

1 **Fine-mapping of nuclear compartments using ultra-deep Hi-C shows that active promoter and**
2 **enhancer elements localize in the active A compartment even when adjacent sequences do not**
3

4 Huiya Gu^{1,12}, Hannah Harris^{2,12}, Moshe Olshansky³, Kiana Mohajeri⁴, Yossi Eliaz¹, Sungjae Kim⁵, Akshay
5 Krishna², Achyuth Kalluchi², Mozes Jacobs⁶, Gesine Cauer⁶, Melanie Pham¹, Suhas Rao¹, Olga
6 Dudchenko¹, Michael H Nichols⁷, Eric S. Davis⁸, Devika Udupa², Victor G. Corces⁷, Douglas H. Phanstiel^{8,9},
7 William Stafford Noble⁶, Jeong-Sun Seo⁵, Michael E. Talkowski^{4,10,11}, Erez Lieberman Aiden^{1*}, and M.
8 Jordan Rowley^{2*}

9 1. Center for Genome Architecture, Department of Molecular and Human Genetics, Baylor College of
10 Medicine, Houston, TX, USA. Center for Theoretical Biological Physics, Rice University, Houston, TX, USA.
11 2. Department of Genetics, Cell Biology and Anatomy, University of Nebraska Medical Center, Omaha,
12 NE, USA.
13 3. Computational Biology and Clinical Informatics, Baker Heart and Diabetes Institute, Melbourne,
14 Victoria, Australia.
15 4. Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA
16 5. Precision Medicine Institute, Seoul, 08511, Republic of Korea.
17 6. Department of Genome Science, University of Washington, Seattle, USA; Paul G. Allen School of
18 Computer science & Engineering, University of Washington, Seattle, USA.
19 7. Department of Human Genetics, Emory University School of Medicine, Atlanta, GA, USA.
20 8. Curriculum in Bioinformatics and Computational Biology, University of North Carolina at Chapel Hill,
21 Chapel Hill, NC, USA.
22 9. Thurston Arthritis Research Center, University of North Carolina, Chapel Hill, NC, USA; Department of
23 Cell Biology and Physiology, University of North Carolina, Chapel Hill, NC, USA.
24 10. Department of Neurology, Harvard Medical School, Boston, MA, USA.
25 11. Program in Medical Population Genetics and Stanley Center for Psychiatric Research, Broad Institute
26 of MIT and Harvard, Cambridge, MA, USA.
27 12. These authors contributed equally to this work.
28
29

30 **Running Title:** Sub-genic discordant compartments

31 ***co-corresponding authors:** ELA: erez@erez.com; MJR: jordan.rowley@unmc.edu

32 **Keywords:** cohesin, CTCF, enhancers, extrusion, nucleus, transcription

33 Abstract

34 Megabase-scale intervals of active, gene-rich and inactive, gene-poor chromatin are known to
35 segregate, forming the A and B compartments. Fine mapping of the contents of these A and B
36 compartments has been hitherto impossible, owing to the extraordinary sequencing depths required to
37 distinguish between the long-range contact patterns of individual loci, and to the computational
38 complexity of the associated calculations. Here, we generate the largest published *in situ* Hi-C map to
39 date, spanning 33 billion contacts. We also develop a computational method, dubbed PCA of Sparse,
40 SUper Massive Matrices (POSSUMM), that is capable of efficiently calculating eigenvectors for sparse
41 matrices with millions of rows and columns. Applying POSSUMM to our Hi-C dataset makes it possible to
42 assign loci to the A and B compartment at 500 bp resolution. We find that loci frequently alternate
43 between compartments as one moves along the contour of the genome, such that the median
44 compartment interval is only 12.5 kb long. Contrary to the findings in coarse-resolution compartment
45 profiles, we find that individual genes are not uniformly positioned in either the A compartment or the B
46 compartment. Instead, essentially all (95%) active gene promoters localize in the A compartment, but
47 the likelihood of localizing in the A compartment declines along the body of active genes, such that the
48 transcriptional termini of long genes (>60 kb) tend to localize in the B compartment. Similarly, nearly all
49 active enhancers elements (95%) localize in the A compartment, even when the flanking sequences are
50 comprised entirely of inactive chromatin and localize in the B compartment. These results are consistent
51 with a model in which DNA-bound regulatory complexes give rise to phase separation at the scale of
52 individual DNA elements.

53

54 Main

55 The nucleus of the human genome is partitioned into distinct spatial compartments, such that
56 stretches of active chromatin tend to lie in one compartment, called the A compartment, and stretches
57 of inactive chromatin tends to lie in the other, called the B compartment¹. Compartmentalization was
58 first identified using Hi-C, a method that relies on DNA-DNA proximity ligation to create maps reflecting
59 the spatial arrangement of the genome¹. Loci in the same spatial compartment exhibit relatively
60 frequent contacts in a Hi-C map, even when they lie far apart along a chromosome, or on entirely
61 different chromosomes^{1,2}. Accurate classification of the resulting genome-wide contact patterns
62 requires a large number of contacts to be characterized at each locus. As such, genome-wide
63 compartment profiles have only been generated, in the past, at resolutions ranging from 40 kb – 1 Mb¹⁻
64 ³. Moreover, extant compartment detection algorithms require operations, such as calculation of
65 principal eigenvectors¹, which are computationally intractable when the underlying matrices have
66 millions of rows and columns – such as high-resolution Hi-C matrices.

67 Although the compartments as a whole are often thought to form as a consequence of phase
68 separation³⁻⁶, the low resolution of compartment profiles has made it difficult to determine the protein
69 mechanisms that underlie this process.

70 Here, we construct an *in situ* Hi-C map in lymphoblastoid cells spanning 42 billion read-pairs and
71 33 billion contacts. This map contains an average of 66,000 contacts for every kilobase of genome
72 sequence. We combine this map with a novel algorithm, dubbed POSSUMM, which greatly accelerates
73 the calculation of the principal eigenvector and the largest eigenvalues of a massive, sparse matrix. This
74 makes it possible to, e.g., calculate the principal eigenvector for correlation matrices containing millions
75 of rows, and billions of nonzero entries. Combining our ultra-deep map with POSSUMM, we find that it
76 is possible to map the contents of the A and B compartments with 500 bp resolution, a 100-fold

77 improvement in resolution. We also show that when we classify loops based on their appearance, at fine
78 resolution, in our ultra-deep map, it becomes possible to distinguish between loops that form by
79 extrusion and those that form via non-extrusion mechanisms.

80

81 **Generation of an ultra-deep in situ Hi-C map in lymphoblastoid cells spanning 33 billion contacts**

82 We produced an ultra-deep Hi-C map using lymphoblastoid cells from a panel of 17 individuals,
83 obtaining over 42 billion PE150 read-pairs. This map was generated by aggregating the results of over
84 150 individual Hi-C experiments. In order to enhance the resolution of the maps, we used a variety of 4-
85 cutter restriction enzymes in the different experiments, thus enhancing the density of cut sites across
86 the genome. Together, these experiments yielded 33 billion contacts after alignment, deduplication, and
87 quality filtering (Table S1). The resulting dataset is far deeper than any prior published Hi-C map. By
88 comparison, the average published Hi-C map contains roughly 300 million contacts; 93% of Hi-C maps in
89 the 4DNucleome database⁷ have less than 1 billion contacts (Fig. S1A, Table S2); and the widely used
90 lymphoblastoid Hi-C map generated in Rao et al. contains 4.9 billion contacts.

91 We generated contact matrices at a series of resolutions as fine as 500 bp. These matrices
92 greatly improved the resolution of all features genome-wide, revealing many additional loops and
93 domains (Fig. 1A). This high coverage also enhanced the long-range plaid pattern indicative of
94 compartments (Fig. 1B, S1B), as well as the corresponding compartment domains observed along the
95 diagonal of the map (Fig. 1C, S1C). Critically, because the number of contacts at every locus was greatly
96 increased, with an average of 66,000 contacts incident on each kilobase of the human genome (Fig. 1B,
97 S1B), we were able to distinguish between loci in the A compartment and loci in the B compartment
98 with much finer resolution.

99

100 **Development of PCA of Sparse, SUPER Massive Matrices (POSSUMM) and its use to create a genome- 101 wide compartment profile with 500bp resolution.**

102 Extant methods for classifying loci into one compartment or the other typically rely on
103 numerical linear algebra to calculate the principal eigenvector (called, in this context, “the A/B
104 compartment eigenvector”) and the smallest eigenvalues of correlation matrices associated with the Hi-
105 C contact matrix. At 100 kb resolution, these matrices typically have thousands of rows and columns and
106 millions of entries, making them tractable using extant numerical algorithms, such as those
107 implemented by Homer⁸, Juicer⁹, and Cooler¹⁰. However, at kilobase resolution or beyond, these
108 matrices have hundreds of thousands of rows and hundreds of billions of entries, making them
109 intractable using the aforementioned tools. For example, computing an eigenvector for chr1 at 500 bp
110 resolution entails generating a matrix with 250 billion entries and performing a calculation that is
111 projected to require >4.6 TB of RAM for >16 years (Fig. S1D).

112 As such, we developed a method, POSSUMM, for calculating the principal eigenvector and the
113 smallest eigenvalues of a matrix. POSSUMM is based on the power method, which repeatedly multiplies
114 a matrix with itself in order to calculate the principal eigenvector (Fig. 1D). However, POSSUMM does not
115 explicitly calculate all of the intermediate matrices required by the power method. Instead, it explicitly
116 calculates only the tiny subset of intermediate values required to obtain the principal eigenvector itself,
117 not requiring dense matrices, which makes it vastly more efficient (Fig. 1D, Fig. S1EF).

118 Using POSSUMM, we assigned loci to the A and B compartment at resolutions up to, and
119 including, 500 bp (Fig. 1B). The calculation of the A/B compartment eigenvector at 500 bp resolution
120 took only 12 minutes, and 13 GB of RAM (Fig. S1D&G). A and B compartments identified by POSSUM
121 accurately detect the segregation of active from inactive chromatin (Fig. S1H-K).

122

123 **The median compartment interval is 12.5 kb long**

124 It is widely thought that compartment intervals (genomic intervals that lie entirely in one
125 compartments) are typically megabases in length and are partitioned into numerous punctate loops and
126 loop domains^{6,11-13}. To explore this phenomenon, we used our fine map of nuclear compartments to
127 examine the frequency with which loci alternate from one compartment to the other. Nearly 99% of
128 compartment intervals were less than 1 Mb in size, and 95% were smaller than 100 kb (Fig. 2A). The
129 median compartment interval was only 12.5 kb, and thousands of compartment intervals were no
130 longer than 5 kb (Fig. S1L). In comparison, the median size of CTCF loops in our map was 360 kb in
131 length, demonstrating that compartment intervals are smaller than individual loops.

132

133 **Kilobase-scale compartment intervals frequently give rise to contact domains**

134 It is well known that long compartment intervals often give rise to contact domains, i.e.,
135 genomic intervals in which all pairs of loci exhibit an enhanced frequency of contact among
136 themselves^{6,14-17} (Fig. 1C). Such contact domains are referred to as compartment domains. We found
137 that even short compartment intervals less than 5 kb frequently give rise to contact domains (Fig. S1M),
138 demonstrating that intervals of chromatin in the same compartment possess the ability to form contact
139 domains regardless of scale.

140

141 **Essentially all active promoter and enhancer elements localize in the A compartment**

142 Next, we compared our fine map of nuclear compartments to ENCODE's catalog of regulatory
143 elements in GM12878 cells. We examined active promoters (defined as 500 bp near the TSS, absence of
144 repressive marks H3K27me3 or H3K9me3, and with ≥ 1 RPKM gene expression in RNA-seq) and found
145 that nearly all lie in the A compartment: out of 9,324 active promoters annotated in GM12878, only 496
146 (5%) were assigned to the B compartment (Fig. 2B - top). We noticed that active promoters in the B
147 compartment had higher values in the principal eigenvector compared to the surrounding regions (Fig.
148 S1N). Indeed, if we use a slightly more stringent threshold (assigning promoters to the B compartment
149 only if the corresponding entry of the principal eigenvector is $<-.001$), we find that only 233 (2.5%) of
150 promoters are assigned to the B compartment. Notably, when 1 Mb resolution compartment profiles
151 are used, the number of active promoters assigned to the B compartment increases 4-fold, to ~21% (Fig.
152 S1O). This is at least in part because the use of coarse resolutions leads to the averaging of interaction
153 profiles from neighboring loci, such that a DNA element in the A compartment might be erroneously
154 assigned to the B compartment if most of the flanking sequence was inactive (Fig. 2C, S2A-G).

155 Similarly, we found that essentially all active proximal enhancers (defined by annotation in
156 DenDB¹⁸, ≤ 10 kb from a TSS, and overlapping H3K27ac but not H3K27me3 & H3K9me3¹⁹) lie in the A
157 compartment (Fig. 2B - middle). Moreover, essentially all active distal enhancers (DenDB¹⁸, >10 kb from
158 a TSS, with H3K27ac, but not H3K27me3 or H3K9me3¹⁹) lie in the A compartment (Fig. 2B - bottom): out
159 of 30,868 active distal enhancers annotated in GM12878, only 1,607 (5%) were assigned to the B

160 compartment. Many of these distal enhancer elements represent small islands of A-compartment
161 chromatin in a sea of inactive, B compartment chromatin (Fig. 2D). This demonstrates that individual
162 DNA elements can escape a neighborhood that is overwhelmingly associated with one compartment in
163 order to localize with a different compartment (g. 1C-E, S2H-I). When 1 Mb resolution compartment
164 profiles are used, the number of active distal enhancers assigned to the B compartment increases 4.6-
165 fold, to 23% (Fig. S2J). Again, this is at least in part because the use of coarse resolutions leads to the
166 averaging of interaction profiles from neighboring loci (Fig. S2H&K). Taken together, we find that
167 essentially all active regulatory elements, including both promoters and enhancers, lie in the A
168 compartment, even when immediately neighboring sequences do not.

169

170 **Many genes exhibit discordant compartmentalization, with the TSS in the A compartment and the TTS** 171 **in the B compartment**

172 When exploring the fine map of nuclear compartmentalization, we noticed many genes where
173 the TSS and TTS localize to opposite compartments (Fig. 3A., see also Fig. 1C,2C,2E). These intra-genic
174 compartmental switches are more easily seen at large genes (Fig. 3B-E, S3AB). We therefore asked if
175 gene size can affect the compartment localization of the TTS. Indeed, average profiles of compartmental
176 status revealed that TSSs were most likely to be in the A compartment (Fig. 3F), but that the likelihood
177 of lying in the A compartment decreases steadily as one examines increasingly distal portions of the
178 gene body, such that the TTSs of large genes are more likely to localize to the B compartment (Fig.
179 3F&G, S3C). This was especially evident if we consider very large genes (Fig. 3E&H, S3D), where the TSS
180 was overwhelmingly in the A compartment, but the TTS was usually in the B compartment.

181 We next asked if genes with discordant compartmentalization (i.e., the TSS was in compartment
182 A, but the TTS was in compartment B) could be explained by different chromatin marks at the TSS vs.
183 TTS. We examined chromatin marks at the TTS in active genes larger than 20 kb and found that
184 diminished levels of active marks at the TTS, specifically RNAPII, H3K4me1, and H3K36me3, were
185 correlated with presence of discordant compartmentalization (Fig. 3H, Fig. S3E). Notably, although
186 repressive chromatin marks are frequently seen at loci in the B compartment, genes with discordant
187 compartmentalization typically lacked such marks at the TTS (Fig. 3H, S3E). We found that chromatin
188 marks at the TSS were not predictive of whether the gene exhibited discordant compartmentalization
189 (Fig. S3E&F).

190 Finally, we sought to determine if discordant compartmentalization was associated with
191 transcriptional pausing as measured by GRO-Seq. We found that elongating genes longer than 20 kb
192 were more likely to exhibit concordant compartmentalization (Fig. 3I), whereas paused genes were
193 more likely to exhibit discordant compartmentalization (Fig. 3J).

194 Taken together, these data support a model where an active TSS localizes to the A compartment
195 but brings with it only a small portion of the gene body, depending on the elongation status (Fig. 3K).

196

197 **Loop extrusion forms diffuse loops, whereas compartmentalization forms punctate loops**

198 We examined loops in our Hi-C dataset. Using SIP²⁰ and HiCCUPS², we identified 32,970 loops.
199 Ninety-one percent of these loops contained a CTCF-bound motif at both anchors, with a strong
200 preference for the convergent orientation (Fig. S4A).

201 Interestingly, when we examined loops at 1 kb resolution, we noticed that the signal is diffuse
202 (Fig. 4A, S4B), indicative of frequent contacts proximal to the CTCF binding sites (Fig. 4B). The elevated
203 contact frequency decays as the distance from the corresponding anchors increases (Fig. 4C, rainbow) (a
204 loss of signal of c.a. -6% from one bin to the next; i.e. -6%/kb compounding). Curiously, this decay rate is
205 much slower than the decay rate reflected by the diagonal of the Hi-C map (Fig. 4C, S4C – expected) (c.a.
206 -28%/kb), which is thought to reflect the properties of the chromatin polymer. The decay was
207 unchanged as a function of loop size or sequencing depth (Fig. S4DE).

208 We wondered whether this slow decay in contact frequency was seen for loops in other species.
209 We therefore examined hundreds of loops observed in a published high-resolution Hi-C map from
210 *Drosophila melanogaster* Kc167 cells at 1 kb resolution^{14,21} (Fig. 4D&E). Interestingly, the loops in
211 *Drosophila* decayed at a rate (c.a. -20%/kb) that matched the diagonal of the *Drosophila* Hi-C map (c.a. -
212 23%/kb) and was much faster than the rate seen for human CTCF-mediated loops (Fig. 4F). This suggests
213 that CTCF loops create interactions between sequences bound by CTCF, as well as interactions between
214 CTCF bound and adjacent sequences. However, in *Drosophila*, Pc loops only create interactions directly
215 between the Pc bound sequences.

216 Finally, we examined loops previously identified in *C. elegans*^{20,22,23}. The loop decay was slower
217 (c.a. -11%/kb) than the decay seen at the diagonal (c.a. -24%/kb) (Fig. 4F, green vs. grey), and was more
218 similar to the rate of decay seen for human CTCF-mediated loops than the one observed for *D.*
219 *melanogaster* loops (Fig. 4F, Fig. S4I).

220 It was notable that the type of decay observed (fast or slow) matched the putative mechanism
221 by which the loops formed. CTCF-mediated loops in human are bound by, and dependent on, the SMC
222 complex cohesin (Fig. S4H), and form by cohesin-mediated extrusion²⁴⁻²⁷. Similarly, the loops in *C.*
223 *elegans* are bound by the SMC complex condensin and we previously suggested that they are formed by
224 condensin-mediated loop extrusion^{20,22,23}. Indeed, the interactions between loop-adjacent sequences
225 are in further support of loop formation by extrusion in *C. elegans*. By contrast, *Drosophila* loops are
226 much less likely to be bound by CTCF, cohesin, condensin, or other extrusion-associated proteins¹⁴.
227 Instead, they are bound by the Polycomb complex, *Pc*, and may form by means other than extrusion²⁸⁻³⁰.

228 These findings suggest that the mechanism of loop formation influences whether loops will be
229 punctate or diffuse, with extrusion-mediated loops forming diffuse peaks and compartmentalization-
230 mediated loops forming more punctate features.

231

232 Diffuse loops enhance the contact frequency of nearby promoter-enhancer interactions

233 Using Fit-Hi-C³¹, we called promoter-enhancer interactions at 1 kb resolution on human chr1.
234 We examined those interactions where both the promoter and enhancer lie within 100 kb of a loop
235 anchor. In some cases, these interactions lie completely inside the loop, but in others they cross the
236 loop anchor. Both cases exhibited strongly enriched contact frequency as compared to enhancer-
237 promoter interactions that are unrelated to CTCF loops, i.e., near permuted random sites (Fig. 4G).
238 These data suggest that CTCF loops enhance the contact frequency of promoter-enhancer interactions,
239 even when both elements lie outside the loop (Fig. 4H). By contrast, in *Drosophila*, Fit-Hi-C interactions
240 between promoters and enhancers tend to be much shorter (Fig. S4J).

241

242 Deletion of CTCF's RNA binding domains leads to more punctate loops

243 Interestingly, we observed some variability in the decay rate for different loops (Fig. S4K). This
244 decay did not correlate strongly with either CTCF motif strength, CTCF ChIP-seq peak strength, or Rad21
245 ChIP-seq peak strength (Fig. S4L-O). Instead, we found that CTCF-mediated loops exhibiting slower decay
246 are associated with higher levels of transcription (Fig. 4I) and chromatin accessibility (Fig. S4P) near the
247 loop anchors. This suggests that nearby transcriptional activity could impact how CTCF interacts with the
248 nearby sequences and / or with the loop extrusion process.

249 The CTCF protein contains 11 zinc finger domains. Recently, it was shown that ZF1 and ZF10 bind
250 to RNA, and that deletion of these two domains causes weakening of loops throughout the genome³².
251 We performed aggregate peak analysis on the published Hi-C in ZF1 and ZF10 mutants³² using “bullseye”
252 plots in order to explore the effect of these deletions on loop decay. Interestingly, we found that loops
253 appeared more punctate in both CTCF RNA binding mutants (Fig. 4J). This effect was especially
254 pronounced in the ZF1 mutant.

255 Taken together, these findings are consistent with a model where CTCF’s RNA-binding domains
256 and the presence of bound RNAs results in more diffuse contacts between loop anchors, and thus to
257 enriched contacts among regulatory elements near the loop.

258

259 Discussion

260 By generating a Hi-C map with extraordinary sequencing depth (33 billion PE, or 9.9 terabases of
261 uniquely mapped sequence), we create the first fine-map of nuclear compartmentalization.

262 Our findings demonstrate that compartment intervals and compartment domains can be far
263 smaller than previously appreciated. This contrasts with the common hierarchical model of chromosome
264 organization in which compartments are partitioned into TADs and loops^{6,11-13}. In fact, our results
265 indicate that compartment intervals can be so small that active DNA elements will localize with the A
266 compartment even when surrounded by inactive chromatin localizing in the B compartment (Fig. 5).

267 Strikingly, we find that essentially all distal enhancer elements lie in the A compartment. This
268 contrasts with earlier work, using coarse-resolution maps of compartmentalization, which only report
269 general enrichment of active distal enhancers in the A compartment, rather than as an absolute
270 characteristic of active enhancers^{33,34}. Similarly, many previous studies have reported a coarse
271 enrichment of active genes in the A compartment⁶, yet we find that essentially all active promoters lie in
272 the A compartment.

273 We also observe that the likelihood that a locus lies inside the A compartment declines as one
274 moves away from the promoter, along the gene body. Interestingly, we observe numerous genes with
275 discordant compartmentalization, where the TSS and TTS tend to be in different compartments. This
276 observation suggests that opposing compartments need not correspond to widely separated locations
277 within the nucleus. For instance, recent work indicates that compartments could be phase-separated
278 droplets³⁵, suggesting that the TSS and TTS of a gene with discordant compartmentalization might be
279 physically proximal within the nucleus, in neighboring A and B droplets (Fig. 5).

280 The finding that active promoters – specifically, active TSSs – are overwhelmingly localized in the
281 A compartments; that TTS compartment status correlates with RNAPII levels at the TTS; and that genes
282 with discordant compartmentalization tend to be transcriptionally paused is consistent with a model in
283 which RNAPII drives localization to the A compartment. Although a recent RNAPII degradation study
284 showed little effect on genome organization, these experiments did not achieve the sequencing depth

285 required to perform fine mapping of nuclear compartmentalization, nor to resolve phenomena such as
286 genes with discordant compartmentalization³⁶. Alternatively, other components of the transcription
287 complex that travel along the gene body during transcription elongation may be responsible for
288 mediating interactions that assign sequences to the A compartment. In future studies, it will be of great
289 interest to examine how RNAPII and other components of the transcription complex impact genome
290 organization at the TSS and TTS separately.

291 We note that our data represent averages within the cellular population, and it is unclear where
292 each component lies during the transcriptional process itself. In the future, fine mapping of nuclear
293 compartments in single cells will be needed to decipher these dynamics. Moreover, we note that our
294 study did not attempt to study subcompartments or models with ≥ 3 distinct compartment states^{2,37},
295 which will be an important topic for future analyses.

296 Our ultra-deep Hi-C map also helped identify interesting properties of chromatin loops. In
297 particular, we observe that CTCF-mediated loops are highly diffuse, more so than would be predicted
298 based on polymer behavior alone (Fig. 5). Interestingly, this diffusivity is observed for loops that form by
299 extrusion, such as loops in human^{2,24-27} and *C. elegans*^{20,22,23}, but is not observed for loops that are
300 believed to form by compartmentalization, such as the numerous *Pc*-associated loops observed in
301 *Drosophila*^{14,21,29,30}. Intriguingly, variations in diffusivity between different loops could explain
302 differences in domains signal (See Supplemental Discussion, Fig. S5).

303 *In vitro* studies have found that large chromatin complexes can impede looping factors^{38,39}, and
304 cohesin was shown to build up near transcriptionally active regions⁴⁰. Yet studies have also reported
305 independence of CTCF loops and transcription^{36,41,42}, bringing the relationship between transcription and
306 CTCF looping in question. Recently, it was shown that CTCF RNA-binding domains, ZF1 and ZF10, are
307 important for looping³². Our finding that loop-decay is altered in CTCF RNA-binding mutants supports
308 the argument that transcription can impact fine-scale chromatin organization in mammals, as does the
309 correlation between TTS compartmental domains and elongation status.

310 Our POSSUMM method, a novel numerical linear algebra algorithm for calculating principal
311 eigenvectors, is now available as part of the Juicer pipeline for Hi-C analysis. Our power analyses suggest
312 that fine mapping of nuclear compartments at sub-kilobase resolution becomes possible for maps
313 containing 7 billion contacts or more (See Supplemental Discussion, Fig. S6&S7). As sequencing costs
314 continue to decline, we expect that fine mapping of nuclear compartments will become increasingly
315 common.

316

317 **Methods**

318 *Library Preparation, Initial Processing, and Quality Metrics*

319 Hi-C libraries were prepared according to the published *in-situ* method². The full map represents a
320 mixture of libraries prepared by digestion of various 4-cutter restriction enzymes, MboI, MseI, and NlaIII.
321 Reads were aligned to the hg19 genome, processed, Knight-Ruiz (KR) normalized using Juicer⁹.
322 Subsampled Hi-C maps were created by uniform random selection of read-pairs from the 33.3 billion Hi-
323 C dataset. We provide a script for subsampling Hi-C data at <https://github.com/JRowleyLab/HiCSampler>.

324 *Compartment Analysis*

325 Compartments were identified using the A/B eigenvector of the Hi-C matrix using POSSUM. POSSUMM
326 can be downloaded from: <https://github.com/aidenlab/EigenVector> and is also now implemented in the
327 ENCODE version of the Juicer pipeline: <https://github.com/ENCODE-DCC/hic-pipeline>.

328 **Introduction to PCA of Sparse, SUper Massive Matrices (POSSUM)**

329 We note that the so-called “A/B compartment eigenvector” is simply the eigenvector of A corresponding
330 to its largest eigenvalue, where X is given by the Hi-C contact matrix. This is equivalent to the first
331 principal component in Principal Component Analysis. We note that in our case, X is a large, sparse
332 matrix, containing millions of rows, millions of columns, and tens of billions of nonzero entries (dubbed
333 a “Sparse, SUper Massive Matrix”).

334 Suppose we seek to calculate the largest eigenpairs, λ_i, v_i of A in this case. Although X is sparse, we note
335 that both Y and A are dense matrices. Unfortunately, storing dense matrices with millions of rows and
336 columns in memory is impossible. Hence we cannot use any method for calculating the eigenvectors of
337 A that would require us to explicitly calculate either Y or A. Similarly, traditional sparse matrix methods
338 for eigendecomposition are not usable here, again because A - the correlation matrix we hope to
339 analyze - is a dense matrix.

340 Therefore, in order to calculate eigenvectors for A, we began by implementing a method that makes it
341 possible to calculate the matrix-vector product Av (where v is an arbitrary vector) using a sparse
342 representation of X, i.e., without explicitly computing either A or Y. See POSSUMM details below for a
343 more complete description.

344 Next, we note that there are many methods for calculating eigenvectors in which the input matrix only
345 appears via a matrix-vector product. These include the Power method, the Lanczos method, and their
346 many variants⁴³. Thus, in principle, any of these methods - for which there are many implementations in
347 Fortran, C, C++, Matlab, and R - can be combined with the sparse Av product calculation described
348 above in order to calculate eigenpairs of A. In practice, methods combining these two approaches are
349 not available.

350 To the best of our knowledge, the sole exception is a method in the R package *irlba*, which was released
351 while this study was being performed. The details of this method are unpublished, but the method itself
352 is available at <https://cran.r-project.org/web/packages/irlba/index.html>. However, *irlba* cannot handle
353 cases where X has more than roughly two billion nonzero entries, which is exceeded in the present case.
354 It also does not enable parallelization, which limits performance in highly demanding settings.

355 POSSUMM combines sparse Av product calculation with the power method, is extremely memory-
356 efficient, and enables parallelization via multi-threading.

357 **POSSUMM Details.**

358 To identify compartments from sparse Hi-C matrices, we began by excluding all rows and columns with 0
359 variance. Let X be a matrix with column vectors $X^{(1)}, \dots, X^{(n)}$. Let $Y^{(i)} = (X^{(i)} - c_i)/\sigma_i$ $1 \leq i \leq n$,
360 where c_i is the mean of X_i and σ_i is its standard deviation. Let $Y = (Y^{(1)}, \dots, Y^{(n)})$ be an $n \times n$ matrix with
361 column vectors. The correlation matrix of X is $A = Y^T Y$ where Y^T is transposed Y. Since A is symmetric
362 and positive semi-definite it has n real eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ and n eigenvectors.
363 v_1, \dots, v_n where $A v_i = \lambda_i v_i$.

364 These eigenvectors are a basis of R^n (i.e., a set of vectors which are independent and span the space) if
365 $\lambda_i \neq \lambda_j$ and $v_i \perp v_j$ (i.e., $v_i^T v_j = 0$). To compute v_1 using the power method (a.k.a power iterations),

366 suppose that $\lambda_1 > \lambda_2$ and let x_0 be any nonzero vector in R^n , we define the recursive relation:
367 $x_{k+1} = Ax_k = A^{k+1}x_0$. We can represent x_0 as $x_0 = a_1v_1 + \dots + a_nv_n$ and therefore $A^kx_0 =$
368 $a_1\lambda_1^k v_1 + \dots + a_n\lambda_n^k v_n = \lambda_1^k(a_1v_1 + a_2\left(\frac{\lambda_2}{\lambda_1}\right)^k v_2 + \dots + a_n\left(\frac{\lambda_n}{\lambda_1}\right)^k v_n)$. Once we have estimates of the
369 eigenvector and the two largest eigenvalues, we can estimate the error given that $\|v - v_1\| \leq$
370 $\frac{\|Av - \lambda_1 v\|}{\|\lambda_1 - \lambda_2\|}$. To find an estimate of λ_2 we know that $v_2 \perp v_1$ and $\|v_1\| = 1$. Let x_0 be any vector and let
371 $x_{k+1} = A(x_k - c_k v_1)$ where $c_k = v_1^T x_k$ (and then $(x_k - c_k v_1) \perp v_1$). If $\lambda^{(k)}_2 = \|Ax_k\|/\|x_k\|$ the
372 using the same argument as before $\lambda^{(k)}_2 \rightarrow \lambda_2$ as $k \rightarrow \infty$. This is true even if $\lambda_2 \approx \lambda_3$ (x_k may not
373 converge to v_2 , but λ_2 will converge to λ_2). In this way we have an estimate of λ_1 and λ_2 and may
374 estimate the error in v . Since $A = Y^T Y$, $Ax = Y^T(Yx) = ((Yx)^T Y)^T$, we do not need to compute A
375 (which has the complexity of $O(n^3)$). We used two matrix vector products at every iteration (which
376 have the complexity of the number of nonzero elements in Y which is at most $O(n)$). Moreover, if X is
377 large a naïve multiplication of a vector by a matrix can still take a long time and storing Y may require a
378 large amount of memory. For example, to store human chr1 at 1 kb resolution (where $n \approx 250000$) 500
379 GB of RAM would be required just to store Y . With sparse implementation we recall that $Y =$
380 $(Y^{(i)}, \dots, Y^{(n)})$ where $Y^{(i)} = \frac{x^{(i)-c_i - c_i}}{\sigma_i} = \frac{x^{(i)}}{\sigma_i} - \frac{c_i}{\sigma_i}$. While $\frac{x^{(i)}}{\sigma_i}$ is sparse, $\frac{x^{(i)}}{\sigma_i} - \frac{c_i}{\sigma_i}$ is not. In lieu of explicit
381 computation, let $\mathbf{1} = (1, 1, \dots, 1)^T$ then $Y^{(i)} = \frac{x^{(i)}}{\sigma_i} - \frac{c_i}{\sigma_i} \mathbf{1}$ and then $Y = XS - \mathbf{1} \cdot \mathbf{1} \cdot r^T$ where
382 $S = [1/\sigma_1 \dots 1/\sigma_n]_n$ and $r = [c_1/\sigma_1 \dots c_n/\sigma_n]^T$ and then $Yx = (X \cdot S)x - \mathbf{1} \cdot r^T \cdot x$. Let $Z = X \cdot S$.
383 Since $r^T x = \sum_{i=1}^n r_i x_i$, $Yx = Zx - (\sum_{i=1}^n x_i r_i) \mathbf{1}$. Since Z is as sparse as X we can do everything with
384 sparse matrices as $x^T Y = x^T Z - (x^T \mathbf{1}) r^T = x^T Z - (\sum_{i=1}^n x_i) r^T$. Projected time and memory usage
385 were calculated by fitting a power decay curve, R^2 of fit = 0.95 for time, and R^2 of fit = 0.98 for memory
386 usage.

387 After compartment calling, chromatin marks were profiled at features that overlap A or B compartments
388 by overlapping with ChIP-seq peaks and by using average signal profiles created by pyBigWig from the
389 deepTools package⁴⁴. ChIP-seq peaks and bigwig files were obtained from the ENCODE Roadmap
390 Epigenomics project⁴⁵. We filtered promoters with bivalent marks as active genes that had 2-fold higher
391 H3K27me3 or H3K9me3 signal compared to the average at promoters. Contiguous compartment domain
392 sizes were calculated by requiring at least two consecutive bins to have the same sign in the
393 eigenvector. To create profiles of A compartmental status along genes, we assigned genes to elongating,
394 mid, and paused. Elongation status was determined by RPKM GRO-seq signal within 250 bp of the TSS
395 compared to the gene body, excluding 500 bp from the TSS. Differences between Promoter – Gene Body
396 GRO-seq signal were ranked and placed into three equal categories considering only genes ≥ 20 kb in
397 size.

398 *Loop Analysis*

399 Loops were identified by HiCCUPS² or SIP²⁰ at multiple resolutions. For HiCCUPS, we used parameters –
400 m 2000 –r 500,1000,5000,10000 –f .05,.05.05.05. For SIP we used an FDR 0.05 at each resolution with
401 the parameters for resolutions of 500 bp; -d 15 –g 3.0; 1 kb: -d 17 –g 2.5; 5 kb: -d 6 –g 1.5; and 10 kb: -d
402 5 –g 1.3. Loops called by both methods were combined by placing all loops into 10 kb bins, and if
403 HiCCUPS and SIP called the same loop within the 10 kb bin, then only one instance of this loop was kept.
404 Loops in subsampled maps were overlapped with loops called in the full 20.3 billion map if the loop was
405 within +/- 25 kb of each other. Overlap of loops with CTCF was done using a published list of CTCF ChIP-
406 seq peaks and motifs². Central 1 kb bins were assigned to those where we could unambiguously assign a
407 CTCF ChIP-seq peak to a unique bin at motifs in convergent orientation. Only loops with unambiguous

408 CTCF assignment were used in decay analysis. Bullseye plots were created using SIPMeta²⁰ and the
409 decay was calculated as the average at each Manhattan distance (ring) moving away from the central
410 bin. These values were plotted as a ratio to the central bin's signal. The central bin of loops called at AUC
411 values were computed using Simpson's rule. Loops were placed into five equally sized categories
412 (quintiles) based on AUC values. AUC values between WT, $\Delta ZF1$, and $\Delta ZF10$ were normalized by the
413 diagonal to account for differences in the expected decay. The decay percentage rate of change listed in
414 the main text was calculated by averaging the number of kb between each 10% loss of signal.

415 Fit-Hi-C³¹ interactions were identified in 1 kb bin-pairs with an FDR 0.05. 3D loop models were created
416 with Pastis⁴⁶ using the raw Hi-C matrix. Models were visualized in ChimeraX⁴⁷.

417 *Comparison with Other Datasets*

418 Hi-C read-pairs from CTCF $\Delta ZF1$, $\Delta ZF1$, and wild-type were downloaded from GSE125595³² and
419 processed with juicer to the mm10 genome. Hi-C maps from the *D. melanogaster* dm6 genome and the
420 *C. elegans* ce10 genome were obtained from our previously published work^{20,21}. Hi-C maps used in our
421 metric comparison are listed in Tables S2 and S3.

422 Enhancers were downloaded from DENdb¹⁸ and active enhancers were defined as those that overlap
423 with H3K27ac ChIP-seq peaks in GM12878. Histone modification ChIP-seq data was obtained from the
424 ENCODE reference epigenome series (ENCSR977QPF) and RNAPII ChIP-seq peaks were combined from
425 RNAPII, RNAPIISer2ph, and RNAPIISer5ph (ENCSR447YYN and ENCSR000DZK)^{19,48}, with overlapping
426 peaks merged into a single peak. GRO-seq data from GM12878 was downloaded from GSM1480326⁴⁹,
427 and chromHMM states for GM12878 were downloaded from the Roadmap Epigenomics Project⁴⁵.

428

429 **Data and Code Availability**

430 Hi-C data can be downloaded from ENCODE Accession: ENCSXXXXX. Our programs for subsampling,
431 noise estimation, and eigenvector calculation on sparse matrices can be downloaded from
432 <https://github.com/JRowleyLab/HiCSampler>, <https://github.com/JRowleyLab/HiCNoiseMeasurer>, and
433 <https://github.com/aidenlab/EigenVector>. These are open source and include source code as well as
434 implementations in python and C++.

435 **References**

- 436 1 Lieberman-Aiden, E. *et al.* Comprehensive Mapping of Long-Range Interactions Reveals Folding
437 Principles of the Human Genome. *Science* **326**, 289-293, doi:10.1126/science.1181369 (2009).
- 438 2 Rao, S. S. P. *et al.* A 3D map of the human genome at kilobase resolution reveals principles of
439 chromatin looping. *Cell* **159**, 1665-1680, doi:10.1016/j.cell.2014.11.021 (2014).
- 440 3 Belaghzal, H. *et al.* Liquid chromatin Hi-C characterizes compartment-dependent chromatin
441 interaction dynamics. *Nat Genet* **53**, 367-378, doi:10.1038/s41588-021-00784-4 (2021).
- 442 4 Falk, M. *et al.* Heterochromatin drives compartmentalization of inverted and conventional
443 nuclei. *Nature* **570**, 395-399, doi:10.1038/s41586-019-1275-3 (2019).
- 444 5 Erdel, F. & Rippe, K. Formation of Chromatin Subcompartments by Phase Separation. *Biophys J*
445 **114**, 2262-2270, doi:10.1016/j.bpj.2018.03.011 (2018).
- 446 6 Rowley, M. J. & Corces, V. G. Organizational principles of 3D genome architecture. *Nat Rev*
447 *Genet* **19**, 789-800, doi:10.1038/s41576-018-0060-8 (2018).
- 448 7 Dekker, J. *et al.* The 4D nucleome project. *Nature* **549**, 219-226, doi:10.1038/nature23884
449 (2017).

- 450 8 Heinz, S. *et al.* Transcription Elongation Can Affect Genome 3D Structure. *Cell* **174**, 1522-1536
451 e1522, doi:10.1016/j.cell.2018.07.047 (2018).
- 452 9 Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C
453 Experiments. *Cell Systems* **3**, 95-98, doi:10.1016/j.cels.2016.07.002 (2016).
- 454 10 Abdennur, N. & Mirny, L. A. Cooler: scalable storage for Hi-C data and other genomically labeled
455 arrays. *Bioinformatics* **36**, 311-316, doi:10.1093/bioinformatics/btz540 (2020).
- 456 11 Szabo, Q., Bantignies, F. & Cavalli, G. Principles of genome folding into topologically associating
457 domains. *Sci Adv* **5**, eaaw1668, doi:10.1126/sciadv.aaw1668 (2019).
- 458 12 Sikorska, N. & Sexton, T. Defining Functionally Relevant Spatial Chromatin Domains: It is a TAD
459 Complicated. *J Mol Biol* **432**, 653-664, doi:10.1016/j.jmb.2019.12.006 (2020).
- 460 13 Ea, V., Baudement, M. O., Lesne, A. & Forne, T. Contribution of Topological Domains and Loop
461 Formation to 3D Chromatin Organization. *Genes (Basel)* **6**, 734-750, doi:10.3390/genes6030734
462 (2015).
- 463 14 Rowley, M. J. *et al.* Condensin II Counteracts Cohesin and RNA Polymerase II in the
464 Establishment of 3D Chromatin Organization. *Cell Rep* **26**, 2890-2903 e2893,
465 doi:10.1016/j.celrep.2019.01.116 (2019).
- 466 15 Rowley, M. J. *et al.* Evolutionarily Conserved Principles Predict 3D Chromatin Organization.
467 *Molecular Cell* **67**, 837-852, doi:10.1016/j.molcel.2017.07.022 (2017).
- 468 16 Dong, P. *et al.* 3D Chromatin Architecture of Large Plant Genomes Determined by Local A/B
469 Compartments. *Mol Plant* **10**, 1497-1509, doi:10.1016/j.molp.2017.11.005 (2017).
- 470 17 Rao, S. *et al.* Cohesin Loss Eliminates All Loop Domains. *Cell* **171**, 305-320, doi:10.1101/139782
471 (2017).
- 472 18 Ashoor, H., Klefogiannis, D., Radovanovic, A. & Bajic, V. B. DENdb: database of integrated
473 human enhancers. *Database (Oxford)* **2015**, doi:10.1093/database/bav085 (2015).
- 474 19 Zhang, J. *et al.* An integrative ENCODE resource for cancer genomics. *Nat Commun* **11**, 3696,
475 doi:10.1038/s41467-020-14743-w (2020).
- 476 20 Rowley, M. J. *et al.* Analysis of Hi-C data using SIP effectively identifies loops in organisms from
477 *C. elegans* to mammals. *Genome Res* **30**, 447-458, doi:10.1101/gr.257832.119 (2020).
- 478 21 Cubañas-Potts, C. *et al.* Different enhancer classes in *Drosophila* bind distinct architectural
479 proteins and mediate unique chromatin interactions and 3D architecture. *Nucleic Acids Research*
480 **45**, 1714-1730, doi:10.1093/nar/gkw1114 (2016).
- 481 22 Anderson, E. C. *et al.* X Chromosome Domain Architecture Regulates *Caenorhabditis elegans*
482 Lifespan but Not Dosage Compensation. *Dev Cell*, doi:10.1016/j.devcel.2019.08.004 (2019).
- 483 23 Jimenez, D. *et al.* Condensin DC spreads linearly and bidirectionally from recruitment sites to
484 create loop-anchored TADs in *C. elegans*. *BioRxiv* (2021).
- 485 24 Davidson, I. F. & Peters, J. M. Genome folding through loop extrusion by SMC complexes. *Nat*
486 *Rev Mol Cell Biol*, doi:10.1038/s41580-021-00349-7 (2021).
- 487 25 Fudenberg, G. *et al.* Formation of Chromosomal Domains by Loop Extrusion. Report No.
488 biorxiv;024620v1, (2015).
- 489 26 Sanborn, A. L. *et al.* Chromatin extrusion explains key features of loop and domain formation in
490 wild-type and engineered genomes. *Proceedings of the National Academy of Sciences of the*
491 *United States of America*, doi:10.1073/pnas.1518552112 (2015).
- 492 27 Nichols, M. H. & Corces, V. G. A CTCF Code for 3D Genome Architecture. *Cell* **162**, 703-705,
493 doi:10.1016/j.cell.2015.07.053 (2015).
- 494 28 Gutierrez-Perez, I. *et al.* Ecdysone-induced 3D chromatin reorganization involves active
495 enhancers bound by Pipsqueak and Polycomb. *Cell Reports* (2019).

- 496 29 Eagen, K. P., Aiden, E. L. & Kornberg, R. D. Polycomb-mediated chromatin loops revealed by a
497 subkilobase-resolution chromatin interaction map. *Proceedings of the National Academy of*
498 *Sciences* **114**, 8764-8769, doi:10.1073/pnas.1701291114 (2017).
- 499 30 Ogiyama, Y., Schuettengruber, B., Papadopoulos, G. L., Chang, J. M. & Cavalli, G. Polycomb-
500 Dependent Chromatin Looping Contributes to Gene Silencing during Drosophila Development.
501 *Mol Cell* **71**, 73-88, doi:10.1016/j.molcel.2018.05.032 (2018).
- 502 31 Ay, F., Bailey, T. L. & Noble, W. S. Statistical confidence estimation for Hi-C data reveals
503 regulatory chromatin contacts. *Genome Res* **24**, 999-1011, doi:10.1101/gr.160374.113 (2014).
- 504 32 Saldana-Meyer, R. *et al.* RNA Interactions Are Essential for CTCF-Mediated Genome
505 Organization. *Mol Cell* **76**, 412-422 e415, doi:10.1016/j.molcel.2019.08.015 (2019).
- 506 33 Vilarrasa-Blasi, R. *et al.* Dynamics of genome architecture and chromatin function during human
507 B cell differentiation and neoplastic transformation. *Nat Commun* **12**, 651, doi:10.1038/s41467-
508 020-20849-y (2021).
- 509 34 Lucic, B. *et al.* Spatially clustered loci with multiple enhancers are frequent targets of HIV-1
510 integration. *Nat Commun* **10**, 4059, doi:10.1038/s41467-019-12046-3 (2019).
- 511 35 Nuebler, J., Fudenberg, G., Imakaev, M., Abdennur, N. & Mirny, L. A. Chromatin organization by
512 an interplay of loop extrusion and compartmental segregation. *Proc Natl Acad Sci U S A* **115**,
513 E6697-E6706, doi:10.1073/pnas.1717730115 (2018).
- 514 36 Jiang, Y. *et al.* Genome-wide analyses of chromatin interactions after the loss of Pol I, Pol II, and
515 Pol III. *Genome Biol* **21**, 158, doi:10.1186/s13059-020-02067-3 (2020).
- 516 37 Nichols, M. H. & Corces, V. G. Principles of 3D compartmentalization of the human genome. *Cell*
517 *Rep* **35**, 109330, doi:10.1016/j.celrep.2021.109330 (2021).
- 518 38 Stigler, J., Çamdere, G. Ö., Koshland, D. E. & Greene, E. C. Single-Molecule Imaging Reveals a
519 Collapsed Conformational State for DNA-Bound Cohesin. *Cell Reports*,
520 doi:10.1016/j.celrep.2016.04.003 (2016).
- 521 39 Davidson, I. F. *et al.* Rapid movement and transcriptional re-localization of human cohesin on
522 DNA. *The EMBO journal* **35**, 2671-2685, doi:10.15252/embj.201695402 (2016).
- 523 40 Busslinger, G. A. *et al.* Cohesin is positioned in mammalian genomes by transcription, CTCF and
524 Wapl. *Nature* **544**, 503-507, doi:10.1038/nature22063 (2017).
- 525 41 You, Q. *et al.* Direct DNA crosslinking with CAP-C uncovers transcription-dependent chromatin
526 organization at high resolution. *Nat Biotechnol*, doi:10.1038/s41587-020-0643-8 (2020).
- 527 42 Vian, L. *et al.* The Energetics and Physiological Impact of Cohesin Extrusion. *Cell* **175**, 292-294,
528 doi:10.1016/j.cell.2018.09.002 (2018).
- 529 43 Baglama, J. & Lothar, R. Augmented Implicitly Restarted Lanczos Bidiagonalization Methods.
530 *Siam J Sci Comput* **27**, 19-42 (2005).
- 531 44 Ramirez, F., Dundar, F., Diehl, S., Gruning, B. A. & Manke, T. deepTools: a flexible platform for
532 exploring deep-sequencing data. *Nucleic Acids Res* **42**, W187-191, doi:10.1093/nar/gku365
533 (2014).
- 534 45 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes.
535 *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).
- 536 46 Cauer, G., Gurkan, Y., Vert, J., Varoquaux, N. & Noble, W. S. Inferring diploid 3D chromatin
537 structures from Hi-C data. *BioRxiv* (2019).
- 538 47 Goddard, T. D. *et al.* UCSF ChimeraX: Meeting modern challenges in visualization and analysis.
539 *Protein Sci* **27**, 14-25, doi:10.1002/pro.3235 (2018).
- 540 48 Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature*
541 **489**, 57-74, doi:10.1038/nature11247 (2012).
- 542 49 Core, L. J. *et al.* Analysis of nascent RNA identifies a unified architecture of initiation regions at
543 mammalian promoters and enhancers. *Nat Genet* **46**, 1311-1320, doi:10.1038/ng.3142 (2014).

544

545 **Acknowledgements**

546 We acknowledge additional members of the ENCODE consortiums Nuclear Architecture Working Group
547 for thought-provoking discussions. Research reported in this publication was supported by the National
548 Institute of General Medical Sciences of the National Institutes of Health (NIH) under Award Numbers
549 T32-GM067553, R35-GM139408, R35-GM128645, R01-MH115957, U24~HG009446, and R00-
550 GM127671. The content is solely the responsibility of the authors and does not necessarily represent
551 the official views of the NIH.

552 **Author Contributions**

553 H.G prepared Hi-C libraries for sequencing with samples prepared by S.K., K.M., M.E.T, and J.S.S. H.G.,
554 H.H., Y.E., A. Krishna, A. Kalluchi, M.P., S.R., O.D, D.U., M.H.N., and E.D. contributed ideas and in
555 performing various quality metrics. M.J. and G.C. created 3D loop models. M.O. created POSSUM.
556 D.H.P., V.G.C, W.S.N, E.L.A., and M.J.R. supervised the work and wrote the manuscript. All other analyses
557 were performed by M.J.R.

558 **Ethics Declarations**

559 We declare that the authors have no competing interests in this work.

560

561 **Figure Legends**

562 **Figure 1.** By combining ultra-deep Hi-C and POSSUMM, we generated a fine map of nuclear
563 compartmentalization achieving 500bp resolution.

564 A) Example locus showing Hi-C signal in 500 bp bins in our full map with 20.3 billion intrachromosomal
565 read-pairs (left) and when read-pairs are subsampled to 1 billion (right).

566 B) Example of compartment interactions in a Hi-C map identified by the eigenvector in 500 bp bins
567 (bottom track). Black track displays transcription measured by GRO-seq.

568 C) Zoomed in view of a compartment domain.

569 D) Overview of the power method and POSSUM for calculating the eigenvector. See Methods for details.

570

571 **Figure 2.** Nearly all active TSSs and Enhancers localize to kilobase-scale A compartments

572 A) Cumulative fraction of compartment domain sizes when identified at 500 bp resolution.

573 B) Percentage of active gene promoters, proximal enhancers, and distal enhancers assigned to A (green)
574 or B (purple) compartment domains when identified by the 500 bp compartment eigenvector.

575 C) Example of small compartment domains only identifiable at high-resolution (red asterisks). Log
576 transformed and distance normalized Hi-C map is shown alongside the eigenvector tracks at various bin
577 sizes.

578 D) Examples an active enhancers denoted by H3K27ac and H3K4me1 signal localizing to the A
579 compartment and surrounded by the B compartment.

580 E) Examples an active promoters denoted by GRO-seq signal localizing to the A compartment and
581 surrounded by the B compartment.

582

583 **Figure 3.** Many genes exhibit discordant compartmentalization.

584 A-E) Examples of genes of various sizes where the TSS is in the A compartment while the TTS is in the B
585 compartment. GRO-seq signal is shown as an indicator of the gene's transcription status.

586 F) Scaled average profiles of the A compartment signal (positive eigenvector) relative to the TSS for
587 short (blue), mid-sized (gold), large (pink), and randomly selected (black) genes.

588 G) Percentage of TTSs that localize to the B compartment for genes of various sizes (left). Diagram of A
589 compartment signal on short and large genes (right).

590 H) ChIP-seq signal at the TTS of discordant A/B genes vs. concordant A/A genes. Genes are sorted by the
591 TTS compartmental signal.

592 I) Scaled average profiles of the A compartment signal (positive eigenvector) relative to the TSS for
593 elongating (blue), mid (red), paused (black), or randomly selected (grey) genes.

594 J) Percentage of TTs that localize to the B compartment for paused, mid, or elongating genes.

595 K) Diagram of TSS and TTS localization to the A compartment depending on gene size and elongation
596 status.

597

598 **Figure 4.** CTCF loop-decay enhances proximal interactions and is dependent on RNA-binding domains.

599 A) Example of broad signal enrichment near CTCF loops when binned at 1 kb.

600 B) Average signal at CTCF loops when binned at 10, 5, or 1 kb, centered on convergent CTCF anchors.

601 C) Average Hi-C signal in 1 kb bins at each radial distance away from the CTCF loop anchors (rainbow).
602 Average signal of the diagonal decay is shown for reference (grey) to estimate interactions due to
603 polymeric distance. AUC=area under the curve.

604 D) Example of punctate signal enrichment at Pc loops in *D. melanogaster* when binned at 1 kb.

605 E) Average signal at *D. melanogaster* Pc loops when binned at 10, 5, or 1 kb.

606 F) Average Hi-C signal in 1 kb bins at each radial distance away from human CTCF loop anchors (blue) vs.
607 *D. melanogaster* Pc loops (orange), and *C. elegans* X-chromosome loops (green). Average signal at the *C.*
608 *elegans* Hi-C diagonal is shown for reference (grey). AUC=area under the curve.

609 G) Enrichment of Fit-Hi-C enhancer-promoter interactions within 100 kb of loops inside the loop (blue)
610 or crossing over loop boundaries (green). Values are shown as enrichment vs random regions of equal
611 size and number as loops.

612 H) Diagram of how CTCF loops can shorten distances between enhancers (orange) and promoters (blue)
613 even when both are located outside of the loop.

614 I) Average GRO-seq signal at CTCF loop anchors and neighboring loci for loops divided into 5 distinct
615 decay categories.

616 J) Average Hi-C signal in WT (left), $\Delta ZF1$ (middle), or $\Delta ZF10$ (right) CTCF mutants at CTCF loops. AUC=area
617 under the curve

618

619 **Figure 5** Sub-genic compartmentalization and diffuse CTCF looping organize the human genome.

620 Diagram depicting localization of active enhancers and TSSs to the A compartment, while TTSs are
621 oriented to the B compartment dependent on transcription elongation status. This sub-genic and precise
622 enhancer compartmentalization combines with diffuse CTCF loops to mediate genome organization.

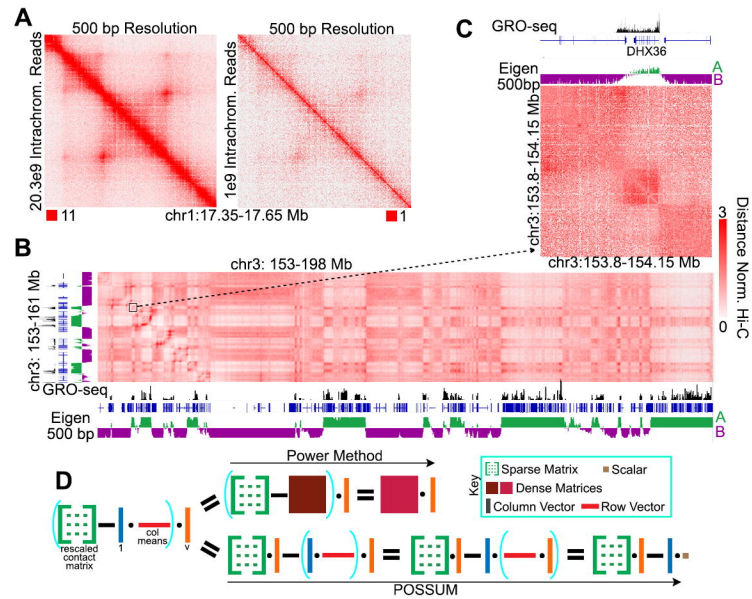


Figure 1

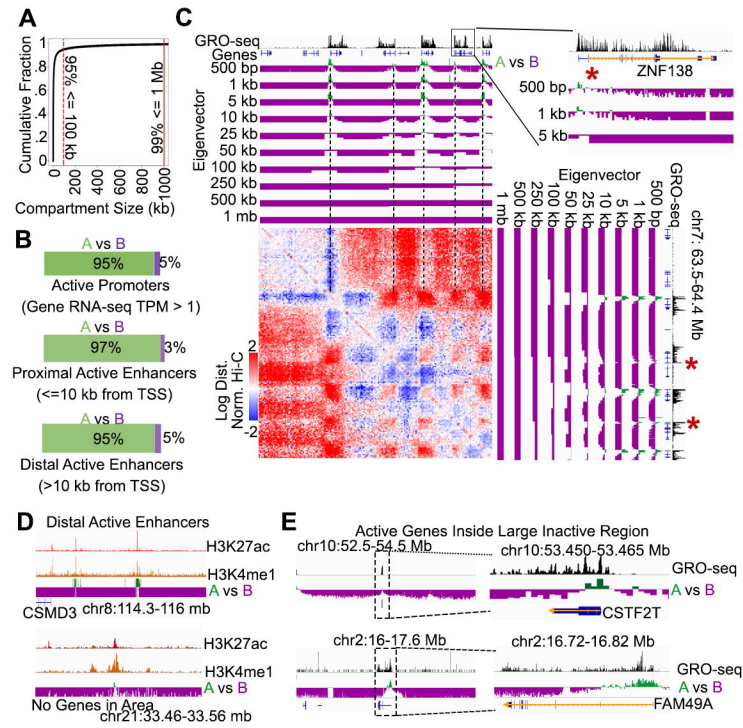


Figure 2

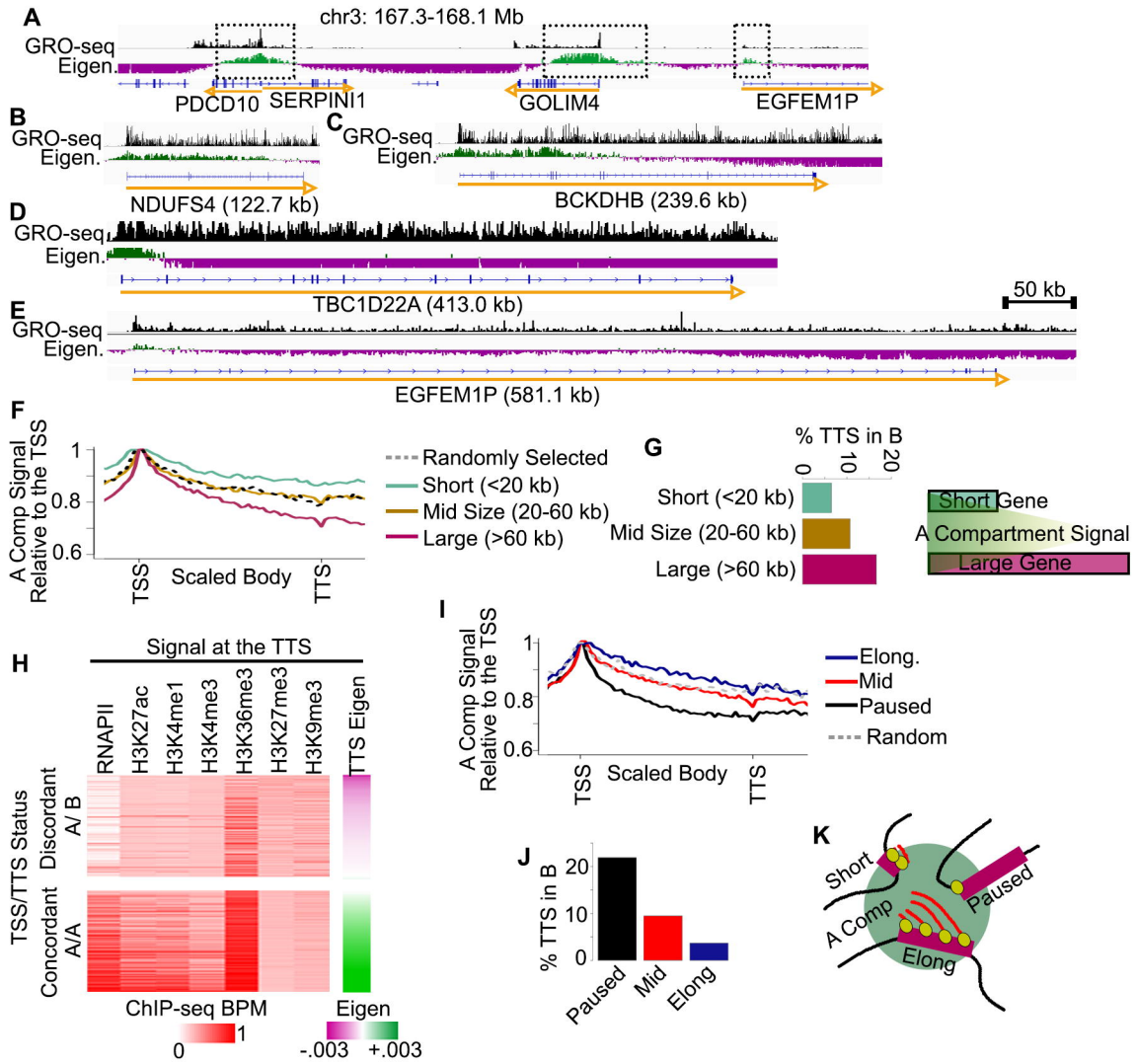


Figure 3

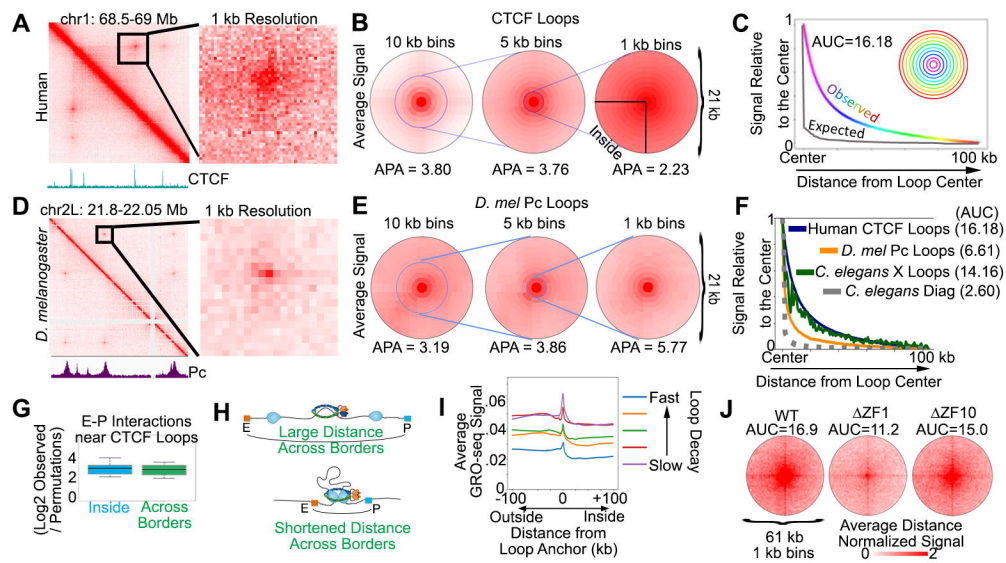


Figure 4

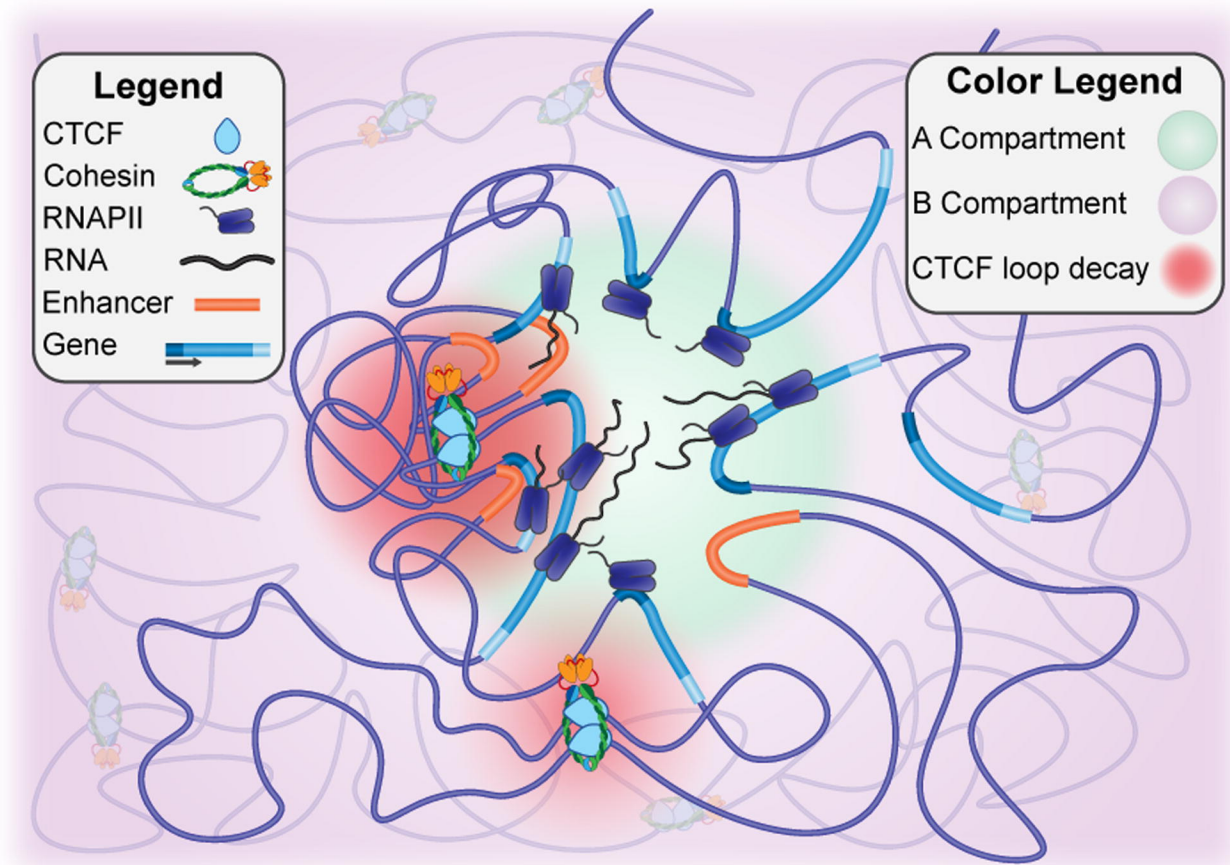


Figure 5