

1 **Predicting CTCF-mediated chromatin interactions by integrating genomic**  
2 **and epigenomic features**

3

4 Yan Kai<sup>1,2,3</sup>, Jaclyn Andricovich<sup>2,3</sup>, Zhouhao Zeng<sup>1</sup>, Jun Zhu<sup>4</sup>, Alexandros Tzatsos<sup>2,3,\*</sup>,  
5 and Weiqun Peng<sup>1,\*</sup>

6 <sup>1</sup>Department of Physics, George Washington University (GWU), Washington DC 20052  
7 and

8 <sup>2</sup>Cancer Epigenetics Laboratory, Department of Anatomy and Regenerative Biology,  
9 GWU, and

10 <sup>3</sup>GWU Cancer Center, GWU School of Medicine and Health Sciences, Washington DC  
11 20052

12 <sup>4</sup>Systems Biology Center, National Heart Lung and Blood Institute, National Institute of  
13 Health, Bethesda, MD 20892

14 \*Correspondence: [atzatsos@gwu.edu](mailto:atzatsos@gwu.edu), [wpeng@gwu.edu](mailto:wpeng@gwu.edu)

15

## Abstract

16 The CCCTC-binding zinc finger protein (CTCF)-mediated network of long-range chromatin  
17 interactions is important for genome organization and function. Although this network has been  
18 considered largely invariant, we found that it exhibits extensive cell-type-specific interactions  
19 that contribute to cell identity. Here we present Lollipop—a machine-learning framework—which  
20 predicts CTCF-mediated long-range interactions using genomic and epigenomic features. Using  
21 ChIA-PET data as benchmark, we demonstrated that Lollipop accurately predicts CTCF-  
22 mediated chromatin interactions both within and across cell-types, and outperforms other  
23 methods based only on CTCF motif orientation. Predictions were confirmed computationally and  
24 experimentally by Chromatin Conformation Capture (3C). Moreover, our approach reveals novel  
25 determinants of CTCF-mediated chromatin wiring, such as gene expression within the loops.  
26 Our study contributes to a better understanding about the underlying principles of CTCF-  
27 mediated chromatin interactions and their impact on gene expression.

## 28 Introduction

29  
30 Higher-order chromatin structure plays a critical role in gene expression and cellular  
31 homeostasis<sup>1, 2, 3, 4, 5, 6, 7</sup>. Genome-wide profiling of long-range interactions in multiple cell-types  
32 revealed that CCCTC-binding factor (CTCF) binds at loop anchors and delimits the boundaries  
33 of Topologically Associating Domains (TADs)<sup>8, 9, 10, 11</sup>, suggesting that CTCF plays a central role  
34 in regulating the organization and function of the 3D genome<sup>12, 13</sup>. Depletion of CTCF revealed  
35 that it is required for chromatin looping between its binding sites and insulation of TADs<sup>14, 15</sup>,  
36 and disruption of individual CTCF binding sites deregulated the expression of surrounding  
37 genes<sup>16, 17, 18, 19</sup>. Mechanistically, many of the CTCF-mediated loops define insulated  
38 neighborhoods that constrain promoter-enhancer interactions<sup>13</sup>, and in some cases CTCF is  
39 directly involved in promoter-enhancer interactions<sup>9, 10, 20</sup>.

40 The CTCF-mediated interaction network has been considered to be largely invariant across cell-  
41 types. However, in studies of individual loci, cell-type-specific CTCF-mediated interactions were  
42 found to be important in gene regulation<sup>17, 21</sup>. Furthermore, CTCF binding sites vary extensively  
43 across cell-types<sup>22, 23</sup>. These findings suggest that the repertoire of CTCF-mediated interactions  
44 can be cell-type specific, and it is necessary to understand the extent and functional role of cell-  
45 type-specific CTCF-mediated loops. If cell-type-specific interactions are prevalent and contribute  
46 to cellular function, it would be inappropriate to use the CTCF-mediated interactome derived  
47 from a different cell-type.

48 CTCF-mediated loops can be mapped through Chromatin Conformation Capture (3C)-based  
49 technologies<sup>2</sup>. Among them, Hi-C<sup>9, 24</sup> provides the most comprehensive coverage for identifying  
50 looping events. However, it requires billions of reads to achieve kilo-base resolution<sup>9</sup>. On the  
51 other hand, Chromatin Interaction Analysis using Paired End Tags (ChIA-PET) increases  
52 resolution by only targeting chromatin interactions associated with a protein of interest<sup>10, 25, 26</sup>.  
53 Recently developed protocols, including Hi-ChIP<sup>27</sup> and PLAC-seq<sup>28</sup>, improved upon ChIA-PET  
54 in sensitivity and cost-effectiveness. Despite recent technical advances, experimental profiling  
55 of CTCF-mediated interactions remains difficult and costly, and few cell-types have been  
56 analyzed<sup>9, 10, 24, 29</sup>. Therefore, computational predictions that take advantage of the routinely  
57 available ChIP-seq and RNA-seq data is a desirable approach to guide the interrogation of the  
58 CTCF-mediated interactome for the cells of interest.

59 Here, we carried out comprehensive analysis of CTCF-mediated chromatin interactions using  
60 ChIA-PET data sets from multiple cell-types. We found that CTCF-mediated loops exhibit  
61 widespread plasticity and the cell-type-specific loops are biologically significant. Motivated by  
62 this observation, we developed Lollipop—a machine-learning framework based on random  
63 forests classifier—to predict the CTCF-mediated interactions using genomic and epigenomic  
64 features. Lollipop significantly outperforms methods based solely on convergent motif  
65 orientation when evaluated both within individual and across different cell-types. Our predictions  
66 were also experimentally confirmed by 3C. Moreover, our approach reveals novel determinants  
67 of CTCF-mediated chromatin wiring, such as gene expression within the loop.

68

## 69 Results

### 70 CTCF-mediated chromatin interactions exhibit cell-type specificity

71 We used the ChIA-PET2 pipeline<sup>30</sup> and analyzed published ChIA-PET data sets from three cell-  
72 lines (**Supplementary Table 1**): GM12878 (lympho-blastoid cells)<sup>10</sup>, HeLa-S3 (cervical  
73 adenocarcinoma cells)<sup>10</sup>, and K562 (chronic myelogenous leukemia cells)<sup>29</sup>. By using false  
74 discovery rate (FDR)  $\leq 0.05$  and paired-end tag (PET) number  $\geq 2$ , we identified 51966, 16783,  
75 13076 high-confidence chromatin loops for GM12878, HeLa, and K562, respectively  
76 (**Supplementary Table 2**). A significant fraction of loops was found to be cell-type-specific  
77 (67.9%, 26.2%, and 21.5% of loops in GM12878, HeLa, and K562, respectively (**Fig. 1a**). It is  
78 worth noting that the higher number of loops and cell-type-specific loops observed in GM12878  
79 may be attributed to the higher sequencing depth of GM12878 ChIA-PET library  
80 (**Supplementary Table 2**).

81 To elucidate what contributes to this plasticity, we compared the CTCF binding sites identified in  
82 ChIA-PET data sets across the three cell-lines. We found that only 36% of CTCF binding sites  
83 are constitutive (i.e., “+++”, **Fig. 1b**), consistent with previous reports<sup>22, 23</sup>. Besides cell-type-  
84 specific binding sites, rewiring of shared binding sites also contributes to the cell-type-specific  
85 loops (**Fig. 1c**).

### 86 Cell-type-specific CTCF-mediated loops contribute to gene regulation

87 Loops shared among different cell-types exhibit significantly higher interaction strength than the  
88 cell-type-specific loops (**Supplementary Fig. 1a**), questioning whether the latter are biologically  
89 relevant. To address this question, we asked whether these loops are involved in gene  
90 regulation. We found that cell-type-specific loops harbor a significantly higher ratio of tandem  
91 CTCF motif orientation compared to shared loops (**Supplementary Fig. 1b**), suggesting their  
92 involvement in gene regulation, considering that tandem loops exhibit more regulatory potential  
93 than convergent ones<sup>10</sup>.

94 Super-enhancers (SEs) are defined as stretches of chromatin that cluster multiple enhancers  
95 decorated with H3K27ac. A recent study revealed that CTCF plays a critical role in the  
96 hierarchical organization of SEs<sup>31</sup>. Considering that SEs play critical roles in cell identity,  
97 development, and cancer<sup>32, 33, 34</sup>, we examined whether they are enriched within cell-type-  
98 specific loops. Disease Ontology analysis using GREAT<sup>35</sup> confirmed that these SEs are linked  
99 with the corresponding disease origin of the three cell-types (**Supplementary Fig. 1c**).  
100 Comparison of SEs in HeLa and K562 identified three sets of SEs: HeLa-specific, common, and  
101 K562-specific. HeLa-specific SEs are significantly enriched within HeLa-specific loops,  
102 compared to common SEs (**Fig. 1d** left panel). Similarly, K562-specific SEs are preferentially  
103 enriched within K562-specific loops compared to common SEs (**Fig. 1d** left panel). The same  
104 conclusion was reached when we compared GM12878 vs HeLa as well as GM12878 vs K562  
105 (**Fig. 1d** central and right panels). Taken together, we found that cell-type-specific SEs are more  
106 likely to be associated with loops specific to that cell-type, suggesting the functional significance  
107 of cell-type-specific loops.

108 Consistently, differentially expressed genes (DEGs) between the three cell types are  
109 significantly associated with cell-type-specific loops (**Supplementary Fig. 1d**). Ingenuity  
110 Pathway Analysis (IPA)<sup>36</sup> revealed that DEGs between HeLa and K562 categorized based on  
111 loop association are enriched in distinct canonical pathways (**Fig. 1e**). Similar results were  
112 obtained in pair-wise comparisons between GM12878 and the other two cell lines



113 **(Supplementary Fig.1e-f)**. For instance, **Fig. 1f** illustrates the loop architecture and epigenomic  
114 features of ROR2, a receptor involved in non-canonical Wnt signaling with a significant role in  
115 human carcinogenesis<sup>37, 38</sup>. ROR2 is highly expressed in K562 compared to HeLa, and these  
116 CTCF-mediated loops are present only in K562. The up-regulation of ROR2 expression is  
117 associated with a concomitant decrease of H3K27me3 and increase in H3K36me3 in the region,  
118 as well as the appearance of a K562-specific SE in the gene body.

119 Altogether, cell-type-specific CTCF-mediated loops are prevalent and may play a significant role  
120 in the transcriptional programs of cell-type-specific genes. Therefore, we sought to develop a  
121 computational approach to infer the CTCF-mediated loops.

## 122 **An ensemble learning method to predict CTCF-mediated loops from genomic and** 123 **epigenomic features**

124 We employed a random forest classifier, a tree-based ensemble learning method, to predict  
125 CTCF-mediated loops. This classification method takes into consideration the complex  
126 interactions among features and is robust against overfitting<sup>39, 40, 41</sup>. The pipeline, named  
127 Lollipop, aims to find an optimized combination of genomic and epigenomic features to  
128 distinguish interacting from non-interacting pairs of CTCF sites. The schema of the pipeline is  
129 shown in **Fig. 2a**. The trained model can be used to predict CTCF-mediated loops in the same  
130 or a different cell-type.

131 For training purposes, the positive and negative loops were derived from ChIA-PET data sets<sup>10,</sup>  
132 <sup>29</sup>. To ensure confident labeling of positive loops, we used stringent criteria (FDR  $\leq$  0.05 and at  
133 least 2 PETs connecting the two anchors). Negative loops were constructed by random pairing  
134 of CTCF binding sites and were 5 times as abundant as the positive loops. Additional rules to  
135 select negative loops included: (a) lack of PET in the ChIA-PET dataset; and (b) absence in the  
136 list of identified interactions from the Hi-C experiments (see methods for details).

137 A total of 77 features were derived from genomic and epigenomic data sets (**Fig. 2a**). Genomic  
138 features include loop length and features defined at the CTCF binding sites, including CTCF  
139 motif orientation, strength, and sequence conservation. We included loop length because it is an  
140 inherent determinant of contact frequency between two genomic regions<sup>42</sup>, and motif orientation  
141 pattern because CTCF anchors preferentially adopt a convergent motif orientation<sup>9</sup>. Epigenomic  
142 features include chromatin accessibility, a variety of histone modifications, and architectural  
143 proteins CTCF and Cohesin (RAD21). For the use of DNase-seq and ChIP-seq data sets, three  
144 types of features were used: (a) local features defined at the anchors, (b) in-between features  
145 defined over the loop region, (c) and flanking features defined over the region from the loop  
146 anchor to the nearest CTCF binding event outside the loop (**Fig. 2b**). The use of the in-between  
147 features was motivated by a recent study<sup>43</sup> showing that signals over the loop regions were  
148 more important in predicting promoter-enhancer interactions than signals at anchors. In addition,  
149 given the insulator role of CTCF, we reasoned that the signals over the flanking regions might  
150 help to distinguish interacting from non-interacting CTCF binding sites. Finally, we also included  
151 gene expression within the looped region as a feature (see methods for details).

## 152 **Assessment of Lollipop's performance within individual cell-types**

153 We employed Receiver Operator Characteristic (ROC) and Precision-Recall (PR) curves with  
154 10-fold cross-validation to assess the performance of Lollipop. To account for possible bias  
155 introduced by random partitioning of training data, we performed 5 iterations for cross-validation  
156 and reported the mean performance. For evaluation of Lollipop's performance, two methods

157 were used for comparison. Both methods are inspired by the finding that the CTCF motifs in  
158 anchors preferentially adopt convergent orientation<sup>9, 10</sup>: (a) The naïve method, which pairs a  
159 CTCF-bound motif that resides on the forward strand to the nearest downstream CTCF-bound  
160 motif that resides on the reverse strand (**Supplementary Fig. 2a**); (b) The Oti method<sup>44</sup>, which  
161 iteratively applies the naïve method to CTCF binding sites selected by different signal intensity  
162 thresholds (see **Supplementary Fig. 2b** for illustration and methods for details). By doing so,  
163 the Oti method identifies more loops than the naïve method and partially recovers the nested  
164 structure of some CTCF-mediated loops.

165 **Fig. 3a-b** show that Lollipop achieved an area under ROC curve (AU-ROC) value of  $\geq 0.97$  and  
166 area under PR curve (AU-PR) value of  $\geq 0.86$  in all cell lines. Compared to other methods,  
167 Lollipop achieved similar or higher precision and superior recall. The latter can be partially  
168 attributed to the failure of naïve and Oti methods to capture tandem loops or loops without  
169 CTCF motif on anchors, which account for a significant fraction of CTCF-mediated loops (64%  
170 for GM12878, 61% for HeLa, 49% for K562). We then independently evaluated Lollipop's  
171 performance on convergent and non-convergent loops. Even on convergent loops, Lollipop  
172 achieved a superior recall score with a precision score comparable those of the naïve and Oti  
173 method (**Fig. 3c**). Furthermore, Lollipop also performed well in the prediction of non-convergent  
174 loops (**Fig. 3d**). In summary, Lollipop can account for the complexity of loop structures by  
175 integrating genomic and epigenomic features and outperforms methods that only consider the  
176 convergent CTCF motif orientation.

### 177 **Feature analysis identified novel determinants of CTCF-mediated chromatin loops**

178 Considering that convergent motif orientation does not suffice to identify CTCF-mediated loops,  
179 we ranked features that significantly improve the performance, by measuring the mean  
180 decrease impurity during training the random forests classifier<sup>45</sup>. We found that the average  
181 binding intensity of CTCF and Cohesin (RAD21) at the loop anchors are the most important  
182 features (**Fig. 4a** and **Supplementary Fig. 3a**), suggesting that sites with stronger CTCF and  
183 Cohesin binding are more likely to become anchors (**Supplementary Fig. 3b**), consistent with  
184 the observation that that these proteins are important for chromatin interactions<sup>14, 15</sup>. In addition,  
185 loop length and motif orientation pattern were amongst the top features, in agreement with  
186 previous results<sup>9, 42</sup>. The list also includes features defined within loop regions, among which  
187 gene expression was of particular interest. Regions inside positive loops exhibit significantly  
188 lower gene expression levels compared to negative loops (**Fig. 4b**). This finding is supported by  
189 similar trends exhibited by histone marks for active gene bodies H3K79me2 and H3K36me3  
190 (**Supplementary Fig. 3c**). Another interesting feature is the standard deviation of CTCF and  
191 Cohesin binding at the anchors (**Fig. 4a**). We therefore examined the relative fluctuation,  
192 defined as standard deviation divided by average intensity, of CTCF and Cohesin on anchor  
193 pairs of the positive and negative loops. As shown in **Fig. 4c** and **Supplementary Fig. 3d**,  
194 anchor-pair CTCF and RAD21 have significantly lower relative fluctuation in positive loops than  
195 in negative loops.

196 While CTCF binding at anchors is clearly critical for looping, formation of a loop requires wiring  
197 (i.e. physical interaction) between specific pair of anchors. We therefore asked what features  
198 contribute to the wiring. To this end, we changed negative loops to be random pairings of actual  
199 anchors, and then reanalyzed feature importance. As shown in **Supplementary Fig. 3e**, length,  
200 motif-orientation and expression are strongly contributing, whereas CTCF and Cohesin binding

201 at anchors become much less important. It is worth noting that more in-between features  
202 showed up in the list, compared to those in **Fig. 4a** and **Supplementary Fig. 3a**.

203 As the features employed are correlated (**Fig. 4d** and **Supplementary Fig. 3f**), the feature  
204 importance scores might be skewed. To validate the ranking of feature importance, we applied  
205 the Recursive Feature Elimination method to evaluate the performance of the recursively  
206 reduced feature set. The results are consistent with the feature ranking from the mean decrease  
207 impurity (**Supplementary Table 3**). Last, performance evaluation under different feature sets  
208 suggests that near-optimal performance can be achieved by using ~16 features (**Fig. 4e**).  
209 These features include those derived from CTCF and RAD21 binding, loop length, CTCF motif  
210 orientation, gene expression, as well as epigenetic features (**Supplementary Table 3**).

### 211 **Assessment of Lollipop's performance across cell-types**

212 Having demonstrated Lollipop's superior performance within individual cell-types, we next used  
213 the model trained in one cell-type to make predictions and assessment in another cell-type (see  
214 methods for details). This is more realistic and challenging, as a large number of CTCF-  
215 mediated loops are cell-type-specific. In all three cell-types Lollipop achieved AU-ROC  $\geq 0.93$   
216 and AU-PR  $\geq 0.79$  (**Fig. 5a-b**), only moderately lower than its performances within individual  
217 cell-types (**Fig. 3a-b**). It is worth noting that Lollipop outperforms motif-orientation based  
218 methods (**Fig. 5a-b**). Given that a loop consists of a pair of anchors and the wiring between  
219 them, we then dissected Lollipop's predictive power on anchors and wiring, respectively. For  
220 assessment of anchor prediction, we evaluated Lollipop by comparing the anchor usage of the  
221 predicted loops with that of loops identified from ChIA-PET in the target cell-type. For  
222 assessment of wiring prediction, we constructed negative loops by random pairing of actual  
223 anchors in the target cell-type (see methods for details). **Fig. 5c-d** show the PR curves  
224 demonstrating that Lollipop performed reasonably well in both, and better in predicting anchors  
225 than in predicting wiring. The results in terms of ROC (**Supplementary Fig. 4a-b**) are consistent  
226 with those in terms of PR.

### 227 **Evaluation of de novo predictions of CTCF-mediated loops**

228 After training Lollipop in individual cell-types, we then applied it to scan the genome of the same  
229 cell-type to make *de novo* genome-wide predictions. Lollipop predicted 67855, 38274, 32237  
230 loops in GM12878, HeLa and K562, respectively. Notably, the number of predicted loops in  
231 GM12878 is much larger than those of the other two cell-types, due to the much larger number  
232 of loops identified by ChIA-PET in GM12878 (see last column of **Supplementary Table 2**).  
233 These loops were used in training the model and thus affect the number of predicted loops.  
234 Indeed, if we down-sample the GM12878 ChIA-PET library to 15% so that the number of called  
235 loops is on par with those in K562 and HeLa (see last column of **Supplementary Table 2**), the  
236 number of predicted loops is comparable to the number of predictions in K562 and HeLa.

237 As shown in **Supplementary Fig. 5a**, a large fraction of the predicted loops (48%, 73% and 77%  
238 for GM12878, HeLa and K562, respectively) was not supported by ChIA-PET under the  
239 stringent criterion of  $FDR \leq 0.05$  and  $PET \geq 2$  used for defining positive loops. However, if we  
240 relaxed the stringency to  $PET \geq 1$  in ChIA-PET, the fraction of predicted loops not supported by  
241 ChIA-PET was significantly reduced, to 24%, 42% and 50% in GM12878, HeLa and K562,  
242 respectively. Similar result can be obtained with the down-sampled GM12878 library  
243 (**Supplementary Fig. 5b**). This observation raises the question of whether the predicted loops  
244 with less or no ChIA-PET support are indeed false positives. To address this question, we

245 carried out the following computational as well as experimental evaluations on those predicted  
246 loops without any ChIA-PET support.

247 First, we used the published Hi-C contact matrices for GM12878 and K562<sup>9</sup> (see methods for  
248 details) to evaluate these loops, and found that they have significantly higher contact  
249 frequencies than pairs of randomly-chosen genomic loci (**Fig. 6a**). For fair comparison, the  
250 control regions were sampled to have a length distribution matching those of the target loops.  
251 Second, we randomly selected two such cases and performed 3C experiments. **Fig. 6b** shows  
252 the sequence of the ligation junctions from the long-range interactions (PRKAG2-KMT2C and  
253 PDE6A-PDGFRB) in HeLa. 3C-qPCR further confirmed the contact frequency of the PRKAG2-  
254 KMT2C loop in respect to neighboring HindIII fragments (**Supplementary Fig. 5d**).

255 Having shown that the predicted loops lacking ChIA-PET support could be real, we sought to  
256 understand why they were not observed in ChIA-PET. To this end, we performed scaling  
257 analysis in the ChIA-PET data of GM12878 cells, which received significantly higher sequencing  
258 coverage than those of K562 and HeLa (**Supplementary Table 2**). Specifically, we used the 15%  
259 down-sampled GM12878 ChIA-PET library to identify loops with the same approach employed  
260 for the full data set, and trained a classifier. We then applied this classifier to make genome-  
261 wide predictions. Of the 33463 predicted loops, 12047 are without any support from the down-  
262 sampled ChIA-PET data set. However, 46% of these loops find support in the full ChIA-PET  
263 library, and 20% of these loops even find significant support (**Fig. 6c**). This down-sampling  
264 process was repeated for 10 times and similar results were obtained (data not shown). Taken  
265 together, the scaling analysis suggests that insufficient sequencing depth contributes to the  
266 presence of predicted loops lacking support in ChIA-PET.

## 267 **Topological properties of CTCF-mediated interaction network and associated** 268 **biological functions**

269 To gain a better understanding of these interactions, we took a systems approach to visualize  
270 and analyze the CTCF-mediated interactions. We constructed the CTCF-mediated interaction  
271 network by denoting the anchors as nodes and the long-range interactions as edges. As  
272 exemplified in **Fig. 7a**, where the interaction network on chromosome 1 (visualized using graph-  
273 tool V2.22, <https://graph-tool.skewed.de>) is shown, the CTCF-mediated interactions form a  
274 disconnected network encompassing many linear-polymer-like components. This is dramatically  
275 different from the RNA-PolIII-mediated interaction network<sup>46</sup>, which is dominated by one scale-  
276 free connected graph<sup>46</sup>. This dramatic difference in topological structure is also manifested in  
277 the degree distributions (**Supplementary Fig. 6**), where the distribution for RNA PolIII exhibits a  
278 fatter tail.

279 It is worth noting that degrees of connections among the anchors vary. We therefore examined  
280 CTCF hubs, anchors involved in multiple interactions. Ranking anchors according to the  
281 degrees of connections, we defined hubs as those among the top 10% anchors and non-hubs  
282 as the bottom 10% (see methods for details), and identified 2914, 2111 and 1843 nodes for  
283 GM12878, HeLa and K562, respectively. Subsequent comparison between hubs and non-hub  
284 nodes revealed that hubs are (a) more conserved across cell-types than non-hubs, likely  
285 because they serve as the structural foci of genome organization in the nucleus, (b)  
286 characterized by significantly higher binding affinity for CTCF and Cohesin (**Fig. 7c**), and (c)  
287 associated with distinct biological functions. Gene ontology analysis<sup>35</sup> showed that the hubs are  
288 preferentially associated with immunology-related functions in GM12878 and K562 cells, but not  
289 in HeLa cells (**Fig. 7d**), consistent with the cellular origin of these cell-lines. For example, the



290 hubs in GM12878 and K562 cells were found to be significantly associated with antigen binding,  
291 and the GM12878 hubs were significantly associated with the MHC (major histocompatibility  
292 complex) protein complex. MHC is a set of cell surface proteins that are essential for immune  
293 system, while MHC class II (MHC-II) genes encode cell-surface glycoproteins that present  
294 antigens to CD4 T cells to initiate and control adaptive immune responses<sup>47</sup>. Our results were  
295 consistent with previous studies<sup>47, 48</sup> which found that CTCF plays an important role in  
296 controlling MHC-II gene expression.

## 297 Discussion

298 Here we showed that CTCF-mediated chromatin interactions exhibit extensive variations across  
299 cell-types. These cell-type-specific interactions are functionally important, as they are linked to  
300 differentially expressed genes and cell-type-specific SEs contributing to cell identity. However,  
301 genome-wide profiling of CTCF-mediated interactions is available in a very limited number of  
302 cell-types and conditions, as experimental approaches remain challenging and costly. Therefore,  
303 we developed Lollipop, a machine-learning framework, to make genome-wide predictions of  
304 CTCF-mediated loops using widely accessible genomic and epigenomic features. Using  
305 computational as well as experimental validations, we demonstrated that Lollipop performed  
306 well within and across cell-types. Analysis of the machine learning model revealed novel  
307 features associated with CTCF-mediated loops, and shed light on the rules underlying CTCF-  
308 mediated chromatin organization.

309 While previous studies focused on the significance of conserved CTCF binding at TAD  
310 boundaries or loop anchors, our study showed a significant proportion of CTCF-mediated  
311 interactions are cell-type-specific. Based on our analysis, both lineage-specific recruitment of  
312 architectural proteins and alternative wiring among available anchor sites contribute to the  
313 establishment of cell-type specificity. Although the process of establishing cell-type-specific is  
314 not well understood, it is conceivable that multiple factors combine to orchestrate a cell-type-  
315 specific chromatin context to promote the formation of a loop.

316 The convergent orientation of CTCF motifs at loop anchors is a prominent feature of CTCF-  
317 mediated interactions<sup>9, 10</sup>, as it is also manifested by our model. However, model comparison  
318 demonstrated that motif orientation alone is limited in its predictive power, and inclusion of other  
319 features significantly improved the performance. Interestingly, we found that features for the  
320 loop regions, which are away from the anchors, contribute significantly to the predictive power,  
321 consistent with findings in enhancer-promoter interaction prediction<sup>43</sup>. Specifically, gene  
322 expression exhibits distinct distributions over positive loop regions compared to negative loops  
323 (**Fig. 4b**, and **Supplementary Fig. 4c**), which may be attributed to the enhancer-blocking role of  
324 CTCF loop anchors.

325 In evaluating our predictions, we showed that false positives could be due to mislabeling in the  
326 testing data. As advances in experimental protocols and continuous decreases in sequencing  
327 cost would result in better training data in reference cell-types, it is likely that the performance of  
328 Lollipop would further improve. Since CTCF plays a major role in defining regulatory domains,  
329 results obtained from our approach can potentially be used as constraints in predicting  
330 enhancer-promoter interactions, which remains a major challenge. Overall, CTCF-mediated  
331 chromatin interactions are critical for genome organization and function, and our study provides  
332 a computational tool for the exploration of the 3D organization of the genome.

## 333 **Materials and Methods**

### 334 **Data availability**

335 GM12878 and HeLa ChIA-PET data were downloaded from Gene Expression Omnibus (GEO)  
336 with accession number GSE72816<sup>10</sup>. K562 ChIA-PET data was downloaded from ENCODE<sup>29</sup>  
337 with accession number ENCLB559JAA. High-resolution genome-wide Hi-C contact matrices  
338 were obtained from GEO with accession number GEO63525<sup>9</sup>. DNase-Seq, ChIP-Seq and RNA-  
339 Seq data were downloaded from ENCODE and were aligned to hg19. The accession numbers  
340 for the data used in this study were summarized in **Supplementary Table 1**.

341 Lollipop is publically available in <https://github.com/ykai16/Lollipop>.

### 342 **Identification of CTCF-mediated loops from CTCF ChIA-PET data**

343 We employed ChIA-PET2 (v0.9.2)<sup>30</sup> to identify CTCF-mediated loops. Briefly, ChIA-PET2  
344 involves linker filtering, PET mapping, PET classification, binding-site identification, and  
345 identification of long-range interactions. In the step of linker filtering, one mismatch was allowed  
346 in identifying reads with linkers. After linker removal, only reads with at least 15 bp in length  
347 were retained for further analysis for GM12878 and HeLa (read length = 150 bp). For K562, the  
348 read length was shorter (36 bp), therefore reads with at least 10 bp in length were retained for  
349 further analysis. In other steps, default values for parameters were used. Only uniquely  
350 mapped reads were kept, and PETs were de-duplicated. Significant loops were identified with a  
351 value of false discovery rate (FDR)  $\leq 0.05$ . We further required that they are supported by at  
352 least two PETs (i.e., IAB  $\geq 2$ ).

353 We only considered long-range interactions whose length are less than 1 million bps (mb), for  
354 two reasons. First, vast majority of loops (93.2% for GM12878, 97.3% for HeLa, 98.1% for K562)  
355 are less than 1mb long. Similar observations were made in <sup>10</sup>. Second, insulated neighborhoods,  
356 the CTCF loops having higher potential in regulation of gene expression, were found to range  
357 from 25 kb to 940 kb<sup>6, 16</sup> (reviewed in <sup>13</sup>).

### 358 **Comparison of CTCF-mediated loops among cell-types (Fig. 1a, Supplementary Fig. 1a-b)**

360 An anchor is considered as shared by two cell-types if the respective genomic regions  
361 delineating this anchor overlap in the two cell-type. A loop is considered as shared by two cell-  
362 types if both anchors are shared by the two cell-types. A loop is considered cell-type specific if  
363 either of the two anchors are cell-type specific. The loops shared by all three cell-types were  
364 defined as GM12878 loops shared by both K562 and HeLa.

### 365 **Analysis of CTCF binding sites in three cell-types (Fig. 1b)**

366 CTCF peaks were determined by MACS2<sup>49</sup> in the ChIA-PET2 pipeline. A binding site was  
367 defined as peak summit +/- 500 bp. The binding sites in the three cell-types were classified into  
368 seven groups according to the overlapping pattern. Binding intensity for each site was  
369 represented by the log<sub>2</sub> (RPKM) value over the summit +/- 2kb region. For each group, the  
370 binding sites were ordered in descending order according to binding intensity in a prioritized  
371 manner. Namely, CTCF binding sites present in GM12878 were ordered by their binding  
372 strengths in GM12878; CTCF binding sites not present in GM12878 were ordered by binding  
373 strengths in HeLa and then in K562 accordingly. Seaborn (V 0.7.1, <http://seaborn.pydata.org>)  
374 was used to generate the heat map.

### 375 **Super-enhancer analysis (Fig. 1d, Supplementary Fig. 1c)**

376 Super-enhancers (SEs) were identified by the Ranking Ordering of Super-Enhancers algorithm  
377 (ROSE<sup>33, 34</sup>), using H3K27ac ChIP-Seq data as input and default parameters. Identified super-  
378 enhancers were then uploaded to Genomic Regions Enrichment of Annotations Tool (GREAT)  
379 V3.0.0<sup>35</sup> for GO analysis (**Supplementary Fig. 1c**). If a SE in one cell-type does not overlap  
380 with any SEs in a different cell-type, it is deemed as a SE specific to that cell-type. Otherwise, it  
381 is called a shared SE. We then counted the number of cell-type specific loops covering each  
382 type of SEs. The comparison between HeLa and K562 is shown in **Fig. 1d**. For comparison  
383 between GM12878 and another cell-type, the GM12878 ChIA-PET data set is first randomly  
384 down-sampled to 15% of the original size so that the number of loops identified matched those  
385 from the ChIA-PET datasets of the other two cell-types (see **Supplementary Table 2**). Then  
386 analysis identical to that in **Fig. 1d** was carried out. The down sampling and follow-up analysis  
387 was repeated 10 times to ensure reproducibility, and standard-deviations were shown in the **Fig.**  
388 **1d**.

### 389 **Analysis of differentially expressed genes and their association with CTCF-** 390 **mediated loops (Fig.1e, Supplementary Fig. 1d, e, f)**

391 Each cell-line has two RNA-Seq replicates. Cufflinks V2.2.1<sup>50</sup> with default parameters (*q-*  
392 *value*=0.05) was used to identify the differentially expressed genes (DEG).

393 For comparison between HeLa and K562, a DEG was deemed to be associated with HeLa-  
394 specific loops if it is within one or more HeLa-specific loops but not within any K562-specific  
395 loops. If a DEG is covered only by one or more shared loops, this DEG is deemed to be  
396 associated with shared loops. Following the criteria described above, we obtained three sets of  
397 DEGs respectively associating with HeLa-specific loops, shared loops, K562-specific loops.  
398 These three sets of DEGs were then subject to GO analysis using 'Ingenuity Pathway Analysis'  
399 <sup>36</sup>. The GO terms whose P-value are no less than 1e-3 in all three gene sets were then removed.  
400 The result is shown in **Fig. 1e**. Color key represents the  $-\log_{10}$  (P-value). For comparison  
401 between GM12878 and another cell-type (**Supplementary Fig. 1 e, f**), the GM12878 ChIA-PET  
402 library is first randomly down-sampled to 15% of the original size so that the number of loops  
403 identified matched those of the ChIA-PET libraries from the other two cell-types.

404 For **Supplementary Fig. 1d**, non-DEG genes were those with the least significant expression  
405 changes as ranked by P-value, with group size matching to that of the corresponding DEG  
406 group.

### 407 **Identification of CTCF motif occurrences**

408 The position frequency matrix of CTCF for human was downloaded from Jaspar 2016  
409 (<http://jaspar.genereg.net>)<sup>51</sup>. CTCF motif occurrences were identified by the FIMO package  
410 (V4.11.1<sup>52</sup>) with the P-value < 1e-5. In total, 110879 motif occurrences were identified.

### 411 **Preparation of training data**

412 Positive loops were identified using ChIA-PET2 pipeline with FDR<=0.05 and IAB >=2, with loop  
413 length restricted to be in the range of 10 kb to 1mb. The choice of the lower limit of 10 kb is  
414 because the ChIA-PET-identified loops with length below 10 kb are likely caused by self-ligation  
415 in library preparation<sup>25</sup>. The reason for the upper limit of 1mb is given above. Negative loops  
416 were constructed by random pairing of CTCF binding sites, with loop length ranging from 10 kb  
417 to 1mb. The number of negative interactions was chosen to be 5 times that of the positive  
418 interactions. To ensure accurate labeling, we further required that the negative loops (1) do not

419 receive any ChIA-PET support; and (2) are not present in the CTCF-mediated interactions  
420 identified from the Hi-C experiments<sup>9</sup>.

### 421 **Feature calculation (Fig. 2a, b)**

422 Genomic features include motif strength, motif orientation, conservation score and loop length.  
423 Motif strength represents how similar the underlying sequence is to the CTCF consensus motif.  
424 The motif strength score was provided by FIMO<sup>52</sup>. The motif strength score of a CTCF binding  
425 site (summit +/- 1000bp) was represented by the strength of the motif occurrence within the site.  
426 If a CTCF binding site have more than one motif occurrences, the highest score was used. If  
427 there is no motif occurrence, 0 would be assigned. The feature of motif orientation was  
428 represented by the following rule: If neither anchor has CTCF motif, we assign a value of 0; If  
429 one anchor has no motif and the other has one or more than one motifs, we assign a value of 1;  
430 If both anchors have one or more motif occurrences, the orientation of each anchor is  
431 determined by the orientation of its strongest motif occurrence. Divergent orientation would be  
432 assigned a value of 2, tandem orientation would be assigned a value of 3, and convergent  
433 orientation would be assigned a value of 4. For conservation, we used the 100 way phastCons  
434 score downloaded from UCSC  
435 (<http://hgdownload.cse.ucsc.edu/goldenpath/hg19/phastCons100way>)<sup>53</sup>. The conservation  
436 score of a CTCF binding site was defined as the mean value of the conservation score of each  
437 nucleotide in the summit +/- 20 bp region.

438 Functional genomic features include chromosome accessibility profiled by DNase-Seq, histone  
439 modifications, CTCF and Cohesin binding profiles profiled by ChIP-Seq, and gene expression  
440 profiled by RNA-Seq. DNase-Seq and ChIP-Seq data were de-duplicated and then subject to  
441 pre-processing to remove noise as follows. For DNase-Seq data, peaks were downloaded from  
442 ENCODE<sup>29</sup>. For ChIP-Seq data, SICER (V1.1)<sup>54</sup> were used to identify enriched regions with  
443 FDR 1e-5. For histone modifications with diffused signal (H3K27me3, H3K36me3, H3K9me3,  
444 H3K79me2), window size = 200 bp, gap size = 600 bp were used. For other ChIP-Seq libraries,  
445 window size = gap size = 200 bp were used. For both DNase-Seq and ChIP-Seq, only reads  
446 located on signal-enriched regions were used for feature calculation. For RNA-Seq data, gene  
447 expressions were calculated using Cufflinks<sup>50</sup> with default parameters. Each dataset was  
448 characterized by three types of features: local features, in-between features and flanking  
449 features, as illustrated in **Fig. 2b**. Local features were defined around anchors, represented by  
450 the signal intensity (RPKM value) over the CTCF summit position +/- 2kb region. In-between  
451 feature is represented by the average signal intensity (RPKM value) over a presumed loop  
452 region. The value of the expression feature is defined as the average FPKM value of the genes  
453 whose promoters are located inside the presumed loop. The flanking features are represented  
454 by the RPKM value over the region from the loop anchor to the nearest CTCF binding event  
455 identified in the CTCF ChIP-Seq.

### 456 **Implementation of the naïve method and the Oti method (Supplementary Fig. 2)**

457 The naïve method is implemented by pairing a CTCF-bound motif that resides on the forward  
458 strand to the nearest downstream CTCF-bound motif that resides on the reverse strand  
459 (**Supplementary Fig. 2a**). The Oti method was introduced in<sup>44</sup>. It ranked all the active motif  
460 sites in terms of CTCF peak strength in descending order. First all active motif sites were used  
461 to construct loops by the naïve method. Then, the same procedure was repeated for the top  
462 80%, top 60%, top 40% and top 20% active motif sites. The loops constructed in different  
463 rounds were then pooled together. The Oti method is illustrated in **Supplementary Fig. 2b**.



### 464 **Performance evaluation within individual cell-types (Fig. 3)**

465 In **Fig. 3c, d**, the performance was evaluated at the looping probability cut-off of 0.5.

### 466 **Evaluation of feature importance (Fig. 4a, d, e and Supplementary Fig. 3a, e, f)**

467 Predictive importance scores of features were obtained from the “feature\_importances  
468 “ attribute of the trained random forest classifier<sup>55</sup>. The ranking of the top 20 features was  
469 visualized in **Fig. 4a** and **Supplementary Fig. 3a**. Pearson correlations of the in-between  
470 features calculated in positive interactions were used to generate the correlation matrix. The  
471 correlation matrix was subject to hierarchical clustering, as shown in **Fig. 4d and**  
472 **Supplementary Fig. 3f**. Recursive Feature Elimination (RFE) method was used to validate the  
473 analysis of the feature importance. After each iteration, model performance was evaluated in  
474 terms of Area Under Receiver Operating Characteristics (AU-ROC) curve and Area Under  
475 Precision Recall (AU-PR) curve. The performance vs. feature number was plotted in **Fig. 4e**.

476 For feature importance analysis of wiring prediction (**Supplementary Fig. 3e**). Negative data  
477 was prepared as follows: the anchors of positive loops were used to construct negative loops by  
478 random pairing. The number of negative loops were set to be 3 times that of positive loops.  
479 Other procedures on construction of negative loops were the same as described in the section  
480 of ‘Preparation of training data’. Positive data remained unchanged.

### 481 **Performance evaluation across cell-types (Fig. 5 and Supplementary Fig. 4)**

482 In the across-cell-type performance evaluation, the model trained in cell-type A was applied to  
483 the cell-type B, using training data prepared in B for evaluation of performance.

484 For evaluation of anchor prediction, the anchors of positive loops in cell-type B were labeled  
485 positive, while the anchors belonging only to negative loops in cell-type B were labeled negative.  
486 The anchors of predicted loop were compared with positive and negative labels for evaluation of  
487 anchor prediction. This evaluation was repeated under different thresholds of looping probability  
488 to generate the PR and ROC curves (**Fig. 5c and Supplementary Fig. 4a**).

489 For evaluation of wiring prediction, the anchors of positive loops in cell-type B were used to  
490 construct negative loops by random pairing. The model trained in cell-type A was then applied  
491 to the training data of cell-type B for evaluation.

### 492 **Computational evaluation of predicted CTCF-mediated loops (Fig. 6a, c and** 493 **Supplementary Fig. 5a, b)**

494 Models trained in a cell-type was used to predict loops genome-widely in the same cell-type.  
495 Predicted loops were then compared with loops identified from ChIA-PET datasets and  
496 categorized into three groups. ‘Significant’ loops denote those supported by ChIA-PET under  
497 the stringent criterion of  $FDR \leq 0.05$  and  $PET > 2$ . ‘With evidence’ loops denote those supported  
498 by ChIA-PET reads but do not meet the stringent criterion mentioned above. ‘No support’ loops  
499 denote those without any support from ChIA-PET. The numbers of loops in each group were  
500 shown in **Supplementary Fig. 5a**.

501 Down sampling of ChIA-PET library in GM12878 cells: The ChIA-PET library was first randomly  
502 down-sampled to 15% of the original size, followed by loop identification using ChIA-PET2 and  
503 preparation of training data. Trained model was used to make genome-wide predictions. The  
504 predicted loops were categorized into three groups by comparing with loop calls using the  
505 down-sampled library, as described above. The result was shown in **Supplementary Fig. 5b**.

506 Evaluation of predicted loops without any ChIA-PET support using Hi-C data (Fig. 6a). 10 kb  
507 resolution Hi-C contact matrices for GM12878 and K562<sup>9</sup> were used for validation. The contact  
508 matrices were normalized by Knight and Ruiz (KR) normalization vector<sup>9</sup>. For each cell-type, we  
509 collected contact frequencies from the contact matrix for those predicted loops without any  
510 ChIA-PET support. As a control, we chose a matching set of random pairs of genomic locations  
511 as anchors with matching length-distribution. We then collected the contact frequencies of this  
512 control set. The two contrasting distributions of contact frequencies are shown. HeLa cell was  
513 not included in this analysis because the Hi-C library and Hi-C derived contact matrix are not  
514 available.

515 Scaling analysis in GM12878 cells. Predicted loops belonging to the 'No support' group in the  
516 down-sampled ChIA-PET library (yellow slice in **Supplementary Fig. 5b**) were compared with  
517 the loops identified using the full GM12878 ChIA-PET library and categorized into three groups,  
518 as shown in **Fig. 6b**.

### 519 **Experimental validation using Chromosome Conformation Capture (3C) (Fig. 6b,** 520 **Supplementary Fig. 5c-d)**

521 The loops used for experimental validation were randomly selected from the loops predicted by  
522 Lollipop but not observed in ChIA-PET, as described above. For the 3C assay, cells were fixed  
523 and nuclei were prepared as in ChIP experiments. Nuclei were resuspended in 500  $\mu$ l 1.2X  
524 CutSmart buffer (NEB) with 14  $\mu$ l 10% SDS, and incubated at 37°C for 1 hour. SDS was  
525 sequestered by the addition of 50  $\mu$ l 20% Triton X-100, and incubated at 37°C for 1 hour. Next,  
526 5-20  $\mu$ l "undigested" was reserved, and 400 U of HindIII was added to the remaining sample  
527 and digested overnight at 37°C with end-over-end rotation. The second day, 5-20  $\mu$ l of  
528 "digested" material was reserved, and 40  $\mu$ l of 20% SDS was added to remaining sample to  
529 inactivate HindIII by incubating at 65°C for 25 minutes. The samples were transferred to 15 mL  
530 conical tubes and diluted with the following 1.15X ligation buffer recipe: 352  $\mu$ l 10X T4 ligase  
531 buffer (NEB), 2.71 ml water, and 187.5  $\mu$ l 20% Triton X-100. Samples were incubated at 37°C  
532 for 1 hour. Next, 5000 U T4 ligase was added, and ligation took place with gentle end-over-end  
533 rotation at 16°C for 4 hours, and then 45 minutes at room temperature. Reverse crosslinking  
534 took place by the addition of 300  $\mu$ g (30  $\mu$ l) Proteinase K at 65°C, overnight. On day three, 300  
535  $\mu$ g RNase-A was added, and samples were placed at 37°C for one hour. To begin DNA  
536 extraction, 4 ml of phenol-chloroform was added, samples were vortexed for a full minute, and  
537 centrifuged at 2,200 x g for 15 minutes. The aqueous phase was collected in a new 50 ml tube  
538 and diluted with an equal volume of water (4 ml) and with 800  $\mu$ l of 2 M sodium acetate pH 5.6;  
539 next 20 ml of ethanol was added, samples were inverted 10 times, and placed at -80°C for 1-4  
540 hours to precipitate the DNA. The samples were centrifuged at 2,200 x g for 45 minutes at 4°C  
541 and washed with 70% ethanol. The 3C libraries were then allowed to dry briefly, without letting  
542 the pellet become dull. The libraries were re-suspended in 100-600  $\mu$ l of 10 mM Tris. The  
543 digestion efficiency, as well as the quality and quantity of 3C libraries, were assessed before  
544 downstream analyses. The Q5 Taq polymerase (NEB) was used for PCR reactions using the  
545 following protocol: 98°C 30 sec, 35 cycles [98°C 10 sec, 70°C 15 sec, 72°C 10 sec], 72°C 2  
546 min. Reactions were run on 2% agarose gels and analyzed using the ImageLab software  
547 (BioRad). Bands were extracted and sequenced (Eurofins) to confirm specificity of primers and  
548 loop identity. Data points plotted in the contact matrix are the averages of duplicates  $\pm$  StDev  
549 from two independent library preparations. Primers were designed using a uni-directional  
550 strategy<sup>56</sup> and used are provided in **Supplementary Table 4**.

## 551 **Analysis of CTCF-mediated interaction network (Fig. 7)**

552 Construction of CTCF-mediated interaction network. We used nodes to represent anchors and  
553 edges to represent loops. Graph-tool (V2.22, <https://graph-tool.skewed.de>) was used for  
554 visualization of networks (**Fig. 7a**). In identification of hubs, anchors were ranked according to  
555 the degree of connection in descending order. Anchors with the same degree of connection  
556 were further ranked according to CTCF binding intensity in descending order. The top 10%  
557 anchors were defined as hubs, while the bottom 10% as non-hubs.

558 Functional enrichment analysis of hubs (Fig. 7d). Hubs were uploaded to GREAT (V3.0.0)<sup>35</sup> for  
559 functional enrichment analysis. The whole set of CTCF anchors was used as background. The  
560 GO terms in 'Molecular Functions' with P-value<1e-4 in each cell-type were shown.

561

## 562 **References**

- 563 1. Bickmore WA. The spatial organization of the human genome. *Annu Rev Genomics Hum Genet*  
564 **14**, 67-84 (2013).
- 565 2. Dekker J, Marti-Renom MA, Mirny LA. Exploring the three-dimensional organization of genomes:  
566 interpreting chromatin interaction data. *Nat Rev Genet* **14**, 390-403 (2013).
- 567 3. Dixon JR, Gorkin DU, Ren B. Chromatin Domains: The Unit of Chromosome Organization. *Mol*  
568 *Cell* **62**, 668-680 (2016).
- 569 4. Dixon JR, *et al.* Chromatin architecture reorganization during stem cell differentiation. *Nature*  
570 **518**, 331-336 (2015).
- 571 5. Gorkin DU, Leung D, Ren B. The 3D genome in transcriptional regulation and pluripotency. *Cell*  
572 *Stem Cell* **14**, 762-775 (2014).
- 573 6. Ji X, *et al.* 3D Chromosome Regulatory Landscape of Human Pluripotent Cells. *Cell Stem Cell* **18**,  
574 262-275 (2016).
- 575 7. Bonev B, Cavalli G. Organization and function of the 3D genome. *Nat Rev Genet* **17**, 772 (2016).
- 576 8. Dixon JR, *et al.* Topological domains in mammalian genomes identified by analysis of chromatin  
577 interactions. *Nature* **485**, 376-380 (2012).
- 578 9. Rao SS, *et al.* A 3D map of the human genome at kilobase resolution reveals principles of  
579 chromatin looping. *Cell* **159**, 1665-1680 (2014).
- 580 10. Tang Z, *et al.* CTCF-Mediated Human 3D Genome Architecture Reveals Chromatin Topology for  
581 Transcription. *Cell* **163**, 1611-1627 (2015).

- 591  
592 11. Huang J, Marco E, Pinello L, Yuan GC. Predicting chromatin organization using histone marks.  
593 *Genome Biol* **16**, 162 (2015).
- 594  
595 12. Ghirlando R, Felsenfeld G. CTCF: making the right connections. *Genes Dev* **30**, 881-891 (2016).
- 596  
597 13. Hnisz D, Day DS, Young RA. Insulated Neighborhoods: Structural and Functional Units of  
598 Mammalian Gene Control. *Cell* **167**, 1188-1200 (2016).
- 599  
600 14. Nora EP, *et al.* Targeted Degradation of CTCF Decouples Local Insulation of Chromosome  
601 Domains from Genomic Compartmentalization. *Cell* **169**, 930-944 e922 (2017).
- 602  
603 15. Zuin J, *et al.* Cohesin and CTCF differentially affect chromatin architecture and gene expression  
604 in human cells. *Proc Natl Acad Sci U S A* **111**, 996-1001 (2014).
- 605  
606 16. Downen JM, *et al.* Control of cell identity genes occurs in insulated neighborhoods in mammalian  
607 chromosomes. *Cell* **159**, 374-387 (2014).
- 608  
609 17. Hanssen LLP, *et al.* Tissue-specific CTCF-cohesin-mediated chromatin architecture delimits  
610 enhancer interactions and function in vivo. *Nat Cell Biol* **19**, 952-961 (2017).
- 611  
612 18. Hnisz D, *et al.* Activation of proto-oncogenes by disruption of chromosome neighborhoods.  
613 *Science* **351**, 1454-1458 (2016).
- 614  
615 19. Narendra V, *et al.* CTCF establishes discrete functional chromatin domains at the Hox clusters  
616 during differentiation. *Science* **347**, 1017-1021 (2015).
- 617  
618 20. Ren G, *et al.* CTCF-Mediated Enhancer-Promoter Interaction Is a Critical Regulator of Cell-to-Cell  
619 Variation of Gene Expression. *Mol Cell* **67**, 1049-1058 e1046 (2017).
- 620  
621 21. Hou C, Dale R, Dean A. Cell type specificity of chromatin organization mediated by CTCF and  
622 cohesin. *Proc Natl Acad Sci U S A* **107**, 3651-3656 (2010).
- 623  
624 22. Cuddapah S, Jothi R, Schones DE, Roh TY, Cui K, Zhao K. Global analysis of the insulator binding  
625 protein CTCF in chromatin barrier regions reveals demarcation of active and repressive domains.  
626 *Genome Res* **19**, 24-32 (2009).
- 627  
628 23. Wang H, *et al.* Widespread plasticity in CTCF occupancy linked to DNA methylation. *Genome Res*  
629 **22**, 1680-1688 (2012).
- 630

- 631 24. Lieberman-Aiden E, *et al.* Comprehensive mapping of long-range interactions reveals folding  
632 principles of the human genome. *Science* **326**, 289-293 (2009).
- 633
- 634 25. Li G, *et al.* Chromatin Interaction Analysis with Paired-End Tag (ChIA-PET) sequencing technology  
635 and application. *BMC Genomics* **15 Suppl 12**, S11 (2014).
- 636
- 637 26. Fullwood MJ, *et al.* An oestrogen-receptor-alpha-bound human chromatin interactome. *Nature*  
638 **462**, 58-64 (2009).
- 639
- 640 27. Mumbach MR, *et al.* HiChIP: efficient and sensitive analysis of protein-directed genome  
641 architecture. *Nat Methods* **13**, 919-922 (2016).
- 642
- 643 28. Fang R, *et al.* Mapping of long-range chromatin interactions by proximity ligation-assisted ChIP-  
644 seq. *Cell Res* **26**, 1345-1348 (2016).
- 645
- 646 29. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**,  
647 57-74 (2012).
- 648
- 649 30. Li G, Chen Y, Snyder MP, Zhang MQ. ChIA-PET2: a versatile and flexible pipeline for ChIA-PET  
650 data analysis. *Nucleic Acids Res* **45**, e4 (2017).
- 651
- 652 31. Jialiang Huang KL, Wenqing Cai, Xin Liu, Yuannyu Zhang, Stuart H.Orkin, Jian Xu, Guo-Cheng  
653 Yuan. Dissecting super-enhancer hierarchy based on chromatin interactions. In:  
654 *bioRxiv*(doi:10.1101/149583) (ed^(eds) (2017).
- 655
- 656 32. Hnisz D, *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934-947  
657 (2013).
- 658
- 659 33. Loven J, *et al.* Selective inhibition of tumor oncogenes by disruption of super-enhancers. *Cell*  
660 **153**, 320-334 (2013).
- 661
- 662 34. Whyte WA, *et al.* Master transcription factors and mediator establish super-enhancers at key  
663 cell identity genes. *Cell* **153**, 307-319 (2013).
- 664
- 665 35. McLean CY, *et al.* GREAT improves functional interpretation of cis-regulatory regions. *Nat*  
666 *Biotechnol* **28**, 495-501 (2010).
- 667
- 668 36. Kramer A, Green J, Pollard J, Jr., Tugendreich S. Causal analysis approaches in Ingenuity Pathway  
669 Analysis. *Bioinformatics* **30**, 523-530 (2014).
- 670

- 671 37. Bolzoni M, *et al.* Myeloma cells inhibit non-canonical wnt co-receptor ror2 expression in human  
672 bone marrow osteoprogenitor cells: effect of wnt5a/ror2 pathway activation on the osteogenic  
673 differentiation impairment induced by myeloma cells. *Leukemia* **27**, 451-463 (2013).
- 674  
675 38. Yuan Y, *et al.* The Wnt5a/Ror2 noncanonical signaling pathway inhibits canonical Wnt signaling  
676 in K562 cells. *Int J Mol Med* **27**, 63-69 (2011).
- 677  
678 39. Ho TK. The random subspace method for constructing decision forests. *Ieee T Pattern Anal* **20**,  
679 832-844 (1998).
- 680  
681 40. Wang HH. Pattern classification with random decision forest. *2012 International Conference on*  
682 *Industrial Control and Electronics Engineering (Icicee)*, 128-130 (2012).
- 683  
684 41. Xue J, Zhao YX. Random-Forests-based phonetic decision trees for conversational speech  
685 recognition. *Int Conf Acoust Spee*, 4169-4172 (2008).
- 686  
687 42. Ay F, Noble WS. Analysis methods for studying the 3D architecture of the genome. *Genome Biol*  
688 **16**, 183 (2015).
- 689  
690 43. Whalen S, Truty RM, Pollard KS. Enhancer-promoter interactions are encoded by complex  
691 genomic signatures on looping chromatin. *Nat Genet* **48**, 488-496 (2016).
- 692  
693 44. Oti M, Falck J, Huynen MA, Zhou H. CTCF-mediated chromatin loops enclose inducible gene  
694 regulatory domains. *BMC Genomics* **17**, 252 (2016).
- 695  
696 45. Kramer O. Scikit-Learn. *Stud Big Data* **20**, 45-53 (2016).
- 697  
698 46. Sandhu KS, *et al.* Large-scale functional organization of long-range chromatin interaction  
699 networks. *Cell Rep* **2**, 1207-1219 (2012).
- 700  
701 47. Majumder P, Gomez JA, Chadwick BP, Boss JM. The insulator factor CTCF controls MHC class II  
702 gene expression and is required for the formation of long-distance chromatin interactions. *J Exp*  
703 *Med* **205**, 785-798 (2008).
- 704  
705 48. Majumder P, Boss JM. CTCF controls expression and chromatin architecture of the human major  
706 histocompatibility complex class II locus. *Mol Cell Biol* **30**, 4211-4223 (2010).
- 707  
708 49. Zhang Y, *et al.* Model-based analysis of ChIP-Seq (MACS). *Genome Biol* **9**, R137 (2008).
- 709



- 710 50. Trapnell C, *et al.* Differential gene and transcript expression analysis of RNA-seq experiments  
711 with TopHat and Cufflinks. *Nat Protoc* **7**, 562-578 (2012).
- 712
- 713 51. Mathelier A, *et al.* JASPAR 2016: a major expansion and update of the open-access database of  
714 transcription factor binding profiles. *Nucleic Acids Res* **44**, D110-115 (2016).
- 715
- 716 52. Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif. *Bioinformatics*  
717 **27**, 1017-1018 (2011).
- 718
- 719 53. Rosenbloom KR, *et al.* The UCSC Genome Browser database: 2015 update. *Nucleic Acids Res* **43**,  
720 D670-681 (2015).
- 721
- 722 54. Zang C, Schones DE, Zeng C, Cui K, Zhao K, Peng W. A clustering approach for identification of  
723 enriched domains from histone modification ChIP-Seq data. *Bioinformatics* **25**, 1952-1958  
724 (2009).
- 725
- 726 55. Pedregosa F, *et al.* Scikit-learn: Machine Learning in Python. *J Mach Learn Res* **12**, 2825-2830  
727 (2011).
- 728
- 729 56. Naumova N, Smith EM, Zhan Y, Dekker J. Analysis of long-range chromatin interactions using  
730 Chromosome Conformation Capture. *Methods* **58**, 192-203 (2012).

731

732

### 733 [Acknowledgement](#)

734 We thank Dr. Michael Beer for valuable discussions and suggestions. We thank Nick Waring,  
735 Stephanie Perkail and Coen Lap for proof-reading and editing. This research was supported by  
736 National Institute of Health (Division of Intramural Research of National Lung and Blood Institute)  
737 to Zhu, a Lemonade Stand Foundation Young Investigator Award, and National Cancer Institute  
738 grants (R00CA158582, R21CA182662, and R03CA212068) to Tzatsos, a George Washington  
739 University Cross-Disciplinary Research Fund to Tzatsos and Peng, and National Institute of  
740 Allergy and Infection Diseases grants (R21 AI113806, R01 AI121080) to Peng.

741

### 742 [Author Contributions](#)

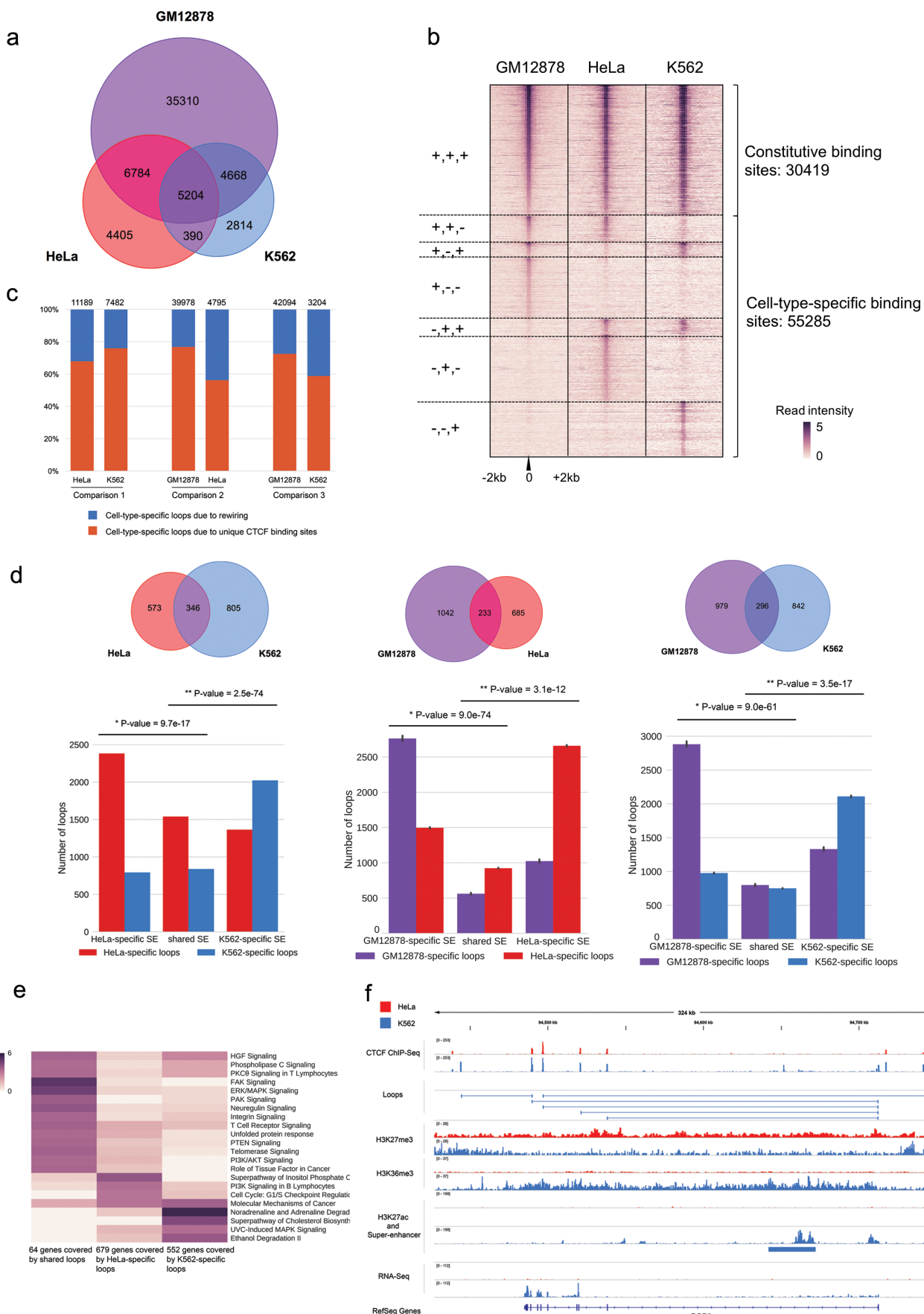
743 Y.K., A.T. and W.P. conceived the project. A.T. and W.P. supervised this study. Y.K. and W.P.  
744 developed the method and analyzed the results. Y.K. wrote the software. J.A. performed the 3C  
745 experiments. Z.Z. and J.Z. contributed to methodology design. Y.K., A.T. and W.P. wrote the  
746 manuscript. All authors discussed the results and commented on the manuscript.

747

### 748 [Competing financial interests](#)

749 The authors declare no competing financial interests.

Figure 1





**Figure 1.** CTCF-mediated loops exhibit cell-type-specificity.

**(a)** Venn diagram of CTCF-mediated loops identified from ChIA-PET experiments in GM12878, HeLa and K562.

**(b)** Heat map of CTCF binding sites in GM12878, HeLa and K562. Each row represents a CTCF binding event identified in ChIA-PET in at least one cell-type. The binding sites are divided into seven groups based on the presence (+) or absence (-) of CTCF binding. Color key shows the log<sub>2</sub>-transformed value of reads per kilobase per million reads (RPKM).

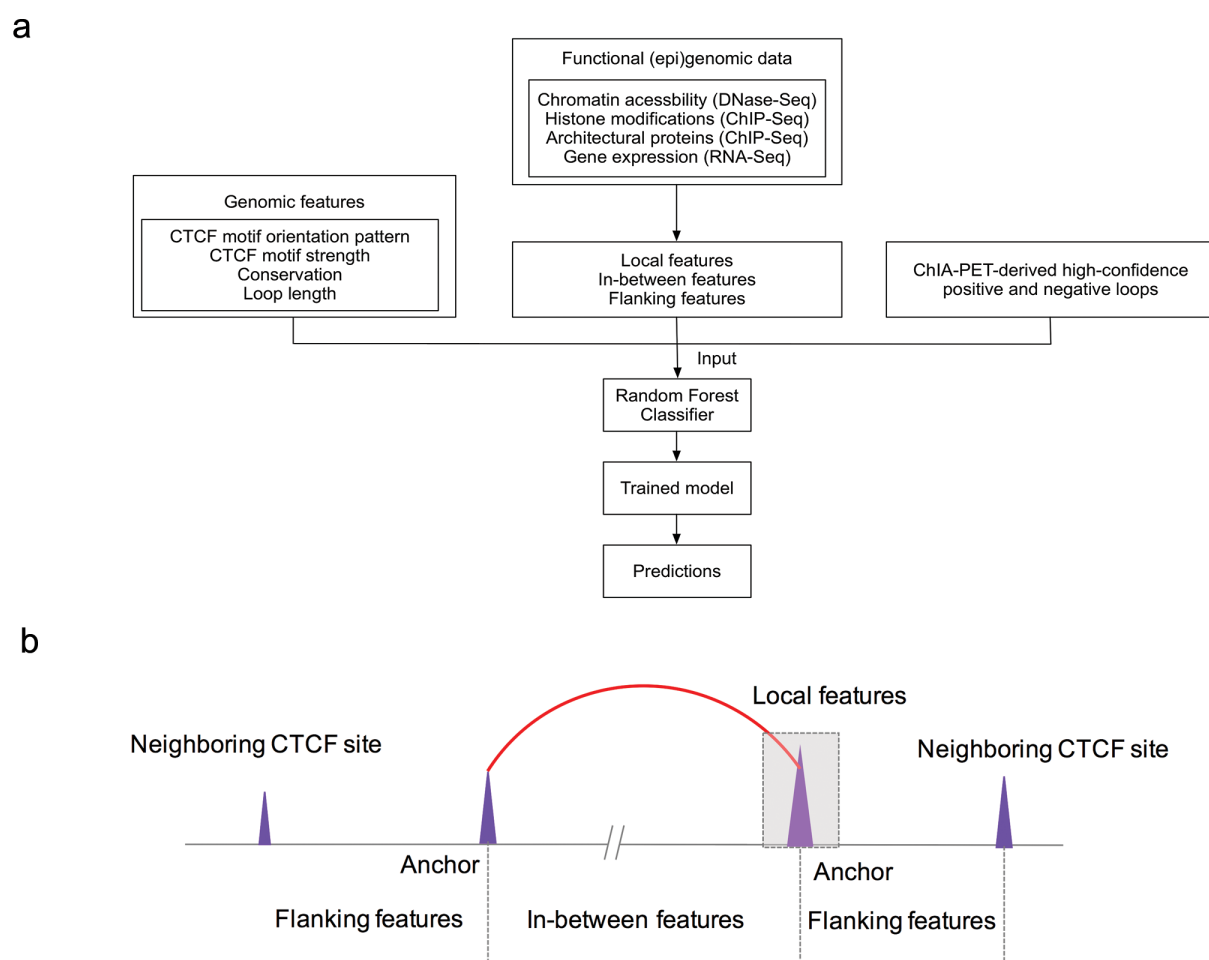
**(c)** Cell-type-specific CTCF binding and rewiring between common CTCF binding sites contribute to cell-type-specific loops.

**(d)** Cell-type-specific SEs are enriched with cell-type-specific loops. Top: Venn diagram of SEs in pairwise comparison of cell types. Bottom: Number of cell-type-specific loops covering cell-type-specific and shared SEs. P-values were calculated by Chi-square test. The GM12878 ChIA-PET dataset was down-sampled to 15% of the original size so that the number of identified loops matched those of the other ChIA-PET datasets. The down sampling and further analysis was repeated 10 times and the standard-deviations were shown.

**(e)** Canonical pathway enrichment analysis of differentially expressed genes associated with K562-specific, HeLa-specific and shared CTCF-mediated loops, respectively. Color represents the -log<sub>10</sub> (P-value).

**(f)** Genome browser snapshot of ROR2 locus. ROR2 is expressed and associated with CTCF-mediated loops in K562 but not in HeLa. Expression of ROR2 in K562 is associated with a concomitant decrease of H3K27me<sub>3</sub> and increase of H3K36me<sub>3</sub> within the gene body, as well as the appearance of a K562-specific SE. The ChIP-Seq and RNA-seq signals are represented in RPKM values.

## Figure 2

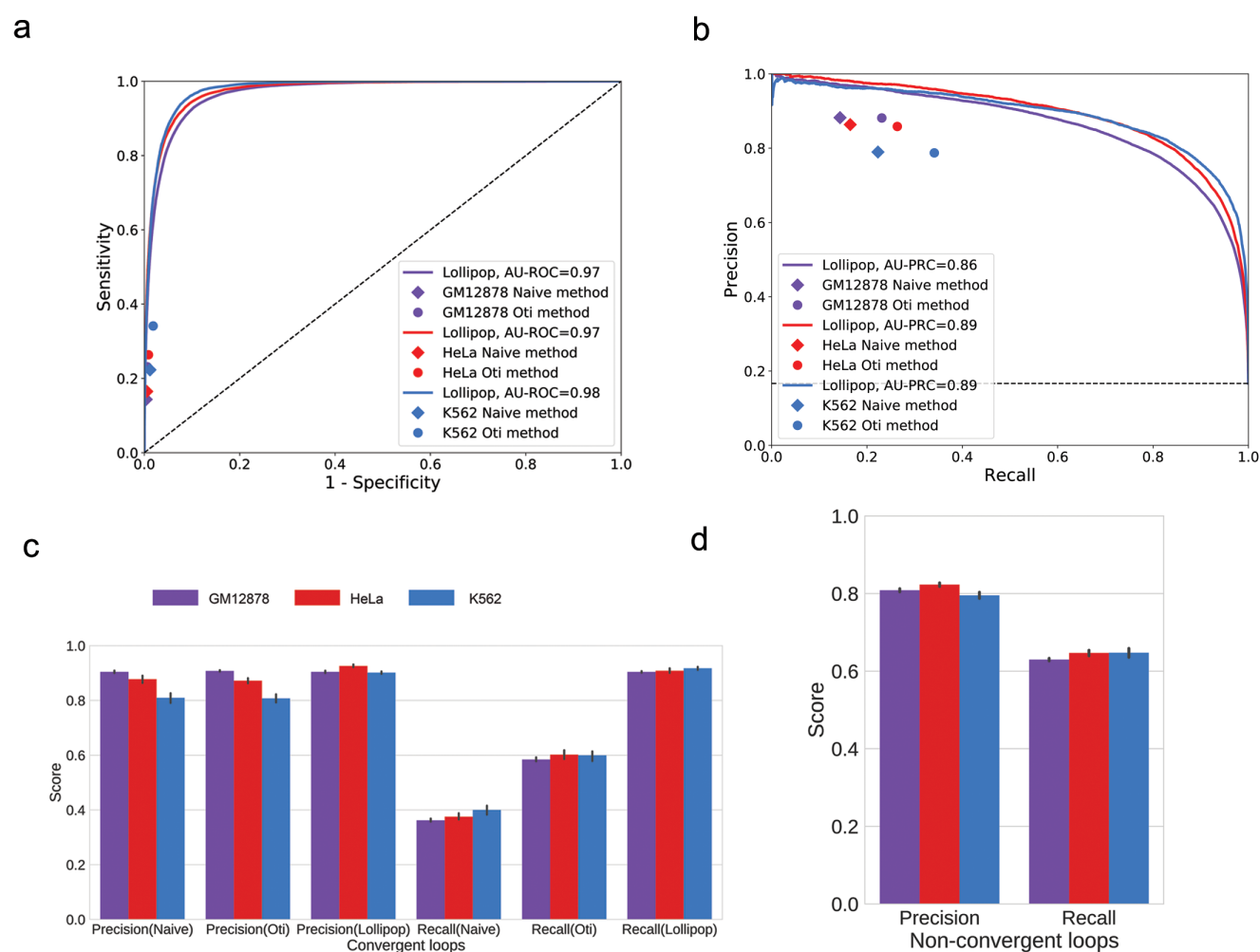


**Figure 2.** Illustration of the Lollipop pipeline and types of features.

(a) Schematic of the Lollipop pipeline. In training data, positive loops were generated from high-confidence interactions identified from ChIA-PET, and negative loops were random pairs of CTCF binding sites without interactions in ChIA-PET or significant contact in Hi-C dataset. A diverse set of features, generated from genomic and epigenomic data, was used to characterize the interactions. A random forests classifier distinguished interacting CTCF binding sites from non-interacting ones. The performance of resulting classifier was then evaluated. Trained model can be used to scan the genome and predict de novo CTCF-mediated loops in the same or a different cell-type.

(b) Illustration of local, in-between, and flanking features.

Figure 3



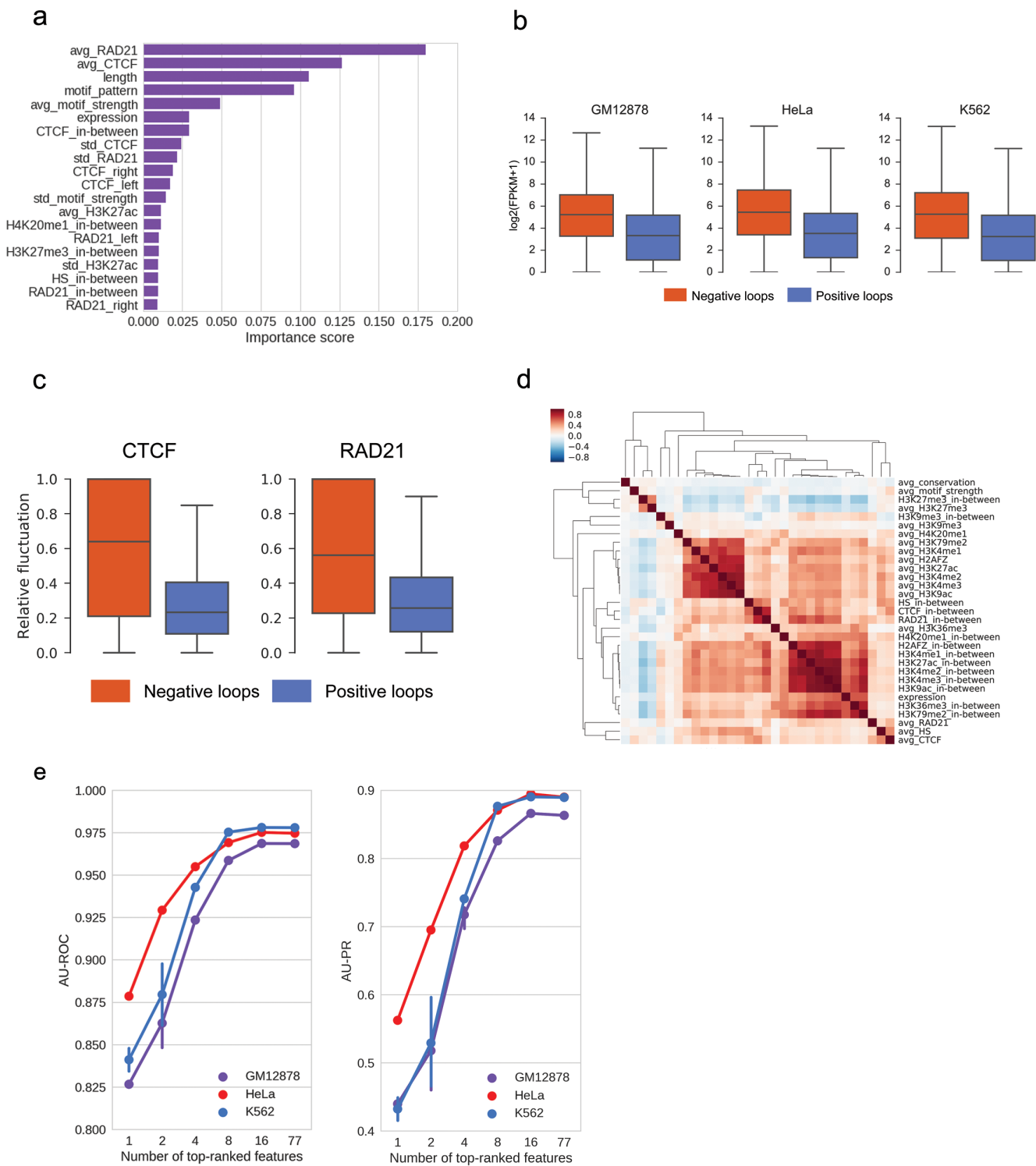
**Figure 3.** Performance evaluation within individual cell types

**(a,b)** Performance evaluation using (A) Receiver Operating Characteristic (ROC) and (B) Precision-Recall (PR) curve. Performance of the naïve and Oti methods are represented by diamonds and circles, respectively. Results in GM12878, HeLa and K562 are shown in purple, red and light blue, respectively.

**(c)** Comparison of the precision and the recall of the three methods in predicting convergent loops.

**(d)** Evaluation of Lollipop's performance on non-convergent loops, which include tandem loops, divergent loops and loops without CTCF motifs in the anchors.

Figure 4



**Figure 4.** Feature analysis identified novel determinants of CTCF-mediated chromatin loops.

(a) Ranking of predictive importance of the top 20 features in the model trained in GM12878 cells. Predictive importance is measured by mean decrease impurity in the training process. 'avg' and 'std' represent the mean and standard deviation of the signal intensity on both anchors. '\_left' and '\_right' represent the flanking features while '\_in-between' is the signal intensity within the loop.

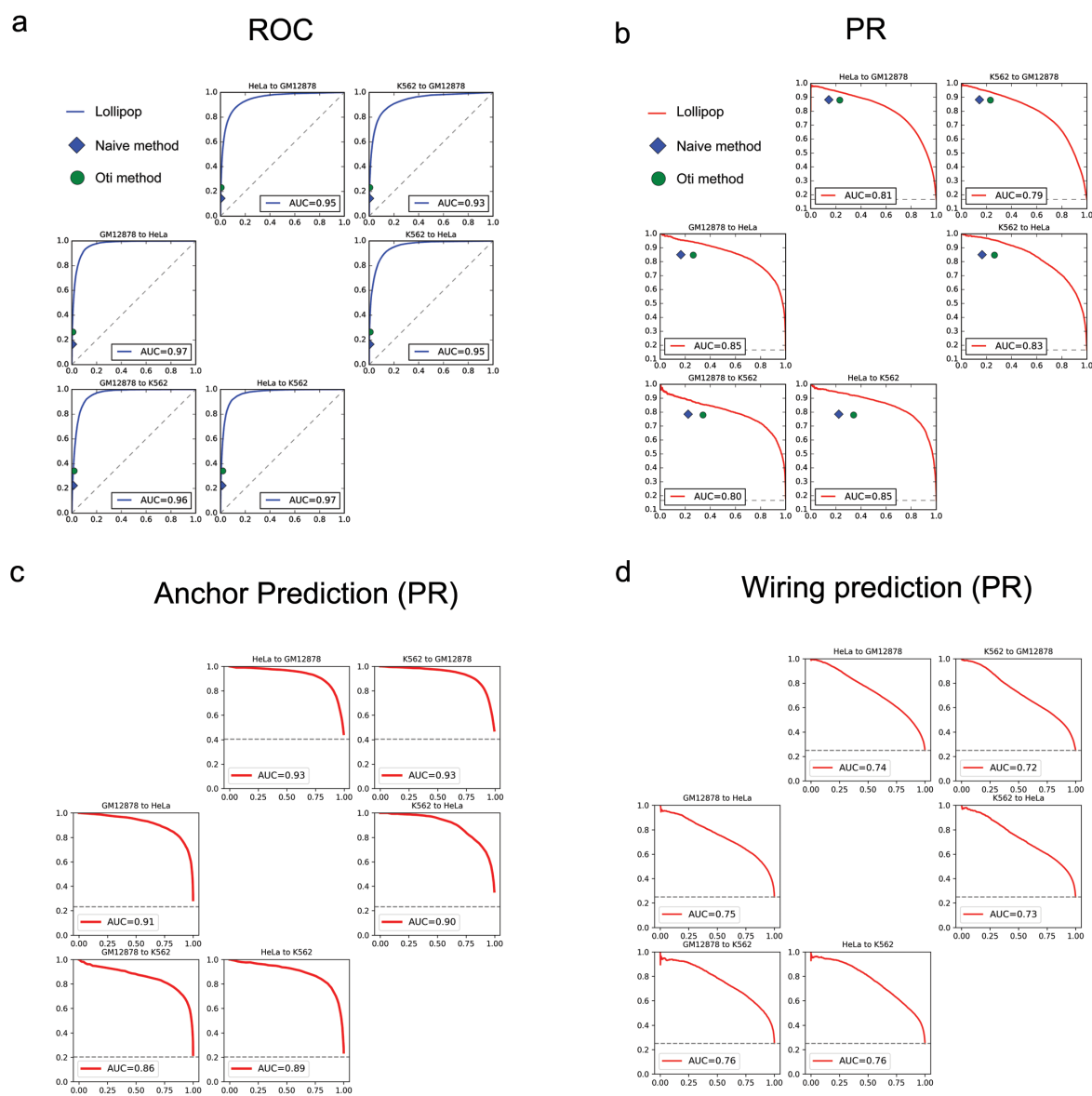
(b) Distributions of average gene expression levels within negative and positive loops. The positive and negative loops were defined in the training data, with those without any promoters inside the loops excluded in this analysis. In all three cases, P-value was  $< 1e-300$  using Mann-Whitney U test.

(c) Distribution of the relative fluctuations of CTCF and RAD21 binding intensities on paired anchors of negative and positive loops in GM12878 cells. Relative fluctuation was defined as the ratio of standard deviation to mean intensity of anchor pairs. In both cases, P-value was  $< 1e-300$  using Mann-Whitney U test.

(d) Heatmap of feature correlations in GM12878. On anchors, active histone marks are highly correlated. Along the loop regions, active histone marks and expression exhibit strong correlation. In addition, RAD21, CTCF and DNase hypersensitive sites are strongly correlated. Spearman's rank correlation and hierarchical clustering were used.

(e) Recursive Feature Elimination analysis on feature reduction. Left: AU-ROC; Right: AU-PR.

Figure 5



**Figure 5.** Assessment of Lollipop's performance across cell-types.

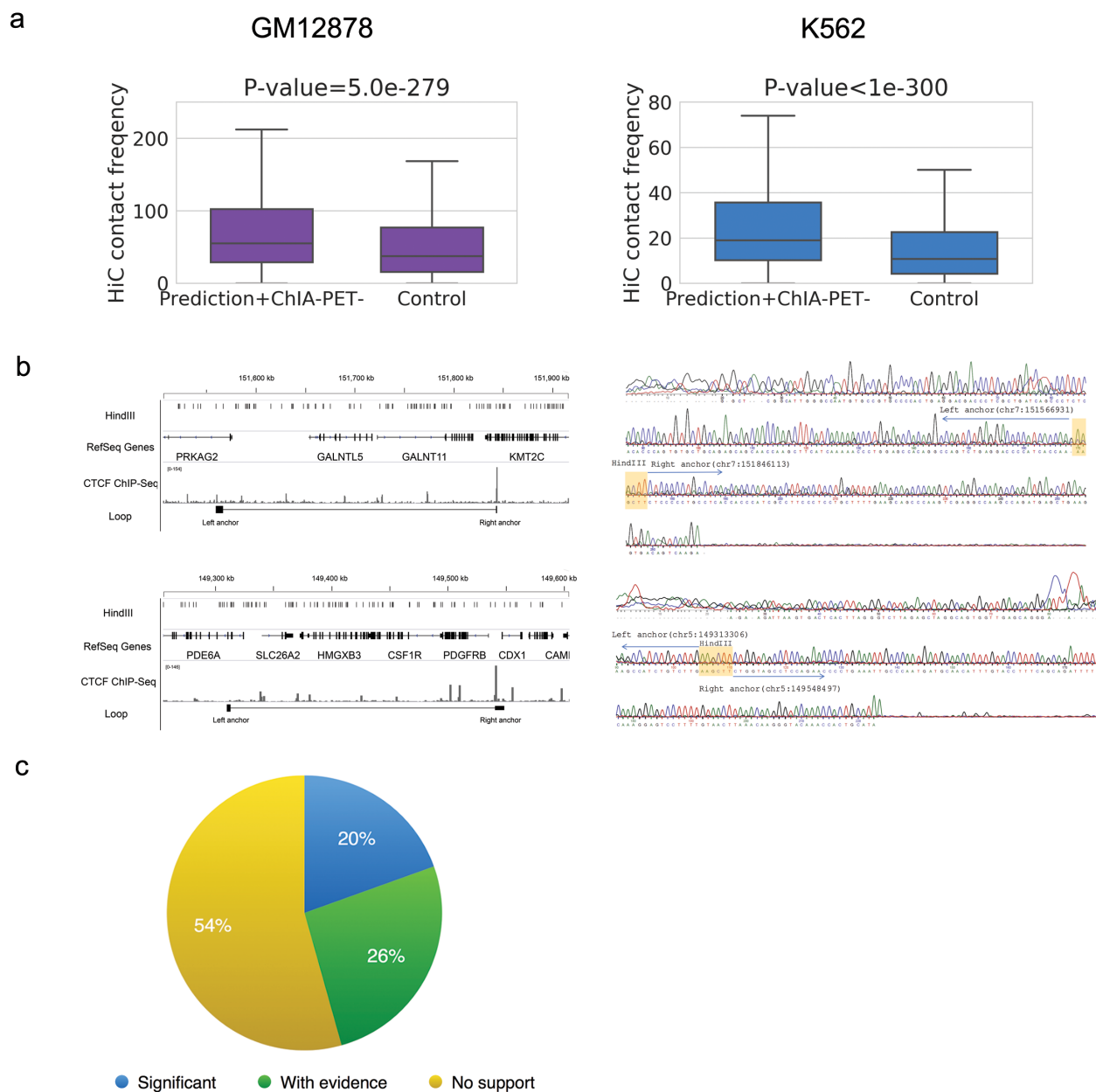
(a, b) Cross-cell-type performance evaluation using (A) ROC and (B) PR curves. In each subplot, 'cell A to cell B', applies the model trained from cell-type A to the data of cell-type B. For comparison, the performance of the naïve and Oti methods in each cell-type were represented by diamonds and circles, respectively.

(c) Performance evaluation of anchor prediction using PR curve.

(d) Performance evaluation of wiring prediction using PR curve.

The dash lines in (a-d) represent baseline performance.

Figure 6



**Figure 6. Validation of predicted CTCF-mediated interactions.**

(a) CTCF-mediated loops predicted by Lollipop but lacking ChIA-PET support exhibit significantly higher contact frequency than background in Hi-C experiments. P-values were calculated using Mann-Whitney U test.

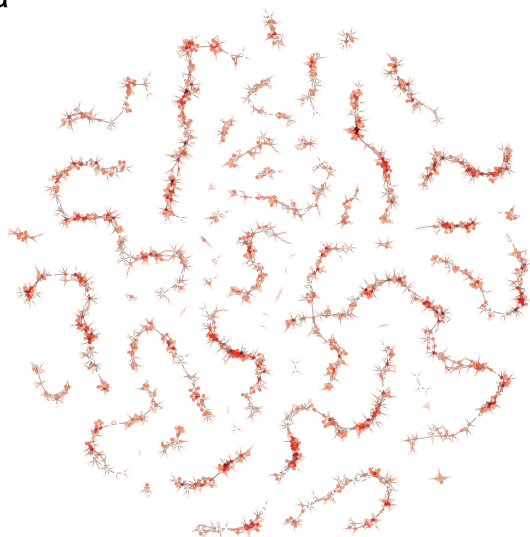
(b) Validation of two loops predicted by Lollipop, but not present in the HeLa ChIA-PET data set. Left: schematic of PRKAG2-KMT2C (chr7:151560677-151843260; top) and PDE6A-PDGFRB loop (chr5:149312517-149547724; bottom). Right: Sanger sequencing confirmation of the ligation junctions. Shaded areas in the right panels indicate the HindIII ligation junctions.

(c) Scaling analysis of loop prediction. Loops predicted using a model trained on the down-sampled (to 15%) GM12878 library, but lacks support in the down-sampled library (i.e., the yellow slice in **Supplementary Fig. 5b**) are evaluated by the full ChIA-PET data. 46% of these loops find support.

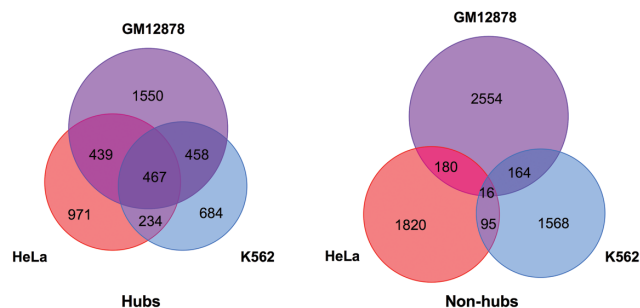


## Figure 7

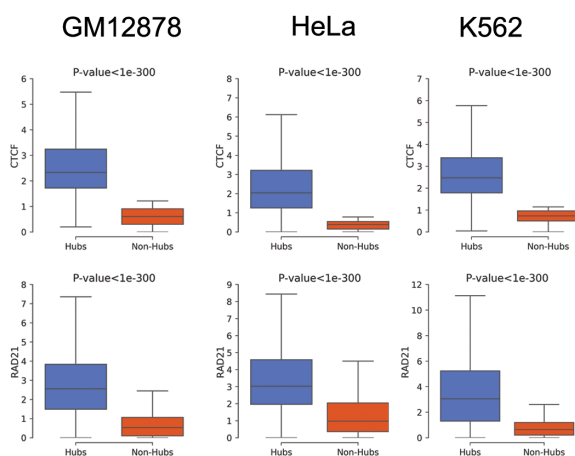
a



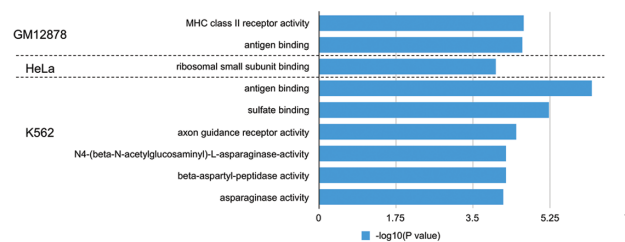
b



c



d



**Figure 7.** Topological properties of the CTCF-mediated interaction network and their association with biological functions.

(a) Visualization of the CTCF-mediated interaction network of chromosome 1 in GM12878 cells. Each node represents an anchor, with color representing degree of connection. Each edge represents an interaction.

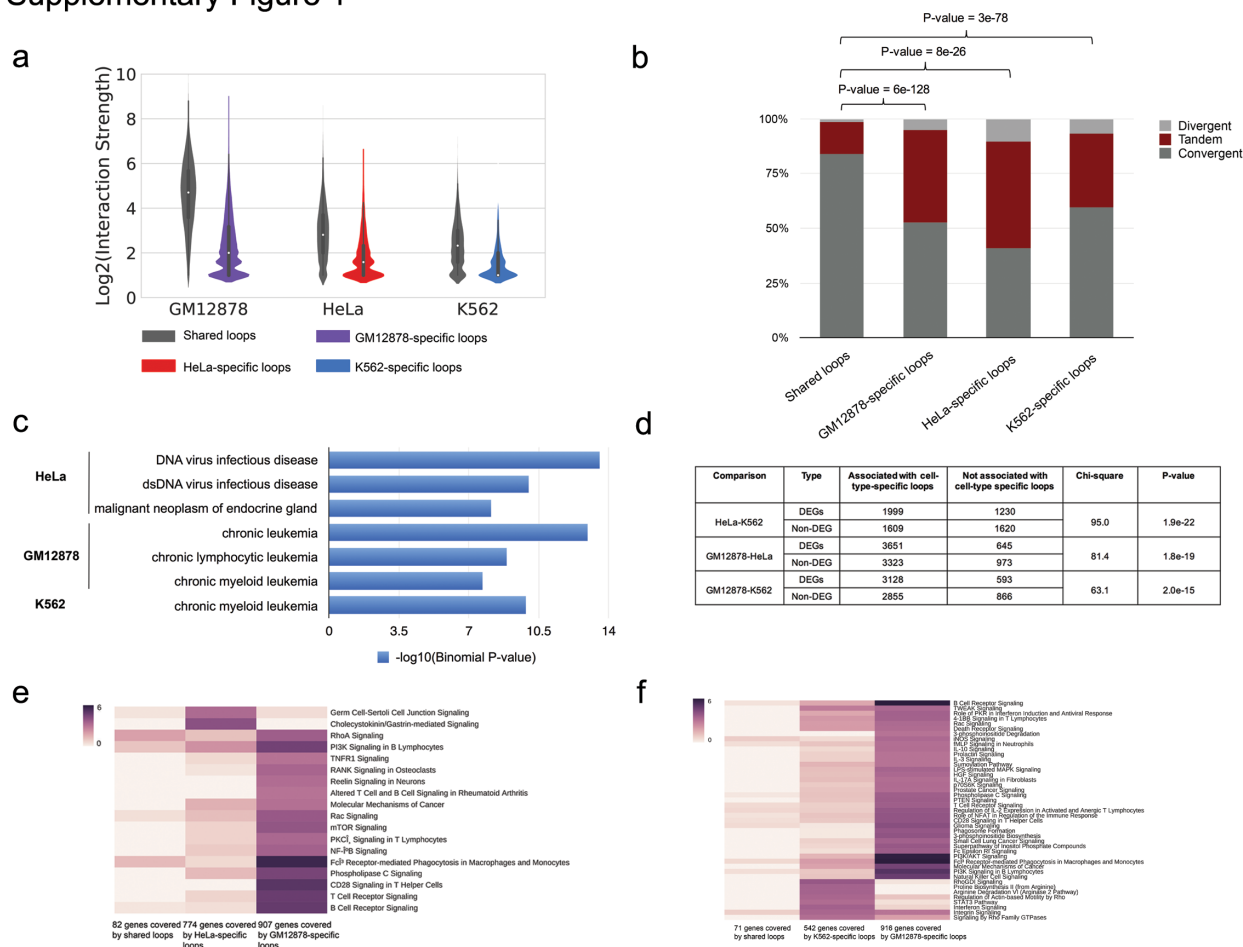
(b) Overlap of predicted hubs and non-hubs among each cell-type. Hubs are more conserved than non-hubs.

(c) Distribution of the binding affinity of architectural proteins, CTCF (top) and RAD21 (bottom), on predicted hubs and non-hubs.

(d) Functional enrichment analysis of hubs using GREAT. The x-axis represents the binomial P-values.



## Supplementary Figure 1



### Supplementary Figure 1. CTCF-mediated loops exhibit cell-type-specificity.

(a) Violin plots show that shared CTCF-mediated loops are stronger than cell-type-specific loops. Interaction strength is defined as the number of Paired-End Tags (PETs) connecting the anchors.

(b) Stacked bar plot comparing the pattern of motif orientation between cell-type-specific and shared loops. The P-values were calculated using Chi-square test.

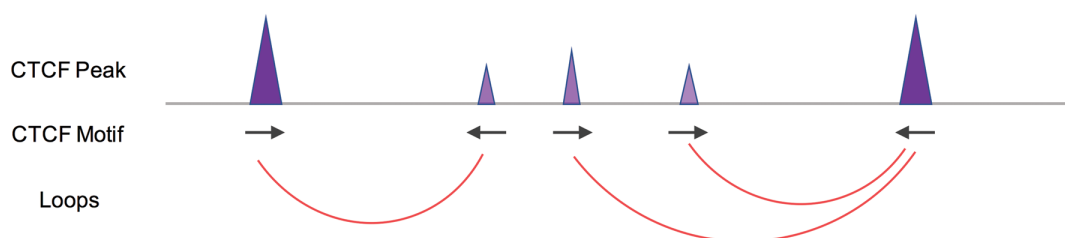
(c) Disease Ontology analysis of SEs using GREAT reveals the disease origin of the three cell-types.

(d) Contingency table for the number of loops associated with DEGs and Non-DEGs among the three cell lines. Pair-wise comparison was shown.

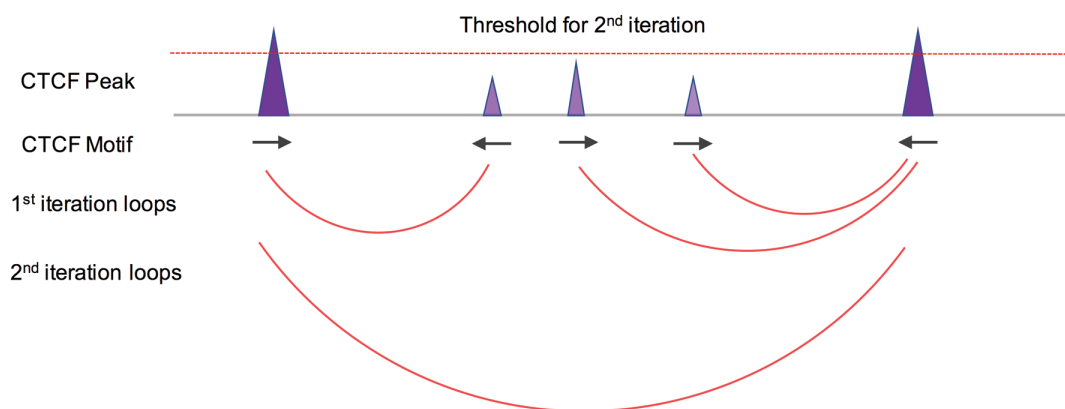
(e-f) Canonical pathway enrichment analysis of DEGs associated with cell-type-specific and shared loops in (e) HeLa-GM12878 and (f) K562-GM12878 comparison. Color represents the  $-\log_{10}$  (P-value).

## Supplementary Figure 2

a



b

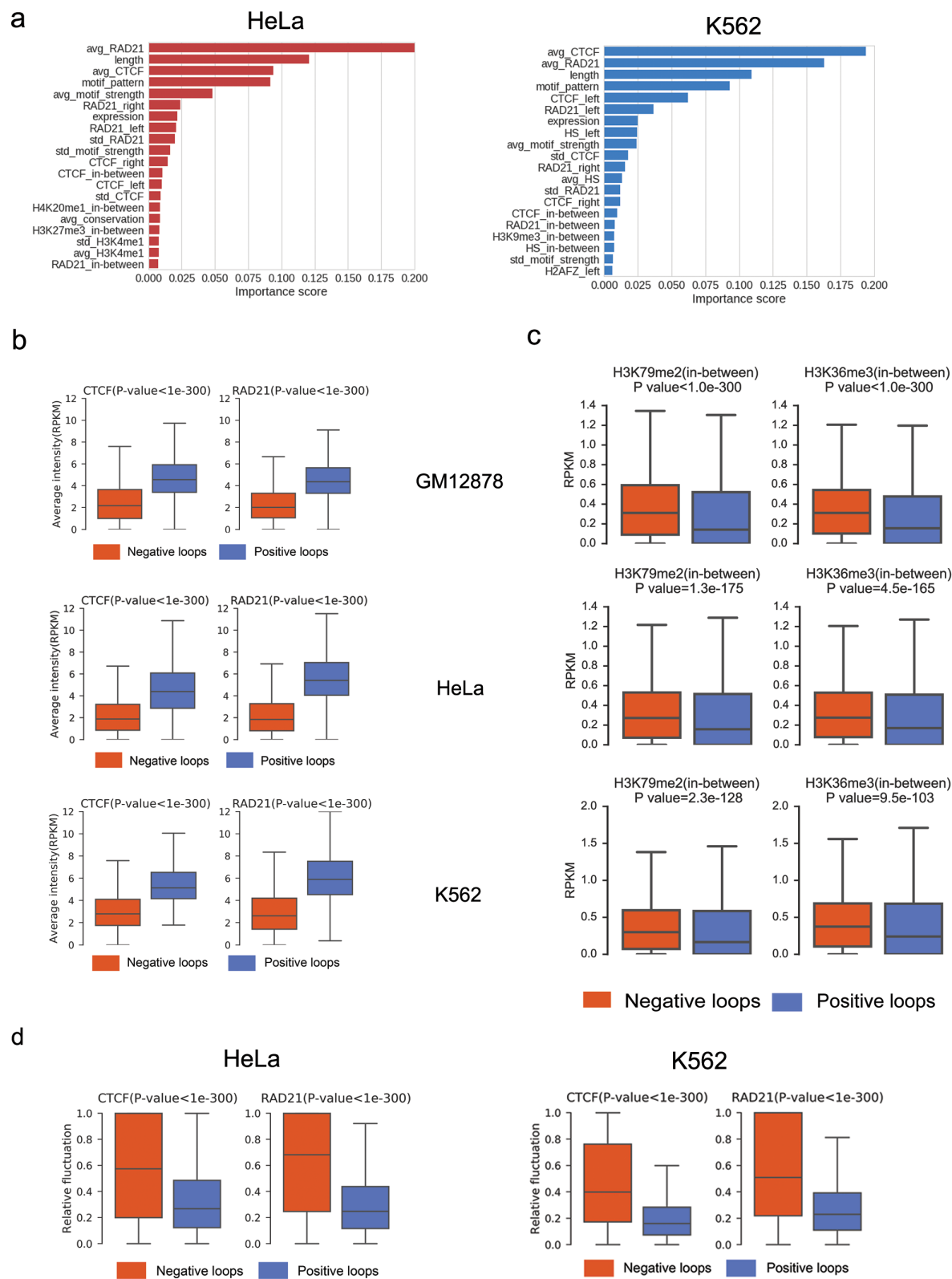


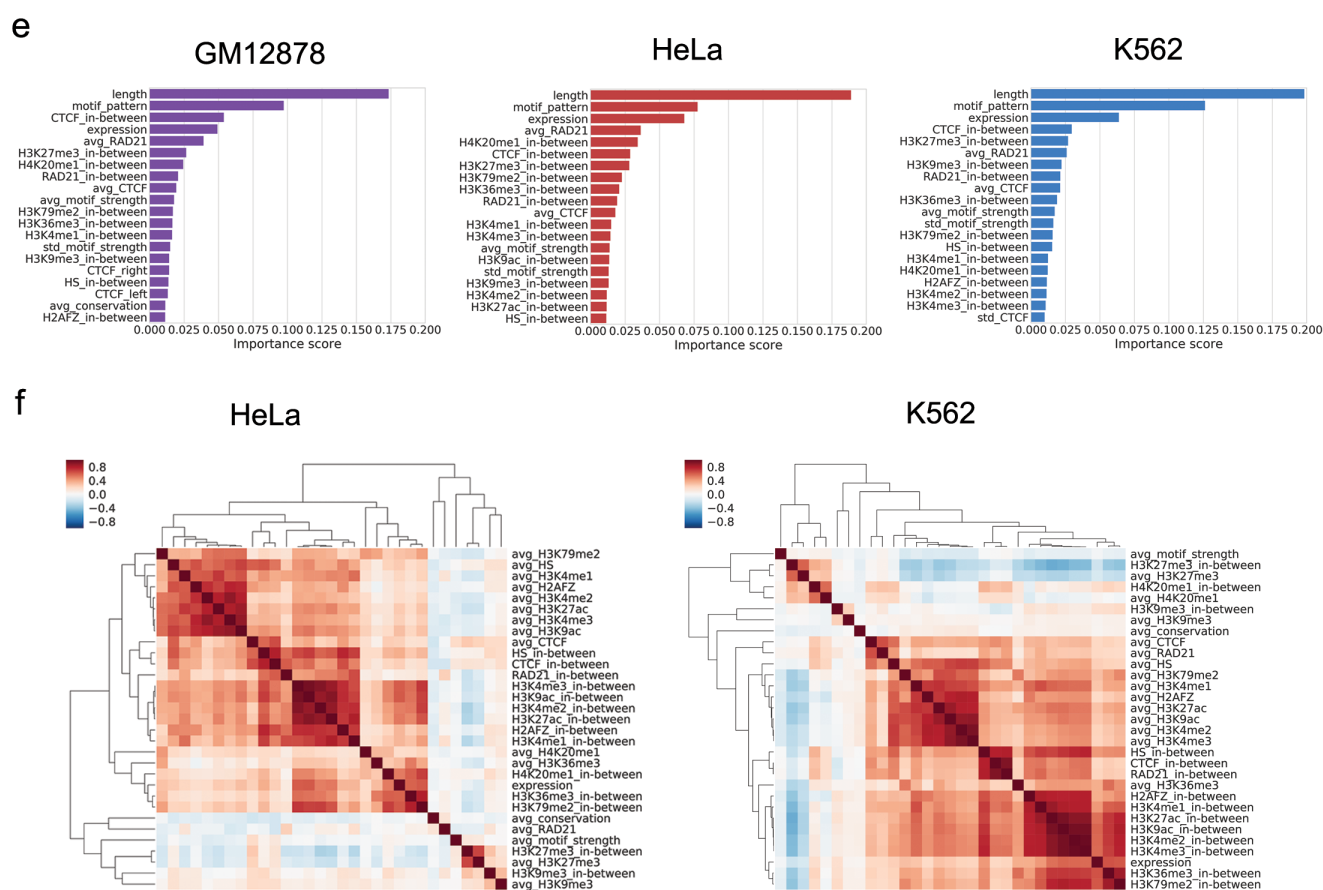
### Supplementary Figure 2. Illustration of the naïve and Oti method.

(a) Illustration of the naïve method. This method pairs a CTCF-bound motif that resides on the forward strand to the nearest downstream CTCF-bound motif that resides on the reverse strand.

(b) Illustration of the Oti method. It constructed loops in iterations by increasing the threshold of CTCF binding intensity. In each iteration, CTCF-bound motifs whose binding intensity are above the threshold were chosen, and naïve method was applied to construct loops. The loops constructed in different iterations were pooled together for the eventual result.

## Supplementary Figure 3





**Supplementary Figure 3.** Results of feature analysis in K562 and HeLa cells are consistent with those in GM12878.

(a) Ranking the predictive importance of the top 20 features in the model trained in HeLa and K562.

(b) Distributions of the average binding intensity of CTCF and RAD21 on anchors in negative and positive loops in the three cell lines. P-values were calculated using Mann-Whitney U test.

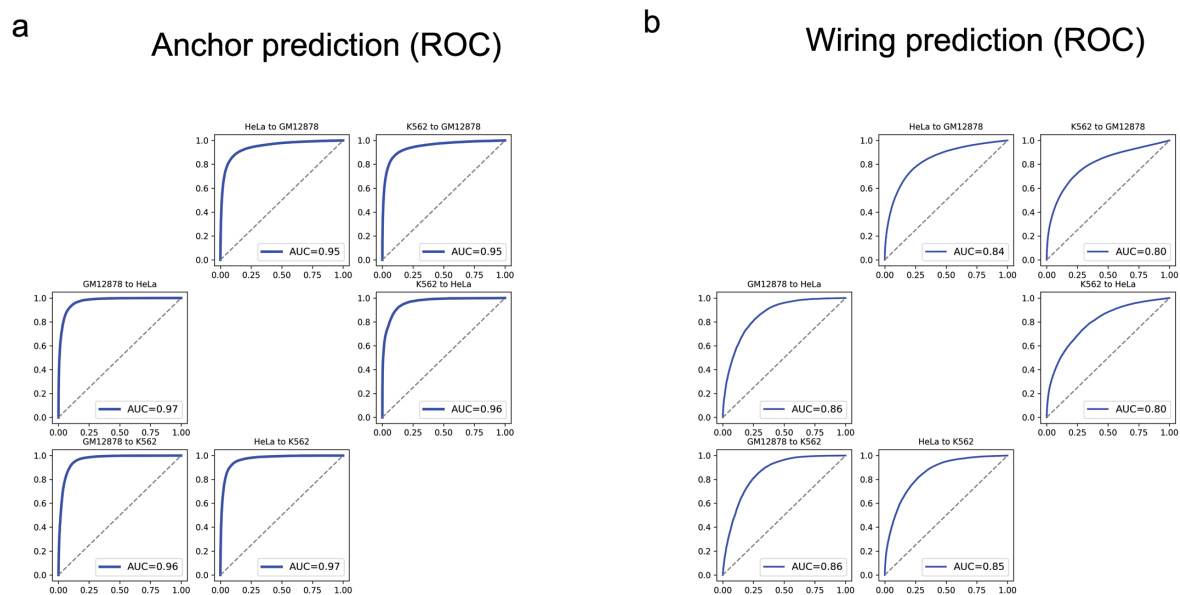
(c) Distributions of the intensity of the indicated histone marks within negative and positive loops.

(d) Distributions of relative fluctuations of CTCF and RAD21's binding intensities on paired anchors of negative and positive loops in HeLa and K562 cells. Relative fluctuation was defined as the ratio of standard deviation to average value.

(e) Ranking the predictive importance of the top 20 features in wiring prediction. The model was trained in the rewiring data (see methods for details) of GM12878, HeLa and K562 cells, respectively.

(f) Heatmaps of feature correlation in HeLa and K562 cells.

## Supplementary Figure 4

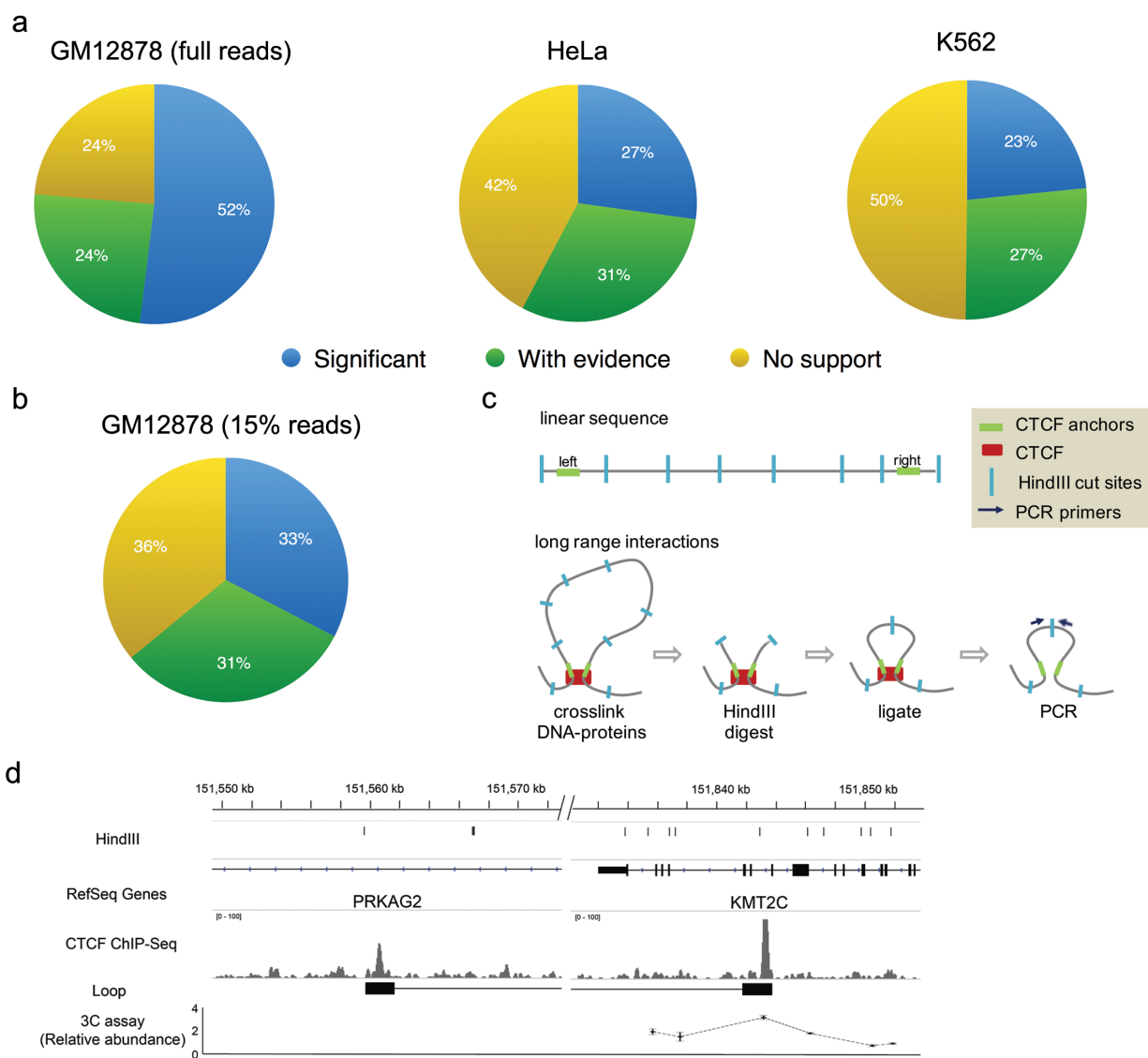


**Supplementary Figure 4.** Assessment of Lollipop's performance across cell-types.

(a) Performance evaluation of anchor prediction using ROC curve.

(b) Performance evaluation of wiring prediction using ROC curve.

## Supplementary Figure 5



### Supplementary Figure 5. Validation of predicted CTCF-mediated interactions

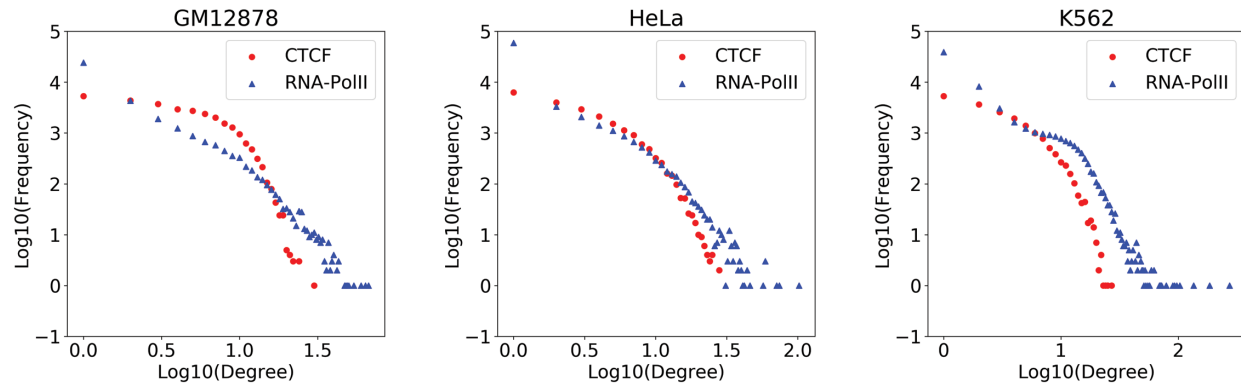
(a) The distribution of de novo predicted loops, as compared to original ChIA-PET data. ‘Significant’ denotes loops with  $FDR \leq 0.05$  and PET number  $\geq 2$  in ChIA-PET. ‘With evidence’ denotes predicted loops with less-significant ChIA-PET support (i.e.,  $FDR > 0.05$  or PET = 1). ‘No support’ denotes predicted loops without any ChIA-PET support.

(b) The distribution of predicted loops using a down-sampled (to 15%) GM12878 library for model building, followed by genome-wide prediction and comparison with ChIA-PET data. The number of loops observed in the downscaled GM12878 library is similar to those of K562 and HeLa (see **Supplementary Table 2**).

(c) Illustration demonstrating the major steps of 3C experiments.

(d) 3C-qPCR analysis shows the relative abundance of PRKAG2 anchor to KMT2C anchor and adjacent HindIII fragments (**Fig. 6b** top panel). Tracks from top to bottom: HindIII cut sites, designed primer for testing interaction, CTCF ChIP-Seq, motif occurrences, and relative quantification of the 3C interaction.

## Supplementary Figure 6



**Supplementary Figure 6.** The connection degree distribution for the CTCF- and RNA-PolIII-mediated interaction network. De novo predictions from Lollipop were used for the CTCF network, whereas loops identified from RNA-PolIII ChIA-PET were used for the RNA-PolIII network.

### Supplementary Table 1: Used data sets

Data	GM12878	HeLa	K562
ChIA-PET	GSE72816 <sup>4</sup>	GSE72816 <sup>4</sup>	ENCLB559JAA <sup>1,2</sup>
Hi-C	GSE63525 <sup>3</sup>		GSE63525 <sup>3</sup>
DNase-Seq	ENCFF000SKV <sup>1,2</sup>	ENCFF000SPJ <sup>1,2</sup>	ENCFF000SVI <sup>1,2</sup>
RNA-Seq	ENCFF000FBU <sup>1,2</sup> ENCFF000FBV <sup>1,2</sup>	ENCFF158RCK <sup>1,2</sup> ENCFF169ZTB <sup>1,2</sup>	GSM765393 <sup>1,2</sup>
ChIP-Seq (CTCF)	ENCFF000ARG <sup>1,2</sup>	ENCFF000BAJ <sup>1,2</sup>	ENCFF000YLT <sup>1,2</sup>
ChIP-Seq (RAD21)	ENCFF000OBV <sup>1,2</sup>	ENCFF000XKH <sup>1,2</sup>	ENCFF084HTD <sup>1,2</sup>
ChIP-Seq (H2AZ)	ENCFF001SUD <sup>1,2</sup>	ENCFF000BAX <sup>1,2</sup>	ENCFF000BWO <sup>1,2</sup>
ChIP-Seq (H3K4me1)	ENCFF000ARY <sup>1,2</sup>	ENCFF000BBA <sup>1,2</sup>	ENCFF000BXK <sup>1,2</sup>
ChIP-Seq (H3K4me2)	ENCFF000ATG <sup>1,2</sup>	ENCFF000BCH <sup>1,2</sup>	ENCFF000BXT <sup>1,2</sup>
ChIP-Seq (H3K4me3)	ENCFF000ATS <sup>1,2</sup>	ENCFF000BCO <sup>1,2</sup>	ENCFF000BXW <sup>1,2</sup>
ChIP-Seq (H3K9ac)	ENCFF000ATY <sup>1,2</sup>	ENCFF000BCW <sup>1,2</sup>	ENCFF000BYK <sup>1,2</sup>
ChIP-Seq (H3K9me3)	ENCFF000AUH <sup>1,2</sup>	ENCFF000BBG <sup>1,2</sup>	ENCFF000BYT <sup>1,2</sup>
ChIP-Seq (H3K27ac)	ENCFF000ASI <sup>1,2</sup>	ENCFF000BBN <sup>1,2</sup>	ENCFF000BWZ <sup>1,2</sup>
ChIP-Seq (H3K27me3)	ENCFF000ASK <sup>1,2</sup>	ENCFF000BBS <sup>1,2</sup>	ENCFF000BXA <sup>1,2</sup>
ChIP-Seq (H3K36me3)	ENCFF000ASX <sup>1,2</sup>	ENCFF000BCC <sup>1,2</sup>	ENCFF000BXE <sup>1,2</sup>
ChIP-Seq (H3K79me2)	ENCFF000ATT <sup>1,2</sup>	ENCFF000BCQ <sup>1,2</sup>	ENCFF000BYC <sup>1,2</sup>
ChIP-Seq (H4K20me1)	ENCFF000AUT <sup>1,2</sup>	ENCFF000BDC <sup>1,2</sup>	ENCFF001QWY <sup>1,2</sup>
ChIP-Seq Input	ENCFF000AQZ <sup>1,2</sup>	ENCFF000BAI <sup>1,2</sup>	ENCFF000BVZ <sup>1,2</sup>
ChIP-Seq Input	ENCFF651WEV <sup>1,2</sup>	ENCFF469INX <sup>1,2</sup>	ENCFF000QEK <sup>1,2</sup>

<sup>1</sup>Consortium, E. P. (2012). "An integrated encyclopedia of DNA elements in the human genome." *Nature* **489**(7414): 57-74.

<sup>2</sup>Sloan, C. A., E. T. Chan, J. M. Davidson, V. S. Malladi, J. S. Strattan, B. C. Hitz, I. Gabdank, A. K. Narayanan, M. Ho, B. T. Lee, L. D. Rowe, T. R. Dreszer, G. Roe, N. R. Podduturi, F. Tanaka, E. L. Hong and J. M. Cherry (2016). "ENCODE data at the ENCODE portal." *Nucleic Acids Res* **44**(D1): D726-732.

<sup>3</sup>Rao, S. S., M. H. Huntley, N. C. Durand, E. K. Stamenova, I. D. Bochkov, J. T. Robinson, A. L. Sanborn, I. Machol, A. D. Omer, E. S. Lander and E. L. Aiden (2014). "A 3D map of the human genome at kilobase resolution reveals principles of chromatin looping." *Cell* **159**(7): 1665-1680.

<sup>4</sup>Tang, Z., O. J. Luo, X. Li, M. Zheng, J. J. Zhu, P. Szalaj, P. Trzaskoma, A. Magalska, J. Wlodarczyk, B. Rusczycki, P. Michalski, E. Piecuch, P. Wang, D. Wang, S. Z. Tian, M. Penrad-Mobayed, L. M. Sachs, X. Ruan, C. L. Wei, E. T. Liu, G. M. Wilczynski, D. Plewczynski, G. Li and Y. Ruan (2015). "CTCF-Mediated Human 3D Genome Architecture Reveals Chromatin Topology for Transcription." *Cell* **163**(7): 1611-1627.



### Supplementary Table 2: Analysis results of ChIA-PET data sets

Cell-type	Raw reads (in million)	Unique PETs (in million)	IAB $\geq$ 2 loops	FDR $\leq$ 0.05 loops	IAB $\geq$ 2 and FDR $\leq$ 0.05 loops
GM12878 (full reads)	680	39.8	93914	73511	51966
GM12878 (15% reads)	102	13.1	37125	22248	15569
HeLa	531	21.1	42430	25047	16783
K562	195	6.6	23884	23377	13076

## Supplementary Table 3: Top-ranked features from the Recursive Feature Elimination analysis

<p>*Numbers inside the parentheses indicate the times of top-ranked feature set appears.          *'avg' and 'std' represent the mean and standard deviation of the signal intensity on both anchors. '_left' and '_right' represent the flanking features while '_in-between' means the signal intensity in the loop region.</p>	
	GM12878
<b>Top 1 feature</b>	avg_RAD21 (5)
<b>Top 2 features</b>	avg_CTCF, avg_RAD21 (3) motif_pattern, avg_RAD21 (2)
<b>Top 4 features</b>	length, motif_pattern, avg_CTCF, avg_RAD21 (4) length, avg_CTCF, CTCF_in-between, avg_RAD21 (1)
<b>Top 8 features</b>	length, motif_pattern, avg_motif_strength, avg_CTCF, std_CTCF, CTCF_in-between, avg_RAD21, expression (5)
<b>Top 16 features</b>	length, motif_pattern, avg_motif_strength, HS_in-between, avg_H3K4me1, avg_H3K27ac, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, CTCF_right, avg_RAD21, std_RAD21, RAD21_in-between, RAD21_left, expression (2) length, motif_pattern, avg_motif_strength, HS_in-between, HS_left, avg_H3K4me1, avg_H3K27ac, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, CTCF_right, avg_RAD21, std_RAD21, RAD21_in-between, expression (3)
	HeLa
<b>Top 1 feature</b>	avg_RAD21 (5)
<b>Top 2 features</b>	length, avg_RAD21 (5)
<b>Top 4 features</b>	length, motif_pattern, avg_CTCF, avg_RAD21 (5)
<b>Top 8 features</b>	length, motif_pattern, avg_motif_strength, avg_CTCF, CTCF_in-between, avg_RAD21, std_RAD21, expression (3) length, motif_pattern, avg_motif_strength, avg_H3K4me1, avg_CTCF, avg_RAD21, std_RAD21, expression (2)
<b>Top 16 features</b>	length, motif_pattern, avg_motif_strength, avg_HS, avg_H3K4me1, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, CTCF_right, avg_RAD21, std_RAD21, RAD21_in-between, RAD21_left, RAD21_right, expression (5)
	K562
<b>Top 1 feature</b>	avg_CTCF (3) avg_RAD21 (2)
<b>Top 2 features</b>	length, avg_CTCF (1) avg_CTCF, avg_RAD21 (3) length, avg_RAD21 (1)
<b>Top 4 features</b>	length, avg_CTCF, CTCF_left, avg_RAD21 (5)
<b>Top 8 features</b>	length, motif_pattern, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, avg_RAD21, expression (3) length, motif_pattern, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, avg_RAD21, RAD21_left (2) length, motif_pattern, HS_left, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, avg_RAD21 (1)
<b>Top 16 features</b>	length, motif_pattern, avg_motif_strength, HS_in-between, HS_left, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, CTCF_right, avg_RAD21, std_RAD21, RAD21_in-between, RAD21_left, RAD21_right, expression (2) length, motif_pattern, avg_motif_strength, std_HS, HS_in-between, HS_left, avg_CTCF, std_CTCF, CTCF_in-between, CTCF_left, CTCF_right, avg_RAD21, std_RAD21, RAD21_in-between, RAD21_left, expression; (3)

#### Supplementary Table 4: Designed primers for 3C validation

Primer Name	Sequence (5' to 3')
KMT2C_U2	FGGAGAGGATGATGGTGCTGTGTAT
KMT2C_U1	CTTGATCGTTTCTCACTCCTTTCA
KMT2C_L	CTTGACTGTCACCTTCAGCTCATC
KMT2C_D1	GACATACCAGAGCAATAACCTGGA
KMT2C_D3	AGCAGCAAATGAATCAGCTCAG
KMT2C_D4	AGTGGTGTCAATGCTGGTTTTTC
KMT2C_R	ATCACTGTCTAGCTGCCCGTTC
PDGFRB_L	TATGCAGTGGTTTGTACCCTTG
PDGFRB_R	GTGGCACCATAATCATCCCTAT