>65</sup>. The SBS parameters used are the same reported in ref <sup>30</sup>, i.e., the beads and binders  
491 interact with an interaction energy  $E_{int}=8.1KbT$  and the binder concentration is high enough to  
492 allow the coil-globule transition ( $c=194nmol/l$  for the *Sox9* WT and similar values for the  
493 duplications). To make MD computation times feasible for the entire chromosome 20, we  
494 considered a coarse-grained version of its SBS polymer, having a 50-fold reduced number of beads.  
495 All the conformations are taken in the equilibrium globular phase. In all the snapshots, beads  
496 coordinates have been interpolated with a smooth third-order polynomial spline curve by using the  
497 POV-RAY<sup>66</sup> software.

### 498 **Characterization of the binding domains arrangement along chromosomes**

499 To study how the different binding domains (colors) span along the genome, we employed two  
500 main measures. The first one, that measures the domain size, is the genomic coverage, i.e., the  
501 fraction of beads of a given color multiplied by the length of the chromosome it belongs to.  
502 Averaging over all the sizes of the domains identified by PRISMR across chromosomes, we find that  
503 the genomic length covered by each domain is on average 3.1 Mb, with a standard deviation of 1.9  
504 Mb, a value close to the mean-size of a TAD. To measure, instead, the range of the interactions due  
505 to a single binding domain, we defined  $r_{int}$  as two times the standard deviation of the center of  
506 mass of that domain. The distribution of  $r_{int}$ ,  $P(r_{int})$ , extends far beyond the size of the single  
507 domain, ranging from a few mega-bases to more than 100 Mb (**Fig. S3A**). To check the statistical  
508 significance of the domains identified by PRISMR, we compared  $P(r_{int})$  with a control model  
509 obtained by randomly bootstrapping the location of our binding sites along the genome, and we  
510 found that the two distributions are significantly different ( $p$ -value<0.001, Wilcoxon's rank sum

511 test). We also found that  $P(r_{int})$  is asymptotically consistent with a power-law scaling, as shown in  
512 **Fig. S3A** where the right-hand side of the distribution is well described by a power-law fit (dotted  
513 red curve in the graph).

514 Another way to test the significance of the binding domains identified by PRISMR is to measure  
515 their mutual overlap<sup>41</sup>, to be compared with the expected level of overlap in the random model of  
516 bootstrapped domains mentioned before. To this aim, given a pair of different domains on a  
517 chromosome, we defined their overlap  $q$  as the sum of products of binding sites occurrences of the  
518 two colors in each genomic window, normalized to have  $q=100\%$  in the case of identical domains  
519 (the cartoon in **Fig. S3B** gives a visual impression of what  $q$  is measuring). We found that the  
520 distribution  $P(q)$  of the overlap of the binding domains predicted by PRISMR is significantly different  
521 ( $p\text{-value}<0.001$ , Wilcoxon's rank sum test) from the one expected in the random control model (red  
522 and blue distributions in **Fig. S3B**, respectively).

### 523 **Epigenetic analysis of the binding domains**

524 To obtain insight into their molecular nature, we analyzed the PRISMR inferred binding domains in  
525 the light of epigenetics data. To this aim, we downloaded from the ENCODE database<sup>49</sup> a set of 5  
526 key histone modifications (H3K4me3, H3K4me1, H3K36me3, H3K9me3 and H3K27me3) in the  
527 human GM12878 cell line. ChIP-Seq signals were binned at 5kb resolution by summing the signal  
528 contained within each 5kb window (using the bedtools map tool from the bedtools<sup>67</sup> software).  
529 After that, to measure the similarity between our binding domains and the histone marks, we  
530 computed Pearson's correlation coefficient between the number of binding sites of each domain  
531 and each histone mark profile. Next, we employed a control model to retain only statistically  
532 significant correlations. To this aim, first, we computed the Pearson correlations between  
533 chromatin mark signals and randomized binding domains signals obtained by bootstrapping their  
534 actual genomic locations; then, we retained as significant only the correlation values above the  
535 95th or below the 5th percentile of the distribution of the random correlations. We then collected  
536 data in a rectangular matrix  $X$ , whose element  $X_{ij}$  is either the significant correlation between the  $i$ -  
537 th binding domains and the  $j$ -th histone mark or zero if the correlation was not significant. Since  
538 each row of  $X$  represents a binding domain's correlation profile with the considered histone  
539 modifications, we refer to them as the epigenomic signature of the binding domain. To find binding  
540 domains with similar epigenomic signatures, we performed a hierarchical clustering analysis on  $X$

541 using the *Python SciPy* clustering package with ‘Euclidean’ distance metric and ‘Ward’ linkage  
542 method. To assess the number of clusters in the hierarchical clustering output, we cut the  
543 dendrogram at different values (ranging from one to the number of binding domains) and  
544 evaluated the Akaike Information Criterion<sup>50</sup> (AIC) as the number of clusters  $k$  is varied. As shown in  
545 **Fig. S6A**, while no sharp transitions are present, the curve has a global minimum at  $k=9$ . We  
546 therefore grouped all the different rows of  $X$  in 9 different classes according to their affinity to each  
547 cluster (**Fig. S6B**). Each of the 9 classes can be characterized by the epigenetic signature of its  
548 centroid, which is the average histone signature of the domains belonging to the given class (**Fig.**  
549 **3A**). To assign biologically meaningful labels to the obtained classification, we looked at the  
550 enrichment of several types of functional annotations. Precisely, we first binned each annotation  
551 track at 5kb resolution, then, for each pair of annotation mark and epigenetic class, we computed  
552 the average of the Pearson correlation values between that mark and the binding domains of that  
553 class (see **Fig. S6C**). The set of functional annotations in GM12878 cell line considered in this study  
554 is taken from ENCODE and include: (1) all remaining available histone modifications; (2)  
555 transcription factors binding sites; (3) DNase hypersensitive sites; (4) replication timing data from  
556 the Repli-seq assay; (5) transcription data from RNA-seq assay (**Fig. S6C**).

557 To further test the association between binding domains and epigenetics, we repeated the above  
558 analysis for the mouse case. Specifically, we computed correlations among the genome-wide  
559 binding domains obtained from Hi-C data in mES cells and a corresponding set of ENCODE histone  
560 modifications in that cell line. As shown in **Fig. S8A-C**, we found an overall similar epigenetic  
561 classification of the binding domains in human and mouse.

## 562 **Characterization of epigenetic classes of binding domains**

563 The genomic coverage of a given epigenetic class has been computed as the fraction of sites of the  
564 binding domains belonging to that class (**Fig. S7A**). To study, instead, how the domains of a given  
565 class are distributed along the chromosomes, we counted, for each class, the number of domains  
566 falling in each chromosome (**Fig. S7B**, dotted lines are the average values). We found that their  
567 distribution is significantly different over the different chromosomes, as measured by the  
568 comparison with a uniform distribution obtained by randomly bootstrapping the domains of a given  
569 class over the chromosomes ( $p$ -value $<0.05$  for each epigenetic class, Kolmogorov-Smirnov test). We  
570 also asked whether the genomic positions of the sites of the different classes (**Fig. 3B**) were

571 correlated with each other. To figure out that, we computed the Pearson correlation between the  
572 genomic location of the sites of all the possible pairs of epigenetic classes, averaged over the  
573 different chromosomes (**Fig. S7C**). We found that classes with similar histone signature correlate  
574 with each other and anti-correlate with classes showing a very different histone pattern.

575 We investigated the impact of the different epigenetic classes on genome architecture by  
576 measuring the effect on contact matrices of the withdrawal of the binding domains belonging to  
577 each class. Precisely, given the list of the binding domains of a class, we replaced their interacting  
578 binding sites with gray, non-interacting elements along each chromosome. We then computed the  
579 PRISMR contact matrices of the modified SBS polymer and measured their correlations  $r$  and  $r'$  with  
580 Hi-C. Finally, we evaluated the variation of the correlation,  $\Delta r$  and  $\Delta r'$ , with respect to the wild-type  
581 model ( $r=0.94$  and  $r'=0.76$ ), averaged over all chromosomes. The variations of  $r$  and  $r'$  obtained are  
582 shown as a function of the genomic coverage of each epigenetic class in **Fig. S7D**.

### 583 **Most abundant and most contributing binding domains to chromatin pairwise contacts**

584 As the different binding domains can overlap with each other, to better visualize their locations  
585 along the genome, we show in **Fig. 4E** (upper bar) the 1<sup>st</sup> most abundant binding domain, i.e. the  
586 one with the largest number of binding sites, per bin. Analogously, **Fig. 4E** (lower bar) shows the 2<sup>nd</sup>  
587 most abundant binding domain per bin. In both cases, to help the visualization, the domains are  
588 colored with their epigenetic class color.

589 The contribution of the different binding domains in forming the interactions between bin pairs is  
590 then highlighted in **Fig. 4F**, where the colors of the most contributing binding domains are shown.  
591 Specifically, for a given pair-wise contact, we defined the contribution of a binding domain to that  
592 contact as the number of pairs of its binding site type between the two considered bins. The  
593 binding domain having the highest number of binding site pairs is the most contributing one and is  
594 colored with the color corresponding to its epigenetic class.

### 595 **Epigenetic linear segmentation model**

596 To obtain a model based exclusively on the interaction among segments with a similar epigenetic  
597 profile, we considered the dataset of five histone modifications discussed in section “*Epigenetic*  
598 *analysis of the binding domains*”. We marked each 5kb genomic window with the z-score value of

599 the signal of each histone mark in that window. Then, we performed a hierarchical clustering  
600 analysis to gather the genomic windows with similar histone profiles in 9 different groups, in order  
601 to match them with the 9 different types of binding domains found above. The obtained linear  
602 segmentation has been employed to define a polymer model for chr.20 with 9 different colors  
603 corresponding to the different linear epigenetic classes (**Fig. 4B**), where interactions can only occur  
604 between same-colored windows. Finally, we derived in-silico the contact map of such a model and  
605 compared it with the corresponding experimental matrix (**Fig. 4A-C** and **Fig. S9A-B**). We found that  
606 the Pearson correlation and distance-corrected Pearson correlation between the matrices are  
607  $r=0.80$  and  $r'=0.21$ .

608 We have also considered an additional model by assigning each of the different binding sites of  
609 chr.20 the color of the epigenetic class it belongs to. We found that this 9 color SBS model, that in  
610 contrast to the linear segmentation model has overlapping binding domains, has correlations  
611  $r=0.89$  and  $r'=0.43$  with Hi-C.

#### 612 **Prediction of *de novo* chromatin structures from epigenetic data by combinatorial barcode**

613 The derived combinatorial code linking 3D conformation to 1D epigenetic signature can be used to  
614 predict *de novo* binding domains in independent chromosomes from epigenetics data only.  
615 Specifically, we used the code derived from the set of even-numbered chromosomes in GM12878  
616 to predict the location of the binding sites along the odd-numbered chromosomes in the same cell  
617 line. To this aim, we partitioned each of their 5kb windows (which is the in situ Hi-C data resolution)  
618 in ten 500-bp sub-windows and binned the signal of the five key histone marks (H3K4me3,  
619 H3K4me1, H3K36me3, H3K9me3 and H3K27me3) in those sub-windows. In this way, we obtained a  
620 state vector for each sub-window, whose components are the histone marks' abundances in that  
621 window. We checked that different sub-windows partitions, ranging from 5 to 20 sub-windows per  
622 bin, led to only marginally different results. To assign each sub-window to a specific binding site  
623 type (in the sense of the SBS model), we compared them with the centroids of the epigenetic  
624 classes of the binding domains of even-numbered chromosomes. Precisely, we computed the  
625 Pearson correlation coefficient between the state vector and each row of the centroid matrix, then  
626 assigned to that sub-window a binding site type corresponding to the epigenetic class with the  
627 highest correlation. Besides, two non-interacting 'gray' beads were added in each sub-window, so  
628 to match the number of beads per 5kb-bin of the PRISMR inferred polymer models. The described

629 procedure results in an SBS polymer with 9 different binding domains for each odd chromosome.  
630 Afterward, we used the SBS model to calculate the predicted polymers' contact matrices and  
631 compared them with the independent Hi-C data (Fig. **S10**). As reflected by the Pearson and distance  
632 corrected Pearson correlations, in all cases, the contact pattern is well described (see for instance  
633 chromosomes 19 and 21 in **Fig. 5**).

634 **REFERENCES**

- 635 1. Misteli, T. Beyond the Sequence: Cellular Organization of Genome Function. *Cell* vol. 128  
636 787–800 (2007).
- 637 2. Bickmore, W. A. & Van Steensel, B. Genome architecture: Domain organization of interphase  
638 chromosomes. *Cell* **152**, 1270–1284 (2013).
- 639 3. Sexton, T. & Cavalli, G. The role of chromosome domains in shaping the functional genome.  
640 *Cell* vol. 160 1049–1059 (2015).
- 641 4. Dekker, J. & Heard, E. Structural and functional diversity of Topologically Associating  
642 Domains. *FEBS Letters* vol. 589 2877–2884 (2015).
- 643 5. Dekker, J. & Mirny, L. The 3D Genome as Moderator of Chromosomal Communication. *Cell*  
644 **164**, 1110–1121 (2016).
- 645 6. Dixon, J. R., Gorkin, D. U. & Ren, B. Chromatin Domains: The Unit of Chromosome  
646 Organization. *Molecular Cell* vol. 62 668–680 (2016).
- 647 7. Lieberman-Aiden, E. *et al.* Comprehensive mapping of long-range interactions reveals folding  
648 principles of the human genome. *Science (80-. )*. **326**, 289–293 (2009).
- 649 8. Beagrie, R. A. *et al.* Complex multi-enhancer contacts captured by genome architecture  
650 mapping. *Nature* **543**, 519–524 (2017).
- 651 9. Quinodoz, S. A. *et al.* Higher-Order Inter-chromosomal Hubs Shape 3D Genome Organization  
652 in the Nucleus. *Cell* **174**, 744-757.e24 (2018).
- 653 10. Boettiger, A. N. *et al.* Super-resolution imaging reveals distinct chromatin folding for different  
654 epigenetic states. *Nature* **529**, 418–22 (2016).
- 655 11. Cattoni, D. I. *et al.* Single-cell absolute contact probability detection reveals chromosomes  
656 are organized by multiple low-frequency yet specific interactions. *Nat. Commun.* (2017)  
657 doi:10.1038/s41467-017-01962-x.
- 658 12. Bintu, B. *et al.* Super-resolution chromatin tracing reveals domains and cooperative  
659 interactions in single cells. *Science (80-. )*. **362**, eaau1783 (2018).
- 660 13. Finn, E. H. *et al.* Extensive Heterogeneity and Intrinsic Variation in Spatial Genome  
661 Organization. *Cell* (2019) doi:10.1016/j.cell.2019.01.020.
- 662 14. Dixon, J. R. *et al.* Topological domains in mammalian genomes identified by analysis of  
663 chromatin interactions. *Nature* **485**, 376–80 (2012).
- 664 15. Nora, E. P. *et al.* Spatial partitioning of the regulatory landscape of the X-inactivation centre.



- 665 *Nature* **485**, 381–385 (2012).
- 666 16. Fraser, J. *et al.* Hierarchical folding and reorganization of chromosomes are linked to  
667 transcriptional changes in cellular differentiation. *Mol. Syst. Biol.* **11**, 852–852 (2015).
- 668 17. Oudelaar, A. M. *et al.* Between form and function: The complexity of genome folding. *Human*  
669 *Molecular Genetics* vol. 26 R208–R215 (2017).
- 670 18. Krijger, P. H. L. & De Laat, W. Regulation of disease-associated gene expression in the 3D  
671 genome. *Nature Reviews Molecular Cell Biology* vol. 17 771–782 (2016).
- 672 19. Spielmann, M., Lupiáñez, D. G. & Mundlos, S. Structural variation in the 3D genome. *Nat.*  
673 *Rev. Genet.* 1–15 (2018) doi:10.1038/s41576-018-0007-0.
- 674 20. Nicodemi, M. & Pombo, A. Models of chromosome structure. *Current Opinion in Cell Biology*  
675 vol. 28 90–95 (2014).
- 676 21. Nicodemi, M. & Prisco, A. Thermodynamic pathways to genome spatial organization in the  
677 cell nucleus. *Biophys. J.* **96**, 2168–2177 (2009).
- 678 22. Sanborn, A. L. *et al.* Chromatin extrusion explains key features of loop and domain formation  
679 in wild-type and engineered genomes. *Proc. Natl. Acad. Sci.* **112**, E6456–E6465 (2015).
- 680 23. Fudenberg, G. *et al.* Formation of Chromosomal Domains by Loop Extrusion. *Cell Rep.* **15**,  
681 2038–2049 (2016).
- 682 24. Brackley, C. A. *et al.* Nonequilibrium Chromosome Looping via Molecular Slip Links. *Phys. Rev.*  
683 *Lett.* **119**, (2017).
- 684 25. Bohn, M. & Heermann, D. W. Diffusion-driven looping provides a consistent framework for  
685 chromatin organization. *PLoS One* **5**, (2010).
- 686 26. Barbieri, M. *et al.* Complexity of chromatin folding is captured by the strings and binders  
687 switch model. *Proc. Natl. Acad. Sci.* **109**, 16173–16178 (2012).
- 688 27. Brackley, C. A., Taylor, S., Papantonis, A., Cook, P. R. & Marenduzzo, D. Nonspecific bridging-  
689 induced attraction drives clustering of DNA-binding proteins and genome organization. *Proc.*  
690 *Natl. Acad. Sci.* **110**, E3605–E3611 (2013).
- 691 28. Jost, D., Carrivain, P., Cavalli, G. & Vaillant, C. Modeling epigenome folding: Formation and  
692 dynamics of topologically associated chromatin domains. *Nucleic Acids Res.* **42**, 9553–9561  
693 (2014).
- 694 29. Brackley, C. A. *et al.* Predicting the three-dimensional folding of cis-regulatory regions in  
695 mammalian genomes using bioinformatic data and polymer models. *Genome Biol.* **17**, 59  
696 (2016).

- 697 30. Chiariello, A. M. A. M., Annunziatella, C., Bianco, S., Esposito, A. & Nicodemi, M. Polymer  
698 physics of chromosome large-scale 3D organisation. *Sci. Rep.* **6**, (2016).
- 699 31. Di Stefano, M., Paulsen, J., Lien, T. G., Hovig, E. & Micheletti, C. Hi-C-constrained physical  
700 models of human chromosomes recover functionally-related properties of genome  
701 organization. *Sci. Rep.* **6**, (2016).
- 702 32. Di Pierro, M., Zhang, B., Aiden, E. L., Wolynes, P. G. & Onuchic, J. N. Transferable model for  
703 chromosome architecture. *Proc. Natl. Acad. Sci.* **113**, 12168–12173 (2016).
- 704 33. Li, Q. *et al.* The three-dimensional genome organization of *Drosophila melanogaster* through  
705 data integration. *Genome Biol.* (2017) doi:10.1186/s13059-017-1264-5.
- 706 34. Serra, F. *et al.* Automatic analysis and 3D-modelling of Hi-C data using TADbit reveals  
707 structural features of the fly chromatin colors. *PLoS Comput. Biol.* (2017)  
708 doi:10.1371/journal.pcbi.1005665.
- 709 35. Nir, G. *et al.* Walking along chromosomes with super-resolution imaging, contact maps, and  
710 integrative modeling. *PLoS Genet.* (2018) doi:10.1371/journal.pgen.1007872.
- 711 36. Lin, D., Bonora, G., Yardimci, G. G. & Noble, W. S. Computational methods for analyzing and  
712 modeling genome structure and organization. *Wiley Interdiscip. Rev. Syst. Biol. Med.* (2018)  
713 doi:10.1002/wsbm.1435.
- 714 37. Chiariello, A. M. *et al.* A Dynamic Folded Hairpin Conformation Is Associated with  $\alpha$ -Globin  
715 Activation in Erythroid Cells. *Cell Rep.* (2020) doi:10.1016/j.celrep.2020.01.044.
- 716 38. Bianco, S. *et al.* Modeling Single-Molecule Conformations of the HoxD Region in Mouse  
717 Embryonic Stem and Cortical Neuronal Cells. *Cell Rep.* **28**, 1574-1583.e4 (2019).
- 718 39. Conte, M. *et al.* Polymer physics indicates chromatin folding variability across single-cells  
719 results from state degeneracy in phase separation. *Nat. Commun.* (2020)  
720 doi:10.1038/s41467-020-17141-4.
- 721 40. Barbieri, M. *et al.* Active and poised promoter states drive folding of the extended HoxB  
722 locus in mouse embryonic stem cells. *Nat. Struct. Mol. Biol.* **24**, 515–524 (2017).
- 723 41. Bianco, S. *et al.* Polymer physics predicts the effects of structural variants on chromatin  
724 architecture. *Nat. Genet.* **50**, 662–667 (2018).
- 725 42. Rao, S. S. P. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles  
726 of chromatin looping. *Cell* **159**, 1665–80 (2014).
- 727 43. Franke, M. *et al.* Formation of new chromatin domains determines pathogenicity of genomic  
728 duplications. *Nature* **538**, 265–269 (2016).

- 729 44. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types.  
730 *Nature* **473**, 43–49 (2011).
- 731 45. Gifford, C. A. *et al.* Transcriptional and epigenetic dynamics during specification of human  
732 embryonic stem cells. *Cell* **153**, 1149–1163 (2013).
- 733 46. Ho, J. W. K. *et al.* Comparative analysis of metazoan chromatin organization. *Nature* **512**,  
734 449–452 (2014).
- 735 47. Javierre, B. M. *et al.* Lineage-Specific Genome Architecture Links Enhancers and Non-coding  
736 Disease Variants to Target Gene Promoters. *Cell* **167**, 1369–1384.e19 (2016).
- 737 48. Yang, T. *et al.* HiCRep: assessing the reproducibility of Hi-C data using a stratum-adjusted  
738 correlation coefficient. *Genome Res.* **27**, 1939–1949 (2017).
- 739 49. Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature*  
740 **489**, 57–74 (2012).
- 741 50. Akaike, H. A New Look at the Statistical Model Identification. *IEEE Trans. Automat. Contr.* **19**,  
742 716–723 (1974).
- 743 51. Kundu, S. *et al.* Polycomb Repressive Complex 1 Generates Discrete Compacted Domains  
744 that Change during Differentiation. *Mol. Cell* **65**, 432–446.e5 (2017).
- 745 52. Tang, Z. *et al.* CTCF-Mediated Human 3D Genome Architecture Reveals Chromatin Topology  
746 for Transcription. *Cell* **163**, 1611–1627 (2015).
- 747 53. Nora, E. P. *et al.* Targeted Degradation of CTCF Decouples Local Insulation of Chromosome  
748 Domains from Genomic Compartmentalization. *Cell* **169**, 930–944.e22 (2017).
- 749 54. Rao, S. S. P. *et al.* Cohesin Loss Eliminates All Loop Domains. *Cell* **171**, 305–320.e24 (2017).
- 750 55. Kubo, N. *et al.* Preservation of Chromatin Organization after Acute Loss of CTCF in Mouse  
751 Embryonic Stem Cells. *bioRxiv* 118737 (2017) doi:10.1101/118737.
- 752 56. Rodríguez-Carballo, E. *et al.* The HoxD cluster is a dynamic and resilient TAD boundary  
753 controlling the segregation of antagonistic regulatory landscapes. *Genes Dev.* **31**, 2264–2281  
754 (2017).
- 755 57. Barutcu, A. R., Maass, P. G., Lewandowski, J. P., Weiner, C. L. & Rinn, J. L. A TAD boundary is  
756 preserved upon deletion of the CTCF-rich Firre locus. *Nat. Commun.* **9**, (2018).
- 757 58. Di Pierro, M., Cheng, R. R., Lieberman Aiden, E., Wolynes, P. G. & Onuchic, J. N. De novo  
758 prediction of human chromosome structures: Epigenetic marking patterns encode genome  
759 architecture. *Proc. Natl. Acad. Sci.* **114**, 12126–12131 (2017).
- 760 59. Yan, J. *et al.* Histone H3 lysine 4 monomethylation modulates long-range chromatin

- 761 interactions at enhancers. *Cell Res.* **28**, 204–220 (2018).
- 762 60. Kirkpatrick, S., Gelatt, C. D. & Vecchi, M. P. Optimization by simulated annealing. *Science* (80-  
763 ). **220**, 671–680 (1983).
- 764 61. Salamon, P., Sibani, P. & Frost, R. *Facts, conjectures, and improvements for simulated*  
765 *annealing*. *SIAM monographs on mathematical modeling and computation* (2002).  
766 doi:10.1137/1.9780898718300.
- 767 62. Python module for Simulated Annealing optimization.
- 768 63. Knight, P. A. & Ruiz, D. A fast algorithm for matrix balancing. *IMA J. Numer. Anal.* **33**, 1029–  
769 1047 (2013).
- 770 64. Plimpton, S. Fast parallel algorithms for short-range molecular dynamics. *J. Comput. Phys.*  
771 **117**, 1–19 (1995).
- 772 65. Kremer, K. & Grest, G. S. Dynamics of entangled linear polymer melts: A molecular-dynamics  
773 simulation. *J. Chem. Phys.* **92**, 5057–5086 (1990).
- 774 66. Persistence of Vision Pty. Ltd. Persistence of Vision Raytracer. (2004).
- 775 67. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing genomic  
776 features. *Bioinformatics* **26**, 841–842 (2010).
- 777 68. Iannone, F. *et al.* CRESCO ENEA HPC clusters: a working example of a multifabric GPFS  
778 Spectrum Scale layout. in *CRESCO ENEA HPC clusters: a working example of a multifabric*  
779 *GPFS Spectrum Scale layout* 1051–1052 (2019).
- 780

781 **ACKNOWLEDGMENTS**

782 M.N. acknowledges support from the NIH grant ID 1U54DK107977-01 and 1UM1HG011585-01, the  
783 EU H2020 Marie Curie ITN n.813282, CINECA ISCRA ID HP10CYFPS5 and HP10CRTY8P, Einstein BIH  
784 Fellowship Award (EVF-BIH-2016 and 2019), Regione Campania SATIN Project 2018-2020. S.B. and  
785 A.M.C. acknowledge support from the CINECA ISCRA grant ID HP10CCZ4KN. We acknowledge  
786 computer resources from INFN, CINECA, ENEA CRESCO/ENEAGRID<sup>68</sup> and *Scope/ReCAS/Ibisco* at the  
787 University of Naples.

788

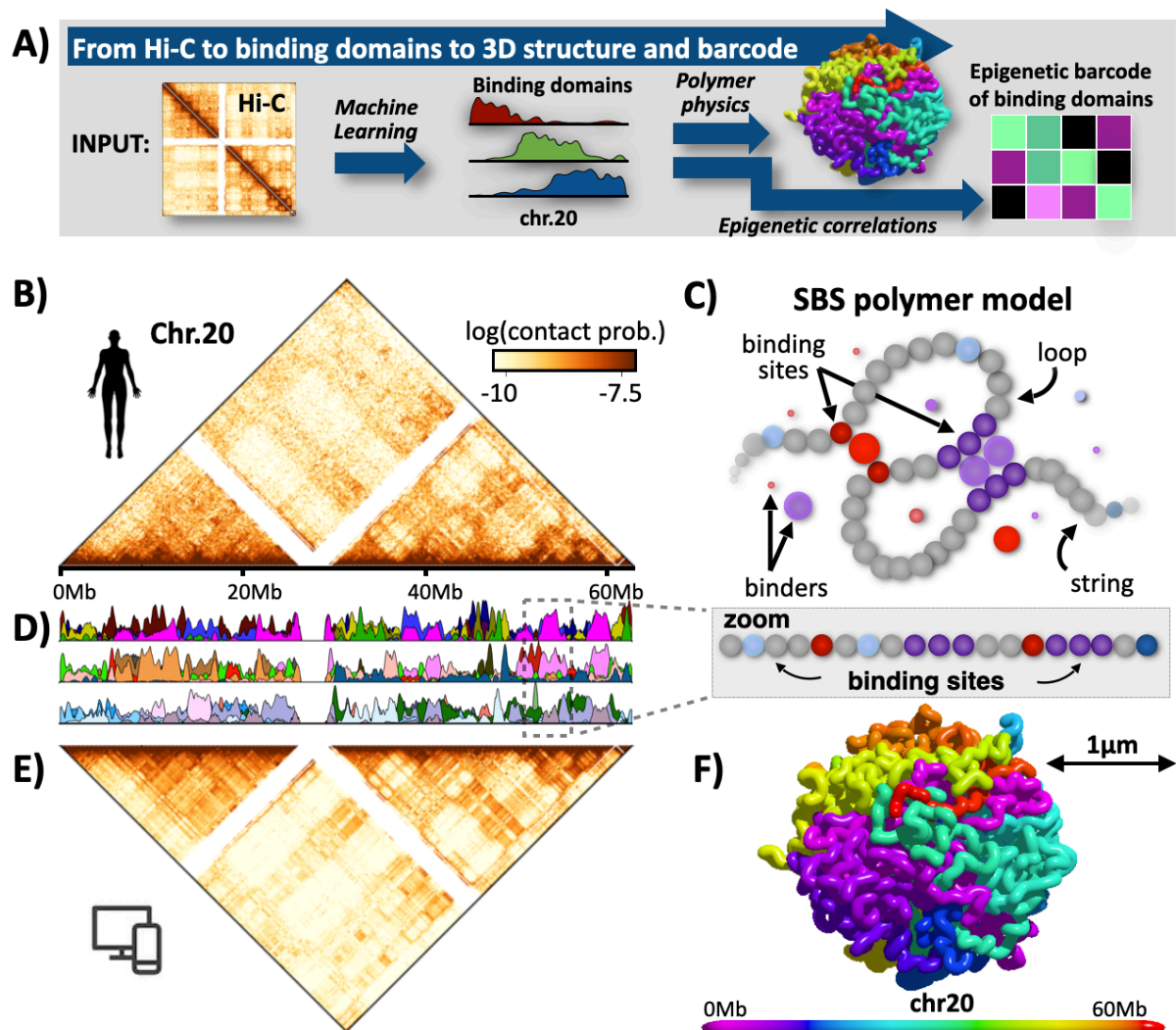
789 **AUTHOR CONTRIBUTIONS**

790 MN designed the project. AE, SB developed modeling. AE, SB, AMC, AA, LF, MC, RC run the  
791 computer simulations and performed analyses. MN, AE, SB, AMC wrote the manuscript.

792 **COMPETING INTERESTS STATEMENT**

793 The authors declare no competing interests.

## FIGURES



**Fig.1 The inferred binding domains explain Hi-C data across chromosomes.**

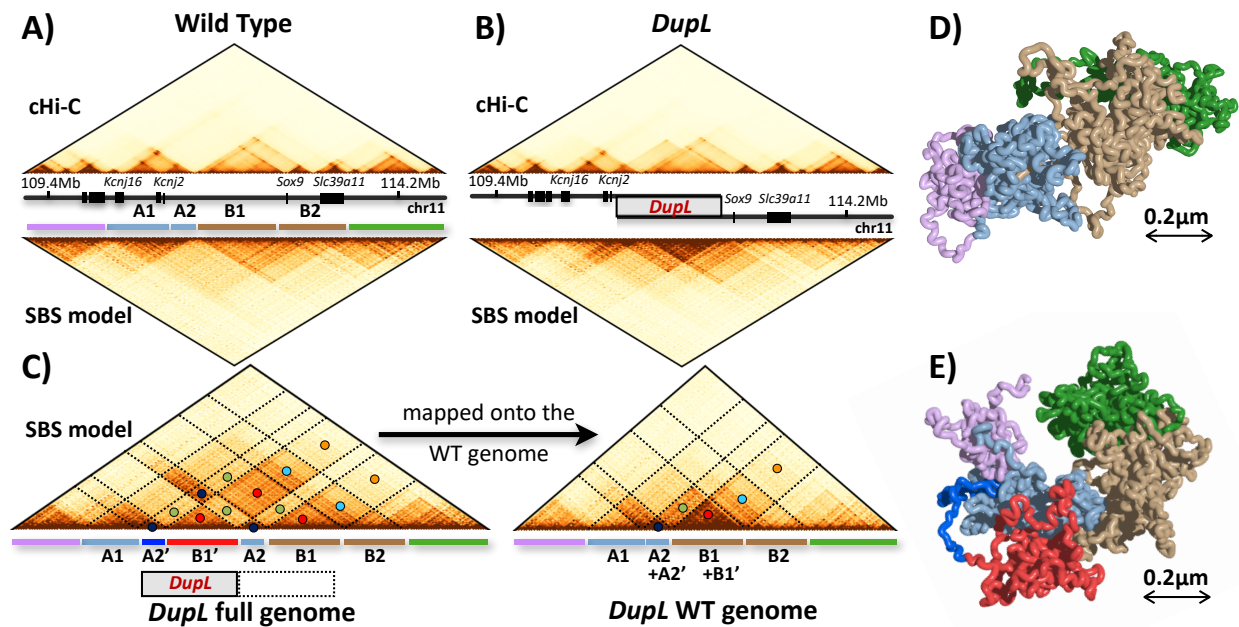
**(A)** Our method combines machine learning and polymer physics to infer from only Hi-C data the genomic location of the minimal set of binding sites required to recapitulate chromatin conformations genome-wide by use of the SBS polymer model of chromatin. Additionally, by correlations with epigenetic data, the inferred binding domains can be assigned a molecular barcode.

**(B)** *In situ* Hi-C data<sup>42</sup> of chromosome 20 at 5kb resolution (log-scale).

**(C)** A scheme of the SBS polymer model of chromatin: it quantifies the scenario where diffusing binders bridge and loop distal cognate binding sites. Each colored bead is a single binding site. The genomic location of the binding sites encodes the 1D information whereby their cognate binders produce the 3D structure via polymer physics.

**(D)** Plots displaying the position and abundance of the different types of binding sites (binding domains) along chromosome 20,

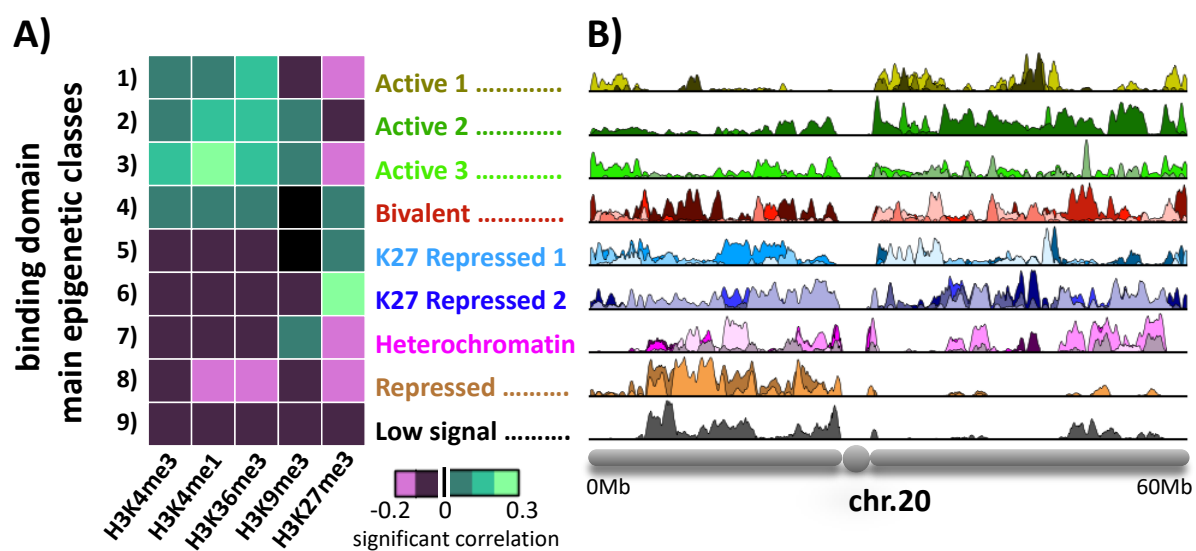
as inferred by our method. For visualization purposes, the different domains, each represented with a different color, are drawn in groups of 10 in different rows. Albeit derived from only Hi-C data, the binding domains have specific correlations each with a set of epigenetic marks, and the colors reflect those associations (see Fig.3). **(E)** The model inferred contact matrix of chromosome 20 has a Pearson, distance-corrected Pearson and stratum adjusted correlation with Hi-C respectively equal to  $r=0.97$ ,  $r'=0.85$ ,  $SCC=0.92$ . Similar results are found across chromosomes (Fig. S1). **(F)** A time snapshot of the 3D structure of the SBS model of chromosome 20.



**Fig.2 The inferred binding domains are validated against mutations at the *Sox9* locus.**

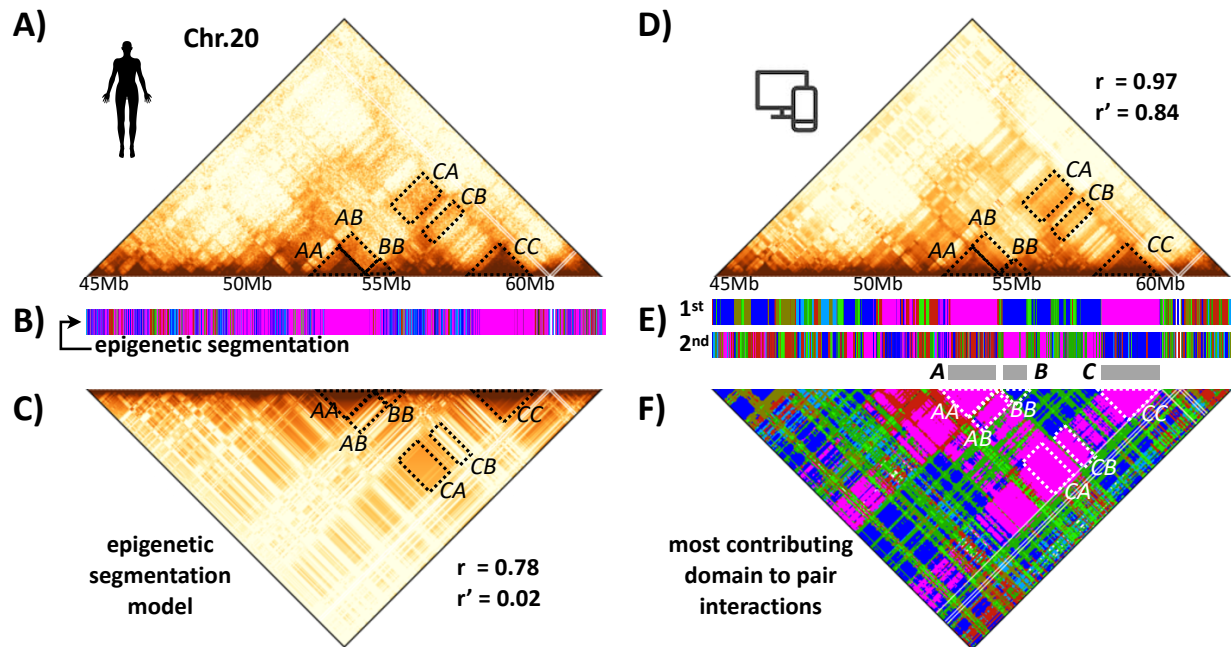
(A) Contact data<sup>43</sup> of the wild type *Sox9* locus from cHi-C experiments in E12.5 limb buds (top) and of the SBS model of the locus in mESC (bottom) have a correlation  $r=0.89$  and  $r'=0.44$ . (B) Based on the WT model, the contact map of a mutant bearing the *DupL* duplication is predicted from only physics (bottom). It has a good correlation ( $r=0.82$ ,  $r'=0.41$ ) with independent *DupL* cHi-C data<sup>43</sup> (top). Model predictions are also validated across the other available *Sox9* mutations (Fig. S3, S4). (C) Mapping the model contacts on the *DupL* full genome clarifies the origin of the associated neo-TAD (red). The colored circles mark corresponding interaction regions as mapped on the WT and *DupL* full genomes. (D)-(E) Snapshots of the model predicted 3D conformation of respectively the WT and *DupL* locus (the color scheme reflects the colored bars in panel A and C) with its neo-TAD. Different mutations result in different 3D structures and distinct enhancer-hijackings, explaining their phenotypes (Fig. S3, S4).





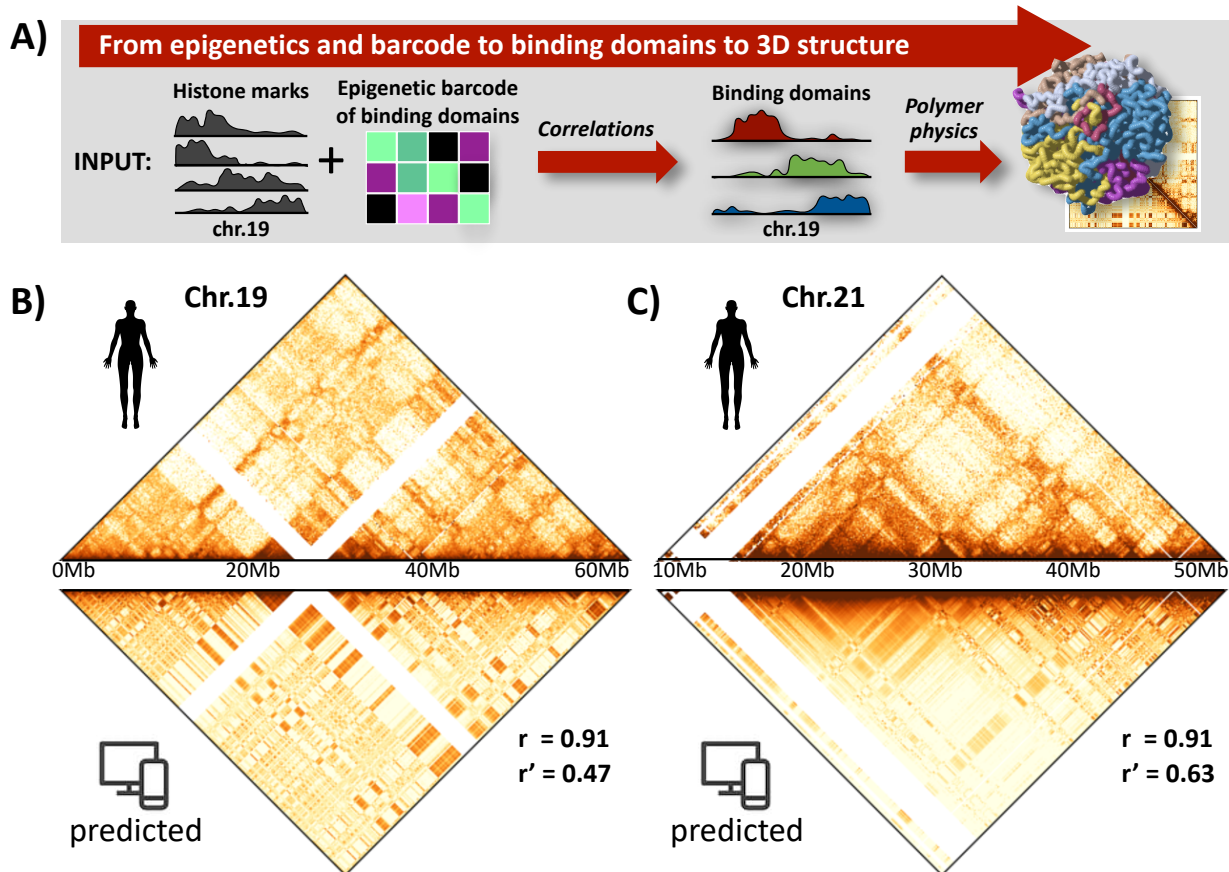
**Fig.3 Epigenetic profiles of the inferred binding domains.**

**(A)** The model binding domains, inferred from Hi-C data only, correlate each with a specific set of epigenetic tracks. They cluster in 9 main classes genome-wide according to their correlations with the shown ENCODE key histone marks (Fig. S6). The epigenetic profile, i.e., the barcode of the centroid of each class is shown in the heat-map. The 9 classes match well chromatin states derived in epigenetic segmentation studies. **(B)** Their abundance along chromosomes is not uniform ( $p$ -value $<0.05$ , Fig.7), as shown here for the binding sites of chromosome 20.



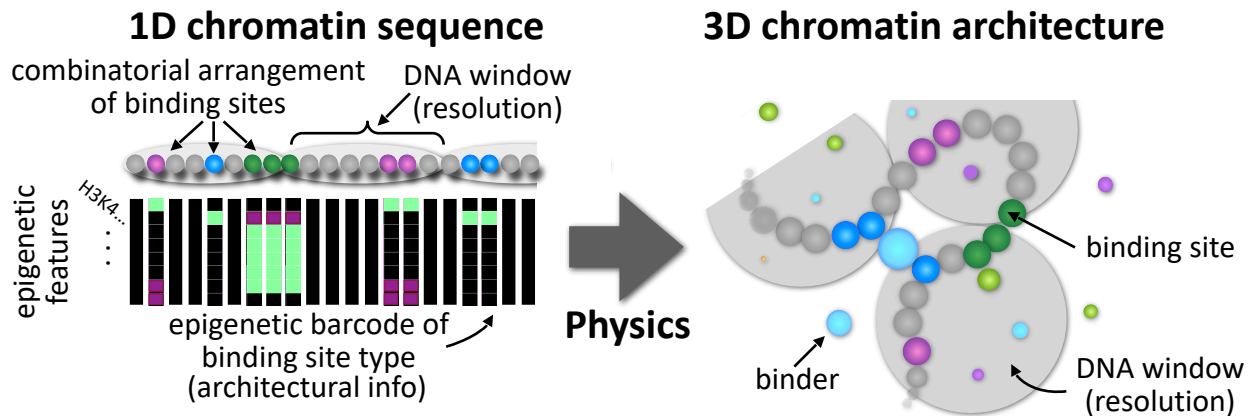
**Fig.4 Chromatin architecture patterns are only partially captured by linear epigenetic segmentation.**

(A) *In situ* Hi-C data<sup>42</sup> (scales as in Fig. 1) of a 20Mb wide region on chr20 in GM12878 and (B) its linear epigenetic segmentation are shown. (C) The contact map of a model based only on homotypic interactions between linear segmented epigenetic domains has a Pearson correlation  $r=0.78$  with Hi-C data. Yet, its distance corrected correlation is much lower,  $r'=0.02$ , returning only a marginal improvement over a control model where each interaction is replaced by the average at the corresponding genomic separation. (D) The contact map of the inferred SBS model of the region has  $r=0.97$  and  $r'=0.84$  with Hi-C data. (E) The PRISMR inferred 1st and 2nd most abundant binding site types of the SBS model of that region are shown. (F) The plot of the SBS most contributing binding domain to each pairwise contact highlights that a combinatorial overlap of different binding site types along the sequence, missing in linear segmentations, is required to capture the complexity and specificity of interaction patterns. For example, interactions (CC) within the TAD in region C are mainly related to binding domains in class 7 (magenta), the most abundant one in C. A and C also interact mainly through class 7, the most abundant in A too. Yet, region B, where class 6 (dark blue) is the most abundant, interacts with C mainly through class 7, its 2<sup>nd</sup> most abundant. Analogously, contacts between A and B originate from different overlapping binding domains in those regions.



**Fig.5 The epigenetic barcode of binding domains predicts chromatin contacts.**

**(A)** In a reverse approach, we correlate the epigenetic profiles of binding domains from even chromosomes with epigenetic signals from odd chromosomes to identify the binding sites of the latter. Next we use the SBS polymer model to predict 3D structures and contact matrices of odd chromosomes to be compared against independent Hi-C data. **(B)** Top: *in situ* Hi-C data<sup>42</sup> (scales as in Fig. 1) of chromosome 19 in GM12878. Bottom: the predicted contact matrix has a correlation, a distance-corrected correlation and a stratum adjusted correlation with Hi-C respectively equal to  $r=0.91$ ,  $r'=0.47$  and  $SCC=0.65$ . **(C)** Top: Hi-C data of chromosome 21. Bottom: the predicted contact matrix has correlations with Hi-C equal to  $r=0.91$ ,  $r'=0.63$  and  $SCC=0.50$ .



**Fig.6 Chromatin 3D architectural information is encrypted in a combinatorial 1D arrangement of epigenetic barcoded sites.**

Our approach infers, from Hi-C data only, the minimal set of binding sites along the 1D genome sequence (left) required to produce, via polymer physics (e.g., interactions with diffusing cognate binding molecules), 3D structures (right) consistent with Hi-C contacts. The inferred binding sites are barcoded by specific epigenetic marks (vertical bars) and fall into epigenetic classes (bead color) well matching functional chromatin states known from linear segmentation studies. However, they have a genomic overlapping, combinatorial organization, lacking in epigenetic segmentations, necessary to explain Hi-C contacts with high accuracy genome-wide. Their epigenetic barcode was shown to predict *de novo* chromatin conformations, e.g., after genetic or epigenetic variations, showing that the inferred combinatorial 1D arrangement of binding sites carry accurate, specific 3D architectural information genome-wide.