



Analysis of Hi-C data using SIP effectively identifies loops in organisms from *C. elegans* to mammals

M Jordan Rowley, Axel Poulet, Michael Nichols, et al.

Genome Res. published online March 3, 2020

Access the most recent version at doi:[10.1101/gr.257832.119](https://doi.org/10.1101/gr.257832.119)

P<P	Published online March 3, 2020 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Analysis of Hi-C data using SIP effectively identifies loops in organisms from *C. elegans* to mammals

M. Jordan Rowley^{1,3,7}, Axel Poulet^{1,3,8}, Michael H. Nichols², Brianna J. Bixler², Adrian L. Sanborn⁵, Elizabeth A. Brouhard⁴, Karen Hermetz², Hannah Linsenbaum², Gyorgyi Csankovszki⁴, Erez Lieberman Aiden^{5,6}, and Victor G. Corces^{2*}

¹These authors contributed equally to this work.

²Department of Human Genetics, Emory University School of Medicine, 655 Michael St., Atlanta, GA 30322, USA

³Department of Biology, Emory University, 1510 Clifton Rd NE, Atlanta, GA 30322, USA

⁴Department of Molecular, Cellular, and Developmental Biology, University of Michigan, Ann Arbor, MI, USA

⁵Center for Genome Architecture, Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX

⁶Center for Theoretical Biological Physics and Department of Computer Science, Rice University, Houston, TX

Current Addresses:

⁷Department of Genetics, Cell Biology and Anatomy, University of Nebraska Medical Center, Omaha, NE

⁸Department of Molecular, Cellular and Developmental Biology, Yale University, 165 Prospect St New Haven, CT 06511, USA

*Corresponding Author:

Victor G. Corces

Email: vgcorces@gmail.com

Phone: 404-727-4250

Running Title: Analysis of CTCF loops

Key Words: CTCF, cohesin, condensin, extrusion, chromatin, dosage compensation, Transcription, HiChIP

32 ABSTRACT

33 Chromatin loops are a major component of 3D nuclear organization, visually apparent as
34 intense point-to-point interactions in Hi-C maps. Identification of these loops is an critical part of
35 most Hi-C analyses. However current methods often miss visually evident CTCF loops in Hi-C
36 datasets from mammals and they completely fail to identify high intensity loops in other
37 organisms. We present SIP, Significant Interaction Peak caller, and SIPMeta, which are
38 platform independent programs to identify and characterize these loops in a time and memory
39 efficient manner. We show that SIP is resistant to noise and sequencing depth and can be used
40 to detect loops that were previously missed in human cells as well as loops in other organisms.
41 SIPMeta corrects for a common visualization artifact by accounting for Manhattan distance to
42 create average plots of Hi-C and HiChIP data. We then demonstrate that the use of SIP and
43 SIPMeta can lead to biological insights by characterizing the contribution of several transcription
44 factors to CTCF loop stability in human cells. We also annotate loops associated with the SMC
45 component of the Dosage Compensation Complex (DCC) in *C. elegans* and demonstrate that
46 loop anchors represent bi-directional blocks for symmetrical loop extrusion. This is in contrast to
47 the asymmetrical extrusion until unidirectional blockage by CTCF that is presumed to occur in
48 mammals. Using HiChIP and multi-way ligation events, we then show that DCC loops form a
49 network of strong interactions that may contribute to X Chromosome-wide condensation in *C.*
50 *elegans* hermaphrodites.

51 (Supplemental material is available for this article)

52 INTRODUCTION

53 High resolution Hi-C in human cells is able to find thousands of strong punctate signals that
54 indicate the presence of loops formed by CTCF sites arranged in a convergent orientation (Rao
55 et al. 2014). Based on this orientation preference, it has been proposed that CTCF loops are
56 formed by a loop extrusion process mediated by cohesin [reviewed by (Rowley and Corces
57 2018)]. Indeed, depletion of cohesin in mammalian cells results in loss of CTCF loops (Rao et
58 al. 2017). However, other transcription factors are also present at CTCF loop anchors and it is
59 unclear whether or not they play a role in loop extrusion or affect the frequency or stability of
60 CTCF loops (Rao et al. 2014).

61 CTCF loops have been identified in mammals but have not been observed in other
62 organisms. For example, *Drosophila* Hi-C maps do not display CTCF loops despite the
63 existence of a conserved *Drosophila* homologue (Rowley et al. 2017). Instead, *Drosophila*
64 contact maps in Kc167 cells contain a few hundred loops that lack CTCF and their formation
65 does not depend on cohesin (Rowley et al. 2019). Many non-vertebrate organisms, including *C.*
66 *elegans*, lack a CTCF homologue (Heger et al. 2012). It is possible that proteins distinct from
67 CTCF are able to form point to point interactions in these organisms, as is the case of
68 *Drosophila*, or to stop the extrusion of SMC complexes to form loops. For example, *C. elegans*
69 hermaphrodites regulate X-Chromosome expression through the use of the DCC complex,
70 which contains a condensin complex presumably able to extrude DNA (Lau and Csankovszki
71 2014). However, although published Hi-C contact maps reveal the presence of large self-
72 interacting domains in the dosage compensated X Chromosome and evidence of loop
73 formation, current algorithms have not been successful at systematically annotating punctate
74 signals corresponding to loops in *C. elegans* (Crane et al. 2015; Anderson et al. 2019). This has
75 made it difficult to fully explore these features in non-mammalian organisms.

76 Here we report a method of loop identification named Significant Interaction Peak caller
77 (SIP) that relies on CPU based image analysis of Hi-C contact maps to find loops. SIP detects
78 additional functionally relevant loops in human cells and can be used to detect loops in a variety

79 of other organisms. We also present a companion tool, SIPMeta, that creates average
80 metaplots of loops. We show that current standard metaplots contain visual biases that SIPMeta
81 corrects. Using SIP and SIPMeta, we test whether several transcription factors, including
82 ZNF143, YY1, RNA Polymerase II, and CTCFL affect the strength of CTCF loops. We then
83 perform Hi-C and DPY-27 HiChIP in *C. elegans* hermaphrodites and show that the X
84 Chromosome contains dozens of high intensity loops configured into a complex network. These
85 loops are associated with condensin I-DCC, suggesting the existence of extrusion-mediated
86 non-CTCF loops. The results suggest the formation of a rosette-like structure that may be
87 responsible for dosage compensation in this organism. Therefore, SIP and SIPMeta represent
88 sensitive and versatile new methods for loop calling and analysis that can lead to the discovery
89 of novel biological information from Hi-C data.

90 RESULTS

91 *SIP Software*

92 Loops present in Hi-C heatmaps appear as intense saturated punctae (Rao et al. 2014). To
93 identify these visibly evident interactions we took advantage of image processing methods to
94 create SIP. SIP includes options to use command line or graphical user interfaces (Fig. 1A). The
95 SIP pipeline (Fig. 1B) reads Hi-C data in either the Juicer .hic format (Durand et al. 2016) or in a
96 bedpe-like format with distance-normalized signal. The genome is analyzed by sliding windows
97 using image processing to identify potential loops, which are then filtered based on several
98 aspects of the matrix. Images undergo a gaussian blur, contrast enhancement, white top-hat,
99 and a Minimum-Maximum Filter. These steps provide a corrected image of the interactions (Fig.
100 1), which is used with a regional maxima detection algorithm to detect a preliminary list of
101 candidate loops.

102 Candidate loops must then pass several filters that utilize the original distance
103 normalized signal in the Hi-C data. First, pixels near unmappable or repetitive regions are
104 removed. Second, to remove isolated pixels representing noise, loops must display a decay
105 such that the central pixel is the highest, followed by decreases at 1 and 2 pixels away from the
106 center. The center must also be 1.2-fold higher than the average of nearby pixels and must
107 pass a Poisson CDF filter such that the probability that the center is higher than other nearby
108 pixels is greater than 0.9. Thus, SIP utilizes the local background to identify loops. While this is
109 useful for identifying punctate signals, other programs that model enrichment over global
110 background can be useful to create large lists of enhancer-promoter interactions (Ay et al.
111 2014). Finally, candidate loops are filtered based upon an empirical FDR calculated as the
112 enrichment of loops vs random sites at equal distances.

113 *Performance of SIP*

114 We tested the performance of SIP on Hi-C data from GM12878 cells containing approximately
115 2.4 billion intra-chromosomal reads (Rao et al. 2014). As a benchmark, we compared the time
116 and memory usage of SIP to other interaction callers designed to identify CTCF loops (Durand
117 et al. 2016; Heinz et al. 2018; Cao et al. 2019). SIP is intended to be used without the need of
118 large computing power, therefore we intentionally limited SIP to one thread and memory usage
119 to 1 GB using java -Xmx1g for both SIP and HiCCUPS on Chromosome 1. In comparison,
120 HOMER and cLoops used 62 and 103 GB respectively for Chromosome 1 (Table S1). Even
121 with these parameters, SIP identified loops 2x, 14x, or even 1,057x faster than HiCCUPS,
122 HOMER, and cLoops, respectively. We then tested SIP on a laptop with a 2-core processor and
123 4 GB of RAM and we were able to call loops in the full dataset, including all chromosomes, at 5
124 kb resolution in 46 min, including dumping data from the Juicer .hic file (Fig. 2A). To allow easier
125 parameter optimization, users have the option of saving these dumped files and rerunning SIP,

126 which took only 31 min. On a Linux machine using 23 cores, we were able to call loops in the
127 full GM12878 dataset in 12 minutes (Fig. 2B). A comparison of memory and time usage by SIP
128 on different systems can be found in Table S2.

129 We compared loops called by SIP to those called by other loop identification programs
130 (Fig. S1A) and found that SIP identified more loops than existing tools (Fig. S1B). An exception
131 was Fit-Hi-C, which was designed to identify enhancer-promoter interactions rather than
132 punctate spots present in Hi-C data and associated with CTCF loops (Ay et al. 2014) (Fig. S1A-
133 B). We compared loops called by SIP vs HiCCUPS and observed good overlap. However, 33%
134 of HiCCUPS loops were not identified by SIP, and 67% of SIP loops were not identified by
135 HiCCUPS (Fig. S1C). These show punctate signal in average metaplots and are therefore likely
136 true loops that each program missed (Fig. S1C). However, when we include cLoops, HOMER,
137 and Fit-Hi-C in this analysis, we found that both SIP and HiCCUPS had more than 95% of loop
138 calls identified in at least one other program (Fig. S1D). In order to compare loops called by
139 each program, we designated loops that were identified by at least two programs as pseudo
140 true positives while loops that were unique to each program were designated as pseudo-false
141 positives. SIP had a low pseudo-false positive rate and low pseudo-false negative rate in
142 comparison to other programs (Fig. S1E). These results indicate that while current loop callers
143 are unable to identify 100% of loops, SIP has an improved detection rate.

144 To further benchmark SIP, we evaluated three different aspects of the program - the
145 ability to accurately capture loops with sparse datasets, the reproducibility of loop calls, and the
146 resistance to noise. To test the ability to identify loops in datasets with fewer sequenced reads,
147 we subsampled a dataset with 2.4 billion intra-chromosomal paired reads and created contact
148 maps with 1 billion, 500 million, 250 million, and 100 million reads. Regardless of the method,
149 lower sequencing depth correlates with a decreased ability to identify loops at 5 kb resolution
150 (Fig. 2C). However, SIP consistently identified a higher percentage of loops from the full dataset
151 than HiCCUPS (Fig. 2C). We also tested whether lower read counts resulted in identification of
152 loops in the subsampled dataset that were not identified in the full dataset i.e. likely false
153 positives. We found that each method identified a low number of potential false positives with no
154 correlation to sequencing depth (Fig. 2D). As a secondary test, we called loops with each
155 method in the 1 billion read dataset but varied FDR parameters. For each FDR parameter
156 tested, the false positive rate was calculated by the number of loops called in the subsampled
157 data that were not called in the full dataset. As expected, both methods displayed increased
158 false positives with decreased FDR stringency. However, at similar false positive rates in the
159 subsampled data, SIP was able to identify approximately twice the number of loops as
160 HiCCUPS (Fig. 2E). Overall, we find that SIP is able to recover a high percentage of loops
161 without increasing the false positive rate using Hi-C datasets with a low number of sequenced
162 reads.

163 In order to determine the reproducibility of loop calls with different datasets, we called
164 loops in Hi-C maps from 8 distinct cell lines (Rao et al. 2014). Each of the 8 datasets have
165 different depths of sequencing and, in general, the number of loops identified approximately
166 matches the number obtained from down-sampled GM12878 datasets (Fig. S1F). In
167 comparison to HiCCUPS, SIP identified a larger number of loops in each dataset that were also
168 present in GM12878 cells (Fig. 2F and Fig. S1G). Loops specific to each dataset display Hi-C
169 signal specific to that dataset (Fig. S1H). This suggests that SIP is able to reproducibly identify
170 loops and that differences in loop calls between Hi-C maps are due to differences in looping. To
171 further estimate the reproducibility of loop calls by SIP, we created distinct Hi-C datasets by
172 random sampling the full dataset down to 1 billion reads in independent iterations to create 10
173 different .hic maps. We then examined how many of the loops in each iteration were the same
174 between datasets. Both SIP and HiCCUPS were able to reproducibly identify loops obtaining on

175 average 91% (SIP) or 86% (HiCCUPS) of loops in each subsampled iteration that were
176 consistent between datasets (Fig. 2G).

177 Next, we evaluated the ability of each method to identify loops in noisy datasets. We
178 created Hi-C maps where noise was simulated by distributing random additional signal within
179 the map (see Methods). We noticed that in maps with 50% additional noise signal, HiCCUPS
180 called a large number of false positives at extreme distances crossing over the entire
181 chromosome (Fig. S2A, blue). These can be easily filtered using a distance cutoff. Thus, to
182 more fairly benchmark SIP and HiCCUPS we only examined loops less than 10 Mb in size. This
183 noise model decreased the original loop signal vs background (Fig. S2B), but both methods
184 recovered a comparable fraction of the original loop calls despite the additional noise (Fig. S2C-
185 D). However, increased noise caused an increase in the false positive rate by HiCCUPS, while
186 SIP remained consistently low (Fig. S2C and S2E). While this noise model is purely artificial and
187 may not recapitulate the true noise in a sample, these results indicate that SIP is at least
188 partially resistant to these variations while HiCCUPS is not.

189 Lastly, we examined the effects of bin size on loop calling by identifying loops at 5, 10,
190 and 25 kb. We found that SIP and HiCCUPS had similar overlaps between these calls, but each
191 program had loops uniquely identified at each resolution (Fig. S2F). Therefore, as in the original
192 HiCCUPS caller, it may be advantageous to call loops at multiple resolutions (Rao et al. 2014).
193 Overall, these results suggest that SIP is memory and time efficient, identifies loops that are
194 semi-resistant to sequencing depth, has high reproducibility, and has high resistance to noise.

195 ***SIP and SIPMeta Allow Identification and Visualization of Loops in Various Organisms***

196 One problem with loop identification in Hi-C datasets has been that the training on one dataset
197 impacts loop calling on other datasets. This is the reason loop identification in *Drosophila* was
198 done using separate custom scripts or by hand (Cubeñas-Potts et al. 2016; Eagen et al. 2017).
199 We used SIP on Hi-C maps of Kc167 cells and identified 143 high intensity loops at 1 kb
200 resolution (Fig. 3A). Visual inspection of these loops shows that they correspond to punctate
201 signal (Fig. 3A). In comparison, other loop calling methods have a tendency to also call
202 interactions near sparse signal that likely corresponds to repetitive regions (Fig. S3A, see also
203 Fig. S2C). We then tested if anchors of loops identified by SIP were enriched in proteins
204 previously found to be important for looping in *Drosophila* (Eagen et al. 2017; Ogiyama et al.
205 2018; Gutierrez-Perez et al. 2019). Indeed Polycomb (Pc) and Pipsqueak (Psq) are enriched on
206 SIP loop anchors (Fig. S3B).

207 A common approach to evaluating loops is through a metaplot analysis that averages the Hi-C
208 signal at loops compared to the surrounding region (Rao et al. 2014; Rowley et al. 2019).
209 Standard metaplots of *Drosophila* loops display central signal enrichment, but with a crosshair
210 like pattern (Fig. 3B left). This could be interpreted as evidence of extrusion, similar to enriched
211 stripes in Hi-C maps of mammals that occur at some CTCF loops due to proximal loading of
212 cohesin (Vian et al. 2018). However, it was previously found that depletion of cohesin or
213 condensin II has no effect on *Drosophila* loop intensity (Rowley et al. 2019), thus it is unlikely
214 that these loops are formed via extrusion. When considering the crosshair pattern in square
215 metaplots, we realized that distance from the loop is different between pixels adjacent
216 horizontally or vertically vs pixels adjacent diagonally. For example, one pixel to the right of the
217 loop is 0 kb away from the left anchor and 1 kb away from the right anchor, equivalent to 1 kb
218 Manhattan distance (Fig. 3B left). However, one pixel diagonally away is 1 kb away from the left
219 anchor and 1 kb away from the right anchor, equivalent to 2 kb Manhattan distance. Thus, the
220 juxtaposition of pixels at different distances likely creates the observed crosshair pattern and
221 could be potentially misinterpreted. To more accurately depict Hi-C signal vs distance from
222 loops and thereby alleviate this common visualization issue, we created SIPMeta, which

223 generates both the standard square metaplots, as well as “bullseye” plots where pixels in each
224 ring represent the same Manhattan distance away from the loop (Fig. 3B right). The bullseye
225 visualization of *Drosophila* loops eliminates the crosshair pattern, demonstrating the potential
226 usefulness and impact of SIPMeta on data interpretation. For comparison we examined a “true”
227 stripe found in human GM12878 cells (Fig. 3C left) and found that the SIPMeta bullseye plot is
228 able to display this stripe (Fig. 3C right). Therefore, SIPMeta can distinguish extrusion-mediated
229 stripes from crosshair patterns, which are due to Euclidean distance effects relative to the loop.

230 We then examined the ability of SIP to identify loops in an organism where they have not
231 previously been characterized. We examined published Hi-C maps in the mosquito *Aedes*
232 *aegypti* (Matthews et al. 2018) and detect visually apparent loops (Fig. S3C). Using SIP we
233 identified 231 high intensity loops that display central enrichment (Fig. S3D). In this case,
234 cLoops was also able to identify these peaks while other programs were not (Fig. S3E). To test
235 whether these loops correlate with the presence of CTCF at anchor sites as is the case in
236 mammals or if they are similar to those found in *Drosophila* cells, we examined the enrichment
237 of CTCF motifs at loop anchors. Unlike human cells, which display a large enrichment of CTCF
238 motifs at loop anchors, we found that *A. aegypti* loop anchors are not enriched in CTCF motifs
239 (Fig. 3D). Therefore, it is likely that *A. aegypti* loops are like those found in *D. melanogaster* and
240 contain proteins other than CTCF at their anchors.

241 **Characterization of SIP Loops in Human Cells**

242 Having found that SIP identifies visually observable loops in *D. melanogaster* and *A. aegypti*, we
243 next evaluated SIP loop calls in GM12878 human cells. SIP identified 13,692 loops in the 5
244 billion Hi-C contacts dataset obtained in GM12878 cells at 5 kb resolution. This is nearly twice
245 (1.93-fold) the 7,101 loops identified by HiCCUPS using default parameters. We found a strong
246 preference for loop anchors containing CTCF peaks assignable to a convergent orientation
247 (10,663, 78%) (Fig. 4A). Compared to other programs, SIP and HiCCUPS detect the highest
248 percentage of loops with convergent CTCF, which is indicative of their ability to identify these
249 features (Fig. S4A). Additionally, we examined CTCF ChIA-PET data (Tang et al. 2015) and
250 found that SIP and HiCCUPS loops had the highest enrichment signal (Fig. S4B). Of the loops
251 identified by SIP, 1,038 and 56 were assigned to tandem and divergent orientations,
252 respectively, whereas 1,935 (14%) did not coincide with detected CTCF peaks in any particular
253 orientation (Fig. 4A). Analysis of metaplots of loops in all categories indicates that convergent
254 loops are strongest, followed by tandem and then unassigned loops (Fig. 4B). We then tested
255 the ability of each loop category to form domains by taking the Z-score values of each ring in the
256 bullseye plot and calculating an Aggregate Domain Analysis (ADA) score from the number of
257 high Z-scores in the bottom right corner compared to the surrounding regions. This Z-score
258 transformation and ADA calculation is included as an option in SIPMeta. Using this method, we
259 found that loops between convergent CTCF sites form the strongest domains and tandem loops
260 contain slightly weaker intra-domain interaction frequencies (Fig. 4C). Loops without identified
261 CTCF peaks do not display an underlying domain (Fig. 4C). We tested whether the absence of
262 an interaction domain is due to loop strength by examining convergent CTCF loops that display
263 weak loop signal. Weak convergent loops do not display domain signal either, indicating that
264 domain formation correlates with loop strength (Conv. Low Fig. 4C).

265 SIP detects 1,935 loops whose anchors seem to lack CTCF bound to its motif. This
266 could be a result of the stringency of CTCF peak calling in ChIP-seq data. For example, an
267 unassigned loop has a strong CTCF site on one anchor but has weak CTCF signal on the other
268 (Fig. 4D). Indeed, 1,572 (81%) of these unassigned loops have an identifiable CTCF ChIP-seq
269 peak on one anchor. Therefore, these loops are either interactions between CTCF and some
270 other protein or loops where the second anchor shows weak CTCF ChIP-seq signal insufficient

271 to call a peak but sufficient to form a loop. We examined CTCF ChIA-PET data (Tang et al.
272 2015) using SIPMeta and found enrichment of signal on convergent and tandem loops (Fig. 4E).
273 Although unassigned loops display weaker CTCF ChIA-PET signal than even weak convergent
274 loops, we still detect enrichment signal in the center compared to the surrounding region (Fig.
275 4E). Therefore, we believe unassigned loops are CTCF loops where one anchor has low levels
276 of CTCF. Thus, SIP and HiCCUPS identify similar types of features, although SIP is able to
277 identify additional CTCF loops.

278 Next we examined loops in published Hi-C data in HCT116 cells before and after
279 cohesin depletion (Rao et al. 2017). Using SIPMeta, we examined changes in loops after Rad21
280 depletion and reintroduction of this protein (Fig. 4F). Loops in each context are lost after Rad21
281 depletion, confirming that loops in mammals generally depend on the presence of cohesin (Rao
282 et al. 2017) (Fig. 4F). We noticed that CTCF loops in tandem orientation or without an
283 assignable CTCF peak are not able to recover as efficiently as convergent CTCF loops (Fig.
284 4F). Measuring APA scores after 180 min of cohesin recovery, convergent loops return to their
285 original enrichment value almost fully (93% of APA score), whereas tandem and unassigned
286 loops recover to only 52% and 20% of their original APA scores, respectively (Fig. 4G). This
287 slow recovery is not due to weaker loop signal in the control, since weak convergent loops
288 recover better than tandem and unassigned loops (Conv. Low, Fig. 4G).

289 CTCF is thought to block extrusion in an orientation-specific manner, so we reasoned
290 that the strength of the motif may determine the strength of the loop. We examined CTCF motif
291 strength vs loop strength and found that they are correlated (Fig. 4H). We also found that
292 convergent loops display stronger CTCF motifs than tandem loops (Fig. 4I). To examine
293 whether motif strength affects loop recovery after cohesin depletion and repletion, we examined
294 convergent loops in the top and bottom 10% of motif strength. We found that convergent loops
295 with weak motifs did not recover as quickly as convergent loops with strong motifs (Fig. 4J).
296 Indeed, convergent loops with weak motifs recovered as slowly as tandem loops. These data fit
297 with a model where strong convergent motifs efficiently dictate the orientation of the CTCF
298 protein, resulting in more robust blockage of extrusion and quick recovery. Based on the results,
299 weak convergent, tandem, or unassignable motifs are less efficient at dictating the orientation of
300 the CTCF protein on chromatin, resulting in less blockage of extrusion and slower recovery.
301 Therefore, we suggest that loops that appear to overlap tandem or no motifs could still be
302 occupied by convergently oriented CTCF proteins.

303 ***Other Transcription Factors affect the strength of CTCF Loops***

304 Although CTCF has a major role in the establishment of loops in mammalian cells, other
305 transcription factors present at loop anchors may affect the frequency of point-to-point
306 interactions causing the formation of these loops. Using SIPMeta, we investigated several
307 transcription factors whose binding sites have been previously shown to be present at CTCF
308 loop anchors, including ZNF143, YY1, CTCFL, and RNA Polymerase II (RNAPII) (Rao et al.
309 2014; Tang et al. 2015). First, we examined loops with high levels of CTCF on both anchors and
310 divided them into those with high or low levels of ZNF143. We find that when CTCF is high,
311 loops with high ZNF143 are stronger than loops with low ZNF143 signal (Fig. 5A top row).
312 Loops with weak CTCF signal are also stronger when ZNF143 signal is high (Fig. 5A bottom
313 row) indicating that the presence of ZNF143 can enhance looping frequency. In contrast, we
314 found no difference in loop signal between those containing high or low YY1 (Fig. 5B). However,
315 we do detect small signal differences on CTCF loops with high or low RNAPII (Fig. 5C).

316 CTCFL binds to the same motif as CTCF (Pugacheva et al. 2015). While GM12878 cells
317 do not express CTCFL, many CTCF motifs in K562 cells display peaks of both CTCFL and
318 CTCF. ChIP Re-ChIP experiments indicate that these sites are often bound by the two proteins

319 at the same time (Pugacheva et al. 2015). We thus hypothesized that the presence of both
320 proteins could affect looping. To test this, we examined published CTCFL ChIP-seq data in
321 K562 cells (Pugacheva et al. 2015) and compared its presence at loop anchors to that of CTCF.
322 While loop anchors preferentially contain strong CTCF peaks, there is equal presence of weak
323 and strong CTCFL peaks at loop anchors (Fig. 5D). We could not identify enough loops with low
324 CTCF and high CTCFL unambiguously (n=2), but for loops with high CTCF we found no
325 difference in signal when CTCFL was high or low. (Fig. 5E). We should note that other
326 programs were unable to categorize loops in this manner (Fig. S5), thus highlighting the
327 differences between loop callers. These results suggest that despite a similar DNA binding
328 domain, CTCFL is unable to form loops. Additionally, our results indicate that CTCFL does not
329 interfere with looping when present at the same location as CTCF.

330 ***The Dosage Compensated X Chromosomes of C. elegans are Organized in a Network of*** 331 ***Loops***

332 Results described above suggest that non-CTCF proteins can alter CTCF loop strength in
333 mammals, and that non-CTCF loops can be observed in Hi-C data from organisms such as *A.*
334 *aegypti* and *D. melanogaster* and *C. elegans*. These observations prompted us to investigate
335 whether SIP is able to detect loops in Hi-C data from organisms where few loops have been
336 previously detected. Previous Hi-C experiments in *C. elegans* embryos identified interaction
337 domains in the X Chromosomes of hermaphrodites (Crane et al. 2015). These X Chromosomes
338 are bound by a condensin I-containing dosage compensation complex (DCC) that remodels X
339 Chromosome topology and downregulates expression of genes chromosome-wide. This finding
340 represents a significant advance in the understanding of the role of 3D chromatin architecture in
341 the organization of dosage compensated chromosomes. Borders separating these domains on
342 the X Chromosome correspond to binding sites of the specialized condensin I-DCC (Crane et al.
343 2015; Anderson et al. 2019). However, it was difficult to determine whether these domains were
344 formed by self-interactions, as is the case in *Drosophila*, or by point-to-point interactions
345 between DCC sites to form loops by loop extrusion similar to those formed by CTCF and
346 cohesin in mammals. To address this question, we performed Hi-C in *C. elegans* hermaphrodite
347 embryos and used SIP to detect punctate signals (Table S3). Recent experiments performed Hi-
348 C in the same *C. elegans* hermaphrodite embryos (Anderson et al. 2019) and thus we combined
349 Hi-C contacts reported by Anderson et al with ours to obtain over 535 million useable Hi-C
350 contacts. We then used SIP with this combined dataset at 5 kb resolution and we were able to
351 identify 41 loops (Fig. 6A). ChIP-seq for the DPY-27 subunit of condensin I-DCC shows the
352 presence of this protein at loop anchors, suggesting its involvement in the establishment of
353 loops in *C. elegans* (Fig. 6A top track). SIP called zero loops in the Hi-C data for the DCC
354 mutant (Anderson et al. 2019) (Fig. 6B) indicating that the establishment of the 41 loops
355 depends on the presence of DCC. To confirm that these loops are associated with condensin I-
356 DCC, we performed HiChIP using a DPY-27 antibody (Fig. 6C). We detect enrichment of DPY-
357 27 HiChIP signal on SIP loops identified by Hi-C, indicating that condensin I-DCC may play a
358 role in the formation of loops in the X Chromosome of *C. elegans* hermaphrodites (Fig. 6D). In
359 comparison to other loop callers, SIP loops have higher overlap with DPY27 HiChIP (Fig. S6A
360 and S6B). Anderson et al found that deletion of eight *rex* sites at borders of domains results in
361 the loss of these domains with no change to gene expression (Anderson et al. 2019). These
362 deletions overlap with some of the loop anchors we detect (Fig. 6E). We examined average loop
363 signal at sites where one anchor overlaps a deletion and found loss of these loops (Fig. 6F top).
364 However, our HiChIP data indicates that there are many DCC-mediated loop anchors that do
365 not overlap these deletions (Fig. 6E). We examined loops where neither anchor overlaps a
366 deleted site and found that these loops are still present and become stronger (compare Fig. 6F
367 bottom and Fig. 6D). All of these loops are dependent on DCC (Fig. 6F right). These additional

368 DCC-dependent loops could explain the interesting observation that disruption of DCC results in
369 X Chromosome decondensation and increased gene expression while deletion of 8 *rex* sites
370 does not (Anderson et al. 2019). A similar model has been proposed by Anderson et al
371 (Anderson et al. 2019).

372 The presence of condensin I-DCC suggests that loops may be formed via extrusion in *C.*
373 *elegans*. We examined motifs at loop anchors and found the MEX motif, which is enriched at
374 DPY-27 peaks (Jans et al. 2009) and is at *rex* sites previously reported to be present at borders
375 of domains (Crane et al. 2015; Anderson et al. 2019). We then tested whether loops occur
376 between convergent MEX motifs as is the case for CTCF loops in mammals, yet we found no
377 bias in motif orientation (Fig. S6C) consistent with what was reported for domain borders
378 (Anderson et al. 2019). Thus, loop anchors in *C. elegans* likely represent bidirectional blocks to
379 extrusion. In support of this notion, loop anchors generally form interactions both upstream and
380 downstream of the anchor (Fig. 6C and 6E). We detect no loops with visually apparent extrusion
381 stripes in the Hi-C data (Fig. 6A). Stripes in mammalian Hi-C maps indicate unidirectional
382 extrusion starting near one anchor (Vian et al. 2018). We ran molecular dynamics polymer
383 simulations of unidirectional extrusion starting near loop anchors and detect strong stripes at
384 loop anchors which is consistent with an asymmetric extrusion model reported for CTCF loops
385 in mammals (Fig. S6D bottom left) (Vian et al. 2018). We then ran polymer simulations of bi-
386 directional extrusion starting randomly and detect less striping and more filled-in domains (Fig.
387 S6D top right). This pattern is more consistent with the absence of stripes at loops associated
388 with the condensin I-DCC in *C. elegans* (Fig. 6A) and therefore suggests that loading is either
389 random or takes place at many sites, rather than just the high affinity *rex* sites.

390 A chromosome-wide view of DPY-27 HiChIP shows a network of loops spanning the X-
391 Chromosome (Fig. 6E). This loop network can also be seen in scaled metaplots of distance
392 normalized Hi-C data corresponding to DPY-27 peaks (Fig. 6G). The formation of a rosette
393 structure by the compensated X Chromosome is supported by viewing this network as a
394 bullseye plot (Fig. 6H). Since Hi-C and HiChIP are performed on a population of cells, the
395 apparent network of DPY-27 loops could either represent multiway interactions occurring in the
396 same cell or individual two-way interactions occurring in different cells. If all the loops are
397 present simultaneously in all cells, the results would suggest that these nested loops can occur
398 between 5 anchors or more (Fig. 6E and 6G). To distinguish between these two possibilities, we
399 examined Hi-C reads containing multiple interacting fragments. Because the Hi-C protocol
400 involves digestion with DpnII followed by ligation and sonication of fragments for library
401 preparation, sequenced reads can contain several ligation events. Therefore, we examined
402 paired-end reads in which we could determine ligations between three different genomic regions
403 (see Methods). For example, DPY-27 loop anchors at chromosomal coordinates 12.35 Mb,
404 13.70 Mb, and 14.52 Mb were ligated together, indicating that these loops occur within the same
405 cell (Fig. 6I). Analysis of multiway Hi-C interactions shows enrichment of contacts among
406 multiple DPY-27 loop anchors (Fig. 6J). Three-way interactions between DPY-27 loop anchors
407 are 2.4-fold higher ($p < 0.05$) than permutations on sets of random loci at similar distances on
408 the X Chromosome. To improve the ability to detect condensin I-DCC-mediated three-way
409 ligations we examined DPY-27 HiChIP data obtained from 250 bp paired-end sequencing
410 (Supplemental Table S4). We then examined three-way ligations in DPY-27 HiChIP data and
411 found enrichment of multiway DPY-27 anchor interactions compared to Hi-C and compared to
412 random regions (Fig. S6E). Additionally, in a metaplot of all possible three-way interactions
413 between DPY-27 loop anchors and the surrounding regions, we found enrichment at DPY-27
414 loop anchors (Fig. 6K). Altogether, our observations suggest that loops identified by SIP in *C.*
415 *elegans* may represent nested interconnected DCC interactions mediated by condensin I-DCC,

416 implying that the dosage compensated X Chromosome of hermaphrodites is organized in a
417 rosette-like structure.

418

419 **DISCUSSION**

420 Hi-C datasets containing billions of contacts have allowed the identification of thousands of
421 loops representing point-to-point interactions between CTCF sites in mammals (Rao et al.
422 2014). However, there are very few methods capable of identifying these loops and sometimes
423 it has been more feasible to annotate loops by eye (Eagen et al. 2017). SIP utilizes image
424 processing and the local background to identify loops. Here we demonstrate the utility of SIP as
425 a loop caller in identifying additional CTCF loops in mammals, non-CTCF loops in *D.*
426 *melanogaster* and *A. aegypti*, and condensin I-DCC loops in *C. elegans*. The high accuracy of
427 SIP in loop identification allows the detection of nearly double the CTCF loops from the same
428 dataset as well as detection of loops in non-mammalian species. With the companion tool
429 SIPMeta, SIP can facilitate discovery of novel aspects of 3-D chromatin architecture. We intend
430 SIP to be easily useable by anyone performing analysis of Hi-C data on a variety of platforms
431 and have given users the ability to alter most parameters to facilitate custom loop calling.

432 While CTCF has been the major focus of studies of loop formation, other chromatin
433 bound factors may also affect this process. For example, non-CTCF loops are evident in *D.*
434 *melanogaster* and we are also able to detect loops in *A. aegypti* in this study using SIP.
435 Although we cannot determine the nature of the proteins forming loops in *A. aegypti*, these
436 loops appear similar to Pc/Psq loops in *D. melanogaster*, and similar proteins are likely involved
437 in their establishment. In mammals, CTCF loop strength may also be affected by other proteins,
438 since cohesin sliding has been shown to be delayed by other DNA bound complexes *in vitro*,
439 including quantum dot labelled catalytically inactive EcoRI and dCas9 (Davidson et al. 2016;
440 Stigler et al. 2016). Thus, DNA-bound proteins at specific sites in the genome may affect the
441 loop extrusion process and thereby affect loop strength. The involvement of ZNF143 in the
442 establishment of CTCF loops (Wen et al. 2018; Jung et al. 2019) is supported by the increased
443 strength of loops detected by SIP when CTCF and ZNF143 are both present at interacting
444 anchors. The inability of CTCFL to form loops has been confirmed by a recent study that
445 indicates CTCFL cannot stop cohesin extrusion (Pugacheva et al. 2020).

446 Using SIP, we find that condensin I-DCC in *C. elegans* forms dozens of loops along the
447 dosage-compensated X Chromosome of hermaphrodites. Previous studies of the structure of
448 metaphase chromosomes in chicken cells found that condensin II creates the axial scaffold
449 while condensin I creates clusters of nested loops (Gibcus et al. 2018). We speculate that during
450 dosage compensation in *C. elegans*, condensin I-DCC may perform a similar function along the
451 X Chromosome, creating a rosette-like structure and thereby compacting the chromosome
452 sufficiently to decrease transcription by two-fold. This hypothesis agrees with microscopy data
453 showing that the X Chromosome of hermaphrodites occupies a smaller nuclear volume than
454 autosomes, and that mutations in DCC result in decreased compaction (Lau et al. 2014).
455 However, deletion of eight of the *rex* sites that are important for the establishment of contact
456 domains observed in Hi-C data reportedly have no effect on chromosome compaction or gene
457 expression (Anderson et al. 2019). While we find that there are more than eight loop anchors
458 and thus not all loops were lost when removing eight of them, it is curious that removal of a
459 portion of loops did not affect gene expression. It has been suggested that loss of each anchor
460 results in continued extrusion until the next anchor (Anderson et al. 2019). This would keep the
461 overall network or rosette-like structure intact but containing larger loops.

462 CTCF loop anchors in mammals often coincide with a stripe of intense interaction signal
463 in Hi-C maps (Vian et al. 2018). This observation prompted the suggestion of a “loop gun”
464 model of extrusion where cohesin is loaded near loop anchors and proceeds asymmetrically
465 until reaching the other anchor. Our study provides *in vivo* evidence of condensin I mediated
466 extrusion in animals. Unlike the formation of CTCF loops in mammals, our analysis suggests an

467 alternate method of extrusion-mediated looping in *C. elegans*. Our results indicate that SMC
468 proteins are loaded randomly on the X Chromosome of *C. elegans* hermaphrodites and extrude
469 until reaching stopping points from either direction, in a fashion similar to what was suggested
470 by Anderson et al (Anderson et al. 2019). This may occur via a single condensin complex
471 undergoing bi-directional extrusion or by unidirectional extrusion of multiple condensin
472 complexes, recreating a bi-directional effect. While condensin complexes from yeast display
473 unidirectional extrusion and compaction *in vitro* (Ganji et al. 2018; Kong et al. 2019), recent
474 work indicates mammalian condensins perform bi-directional extrusion (Kong et al. 2019).
475 Studies performed *in vitro* suggest that condensin complexes can pass over each other during
476 extrusion and thereby form a z-loop structure (Kim et al. 2019). Therefore, if condensin I-DCC
477 moves unidirectionally, the interaction-enriched interior of the domain may form by randomly
478 placed z-loops in the population of cells. In either case, looping both upstream and downstream
479 of each anchor indicates that extrusion is blocked from both sides. In a simplified three-anchor
480 example where each anchor is a bi-directional block, extrusion on both sides of the middle
481 anchor would naturally cause all three anchors to come into close proximity (Fig. 7A). This
482 indicates that these bi-directional blocks could create the network of interactions observed by
483 Hi-C and HiChIP to form the axis of a rosette like structure (Fig. 7B).

484 SIP and SIPMeta greatly facilitate the analysis of Hi-C data and the detection of point-to-
485 point interactions. Loops formed by these interactions represent an important aspect of the
486 three-dimensional organization of the mammalian genome. Our evaluation indicates that
487 sequencing depth is an important factor in loop calling, thus methods that increase signal by
488 data imputation may become valuable tools (Zhang et al. 2018). As sequencing costs decrease
489 and high resolution Hi-C datasets become standard, memory and time inefficient methods will
490 perform worse as the matrix processing becomes more complex. However, using images
491 instead of matrices along with image processing should limit the increased memory costs
492 associated with deeper sequencing. The ability to call loops and quantitatively measure their
493 strength using SIP will facilitate the discovery of the biological significance of 3D nuclear
494 architecture.

495

496

497 **METHODS**498 ***The SIP program and the loop calling process***

499 SIP retrieves raw Hi-C signal stored in Juicer .hic files using Juicer Tools (Durand et al. 2016) at
 500 the resolution and with the normalization scheme chosen by the user. The genome is analyzed
 501 by sliding windows, the size of which depends on the resolution and matrix size specified by the
 502 user. For example, we used 5 kb resolution data with KR normalization and a matrix size of
 503 2000 for all experiments involving GM12878 or HCT116 cells. This creates 10 Mb snapshots
 504 sliding over 5 Mb each step. Observed-expected (o-e) values are used to create images. Later,
 505 in postfiltering, retrieved data is distance normalized by the formula $value_{normalized} = 1 +$
 506 $((value - expected)/expected + 1)$, which is used to compute the central loop value. SIP then
 507 uses image processing methods to create a list of candidate loops that will be filtered later.
 508 Because even with o-e normalized values the diagonal represents extremes in the data, outliers
 509 within 2 bins along the diagonal are removed if the value is higher than the average + 1 std dev
 510 of the image signal. The first image processing step utilizes gaussian blurring to smooth the Hi-
 511 C signal in order to avoid detection of outlier pixel signals. Afterwards, contrast enhancement is
 512 used to increase the contrast between the background and the signal of interest (Schneider et
 513 al. 2012). White top-hat, a mathematical morphology method from the MorpholibJ plugin, is
 514 used to homogenize the background and make bright structures easier to detect (Legland et al.
 515 2016). Because loops appear as bright punctate signal in images, we use this top-hat method to
 516 transform the grey-scale intensity values of each Hi-C image, causing the bright structures to
 517 have increased contrast from the background. The last step uses a Minimum and Maximum
 518 Filter (Schneider et al. 2012) combination to remove isolated pixels and further homogenize the
 519 background. These steps provide a corrected image of the interactions (Fig. 1). This corrected
 520 image is then used with the regional maxima detection algorithm available from ImageJ
 521 (Schneider et al. 2012) to detect a long list of candidate loops.

522 Candidate loops must then pass several filters that utilize the distance normalized signal
 523 from the original matrix before image processing (Fig. 1). The first step is to exclude pixels near
 524 columns and rows with insufficient data; the default is to filter any with ≥ 6 pixels with zero
 525 values in the surrounding 24-pixel neighborhood. The second filter is to remove pixels without
 526 increased interactions compared to the surrounding 8-pixel neighborhood and the 24-pixel
 527 neighborhood. To remove isolated enriched pixels, loops must display decay between the
 528 central pixel and the surrounding neighborhood pixels. Candidate loops are then filtered so that
 529 the center pixel's KR value $\geq .3$ and > 1.2 -fold higher than nearby pixels (PA score). Loops are
 530 then filtered such that the probability that the Poisson CDF function of the center pixel being
 531 higher than the nearby pixels is greater than 0.9. Finally, candidate loops are filtered if their PA
 532 score is lower than the PA scores of a top percentage of random sites. This percentage is
 533 specified by the user.

534 Parameters for SIP loop calling in *D. melanogaster* used a threshold of 6000, with -
 535 nbZero 10, matrix size 500, resolution 1 kb, -d 20, -fdr 0.05, and -isDroso true. Analysis of *A.*
 536 *aegypti* Hi-C data was performed using parameters -g 1.5, -mat 2000, -d 5 -res 5 kb, -t 5000, -
 537 nbZero 4, -fdr 0.05, and -isDroso true. Human Hi-C maps were originally published with
 538 genome build hg19, and we ensured that all data used was mapped to the same genome build.
 539 Remapping everything to GRCh38 will not alter these results. CTCF motif enrichment was
 540 calculate in 5 kb windows using the formula: $Log_2 \left(\frac{ObservedOverlap}{ObservedNonoverlap} \right) / \left(\frac{ExpectedOverlap}{ExpectedNonoverlap} \right)$.
 541 Expected values were derived using random loci.

542 ***Choosing Parameters***

543 SIP was designed for quick and memory efficient loop calling so that loops can be
544 visually inspected for parameter optimization. While we have listed the specific parameters that
545 we used for each map, we recommend users to optimize loop calls using their own criteria. We
546 recommend calling loops at 5 kb, but if sequencing depth is limited, it may be advantageous to
547 call loops at 5 kb, 10 kb, and 25 kb. We suggest using KR normalization (Rao et al. 2014), but
548 depending on sequencing depth, this normalization scheme may not be available for all
549 chromosomes. In this case, other normalization schemes included in the Juicer tool set such as
550 VC_SQRT (Durand et al. 2016) are acceptable alternatives. Users may also wish to alter the `-d`
551 option, which removes signal near the diagonal, depending on the diagonal signal specific to the
552 Hi-C map or especially if using resolutions other 5 kb. For example, the default `-d 6` will remove
553 interactions at less than 30 kb at 5 kb resolution, but 60 kb at 10 kb resolution.

554 The parameters we recommend altering when optimizing loop calls are the `-g` and `-fdr`
555 options. Raising `-g` will increase the blur for the initial loop calls thereby filtering out more
556 isolated speckles that are potentially not true loops. However, this can also blur actual looping
557 signal. Because loops appear more punctate at 10 kb and 25 kb, we suggest decreasing `-g` to
558 reduce the blur and thereby identify more speckled signal. Altering `-fdr` will change how many
559 loops pass the second filter. As SIP is processing, it outputs the number of loops identified
560 before `fdr` filtering so that users can determine how many spots identified by the first pass are
561 filtered by the second pass. This can serve as a gauge for altering the `fdr` parameter. The final
562 parameter we recommend changing is `-nbZero` which filters pixels near areas with low
563 coverage. If loops are erroneously identified near repetitive regions, we suggest increasing `-`
564 `nbZero`. We recommend optimizing the SIP parameters by visual inspection of a single
565 chromosome first, and then using those parameters for all the chromosomes.

566 ***Performance Testing***

567 Comparison of loops between datasets containing various numbers of Hi-C contacts was
568 performed by random picking intra-chromosomal reads. Bootstrapping was performed by down
569 sampling the full dataset to 1 billion intrachromosomal reads 10 different times. Noise was
570 simulated to follow the same distance decay as the Hi-C data. Additional details can be found in
571 the Supplemental Material. Recovery rates were calculated by the number of loops obtained in
572 the down-sampled or noise-added datasets that were within two pixels of the loops identified in
573 the full dataset. False positive rates were calculated under the assumption that loops called in
574 the down-sampled or noise-added datasets that do not overlap with loops in the full dataset are
575 false.

576 ***SIPMeta***

577 SIPMeta is implemented in java and includes a choice between command line options or a
578 graphical user interface (GUI). SIPMeta first reads a loop file from which the bin size (resolution)
579 is inferred. If the images are not present in the input directory, SIPMeta makes images from the
580 SIP-derived bedpe file corresponding to distance normalized signals. Alternatively, users can
581 specify a .hic file from which values are retrieved. Then SIPMeta examines all signal within a
582 specified distance surrounding the loop, computes an APA score as previously described (Rao
583 et al. 2014), and outputs a matrix of averaged values. This matrix can be run through
584 `bullseye.py` to obtain both square and bullseye plots.

585 The bullseye transformation of a heatmap is a visualization technique intended to more
586 accurately represent the secondary interactions around a strong loop in the genome. The plot is
587 a simple transformation of the rectangular heatmap such that each bin's Euclidean distance to
588 the center now directly corresponds to its Manhattan distance in the original map. Each ring in
589 the bullseye plot has segments corresponding to the $4 \times N$ bins with a Manhattan distance of N

590 from the central bin. Each bin in a ring takes up exactly the same angular area and they are
591 evenly distributed around the circle. Although this represents some distortion from their actual
592 angles in the original plot, this creates the same visual area for each bin. Z-score transformation
593 is done for each ring separately and the ADA score is obtained by percentage of Z-scores > 1 in
594 the bottom right quarter vs the total plot.

595 **Contribution of Transcription Factors to Loop Strength**

596 Overlaps between transcription factors and loops anchor sites were assigned if ChIP-seq peaks
597 were within two pixels of the loop anchors. Loops categorized as overlapping high ChIP-seq
598 peaks were those where both anchors overlap peaks in the top quartile of ChIP-seq signal.
599 Loops where anchors do not overlap peaks in either of the top two quartiles were categorized as
600 low. Loops were also separated into 10 equal categories based on motif scores overlapping the
601 two anchors. Because the purpose of the test is to approximate the role of the motif in loop
602 strength, the lower of the two motif values corresponding to the two anchors was assigned as
603 the motif score.

604 **Hi-C and HiChIP in *C. elegans***

605 Hi-C and HiChIP libraries were prepared as previously described (Crane et al. 2015; Rowley et
606 al. 2019); details can be found in the Supplemental Material. Two biological replicates were
607 obtained for Hi-C and HiChIP experiments and processed using Juicer (Durand et al. 2016) with
608 the ce10 genome. Mapping statistics can be found in Supplemental Tables S1 and S2.

609 Loops in *C. elegans* Hi-C were identified using SIP with parameters -g 1.5, -d 5, -fd
610 0.05, -res 5000 -mat 500. Because DPY-27 HiChIP data showed enrichment across the X
611 Chromosome, we identified potential anchors by peaks in the coverage normalization vector.
612 Network metaplots were done by taking every anchor with the closest five others and scaling
613 the region in between each anchor as well as the same distance upstream and downstream of
614 the first and last anchors. Median Hi-C or HiChIP signal was profiled within these regions.
615 Bullseye plots were generated from this scaled matrix re-centered on the a2-a4 matrix
616 coordinate. Polymer simulations were performed as described (Vian et al. 2018).

617 Three-way interactions were obtained by scanning Hi-C or HiChIP FASTQ files for the
618 ligation sequence GATCGATC and mapping each side to the ce10 genome using Bowtie 2
619 (Langmead and Salzberg 2012). Paired-end reads with at least three different sections mapping
620 to different genomic locations at least 50 kb apart were used in downstream analysis. Overlaps
621 were done with all possible three-way combinations of DPY-27 loop anchors in 10 kb bins or
622 with the same number of random 10 kb bins following the same distance distribution. p-values
623 were derived from Monte-Carlo permutations.

624 **DATA ACCESS**

625 All raw and processed sequencing data generated in this study is available at the NCBI Gene
626 Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>) under accession number
627 GSE132640. The latest release of SIP can be obtained from
628 <https://github.com/PouletAxel/SIP/releases> and SIPMeta from
629 <https://github.com/PouletAxel/SIPMeta/releases> including usage documentation and a separate
630 script for the bullseye transformation of matrices. Source code for SIP and SIPMeta, including
631 bullseye.py, is also provided as Supplemental Code. Any issues with the program can be
632 reported on GitHub or directly via e-mail.

633 **ACKNOWLEDGMENTS**

634 We would like to thank the HudsonAlpha Institute for Biotechnology Genomic Services Lab for
635 their help in Illumina Sequencing. We also thank Drs. William Noble and Doug Phanstiel for

636 constructive discussions during the optimization of SIP. This work was supported by NIH
 637 Pathway to Independence Award K99/R00 GM127671 (M.J.R.) and U.S. Public Health Service
 638 Award R01 GM035463 (V.G.C.) from the National Institutes of Health. BJB was supported by
 639 NIH T32 GM008490. The content is solely the responsibility of the authors and does not
 640 necessarily represent the official views of the National Institutes of Health.

641 DISCLOSURE DECLARATION

642 The authors declare no competing interests.

643 FIGURE LEGENDS

644 **Figure 1.** Overview of SIP. (A) The graphical user interface provides options to specify a Juicer
 645 derived .hic file or processed data, the output directory, and chromosome size file. It also allows
 646 adjustment of the various parameters shown. (B) SIP uses image-based detection to create a
 647 long list of candidate loops, which is further filtered based on properties of the distance
 648 normalized matrix.

649 **Figure 2.** Performance of SIP. (A) CPU usage (orange) and GC - garbage collection (blue) over
 650 time using 2 cores during SIP loop calling. (B) Memory usage of SIP during loop calling. (C)
 651 Fraction of loops called using the full dataset recovered by SIP (green) or by HiCCUPS (blue) in
 652 data down-sampled to different sequencing depths. (D) Ratio of false positives (loops not
 653 identified in the full dataset) vs loops recovered by SIP (green) or HiCCUPS (blue) in down-
 654 sampled data. (E) Number of loops identified in down-sampled data (y-axis) for SIP (green) and
 655 HiCCUPS (blue) when parameters were adjusted to give the same false-positive/recovery rate
 656 (x-axis). (F) Percentage of loops identified by SIP (green) or HiCCUPS (blue) in GM12878 cells
 657 that were identified in a different cell type. (G) Percentage of loops identified in each
 658 permutation down-sampling data (purple) vs new loops (i.e. false positives) (teal). Bars
 659 represent averages of 10 permutations with error bars representing standard deviation.

660 **Figure 3.** SIP and SIPMeta can be used to Detect and Analyze Loops in Different Species. (A)
 661 Example locus for SIP loops detected in Hi-C for *D. melanogaster*. (B) Left: Metaplot of SIP
 662 loops illustrating that the Manhattan distance between a and b is different with respect to loops
 663 but is visually depicted as the same distance in square metaplots. L = left anchor, R = right
 664 anchor. Right: Metaplot of SIP loops illustrating the bullseye transformation performed by
 665 SIPMeta. (C) Left: Example locus in GM12878 cells for a stripe detected in Hi-C maps of
 666 mammals. Right: SIPMeta plots displaying how stripes appear in square vs bullseye plots. (D)
 667 Enrichment of CTCF motifs on loop anchors in *D. melanogaster*, *A. aegypti*, and *H. sapiens*.

668 **Figure 4.** Characterization of CTCF Loops with SIPMeta. (A) Number of loops in GM12878 cells
 669 at 5 kb resolution identified by SIP (green) or by HiCCUPS (blue) corresponding to CTCF sites
 670 in convergent, tandem, or in orientations that could not be assigned. The number of loops in
 671 divergent orientation were negligible and could not be depicted. (B,C) SIPMeta bullseye plots
 672 (B) or Z-score plots (C) of SIP loops in categories based on CTCF motif orientation.
 673 APA=Aggregate Peak Analysis. ADA=Aggregate Domain Analysis. (D) Example of a SIP loop
 674 unassignable to any CTCF orientation. CTCF ChIP-seq signal is shown below. Arrows indicate
 675 loop anchors. (E) SIPMeta bullseye plots of CTCF HiChIP data for SIP loops in categories
 676 based on CTCF motif orientation. (F) Metaplots for SIP loops in each CTCF category in control,
 677 cohesin depletion, and recovery in Hi-C data obtained in HCT-116 cells. (G) APA score changes
 678 after cohesin recovery relative to control and cohesin depletion Hi-C. (H) Average APA scores
 679 of convergent CTCF loops divided into 10 equal categories based on the strength of motifs
 680 found on loop anchors. Error bars indicate standard deviation. (I) CTCF motif scores on
 681 convergent vs tandem loops. (J) APA score changes after cohesin recovery relative to control

682 and cohesin depletion Hi-C for convergent loops with the strongest (red) and weakest (blue)
683 10% motif scores.

684 **Figure 5.** Contribution of Transcription Factors to CTCF Loops. (A-C) SIPMeta bullseye plots for
685 loops with either high or low CTCF ChIP-seq signal on both anchors and either high or low
686 ZNF143 (A), YY1 (B), or RNAPII (C) ChIP-seq signal on both anchors. Scores represent
687 average distance normalized Hi-C signal of the loop at the center of the bullseye plot. (D)
688 Number ChIP-seq peaks for each strength quartile that overlaps loop anchors in K562 cells. (E)
689 SIPMeta bullseye plots for loops with either high or low CTCF and CTCFL ChIP-seq signal on
690 both anchors. N.A. indicates an insufficient number of unambiguous loops in this category.
691 Scores represent average distance normalized Hi-C signal of the loop at the center of the
692 bullseye plot.

693 **Figure 6.** A Network of condensin I-DCC Loops in *C. elegans*. (A) Hi-C contact map showing
694 domains and loops on the X Chromosome of *C. elegans* hermaphrodites. Black squares with
695 arrows depict loops called by SIP. Top track displays the DPY-27 ChIP-seq signal across the
696 region. (B) Hi-C in the DCC mutant *sdc-2* (*y93, RNAi*) from Anderson et al (Anderson et al.
697 2019) showing lack of loops on the X Chromosome. (C) DPY-27 HiChIP contact map depicting
698 the enrichment of loops. Top track displays the DPY-27 ChIP-seq signal across the region. (D)
699 SIPMeta bullseye plot of *C. elegans* SIP loops on the X Chromosome displaying the average
700 wild-type Hi-C (top left), or DCC mutant Hi-C (top right), or DPY-27 HiChIP (bottom) signal. (E)
701 X Chromosome-wide view of DPY-27 HiChIP signal. Bottom track displays the DPY-27 ChIP-
702 seq signal across the region. Xs indicate sites deleted in Anderson et al. (F) SIPMeta bullseye
703 plots of Hi-C data after eight *rex* site deletions (left) or after mutation of the DCC (right). SIP
704 loops that overlap with deleted *rex* sites (top) or do not overlap with deleted *rex* sites (bottom).
705 (G) Scaled metaplots of interactions between every DPY-27 loop anchor with its closest four
706 others shown by Hi-C (top right) and by DPY-27 HiChIP (bottom left). The top and side tracks
707 depict the median DPY-27 ChIP-seq signal. Blue circle indicates the point chosen as the center
708 of bullseye plots shown later. (H) SIPMeta bullseye plots for DPY-27 HiChIP (top) and Hi-C
709 (bottom) centered on the interaction between a2-a4. (I) DPY-27 HiChIP contact map depicting a
710 network of two-way interactions between three anchors found to participate in three-way
711 interactions. Bottom track shows DPY-27 ChIP-seq signal. 3D scatterplot of three-way
712 interactions for the chromosome coordinates shown. (J) Number of three-way interactions
713 discovered by Hi-C connecting DPY-27 loop anchors or the average of permutations using an
714 equal number of random regions on the X Chromosome. * indicates $p < .05$ Monte-Carlo
715 permutation test. (K) Profile of three-way interactions across all possible three-way DPY-27 loop
716 anchor connections.

717 **Figure 7.** Model of Condensin Extrusion in *C. elegans*. (A) Extrusion in the X Chromosome of
718 *C. elegans* likely begins at random locations and proceeds until blocked. Depicted here is bi-
719 directional extrusion through one complex, but unidirectional extrusion through multiple
720 complexes is also possible. (B) Loop anchors for DCC in *C. elegans* represent bi-directional
721 blocks resulting in proximity of each anchor with every other.

722

723 REFERENCES

- 724 Anderson EC, Frankino PA, Higuchi-Sanabria R, Yang Q, Bian Q, Podshivalova K, Shin A,
725 Kenyon C, Dillin A, Meyer BJ. 2019. X Chromosome Domain Architecture Regulates
726 *Caenorhabditis elegans* Lifespan but Not Dosage Compensation. *Dev Cell* **51**: 192-207.
- 727 Ay F, Bailey TL, Noble WS. 2014. Statistical confidence estimation for Hi-C data reveals
728 regulatory chromatin contacts. *Genome Res* **24**: 999-1011.
- 729 Cao Y, Chen Z, Chen X, Ai D, Chen G, McDermott J, Huang Y, Xiaoxiao G, Han JJ. 2019.
730 Accurate loop calling for 3D genomic data with cLoops. *Bioinformatics*
731 doi:10.1093/bioinformatics/btz651.
- 732 Crane E, Bian Q, McCord RP, Lajoie BR, Wheeler BS, Ralston EJ, Uzawa S, Dekker J, Meyer
733 BJ. 2015. Condensin-driven remodelling of X chromosome topology during dosage
734 compensation. *Nature* **523**: 240-244.
- 735 Cubeñas-Potts C, Rowley MJ, Lyu X, Li G, Lei EP, Corces VG. 2016. Different enhancer
736 classes in *Drosophila* bind distinct architectural proteins and mediate unique chromatin
737 interactions and 3D architecture. *Nucleic Acids Research* **45**: 1714-1730.
- 738 Davidson IF, Goetz D, Zaczek MP, Molodtsov MI, Huis In 't Veld PJ, Weissmann F, Litos G,
739 Cisneros DA, Ocampo-Hafalla M, Ladurner R et al. 2016. Rapid movement and
740 transcriptional re-localization of human cohesin on DNA. *The EMBO journal* **35**: 2671-
741 2685.
- 742 Durand NC, Shamim MS, Machol I, Rao SSP, Huntley MH, Lander ES, Aiden EL. 2016. Juicer
743 Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell*
744 *Systems* **3**: 95-98.
- 745 Eagen KP, Aiden EL, Kornberg RD. 2017. Polycomb-mediated chromatin loops revealed by a
746 subkilobase-resolution chromatin interaction map. *Proceedings of the National Academy*
747 *of Sciences* **114**: 8764-8769.
- 748 Ganji M, Shaltiel IA, Bisht S, Kim E, Kalichava A, Haering CH, Dekker C. 2018. Real-time
749 imaging of DNA loop extrusion by condensin. *Science* **360**: 102-105.
- 750 Gibcus JH, Samejima K, Goloborodko A, Samejima I, Naumova N, Nuebler J, Kanemaki MT,
751 Xie L, Paulson JR, Earnshaw WC et al. 2018. A pathway for mitotic chromosome
752 formation. *Science* doi:10.1126/science.aao6135: eao6135.
- 753 Gutierrez-Perez I, Rowley MJ, Lyu X, Valadez-Graham V, Vallejo DM, Ballesta-Illan E, Lopez-
754 Atalaya JP, Kremisky I, Caparros E, Corces VG et al. 2019. Ecdysone-induced 3D
755 chromatin reorganization involves active enhancers bound by Pipsqueak and Polycomb.
756 *Cell Reports* **28**: 2715-2727.
- 757 Heger P, Marin B, Bartkuhn M, Schierenberg E, Wiehe T. 2012. The chromatin insulator CTCF
758 and the emergence of metazoan diversity. *Proceedings of the National Academy of*
759 *Sciences* **109**: 17507-17512.
- 760 Heinz S, Texari L, Hayes MGB, Urbanowski M, Chang MW, Givarkes N, Rialdi A, White KM,
761 Albrecht RA, Pache L et al. 2018. Transcription Elongation Can Affect Genome 3D
762 Structure. *Cell* **174**: 1522-1536 e1522.
- 763 Jans J, Gladden JM, Ralston EJ, Pickle CS, Michel AH, Pferdehirt RR, Eisen MB, Meyer BJ.
764 2009. A condensin-like dosage compensation complex acts at a distance to control
765 expression throughout the genome. *Genes Dev* **23**: 602-618.
- 766 Jung YH, Kremisky I, Gold HB, Rowley MJ, Punyawai K, Buonanotte A, Lyu X, Bixler BJ, Chan
767 AWS, Corces VG. 2019. Maintenance of CTCF- and Transcription Factor-Mediated
768 Interactions from the Gametes to the Early Mouse Embryo. *Mol Cell* **75**: 154-171.
- 769 Kim E, Kerssemakers J, Shaltiel IA, Haering CH, Dekker C. 2019. DNA-loop extruding
770 condensin complexes can traverse one another. *bioRxiv*
771 doi:<https://doi.org/10.1101/682864>

- 772 Kong M, Cutts E, Pan D, Beuron F, Thangavelu K, Xue C, Morris E, Musacchio A, Vannini A,
773 Greene EC. 2019. Human condensin I and II drive extensive ATP-dependent
774 compaction of nucleosome-bound DNA. *BioRxiv* doi:<https://doi.org/10.1101/683540>
775 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:
776 357-359.
- 777 Lau AC, Csankovszki G. 2014. Condensin-mediated chromosome organization and gene
778 regulation. *Front Genet* **5**: 473.
- 779 Lau AC, Nabeshima K, Csankovszki G. 2014. The *C. elegans* dosage compensation complex
780 mediates interphase X chromosome compaction. *Epigenetics Chromatin* **7**: 31.
- 781 Legland D, Arganda-Carreras I, Andrey P. 2016. MorphoLibJ: integrated library and plugins for
782 mathematical morphology with ImageJ. *Bioinformatics* **32**: 3532-3534.
- 783 Matthews BJ, Dudchenko O, Kingan SB, Koren S, Antoshechkin I, Crawford JE, Glassford WJ,
784 Herre M, Redmond SN, Rose NH et al. 2018. Improved reference genome of *Aedes*
785 *aegypti* informs arbovirus vector control. *Nature* **563**: 501-507.
- 786 Ogiyama Y, Schuettengruber B, Papadopoulos GL, Chang JM, Cavalli G. 2018. Polycomb-
787 Dependent Chromatin Looping Contributes to Gene Silencing during *Drosophila*
788 Development. *Mol Cell* **71**: 73-88.
- 789 Pugacheva EM, Kubo N, Loukinov D, Tajmul M, Kang S, Kovalchuk AL, Strunnikov AV, Zentner
790 GE, Ren B, Lobanenko VV. 2020. CTCF mediates chromatin looping via N-terminal
791 domain-dependent cohesin retention. *Proc Natl Acad Sci U S A* **117**: 2020-2031.
- 792 Pugacheva EM, Rivero-Hinojosa S, Espinoza CA, Méndez-Catalá CF, Kang S, Suzuki T,
793 Kosaka-Suzuki N, Robinson S, Nagarajan V, Ye Z et al. 2015. Comparative analyses of
794 CTCF and BORIS occupancies uncover two distinct classes of CTCF binding genomic
795 regions. *Genome Biology* **16**, 161 [https://doi-](https://doi-org.proxy.library.emory.edu/10.1186/s13059-015-0736-8)
796 [org.proxy.library.emory.edu/10.1186/s13059-015-0736-8](https://doi-org.proxy.library.emory.edu/10.1186/s13059-015-0736-8).
- 797 Rao S, Huang S-C, Glenn St. Hilaire B, Engreitz JM, Perez EM, Kieffer-Kwon K-R, Sanborn AL,
798 Johnstone SE, Bochkov ID, Huang X et al. 2017. Cohesin Loss Eliminates All Loop
799 Domains. *Cell* **171**: 305-320.
- 800 Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL,
801 Machol I, Omer AD, Lander ES et al. 2014. A 3D map of the human genome at kilobase
802 resolution reveals principles of chromatin looping. *Cell* **159**: 1665-1680.
- 803 Rowley MJ, Corces VG. 2018. Organizational principles of 3D genome architecture. *Nat Rev*
804 *Genet* **19**: 789-800.
- 805 Rowley MJ, Lyu X, Rana V, Ando-Kuri M, Karns R, Bosco G, Corces VG. 2019. Condensin II
806 Counteracts Cohesin and RNA Polymerase II in the Establishment of 3D Chromatin
807 Organization. *Cell Rep* **26**: 2890-2903.
- 808 Rowley MJ, Nichols MH, Lyu X, Ando-Kuri M, Rivera ISM, Hermetz K, Wang P, Ruan Y, Corces
809 VG. 2017. Evolutionarily Conserved Principles Predict 3D Chromatin Organization.
810 *Molecular Cell* **67**: 837-852.
- 811 Schneider CA, Rasband WS, Eliceiri KW. 2012. NIH Image to ImageJ: 25 years of image
812 analysis. *Nat Methods* **9**: 671-675.
- 813 Stigler J, Çamdere GÖ, Koshland DE, Greene EC. 2016. Single-Molecule Imaging Reveals a
814 Collapsed Conformational State for DNA-Bound Cohesin. *Cell Reports* **15**: 988-998
815 doi:10.1016/j.celrep.2016.04.003.
- 816 Tang Z, Luo Oscar J, Li X, Zheng M, Zhu Jacqueline J, Szalaj P, Trzaskoma P, Magalska A,
817 Włodarczyk J, Rusczycki B et al. 2015. CTCF-Mediated Human 3D Genome
818 Architecture Reveals Chromatin Topology for Transcription. *Cell* **163**: 1611-1627.
- 819 Vian L, Pekowska A, Rao SSP, Kieffer-Kwon KR, Jung S, Baranello L, Huang SC, El Khattabi L,
820 Dose M, Pruett N et al. 2018. The Energetics and Physiological Impact of Cohesin
821 Extrusion. *Cell* **175**: 292-294.

- 822 Wen Z, Huang ZT, Zhang R, Peng C. 2018. ZNF143 is a regulator of chromatin loop. *Cell Biol*
823 *Toxicol* **34**: 471-478.
- 824 Zhang Y, An L, Xu J, Zhang B, Zheng WJ, Hu M, Tang J, Yue F. 2018. Enhancing Hi-C data
825 resolution with deep convolutional neural network HiCPlus. *Nat Commun* **9**: 750 doi:
826 10.1038/s41467-018-03113-2.
- 827

A

Program choice: hic processed

Data type: hic HiChIP

Juicer Tools

Normalization scheme (prefers KR): NONE KR VC VC SQRT

Data and Output directories:

Data (hic or SiP)

Output directory

Matrix parameters:

Matrix size (in bins): resolution*2

Diag size (in bins): resolution*5

Resolution (in bases):

Multi resolution loop calling:

resolution*2

resolution*5

Chromosome size file (same chr names as in .hic file):

Chr size file

Image processing parameters:

Gaussian filter: Threshold for maxima detection:

Maximum filter: Minimum filter:

% of saturated pixel: Empirical FDR:

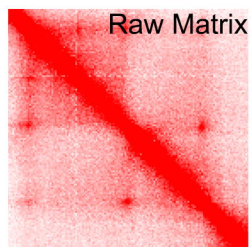
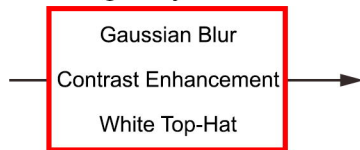
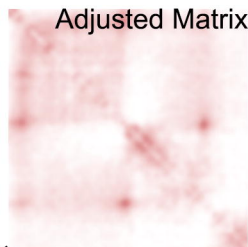
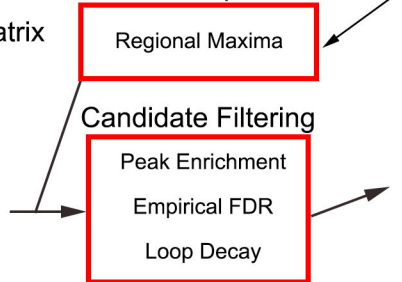
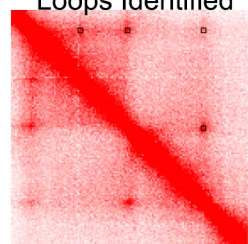
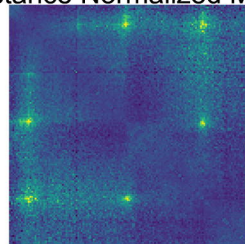
Number of zeros allowed in the 24 surrounding pixels:

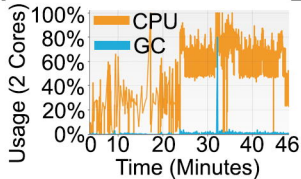
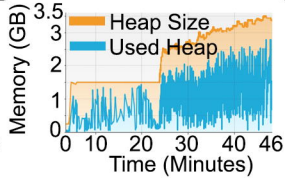
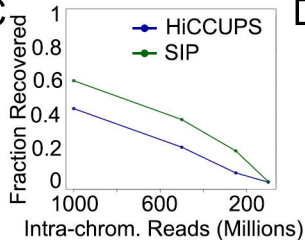
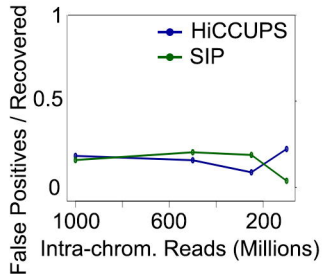
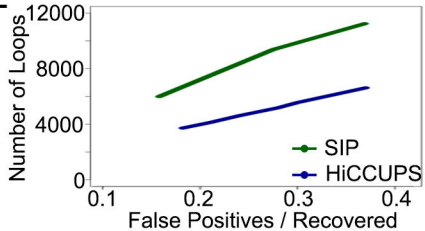
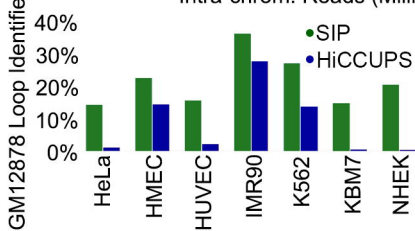
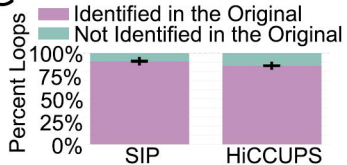
If is droso or like droso HiC map:

Is Droso

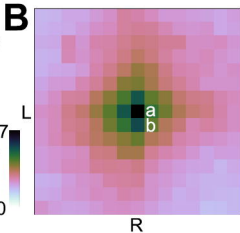
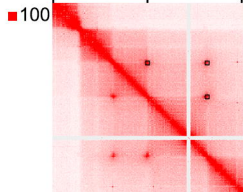
Number of CPU: Delete tif files

Quit Start

B**Image Adjustment****Adjusted Matrix****Candidate Loop Detection****Loops Identified****Distance Normalized Matrix**

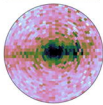
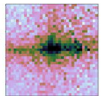
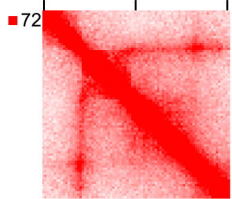
A**B****C****D****E****F****G**

A Chr 2L *D. melanogaster*
 21.7 Mb 21.9 Mb 22.1 Mb



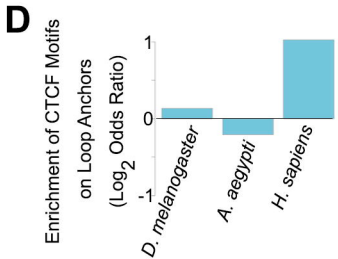
a. 1 kb from R + 0 kb from L = 1 kb from Loop
 b. 1 kb from R + 1 kb from L = 2 kb from Loop

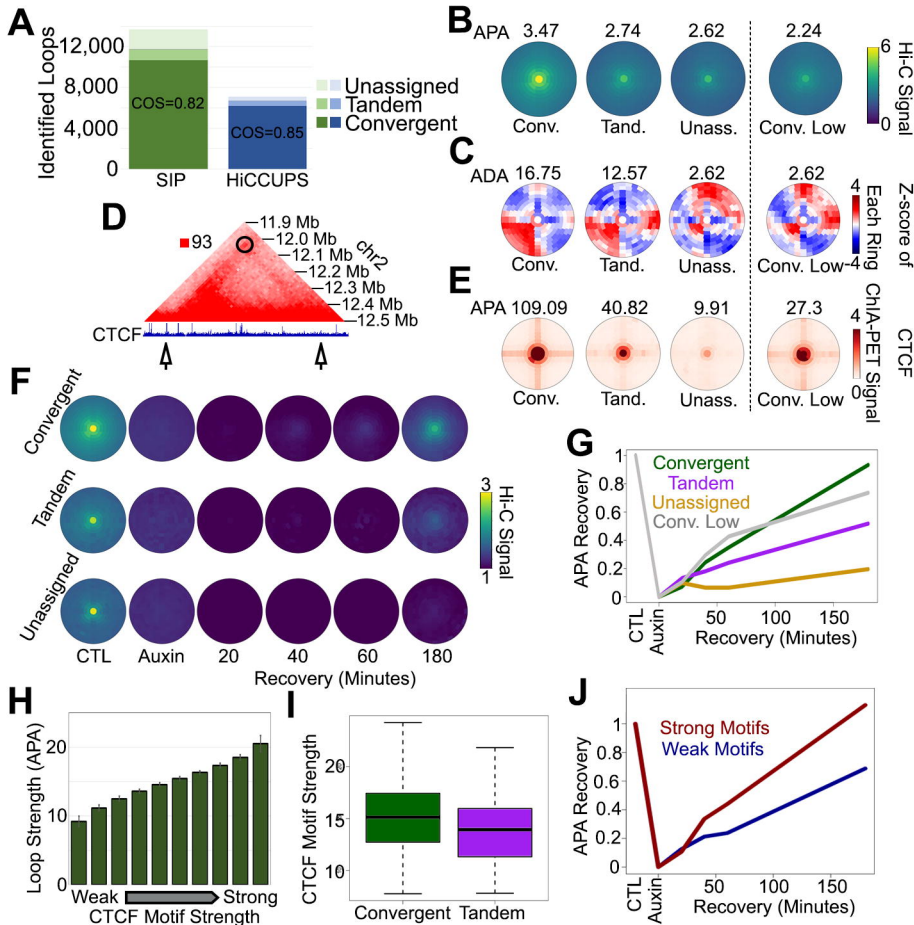
C Chr 2 *H. sapiens*
 223.1 Mb 223.3 Mb 223.5 Mb

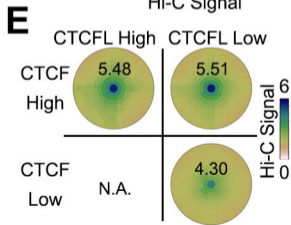
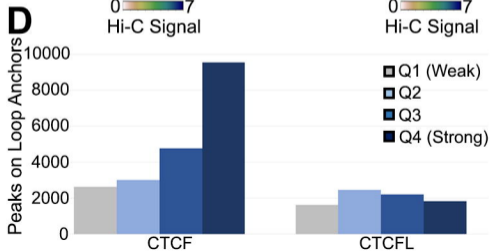
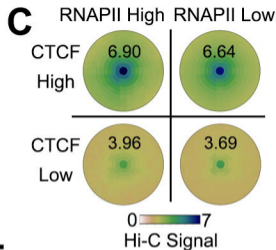
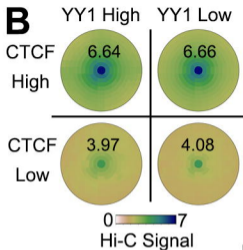
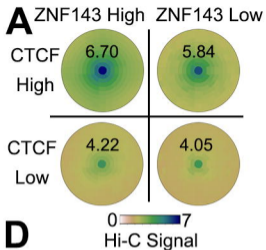


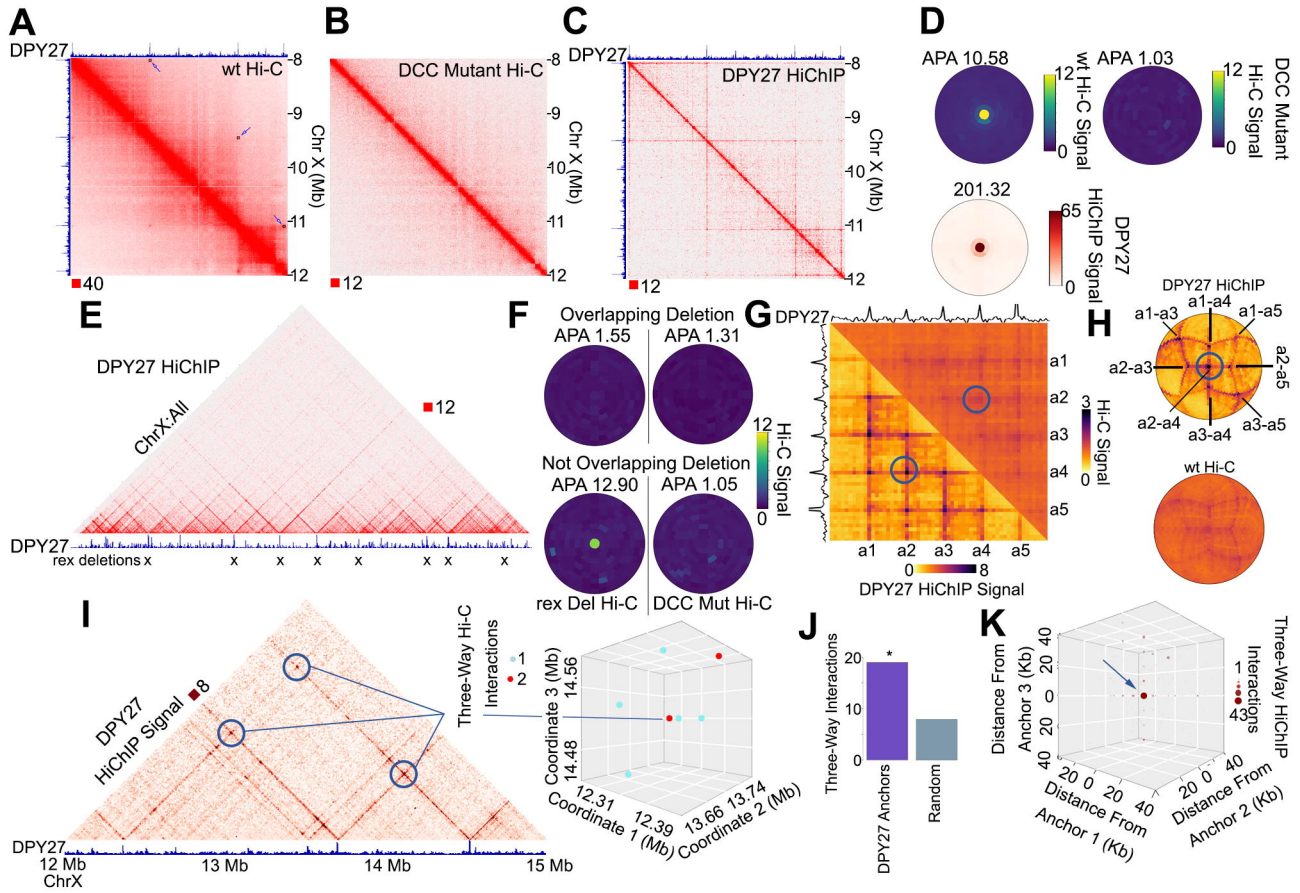
Hi-C Signal

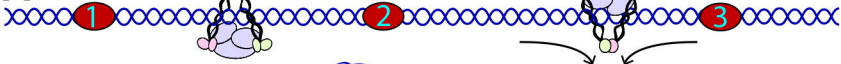
0 6









A**B**