

DiffGR: Detecting Differentially Interacting Genomic Regions from Hi-C Contact Maps

Huiling Liu and Wenxiu Ma*

Department of Statistics, University of California Riverside, Riverside, CA 92521, USA

*Corresponding author: wenxiu.ma@ucr.edu

Abstract

Recent advances in Hi-C techniques have allowed us to map genome-wide chromatin interactions and uncover higher-order chromatin structures, thereby shedding light on the principles of genome architecture and functions. However, statistical methods for detecting changes in large-scale chromatin organization such as topologically-associating domains (TADs) are still lacking. We proposed a new statistical method, DiffGR, for detecting differentially interacting genomic regions at the TAD level between Hi-C contact maps. We utilized the stratum-adjusted correlation coefficient to measure similarity of local TAD regions. We then developed a non-parametric approach to identify statistically significant changes of genomic interacting regions. Through simulation studies, we demonstrated that DiffGR can robustly and effectively discover differential genomic regions under various conditions. Furthermore, we successfully revealed cell type-specific changes in genomic interacting regions in both human and mouse Hi-C datasets, and illustrated that DiffGR yielded consistent and advantageous results compared with state-of-the-art differential TAD detection methods. The DiffGR R code is published under the GNU GPL ≥ 2 license and is publicly available at <https://github.com/wmalab/DiffGR>.

21 1 Introduction

22 Recent developments of chromatin conformation capture (3C)-based techniques—including 4C [1],
23 5C [2], Hi-C [3–5], ChIA-PET [6], and Hi-ChIP [7]—have allowed high-throughput characterization
24 of pairwise chromatin interactions in the cell nucleus, and provided an unprecedented opportunity to
25 investigate the three-dimensional (3D) chromatin structures and to elucidate their roles in nuclear
26 organization and gene expression regulation. Among these techniques, Hi-C and its variants [8–10]
27 are of particular interest because of their ability to map chromatin interactions at a genome-wide
28 scale.

29 A Hi-C experiment yields a symmetric contact matrix in which each entry represents the chro-
30 matin contact frequency between the corresponding pair of genomic loci. A particularly important
31 characteristic of Hi-C contact matrices is the presence of the topologically-associating domains
32 (TADs), which are functional units of chromatin with higher tendency of intra-domain interactions
33 [11]. TADs are largely conserved across cell types and species. Moreover, CTCF and other chro-
34 matin binding proteins are enriched at the TAD boundaries, indicating that TAD boundary regions
35 form chromatin loops and play an essential role in gene expression regulation [11, 12].

36 Several computational methods have been developed to detect TADs in Hi-C contact maps.
37 These methods can be categorized into two groups: one-dimensional (1D) statistic-based methods
38 and two-dimensional (2D) contact matrix-based methods [13]. Of these, 1D statistic-based methods
39 often take a sliding window approach along the diagonal of Hi-C contact matrix and compute a
40 1D statistic for each diagonal bin to detect TADs and/or TAD boundaries. For instance, Dixon
41 et al. [11] introduced a statistic named directionality index (DI) to quantify whether a genomic
42 locus preferentially interacts with upstream or downstream loci and developed a hidden Markov
43 model to call TADs from DIs. Later, Crane et al. [14] proposed a novel TAD detection method,
44 which computes an insulation score (IS) for each genomic bin by aggregating chromatin interac-
45 tions within a square sliding through the diagonal and then searches for the minima along the
46 IS profile as TAD boundaries. Unlike the 1D statistic-based methods which calculate statistics
47 using local information, the 2D contact matrix-based methods utilize global information on the
48 contact matrix to capture TAD structures. For example, the Armatus algorithm [15] identifies
49 consistent TAD patterns across different resolutions by maximizing a quality scoring function of
50 domain partition using dynamic programming. In addition, Lévy-Leduc et al. [16] proposed a TAD
51 boundary detection method named HiCseg, which performs a 2D block-wise segmentation via a
52 maximum likelihood approach to partition each chromosome into its constituent TADs. Recently,
53 several review papers have quantitatively compared the performances of the aforementioned TAD-
54 calling methods and demonstrated that HiCseg detects a stable number of TADs against changes
55 of sequencing coverage and maintains the highest reproducibility among Hi-C replicates across all
56 resolutions when compared with other TAD-calling methods [17, 18].

57 With the fast accumulation of Hi-C datasets, there has been a growing interest in performing
58 differential analysis of Hi-C contact matrices. To date, several computational tools have been
59 developed for comparative Hi-C analysis, but the majority of them focused on the identification
60 of differential chromatin interactions (DCIs), which represent different chromatin looping events
61 between two Hi-C contact maps. In early studies, the most common strategy for DCI detection was
62 to use the fold change values between two Hi-C contact maps. For instance, Wang et al. [19] used a
63 simple fold-change strategy to detect the influence of estrogen treatment on chromatin interactions
64 in MCF-7 Hi-C samples. Additionally, Dixon et al. [20] utilized the fold change values of chromatin
65 interactions to train a random forest model to discover the epigenetic signals that were more

66 predictive of changes in interaction frequencies. In addition to these fold change-based approaches,
67 another commonly utilized method for detecting DCIs was the binomial model implemented by the
68 HOMER software [21]. In contrast, in more recent studies, count-based statistical methods, such
69 as edgeR [22] and DESeq [23], have been adopted to identify pairwise chromatin interactions that
70 show significant changes in contact frequencies. Among them, Lun and Smyth [24] presented a tool
71 named diffHic for rigorous detection of differential interactions by leveraging the generalized linear
72 model (negative binomial regression) of edgeR, and demonstrated that edgeR outperformed the
73 binomial model. Later, Stansfield et al. [25] introduced MD normalization and performed Z-tests
74 to detect statistically significant DCIs. While all these methods assumed independence among
75 pairwise interactions, which holds true only in coarse-resolution Hi-C maps, Djekidel et al. [26]
76 presented a novel method, named FIND, that takes into account the dependency of adjacent loci
77 at finer resolutions. Briefly, FIND utilizes a spatial Poisson process model to detect DCIs that
78 show significant changes in interaction frequencies of both themselves and their neighborhood bins.
79 Lastly, Cook et al. [27] introduced ACCOST to identify differential chromatin contacts by extending
80 the DESeq model used in RNA-seq analysis and repurposing the “size factor” to account for the
81 notable genomic-distance effect in Hi-C contact matrices.

82 In the cell nucleus, chromatin is organized at multiple levels, ranging from active and inactive
83 chromosomal compartments and sub-compartments (on a multi-Mb scale) [3, 9], TADs (0.5–2 Mb
84 on average) [11], to fine-scale chromatin interacting loops [8, 9]. Chromatin structures also exhibit
85 multi-scale differences among different cell types in their compartments, TADs, and chromatin
86 loops. Among these, changes in TAD organizations are of particular interest as TADs are strongly
87 linked to cell type-specific gene expression [11]. For example, Taberlay et al. [28] have shown that
88 genomic rearrangements in cancer cells are partly guided by changes in higher-order chromatin
89 structures, such as TADs. They discovered that some large TADs in normal cells are further
90 segmented into several smaller TADs in cancer cells, and these changes are tightly correlated with
91 oncogene expression levels. Current differential analyses of TAD structures between different cell
92 types and conditions are limited to the detection of TAD boundary changes. Recently, Chen et al.
93 [13] proposed a TAD boundary detection approach named HiCDB, which is constructed based
94 on local measures of relative insulation and multi-scale aggregation. In addition to calling TAD
95 boundaries in single Hi-C sample, HiCDB also provides differential TAD boundary detection using
96 the average values of relative insulation across multiple samples. Later, Cresswell and Dozmorov
97 [29] developed TADCompare, which uses a spectral clustering-derived metric named eigenvector gap
98 to identify differential and consensus TAD boundaries and track TAD boundary changes over time.
99 Lastly, TADreg [30] introduced a versatile regression framework which generalizes the insulation
100 score by estimating the relative insulating effects of genomic loci and adding a sparsity constraint.
101 The TADreg framework was designed for TAD boundary detection, but also allowed differential
102 TAD analysis across various conditions. The HiCDB, TADCompare and TADreg methods focused
103 on detecting changes in TAD boundaries rather than changes in chromatin organization within
104 TADs. However, differential TAD boundaries do not necessarily indicate differential chromatin
105 conformation within those regions. First, Hi-C contact matrices are often sparse and noisy, which
106 might lead to unstable detection of TAD boundaries. Second, chromatin interactions within a
107 TAD could be strengthened or weakened in another Hi-C sample, which would suggest different
108 patterns of chromatin organization within the same TAD region. Unfortunately, few methods have
109 been developed to detect differential TAD regions instead of boundaries. Recently, the Hi-C pre-
110 processing and analysis tool HiCExplorer [31–33] expanded its functions to capture differential TAD
111 regions by comparing the precomputed TAD regions on the target Hi-C map with the same regions
112 on the control map by accounting for the information in both intra-TAD and inter-TAD regions.

113 However, such comparison was only limited to the precomputed genomic regions in only one of
114 the Hi-C conditions. Thus, appropriate statistical methods for detecting differentially interacting
115 regions by considering TAD regions across both conditions are still lacking.

116 To tackle this problem, we developed a novel statistical method, DiffGR, for detecting differ-
117 ential genomic regions at TAD level between two Hi-C contact maps. Briefly, DiffGR utilizes the
118 stratum-adjusted correlation coefficient (SCC), which effectively eliminates the genomic-distance
119 effect in Hi-C data, to measure the similarity of local genomic regions between two contact matri-
120 ces. Subsequently, DiffGR applies a nonparametric permutation test on those SCC values to detect
121 genomic regions with statistically significant differential interactions. We demonstrated, through
122 simulation studies and real data analyses, that DiffGR can effectively and robustly identify differ-
123 entially interacting genomic regions at TAD level.

124 2 Methods

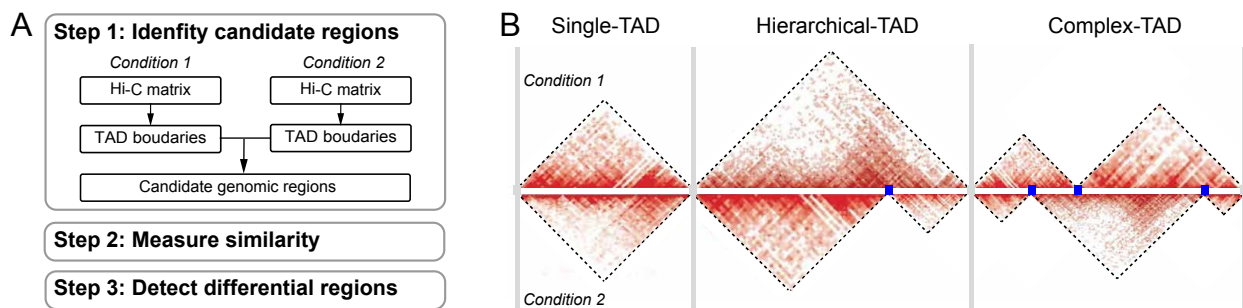


Figure 1: **Overview of DiffGR.** A. Workflow of the DiffGR algorithm. B. Illustration of three candidate types of differential genomic regions. The gray vertical bars represent the common TAD boundaries between two conditions, which partition the genome into three types of candidate regions. The blue points stand for unique TAD boundaries in only one of the two conditions.

125

126 The DiffGR method detects differentially interacting genomic regions in three steps, as shown in
127 Figure 1A and described below in Sections 2.1-2.3. In addition, the simulation settings are outlined
128 in Section 2.4 and real data preprocessing and analyses are described in Section 2.5.

129 2.1 Identifying candidate genomic regions

130 Suppose we have two sets of Hi-C data and their corresponding contact frequency matrices as
131 the input. First, we detect the TAD boundaries in each Hi-C data, separately. Specifically, we
132 apply HiCseg [16] to the raw contact matrices and obtain the corresponding TAD boundaries.
133 Note that in this step one can change HiCseg with any other TAD caller, whose detected TADs
134 satisfy the non-overlapping and continuous properties. We choose HiCseg because it has been
135 shown that HiCseg produces more robust and reliable TAD boundaries than other TAD-calling
136 methods [17, 34]. We next combine the TAD boundaries from both Hi-C contact maps to identify
137 the candidate genomic regions for subsequent analyses. TAD boundaries within two-bin distance
138 are considered to be a common boundary shared by both Hi-C datasets and replaced by the middle

139 bin locus. We then partition the genome into non-overlapping candidate regions using the common
140 TAD boundaries, and categorize these candidate regions into the following three groups: (1) single-
141 TAD candidate regions, (2) hierarchical-TAD candidate regions, and (3) complex-TAD candidate
142 regions, as illustrated in Figure 1B.

143 We expected different patterns of differential features in these three kinds of candidate genomic
144 regions. As to the differential single-TAD region, we would expect strength changes occurred in such
145 areas. For differential hierarchical-TAD regions, one large interacting domain could be evidently
146 split into two or more sub-domains, or vice versa, boundaries between TADs disappeared and thus
147 the corresponding domains merged in one of the contact maps. Lastly, domains might be split,
148 merged, or shifted in a more complicated manner thereby constructing an entirely new structure,
149 which would be defined as differential complex-TAD regions. Unlike differential single-TAD regions,
150 the differential hierarchical-TAD and complex-TAD regions represent more disruptive changes in
151 the 3D structure of the chromatin.

152 2.2 Measuring similarity of candidate regions between two Hi-C contact maps

153 In the second step, we evaluate the similarity of each candidate region between the two samples.
154 Suppose a candidate genomic region is bounded by two common TAD boundaries shared by both
155 Hi-C maps, and contains k unique TAD boundaries in either one of the two Hi-C maps (shown as
156 blue points in Figure 1B). In the single-TAD candidate region, $k = 0$; in the hierarchical-TAD or
157 complex-TAD candidate regions, $k \geq 1$. For each candidate region, we consider all $\binom{k+2}{2}$ possible
158 (sub)TADs, which are separated by any pair of TAD boundaries within that region, as potential
159 differential TADs. For each potential differential TAD, we calculate the stratum-adjusted correla-
160 tion coefficient (SCC) [35] rather than the standard Pearson or Spearman correlation coefficients
161 (CCs) to measure the similarity of intra-TAD chromatin interactions between two Hi-C samples.
162 The advantages of using SCC instead of standard CCs are shown in Supplementary Results in File
163 S1.

164 The SCC metric was introduced by Yang et al. [35] as a measure of similarity and reproducibility
165 between two Hi-C contact matrices. To account for the pronounced distance-dependence effect
166 in Hi-C contact maps, chromatin contacts are first stratified into K stratum according to the
167 genomic distances of the contacting loci pairs, and the correlation coefficients of contacts within
168 each stratum are calculated between two samples. These stratum-specific correlation coefficients
169 are then aggregated to compute the SCC value using a weighted average approach, where the
170 weights are derived from the Cochran-Mantel-Haenszel (CMH) statistic [36]. That is, the SCC ρ
171 is calculated as

$$\rho = \sum_{k=1}^K \left(\frac{N_k r_{2k}}{\sum_{k=1}^K N_k r_{2k}} \right) \rho_k,$$

172 where N_k is the number of elements in the k -th stratum, r_{2k} is the product of standard deviations
173 of the elements in the k -th stratum of both samples, and ρ_k denotes the correlation coefficient of
174 the k -th stratum between two samples.

175 The original SCC metric is computed using the intra-chromosomal contact matrices with a
176 predefined genomic distance limit. The resulting value has a range of $[-1, 1]$ and can be interpreted
177 in a way similar to the standard correlation coefficient. Here we use SCC as a local similarity
178 measurement to evaluate each potential differential TAD between two Hi-C samples. In the SCC

179 calculation, an upper limit of genomic distance is set to 10 Mb because TADs are commonly
180 smaller than 10 Mb and distal interactions over a genomic distance larger than 10 Mb are often
181 sparse and highly stochastic. In addition, as the sparsity of Hi-C matrices might affect the precision
182 of SCC values, the loci pairs with zero contact frequencies in both samples are excluded from the
183 calculation.

184 Hi-C contact maps are often sparse due to sequencing coverage limits and contain various
185 systematic biases. To solve these issues, when preprocessing the Hi-C contact matrices, we first
186 smooth each contact map by a 2D mean filter [35], which substitutes the contact count observed
187 between each bin pair by the average contact count in its neighborhood. This smoothing process
188 improves the contiguity of the TAD regions with elevated contact frequencies, thereby enhancing
189 the domain structures. Next, we utilize the Knight-Ruiz (KR) normalization [37] on the smoothed
190 matrices to remove potential biases.

191 **2.3 Detecting statistically significant differential regions**

192 In the third step, we identify differential genomic regions by first finding differential TADs within
193 these candidate regions. In each candidate genomic region, we calculate the SCC values for all
194 potential differential TADs as described above. Then we develop a nonparametric permutation test
195 to estimate the p -values for these local SCC values. Additionally, we propose a quantile regression
196 strategy to speed up the permutation test (see details in Supplementary Method in File S1). Finally,
197 we consider a candidate region to be a differentially interacting genomic region, if at least one TAD
198 within that region exhibits a statistically significant difference between the two samples and the size
199 of the largest differential TAD meeting this criterion is greater than one third of the length of the
200 entire candidate region. The longest differential TADs within the detected differentially interacting
201 genomic regions are defined as the noticeable differential areas.

202 Specifically, we perform the following nonparametric permutation test for each unique TAD
203 size, as the local SCC values are calculated for all potential differential TADs of various sizes.

204 Suppose s is a potential differential TAD whose length is l_s and SCC value between two Hi-C
205 samples is ρ_s . To assess the statistical significance of the observed SCC value ρ_s , the null distribution
206 of SCC values for TADs of the same size is estimated via the following permutation procedure. To
207 generate a random TAD with length l_s , we first randomly select l_s positions from main diagonal of
208 Hi-C contact matrix, then $l_s - 1$ position from the first off-diagonal, ..., and lastly 1 position from
209 the $(l_s - 1)$ -th off-diagonal. We subsequently extract contact counts of these randomly selected
210 positions from the two Hi-C contact matrices to construct the permuted TAD pair and calculate
211 its SCC value. We repeat the above random TAD generation step N times ($N = 2000$) and obtain
212 the corresponding SCC values $\{\rho_i^{l_s}\}$, $i = 1, \dots, N$. Then the p -value of the observed SCC value ρ_s
213 can be computed as:

$$p_s = \frac{\sum_{i=1}^N I(\rho_i^{l_s} < \rho_s)}{N},$$

214 where $I(\cdot)$ is the indicator function. Lastly, we compare the p -values with a pre-defined significance
215 level α (by default $\alpha = 0.05$) to determine differential TADs meeting the significance thresh-
216 old. Note that the permutation framework accounts for the multiple testing correction using the
217 Benjamini-Hochberg procedure [38].

218 One potential issue of this permutation framework is the false detection of significantly differen-

219 tial TADs when the two samples are highly similar (e.g., biological replicates from same experiment).
220 This is because the high similarity between biological replicates would lead to high SCC values of
221 the corresponding random TAD patterns. As a result, some non-differential TADs with relatively
222 low SCC values would be falsely detected as differential ones. In order to reduce the number of false
223 positives, we provide an option to filter the p-values p_s by an empirical or automatically calculated
224 threshold. This optional filtering step allows us to pre-specify the meaningful SCC between the
225 two Hi-C datasets that should be reached in order to call a differential TAD truly significant.

$$p_s^{adj} = \begin{cases} 0.5 & \text{if } p_s < \alpha \text{ and } \rho_s > \theta \\ p_s & \text{otherwise} \end{cases}$$

226 The threshold θ can normally be defined as 0.85, which corresponds to a clear margin separating
227 non-replicates from biological/pseudo-replicates in the whole-chromosome similarity comparison
228 between multiple cell lines [39]. Alternatively, θ can be calculated automatically as $\theta = \frac{\rho_{nr}^{ls} + \rho_{br}^{ls}}{2}$,
229 where ρ_{nr}^{ls} represents the mean α quantile of SCCs between non-replicate data and ρ_{br}^{ls} is the mean
230 α quantile of SCCs between their corresponding biological/pseudo-replicate data. Here, we call
231 matrices from different cell lines as non-replicates, matrices from the same cell type as biological
232 replicates, and matrices sampled from pooled biological replicates as pseudo-replicates.

233 2.4 Simulation settings

234 To evaluate the performance of the DiffGR method, we conducted a series of simulation experi-
235 ments by varying the proportion of altered TADs, proportion of TAD alternation, noise level, and
236 sequencing coverage level. Specifically, we utilized the published chromosome 1 contact matrix of
237 K562 cells at 50-kb resolution [9] as the original Hi-C data and simulated the altered Hi-C contact
238 matrices as described below.

239 2.4.1 Single-TAD alternation

240 Since TADs are conserved genomic patterns and TAD boundaries are relatively stable across cell
241 types and even across species [11], our simulations primarily focused on the scenarios of single-TAD
242 alternations. Suppose we had an original Hi-C contact matrix M and its identified TAD boundaries.
243 Each of our simulated Hi-C matrices contained two components: the signal matrix S and the noise
244 matrix N , with a certain signal-to-noise ratio.

245 First, to construct the signal matrix S , we randomly selected a subset of TADs from original
246 contact matrix to serve as the true differential TADs. Then we replaced a certain portion of
247 contact counts in each selected TAD by randomly sampling contact counts from the corresponding
248 diagonals of the contact matrix. Second, we simulated the noise matrix N which represents the
249 random ligation events in Hi-C experiments. Briefly, we generated these contacts by randomly
250 choosing two bins, i and j , and adding one to the entry N_{ij} in the noise matrix. The probability
251 of sampling each bin in the bin pair was set proportional to the marginal count of that bin in the
252 original matrix. The sampling process was repeated C times, where C was the total number of
253 contacts in the original Hi-C contact matrix M . The resulting random ligation noise matrix N
254 contained the same number of contacts as the original contact matrix M .

255 To summarize, we had the following parameters in our single-TAD simulations.

- 256 • proportion of altered TADs. Using HiCseg, we detected 189 TADs with a mean size of 1.2 Mb
257 in the original K562 chromosome 1 contact matrix (Supplementary Figure S1). By default,
258 we set the proportion of altered TADs to be 50%, which can vary from 20% to 70%.
- 259 • proportion of TAD alternation. In the default setting, we substituted all contact counts in
260 the selected TADs by random counts permuted from the matching diagonals in Hi-C maps.
261 To reduce the degree of intra-TAD alternation, we gradually decreased the proportion of
262 randomly substituted intra-TAD contacts from 100% to 10%.
- 263 • noise level, i.e., the ratio between the noise and signal matrices. The noise level was set to
264 10% by default, and varied from 1% to 80%.

265 For each simulation parameter setting, we generated 100 altered Hi-C contact matrices to com-
266 pare against the original contact matrix. To evaluate the accuracy of the detection results, we used
267 the false detection rate which defines as inaccurate percentage and is computed as $1 - Accuracy =$
268 $\frac{FP+FN}{N}$, where FP denotes the falsely detected differential regions, FN represents the the falsely
269 detected non-differential regions, and N is the total number of candidate regions being tested.

270 **2.4.2 Hierarchical-TAD alternation**

271 In addition to single-TAD alternation, we also simulated the alternation pattern of hierarchical
272 TADs. We randomly selected 50% of the large TADs whose size was greater than 10 bins in the
273 signal matrix to serve as the true differential TADs. For each of the selected large TAD, we chose
274 a random subTAD boundary to split it into two smaller subTADs (each with size > 5 bins). We
275 then replaced all inter-subTAD contact counts by randomly sampled counts in Hi-C maps. Next,
276 we validated the performance of DiffGR under the hierarchical-TAD condition with respect to
277 different noise levels similar to the single-TAD simulations. Because the complex-TAD condition
278 has complicated TAD boundaries between two samples and occurs less frequently in real data, we
279 did not generate simulation data for this condition.

280 **2.4.3 Simulating low-coverage contact matrices**

281 Low sequencing depth of Hi-C experiments would lead to low-coverage and sparse contact matrices,
282 thus it could potentially affect the performance of the detection of differentially interacting regions.
283 To simulate low-coverage contact matrices, we started with a deep-sequenced Hi-C contact map
284 obtained from human GM12878 cells [9], and down-sampled the contact counts to generate lower-
285 coverage matrices. Specifically, for each non-zero contact count M_{ij} in the original matrix, we
286 assumed that the simulated contact count follows a binomial distribution $M'_{ij} \sim \text{Binomial}(M_{ij}, p)$,
287 where the binomial parameter $p = \{0.2, 0.4, 0.6, 0.8, 1.0\}$ represents the relative coverage level of
288 the down-sampled contact matrix M' . In addition, 10% noise were added to the down-sampled
289 matrices.

290 **2.5 Real data preprocessing steps**

291 In our real data analysis, we used two published Hi-C datasets by Rao et al. [9] (GEO accession
292 GSE63525) and Dixon et al. [11](GEO accession GSE35156). The Rao et al. [9] dataset include five
293 human cell types: B-lymphoblastoid cells (GM12878), mammary epithelial cells (HMEC), umbilical

294 vein endothelial cells (HUVEC), erythrocytic leukemia cells (K562), and epidermal keratinocytes
295 (NHEK). The GM12878 dataset contains two replicates, which were also pooled together in cell
296 type-specific comparison. The Dixon et al. [11] dataset are from mouse embryonic stem (ES)
297 and cortex cells. Two replicates from mouse ES cells were merged together in cell type-specific
298 comparison. We applied DiffGR to detect differential genomic regions between each pair of cell
299 types at 25-kb, 50-kb, and 100-kb resolutions. Since some of these Hi-C datasets were not deeply
300 sequenced, the local variations introduced by low sequencing coverage made it challenging to capture
301 large domain structures, especially in fine-resolution analyses. Therefore, to enhance the domain
302 structures, all contact matrices were first preprocessed by a 2D mean filter smoothing and then
303 normalized by the KR method to eliminate potential biases.

304 In addition to Hi-C contact maps, ChIP-seq and RNA-seq data from the same cell lines were also
305 included in real data analyses. For ChIP-seq analysis, CTCF and histone modification (H3K4me1,
306 H3K4me2, H3K27me3, and H3K36me3) datasets from five human cell lines in Rao et al. [9], and
307 CTCF, Polr2a, and histone modification (H3K4me1, H3K4me3, and H3K27ac) datasets from mouse
308 cell lines in Dixon et al. [11] were obtained from the ENCODE project [40, 41] ([https://www.
309 encodeproject.org/](https://www.encodeproject.org/)). The ChIP-seq peak files were in narrowpeak/broadpeak BED format. The
310 ChIP-seq peaks were aggregated into fixed-size bins with the same resolution as the Hi-C data, and
311 the bin-wise peak counts were normalized by the total number of peaks in each ChIP-seq dataset.
312 The absolute mean differences of the normalized bin-wise peak counts were calculated for each pair
313 of cell lines for the subsequent analyses. In addition, RNA-seq datasets were also obtained from
314 the ENCODE project [41] for human GM12878 and K562 cells (GEO accession GSE78552 and
315 GSE78625) in read count format, and for mouse ES and cortex cells (GEO accession GSM723776
316 and GSM723769) in FPKM format.

317 **3 Results**

318 **3.1 DiffGR accurately detected single-TAD differences in simulated datasets**

319 To validate the accuracy and efficiency of our DiffGR method, we first generated pairs of original
320 and simulated Hi-C contact matrices, where a given proportion of TADs in the simulated contact
321 matrices were altered (see Methods). We used the intra-chromosomal contact matrix of chromosome
322 1 in K562 cells at 50-kb resolution to serve as the original contact matrix. At the default setting,
323 we altered 50% of the original TADs by completely replacing the intra-TAD contact counts by
324 randomly sampled counts outside the TAD regions. In addition, we added 10% random-ligation
325 noise into the altered contact matrices.

326 We first simulated Hi-C matrices with various proportions of altered TADs (20%, 30%, 40%,
327 50%, 60%, and 70%). With each proportion setting, we completely mutated the intra-TAD counts
328 and added 10% noise, and repeated this simulation procedure 100 times. As expected, the perfor-
329 mance of the DiffGR method depended on the proportion of altered TADs. As shown in Figure 2A
330 and Supplementary Table S1, when the proportion of altered TADs changed from 20% to 70%,
331 the false detection rate increased from 0.01 to 0.21. One possible explanation of this observed
332 trend is that when the majority of TADs were altered, the large differences between the original
333 and altered matrices would affect the permutation test and therefore lead to inaccurate detection.
334 However, differential TADs rarely exist in large proportion in real data. The false detection rates
335 of our method remained below 0.07 when the proportion of altered TADs was smaller than or
336 equal to 50%, which demonstrated that our method can accurately and reliably detect single-TAD

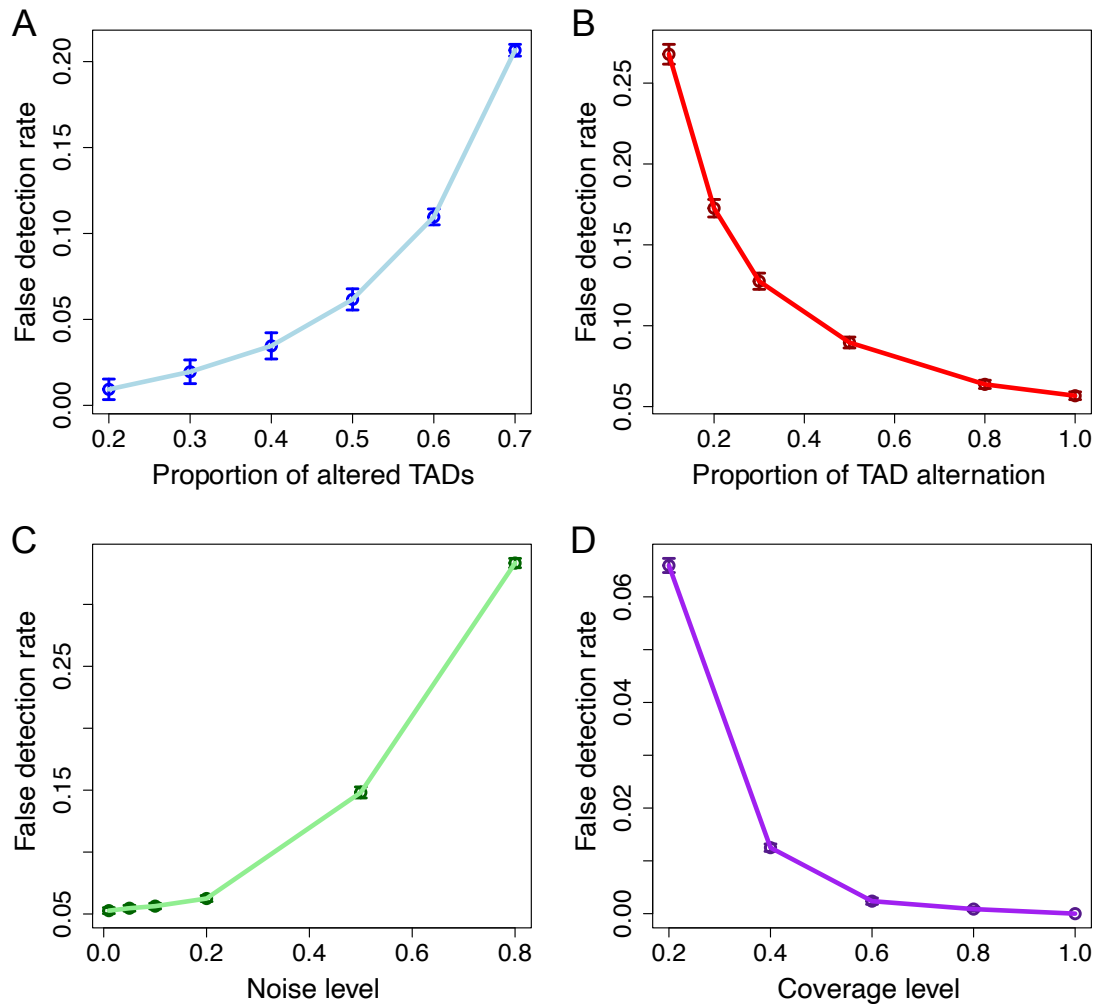


Figure 2: **Performance of single-TAD simulations.**

The curves display the mean false detection rates at different levels of A. proportion of altered TADs, B. proportion of TAD alternation, C. noise, and D. sequencing coverage. Vertical bars represent 95% confidence intervals.

337 differences under these conditions.

338 In the default simulation setting, we completely altered the selected TADs by substituting all
339 intra-TAD contact counts by randomly sampled counts from the matching diagonals outside the
340 TADs. To investigate the influence of the degree of TAD alternation on the DiffGR performance,
341 we generated a series of simulated contact matrices, in which half of original TADs were altered
342 and the proportion of intra-TAD alternation varied from 10%, to 20%, 30%, 50%, 80%, and 100%.
343 In theory, TADs with higher degrees of alternation are easier to identify, whereas TADs with minor
344 changes remain difficult to be detected. As illustrated in Figure 2B and Supplementary Table S2,
345 the performance of DiffGR improved resulting in higher accuracy as the percentage of randomly
346 substituted counts in altered TADs increased. Even with the most challenging case where only
347 10% of the intra-TAD counts were altered, the accuracy of our method was 0.73, suggesting that
348 DiffGR can effectively detect subtle TAD differences.

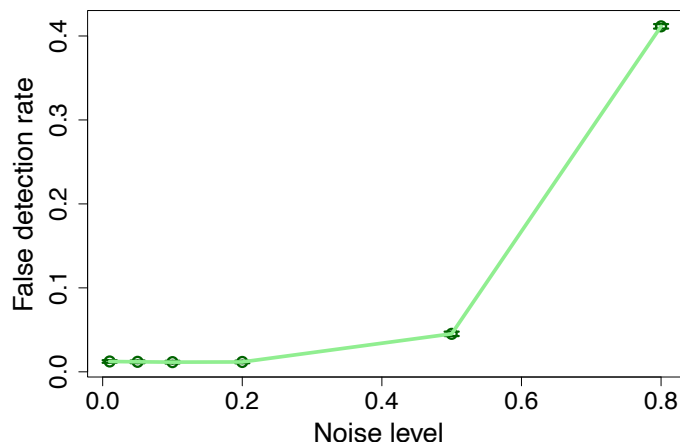


Figure 3: **Performance of hierarchical-TAD simulations.** The curve shows the mean false detection rates at various noise levels. Vertical bars represent 95% confidence intervals.

3.2 DiffGR performed stably against changes in noise and coverage levels

Next we sought to evaluate the robustness of our method under various noise levels and sequencing coverage conditions.

In the earlier simulations, we added 10% noise to the simulated differential contact matrices. To evaluate the performance of our method under different noise levels, we fixed the proportion of altered TADs at 50% and the proportion of intra-TAD alternation at 100%, and simulated the differential contact matrices with a wide range of noise levels (1%, 5%, 10%, 20%, 50%, and 80%). Intuitively, a good detection method should easily discover the differential regions in the less noisy matrices, and it becomes more challenging to detect the differential regions in the noisier cases. Our results demonstrated that DiffGR was able to correctly rank the simulated datasets. We observed a monotonic increasing trend of the false detection rate and a decreasing tendency of other precision measures as the noise levels raised (Figure 2C and Supplementary Table S3). With moderate noise levels that were not greater than 20%, the accuracy of DiffGR remained above 0.93, indicating that our method can correctly detect differential TAD regions in such noisy cases.

The sequencing coverage of the Hi-C contact maps is another major factor that could affect the performance of our method. Considering two Hi-C replicates that have the same underlying TAD structures but different sequencing coverage levels, we questioned whether our DiffGR method can correctly categorize them as non-differential. In other words, we intended to estimate the false positive rates caused by low-coverage and sparse Hi-C data. To directly investigate the influence of the sequencing coverage on the detection of differential regions, we utilized the GM12878 chromosome 1 contact matrix as the original matrix, and generated a series of down-sampled contact matrices with lower coverage levels (20%, 40%, 60%, 80%, and 100%). Figure 2D and Supplementary Table S4 show that the average false detection rates remained below 0.05 for most coverage levels, except for the lowest coverage level of 20%, demonstrating the robustness of our DiffGR method under low-coverage conditions.

3.3 DiffGR successfully detected hierarchical-TAD changes

In addition to single-TAD differences, hierarchical-TAD changes also exist in some genomic regions between different cell types. In these regions, one of the Hi-C contact maps exhibits a single dominant TAD structure, while the other Hi-C contact map presents two or more subTADs separated by additional boundaries in between. Hierarchical TADs are computationally challenging to detect. Although the two Hi-C maps have different TAD boundaries, the chromatin interaction patterns within the subTADs could be very similar. Consequently, the correlation coefficients (CCs) for the strata with small genomic distances might still remain high between two contact maps. In addition, as the genomic distance increases, the weight of the corresponding stratum in the SCC calculation gradually declines. As a result, the SCC values are primarily contributed by CC values from strata with smaller genomic distances, which makes it difficult to detect differential regions in the hierarchical-TAD cases.

To evaluate the performance of DiffGR in this more challenging situation, we simulated contact matrices containing hierarchical-TAD structures with respect to varying noise levels (see Methods) and then computed the false detection rate in a similar manner as in the single-TAD simulations. As demonstrated in Figure 3 and Supplementary Table S5, the trend of the false detection rates and other measure statistics across various noise levels under the hierarchical-TAD setting was similar to the pattern observed in the single-TAD case (Figure 2C and Supplementary Table S3). Furthermore, the false detection rates remained lower than 0.05 when the noise level was within 50%. Taken together, these results indicated that DiffGR can reliably detect the differentially interacting genomic regions with hierarchical-TAD patterns.

3.4 DiffGR revealed cell type-specific genomic interacting regions

Besides validating our method on simulated datasets, we further applied DiffGR to detect cell type-specific differences in five human cell types (GM12878, HMEC, HUVEC, K562, and NHEK) [9] and in two mouse cell types (ES and cortex cells) [11]. In total, we conducted two comparisons between biological replicates in human GM12878 and mouse ES cells, and eleven pairwise comparisons between different cell types (ten pairs among five human cell types and one pair between two mouse cell types). In each pairwise comparison, we first applied HiCseg to identify TAD boundaries from the 50-kb contact matrix for each data and then partitioned the genome into three types of candidate regions: single-TAD candidate regions, hierarchical-TAD candidate regions, and complex-TAD candidate regions. Statistically significant differential genomic regions were identified between each comparison with FDR cutoff 0.05.

We first sought to evaluate the performance of our method on biological replicates of Hi-C data. Previous studies have shown that the high degree of similarity between biological replicates and dominant consistence between TAD boundaries in replicate data [9, 11, 39]. For the comparison between human GM12878 replicates, consistent with our expectations, the majority (89.55%) of the 2325 candidate genomic across the genome regions belonged to single-TAD type and very few (2.45%) candidate genomic regions were detected as differential by our method (Supplementary Figure S2). Specifically, only 1.97% of single-TADs were identified as differential, whereas 6.17% and 4.94% were detected in hierarchical-TAD and complex-TAD cases respectively. Similar results were also witnessed in the comparison between replicates in mouse ES cells: 83.42% candidate genomic regions were classified as single-TAD type and few (6.02%) were identified as differential (Supplementary Table S6). Overall, our DiffGR results confirmed that these biological replicates

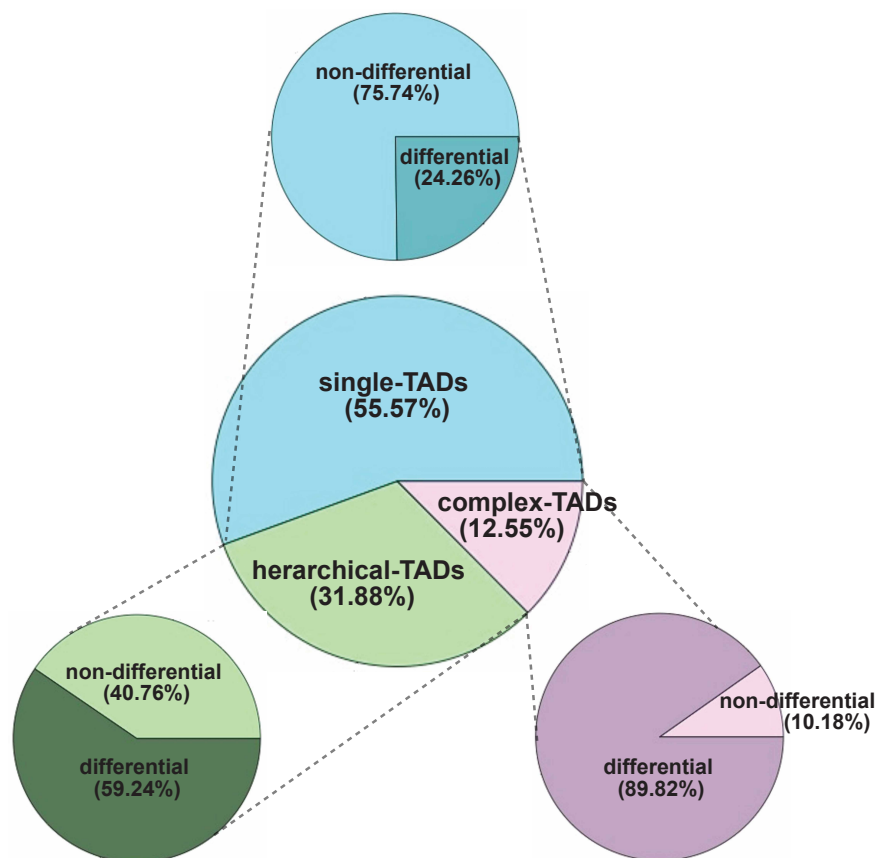


Figure 4: **Piecharts of DiffGR results obtained from human Hi-C datasets.**

The center piechart presents the proportions of three categories of candidate regions. The three outer piecharts display the proportions of DiffGR-detected differential genomic regions, one for each candidate category.

417 displayed highly consistent chromatin structures with minor biological variations.

418 Next, we applied DiffGR to detect cell-type-specific differences and the results are illustrated in
419 Figure 4 and Supplementary Table S7. For the ten pairwise comparisons among human cell types,
420 55.57% of all candidate genomic regions belonged to the single-TAD category (consistent with pre-
421 vious observations indicating that TAD boundaries are stable across cell types [11]), 31.88% to
422 the hierarchical-TAD category, and 12.55% to the complex-TAD category. Our DiffGR analyses
423 showed that only 24.26% of the single-TAD candidate regions showed statistically significant dif-
424 ferences between two samples; 59.24% of the hierarchical-TAD candidate regions were determined
425 to be differential; while the differential proportion of the complex-TAD category was as high as
426 89.82%. In addition, we found that the proportion of detected differential regions varied largely
427 across chromosomes, ranging from 0.14 to 0.76 (Supplementary Figure S3). For the comparison
428 between mouse ES and cortex cells, 20.22% of the candidate genomic regions in the single-TAD cat-
429 egory were identified as differential, while the proportion increased to 75.94% in the complex-TAD
430 category. These observations indicated that candidate genomic regions with more distinct patterns
431 of TAD boundaries are more likely to be detected as differential between two Hi-C samples.

432 In addition to partitioning the genome at 50-kb resolution, we also performed differential analy-
433 ses on the five human Hi-C datasets at 25-kb and 100-kb resolutions, separately. We calculated the

434 overlapping rate (that is, the proportion of the genome that was classified into the same differential
435 or non-differential status) between different resolutions. Overall, we observed a high consistency
436 between the detected differential regions across different resolutions, where the overlapping rate was
437 0.9856 between the detection results at 50-kb and 100-kb resolutions, and 0.9480 between those at
438 25-kb and 50-kb resolutions. These results demonstrated that DiffGR can robustly and consistently
439 detect cell type-specific differential genomic regions across various resolutions.

440 3.5 Changes in CTCF and histone modification patterns were consistent with 441 DiffGR detection results

Table 1: Agreements between ChIP-seq data and DiffGR-detected differential genomic regions in human Hi-C datasets.

	100 kb	50 kb	25 kb
CTCF	76 (34.55%)	124 (56.36%)	142 (64.55%)
H3K4me1	57 (25.91%)	110 (50.00%)	136 (61.82%)
H3K4me2	56 (25.45%)	91 (41.36%)	116 (52.73%)
H3K27me3	53 (24.09%)	86 (39.09%)	114 (51.82%)
H3K36me3	36 (16.36%)	72 (32.73%)	110 (50.00%)

Note: A total of 220 *t*-tests (10 pairwise comparisons between five human cell types, 22 chromosome-wide tests for each comparison) were conducted. If the mean absolute differences of a ChIP-seq signal at the TAD boundaries in the differential regions were significantly higher than those in non-differential regions, the results were labeled significant consistent. The counts and percentages of significant consistent results were reported for each ChIP-seq dataset at each resolution.

442 As there is no ground truth of differential chromatin interacting regions in real data, we sought to
443 evaluate the performance of our method by investigating the association between the changes in
444 1D epigenomic features and 3D genomic interaction regions. The chromatin architectural protein
445 CTCF plays an essential role in establishing higher-order chromatin structures such as TADs. In
446 addition, it has been shown that transcription factors and histone marks are enriched or depleted at
447 TAD boundaries. Therefore, we hypothesized that differential bindings of transcription factors such
448 as CTCF and histone modifications would also be present at the TAD boundaries in differential
449 genomic interacting regions.

450 To test this hypothesis, we first combined TAD boundaries from both Hi-C datasets and clas-
451 sified them into two categories: those within the DiffGR-detected differential regions and those
452 outside the differential regions. We then utilized the ChIP-seq datasets of transcription factors like
453 CTCF and histone modifications from the ENCODE project [40]. For each ChIP-seq dataset, we
454 calculated the mean absolute difference of ChIP-seq peaks between the two cell types within the
455 neighborhood (± 1 bin) of each TAD boundary. We expected that if two cell lines have highly
456 different chromatin structures in certain genomic regions, different patterns of CTCF bindings and
457 histone modifications in these regions would be observed. Therefore, we performed the following
458 *t*-test for each ChIP-seq dataset using the DiffGR detection results. In each chromosome, we eval-
459 uated whether the mean absolute differences of the ChIP-seq signal at the TAD boundaries in
460 differential regions were significantly different from those in non-differential regions. If the ChIP-
461 seq signal differences at the TAD boundaries in differential regions were significantly higher (with
462 a significant level 0.1) than those in non-differential regions, we considered the ChIP-seq changes
463 to be consistent with our DiffGR differential detection results.

464 Table 1 and Supplementary Table S8 summarize the ChIP-seq analyses on the DiffGR detection
465 results obtained from five human Hi-C datasets [9] and two mouse Hi-C datasets [11]. For each
466 human ChIP-seq dataset, we performed 220 *t*-tests (ten pairwise comparisons between cell types,
467 22 chromosome-wide tests one for each autosome) at 100-kb, 50-kb, and 25-kb resolutions; for
468 each mouse ChIP-seq dataset, we conducted 19 *t*-tests one for each autosome at 50-kb resolution.
469 Overall, DiffGR-detected differential genomic regions were supported by 1D epigenomic features in
470 both human and mouse data. Furthermore, we observed that the agreement between the changes
471 in ChIP-seq signal and chromatin structures was improved in finer-resolution analyses. As shown in
472 Table 1, 76 out of 220 (34.55%) tests showed significantly higher absolute mean differences of CTCF
473 values at the TAD boundaries in DiffGR-detected differential genomic regions than those in non-
474 differential regions at 100-kb resolution. Whereas in the results at 25-kb resolution, 142 (64.55%)
475 tests exhibited significantly larger changes in CTCF bindings in differential regions than non-
476 differential ones. In addition, the histone modification datasets showed similar results in agreement
477 with the detection results of differentially interacting regions in Hi-C contact maps. At 25-kb
478 resolution, the majority of the *t*-tests showed significantly larger changes of ChIP-seq signal in
479 differentially interacting regions for all four histone modification datasets, including H3K4me1,
480 H3K4me2, H3K27me3, and H3K26me3. Collectively, these results indicated that the changes in
481 CTCF bindings and histone modifications were in good agreements with the differences in genomic
482 interacting regions. Furthermore, at finer resolution our DiffGR method produced more accurate
483 identification of differentially interacting genomic regions in higher agreement with the CTCF and
484 histone modification data.

485 We would like to point out that for those cases where the changes in CTCF or histone mod-
486 ifications are not in significant agreement with the detection results of differentially interacting
487 genomic regions, it does not necessarily suggest that these epigenomic features are inconsistent
488 with 3D genome organization nor DiffGR detection results are inaccurate. Due to the resolution
489 limit of Hi-C contact maps, the boundaries of differential regions are usually identified with a res-
490 olution of tens to hundreds of kilobases. Aggregating ChIP-seq data with such a large bin size
491 dilutes the signal, thereby yielding less statistical power to detect significant changes. Moreover,
492 CTCF and histone modifications play fundamental roles in regulating chromatin structures and
493 gene expression; their functions are not limited to TAD formations. Therefore, changes in CTCF
494 bindings or histone modifications exist in many genomic loci other than TAD boundaries, thus may
495 not be represented in our analyses.

496 **3.6 Differential RNA-seq analysis results were consistent with DiffGR detection**

497 In addition to investigating the changes in 1D epigenomic features, we further studied the relation-
498 ship between quantitative changes in gene expression levels and 3D genomic interaction regions.
499 Previous studies have showed that topological changes of 3D genome organization have a large ef-
500 fect on the cross-talk between enhancers and promoters therefore can alter gene expression [9, 20].
501 Thus, we expected to observe an enrichment of differential expressed genes in DiffGR-detected
502 differential genomic regions.

503 To evaluate this assumption, we first detected significant changes in gene expression levels be-
504 tween human GM12878 and K562 cells using DESeq2 [23] and those between mouse ES and cortex
505 cells using ballgown [42]. Then we calculated the percentage of differentially expressed genes that
506 were located inside the DiffGR-identified differential genomic regions. To calculate the enrichment
507 of differentially expressed genes, we randomly chose a set of genes, whose number is equivalent

Table 2: **Functional enrichment of differentially expressed genes located in differential genomic regions between GM12878 and K562.**

GO Term	<i>P</i> -value
GO:0002376 immune system process	1.7E-9
GO:0050776 regulation of immune response	5.9E-8
GO:0002757 immune response-activating signal transduction	7.8E-8
GO:0002682 regulation of response to stress	2.2E-7
GO:0080134 regulation of immune system process	2.7E-7
GO:0045321 leukocyte activation	2.8E-7

Note: Top 2000 differentially expressed genes located within differential genomic regions at 25-kb resolution were utilized in GO enrichment analysis.

508 to the number of the DESeq2-detected differentially expressed genes, with 200 times, computed
509 their corresponding proportions located in differential genomic regions, and then performed *t*-test
510 for comparison. In summary, a total number of 8781 differentially expressed genes were detected
511 between human GM12878 and K562 cells and 79.54% of them were located in DiffGR-detected
512 differential genomic regions (p -value = 3.72×10^{-5} , permutation test); whereas 2124 differen-
513 tially expressed genes were identified between mouse ES and cortex cells and 61.66% were within
514 DiffGR-detected differential genomic regions (p -value < 2.2×10^{-16}). Taken together, these results
515 demonstrated that the changes of gene expression in RNA-seq data were highly consistent with the
516 DiffGR detection results.

517 To further explore the potential functional roles of the differentially expressed genes located in
518 differential genomic regions, we performed Gene Ontology (GO) enrichment analysis on the top
519 2000 genes using DAVID [43]. As show in Table 2, we observed a high enrichment of GO terms
520 related to the immune responses, which is consistent with the immunological nature of GM12878
521 lymphoblastoid B-cells.

522 **3.7 DiffGR detection was supported by differential chromatin interactions**

523 Several Hi-C comparative studies have demonstrated that the majority of the chromatin struc-
524 tural changes tend to couple with the formation/disappearance of topologically associated domains
525 (TADs) [9, 20], implying that changes in Hi-C interaction counts are likely to be observed within
526 genomic regions at TAD level. Hence, we checked differential chromatin interactions (DCIs) be-
527 tween GM12878 and K562 cells at 50-kb resolution by FIND [26] and compared FIND results with
528 our DiffGR results. As shown in Figure 5, the percentages of DCIs detected by FIND located
529 within candidate genomic regions were dominant in the majority of chromosomes and with 55.43%
530 across the whole genome. In addition, 82.80% of the DCIs located in candidate genomic regions
531 are classified into differential regions, demonstrating that DiffGR effectively detected the regions
532 with significant changes in chromatin contacts.

533 **3.8 Performance and comparison with state-of-the-art differential TAD detec- 534 tion tools**

535 Next, we compared the DiffGR results with three differential TAD boundaries detection methods
536 (HiCDB [13], TADCompare [29], and TADreg [30]) and one differential TAD regions detection

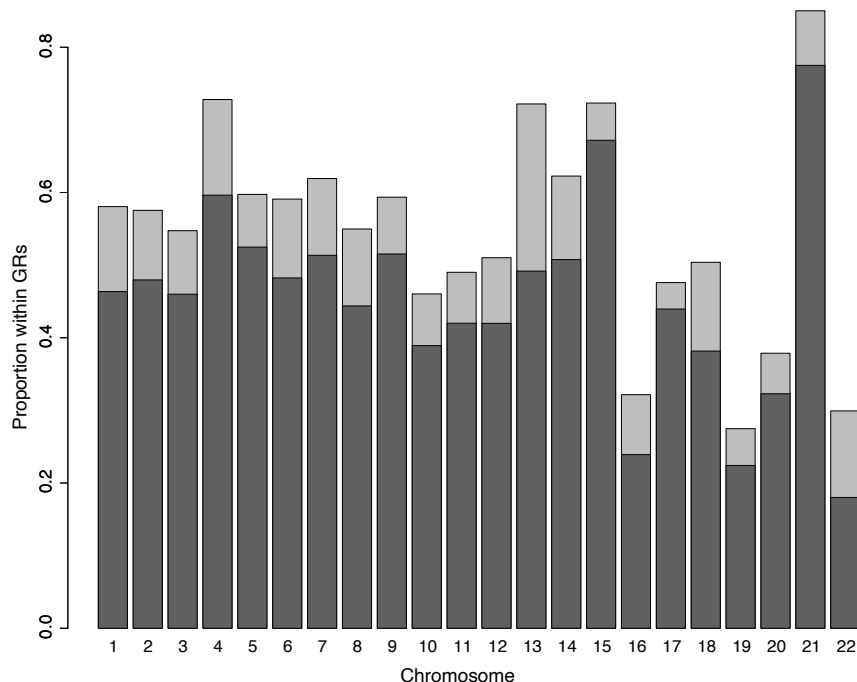


Figure 5: **Comparison between FIND and DiffGR.** Bar chart of the proportions of FIND-detected DCIs located in candidate genomic regions (GRs) and differential GRs for all autosomes between GM12878 and K562. The light gray bars denote the proportions of DCIs located in candidate GRs; the dark gray bars represent the proportions of DCIs located in differential GRs.

Table 3: **Differential TAD boundaries detected by TADcompare in DiffGR-detected differential genomic regions showed higher agreement with ChIP-seq signals than those in non-differential regions in human Hi-C datasets.**

	consistent	significant consistent
CTCF	155 (70.45%)	98 (44.55%)
H3K4me1	145 (65.91%)	89 (40.45%)
H3K4me2	133 (60.45%)	79 (45.91%)
H3K27me3	146 (66.36%)	76 (34.55%)
H3K36me3	127 (57.73%)	51 (23.18%)

Note: A total of 220 tests (10 pairwise comparisons between five human cell types, 22 chromosome-wide tests for each comparison) were conducted. If the mean absolute differences of a ChIP-seq signal at the TADcompare-identified differential TAD boundaries in the differential genomic regions were higher (or significantly higher based on *t*-test) than those in non-differential regions, the results were labeled consistent (or significantly consistent). The counts and percentages of consistent and significant consistent results were reported for each ChIP-seq dataset.

537 tool (provided by HiCExplorer [31–33]) on the five human Hi-C datasets by Rao et al. [9] and the
 538 two mouse datasets by Dixon et al. [11] at 50-kb resolution. Overall, the differential TAD bound-
 539 aries/regions identified by HiCDB, TADCompare, TADreg, or HiCExplorer were highly concordant
 540 with DiffGR-detected differentially interacting genomic regions. Notably, 73.86% of the HiCDB-
 541 detected, 76.25% of the TADCompare-detected, and 71.90% of the TADreg-detected differential
 542 TAD boundaries displayed consistent results with our DiffGR detection in the human datasets. In
 543 addition, highly concordant rates were also witnessed in the mouse dataset with 59.56%, 62.01%,

544 and 60.32% consistency rate with HiCDB, TADCompare, and TADreg, respectively. Furthermore,
545 60.62% of the 2877 HiCExplorer-identified differential regions from the five human cell lines over-
546 lapped with DiffGR-detected differential regions.

547 To investigate the advantages of DiffGR over TADCompare, we further performed tests on
548 changes in CTCF and histone modification patterns for the TADCompare-detected differential
549 TAD boundaries within DiffGR-detected differential and non-differential genomic regions in hu-
550 man datasets. From Table 3, we observed that 155 out of 220 (70.45%) contrasts showed higher
551 absolute mean differences of CTCF values at TADCompare-detected differential TAD boundaries
552 in DiffGR-detected differential genomic regions than those in non-differential regions. In addition
553 98 (44.55%) CTCF tests exhibited significantly larger changes of CTCF bindings with a significant
554 level of 0.1 at differential TAD boundaries in differential regions than those in non-differential re-
555 gions. Furthermore, the histone modification datasets (including H3K4me1, H3K4me2, H3K27me3,
556 and H3K36me3) showed similar results that were in agreement with the advantageous results of dif-
557 ferential TAD boundaries in differentially interacting regions. Collectively, these results indicated
558 that DiffGR-detected differential genomic regions had a better agreement with 1D epigenomic fea-
559 tures than TADCompare-detected differential TAD bounds.

560 4 Discussion and Conclusions

561 With the fast accumulation of Hi-C datasets, there has been a dramatically increasing interest in
562 comparative analysis of Hi-C contact maps. However, most existing methods for comparative Hi-C
563 analysis focused on the identification of differential chromatin interactions, while few studies ad-
564 dressed the detection of differential chromatin organization at TAD scale. To tackle this problem,
565 we developed a novel method, DiffGR, for calling differentially interacting genomic regions between
566 two Hi-C contact maps. Taking genomic distance features of Hi-C data into consideration, our algo-
567 rithm utilized the SCC metric instead of the standard Pearson CC to measure the similarity of local
568 genomic regions between Hi-C contact maps. Furthermore, we proposed a nonparametric permuta-
569 tion test to assess the statistical significance of the local SCC values. In contrast to the parametric
570 approaches that were used by most Hi-C data analysis methods, our nonparametric approach does
571 not have a set of predefined assumptions about the nature of the null distribution and, therefore,
572 is more robust and can be applied to more diverse data from real cases. Additionally, we utilized a
573 non-parametric smoothing spline regression to speed up the permutation test and showed that the
574 speed-up algorithm can steadily produce consistent outputs. Through empirical evaluations, we
575 have demonstrated that DiffGR can effectively discover differential regions in both simulated data
576 and real Hi-C data from different cell types. That is, DiffGR produced robust and stable detection
577 results under various noise and coverage levels in simulated data; DiffGR detection results in real
578 data were effectively validated by the ChIP-seq and RNA-seq data; DiffGR produced consistent
579 and advantageous results compared with state-of-the-art differential TAD boundaries/regions de-
580 tection tools. To summarize, DiffGR provides a statistically rigorous method for the detection of
581 differentially interacting genomic regions in Hi-C contact maps from different cells and conditions,
582 therefore would facilitate the investigation of their biological functions.

583 We envision a few possible extensions and future directions based on this work. First, our
584 method performs pairwise comparison between Hi-C contact maps. One potential future direction is
585 to design a more general statistical framework for differential analyses among three or more samples.
586 Then we could further assign the differentially interacting genomic regions to cell type-specific or
587 condition-specific changing areas. Second, we currently pool biological replicates together in our

588 analyses. Extending DiffGR to incorporate multiple biological replicates to detect reproducible
589 differences would enhance the reliability of the detection results. Third, in our algorithm, we use
590 the shared TAD boundaries between two samples to segment the genome into candidate genomic
591 regions and then detect differential regions. Recently, the notion of TADs being highly conserved
592 across cell types has been questioned [44, 45]. Therefore, a more general approach to define and
593 classify the candidate genomic regions would be beneficial to better characterize the variability of
594 chromatin interactions between different conditions. Lastly, our method is specifically designed
595 for bulk Hi-C data. Given the high sparsity and variability of single-cell Hi-C contact matrices,
596 identifying differential genomic regions at single-cell level remains a significant challenge.

597 **5 Code Availability**

598 The DiffGR R Code (both algorithm and simulation) is publicly available at [https://github.com/](https://github.com/wmalab/DiffGR)
599 [wmalab/DiffGR](https://github.com/wmalab/DiffGR) under the GNU GPL ≥ 2 license. The source code is also available at BioCode
600 <https://ngdc.cncb.ac.cn/biocode/tools/BT007313>.

601 **6 CRediT author statement**

602 **Huiling Liu:** Conceptualization, Methodology, Software, Formal analysis, Writing - Original Draft.

603 **Wenxiu Ma:** Conceptualization, Supervision, Writing - Review & Editing, Funding acquisition.

604 All authors read and approved the final manuscript.

605 **7 Competing interests**

606 The authors have declared that no competing interests exist.

607 **8 Acknowledgments**

608 The authors would like to thank Tiantian Ye, Yangyang Hu, and Luke Klein for helpful discussions
609 and feedback, and the editor and reviewers for their insightful comments and suggestions. This
610 work was supported by the National Science Foundation [DBI-1751317]; and the National Institute
611 of Health [R35GM133678].

612 **9 ORCID**

613 0000-0002-4671-1006 (Huiling Liu)

614 0000-0003-4097-1621 (Wenxiu Ma)

615 References

- 616 [1] Marieke Simonis, Petra Klous, Erik Splinter, Yuri Moshkin, Rob Willemsen, Elzo De Wit, Bas
617 Van Steensel, and Wouter De Laat. Nuclear organization of active and inactive chromatin
618 domains uncovered by chromosome conformation capture–on-chip (4c). *Nature genetics*, 38
619 (11):1348–1354, 2006.
- 620 [2] Josée Dostie, Todd A Richmond, Ramy A Arnaout, Rebecca R Selzer, William L Lee, Tracey A
621 Honan, Eric D Rubio, Anton Krumm, Justin Lamb, Chad Nusbaum, et al. Chromosome
622 conformation capture carbon copy (5c): a massively parallel solution for mapping interactions
623 between genomic elements. *Genome research*, 16(10):1299–1309, 2006.
- 624 [3] Erez Lieberman-Aiden, Nynke L Van Berkum, Louise Williams, Maxim Imakaev, Tobias
625 Ragoczy, Agnes Telling, Ido Amit, Bryan R Lajoie, Peter J Sabo, Michael O Dorschner, et al.
626 Comprehensive mapping of long-range interactions reveals folding principles of the human
627 genome. *science*, 326(5950):289–293, 2009.
- 628 [4] Z. Duan, M. Andronescu, K. Schutz, S. McIlwain, Y. J. Kim, C. Lee, J. Shendure, S. Fields,
629 C. A. Blau, and W. S. Noble. A three-dimensional model of the yeast genome. *Nature*, 465
630 (7296):363–367, 2010.
- 631 [5] R. Kalhor, H. Tjong, N. Jayathilaka, F. Alber, and L. Chen. Genome architectures revealed by
632 tethered chromosome conformation capture and population-based modeling. *Nature biotech-*
633 *nology*, 30(1):90–98, 2012.
- 634 [6] Guoliang Li, Melissa J Fullwood, Han Xu, Fabianus Hendriyan Mulawadi, Stoyan Velkov,
635 Vinsensius Vega, Pramila Nuwantha Ariyaratne, Yusoff Bin Mohamed, Hong-Sain Ooi, Chan-
636 dana Tennakoon, et al. Chia-pet tool for comprehensive chromatin interaction analysis with
637 paired-end tag sequencing. *Genome biology*, 11(2):R22, 2010.
- 638 [7] Maxwell R Mumbach, Adam J Rubin, Ryan A Flynn, Chao Dai, Paul A Khavari, William J
639 Greenleaf, and Howard Y Chang. Hichip: efficient and sensitive analysis of protein-directed
640 genome architecture. *Nature methods*, 13(11):919–922, 2016.
- 641 [8] Wenxiu Ma, Ferhat Ay, Choli Lee, Gunhan Gulsoy, Xinxian Deng, Savannah Cook, Jennifer
642 Hesson, Christopher Cavanaugh, Carol B Ware, Anton Krumm, et al. Fine-scale chromatin
643 interaction maps reveal the cis-regulatory landscape of human lincRNA genes. *Nature methods*,
644 12(1):71–78, 2015.
- 645 [9] Suhas SP Rao, Miriam H Huntley, Neva C Durand, Elena K Stamenova, Ivan D Bochkov,
646 James T Robinson, Adrian L Sanborn, Ido Machol, Arina D Omer, Eric S Lander, et al. A 3d
647 map of the human genome at kilobase resolution reveals principles of chromatin looping. *Cell*,
648 159(7):1665–1680, 2014.
- 649 [10] V. Ramani, D. A. Cusanovich, R. J. Hause, W. Ma, R. Qiu, X. Deng, C. A. Blau, C. M.
650 Disteche, W. S. Noble, J. Shendure, and Z. Duan. Mapping 3D genome architecture through
651 in situ DNase Hi-C. *Nature protocols*, 11(11):2104–2121, 2016.
- 652 [11] Jesse R Dixon, Siddarth Selvaraj, Feng Yue, Audrey Kim, Yan Li, Yin Shen, Ming Hu, Jun S
653 Liu, and Bing Ren. Topological domains in mammalian genomes identified by analysis of
654 chromatin interactions. *Nature*, 485(7398):376–380, 2012.

- 655 [12] Geoffrey Fudenberg, Maxim Imakaev, Carolyn Lu, Anton Goloborodko, Nezar Abdennur, and
656 Leonid A Mirny. Formation of chromosomal domains by loop extrusion. *Cell reports*, 15(9):
657 2038–2049, 2016.
- 658 [13] Fengling Chen, Guipeng Li, Michael Q Zhang, and Yang Chen. Hicdb: a sensitive and robust
659 method for detecting contact domain boundaries. *Nucleic acids research*, 46(21):11239–11250,
660 2018.
- 661 [14] Emily Crane, Qian Bian, Rachel Patton McCord, Bryan R Lajoie, Bayly S Wheeler, Edward J
662 Ralston, Satoru Uzawa, Job Dekker, and Barbara J Meyer. Condensin-driven remodelling of
663 x chromosome topology during dosage compensation. *Nature*, 523(7559):240, 2015.
- 664 [15] Darya Filippova, Rob Patro, Geet Duggal, and Carl Kingsford. Identification of alternative
665 topological domains in chromatin. *Algorithms for Molecular Biology*, 9(1):14, 2014.
- 666 [16] Celine Lévy-Leduc, Maud Delattre, Tristan Mary-Huard, and Stephane Robin. Two-
667 dimensional segmentation for analyzing hi-c data. *Bioinformatics*, 30(17):i386–i392, 2014.
- 668 [17] Mattia Forcato, Chiara Nicoletti, Koustav Pal, Carmen Maria Livi, Francesco Ferrari, and
669 Silvio Bicciato. Comparison of computational methods for hi-c data analysis. *Nature methods*,
670 14(7):679, 2017.
- 671 [18] Zhijun Han and Gang Wei. Computational tools for hi-c data analysis. *Quantitative Biology*,
672 5(3):215–225, 2017.
- 673 [19] Junbai Wang, Xun Lan, Pei-Yin Hsu, Hang-Kai Hsu, Kun Huang, Jeffrey Parvin, Tim HM
674 Huang, and Victor X Jin. Genome-wide analysis uncovers high frequency, strong differen-
675 tial chromosomal interactions and their associated epigenetic patterns in e2-mediated gene
676 regulation. *BMC genomics*, 14(1):70, 2013.
- 677 [20] Jesse R Dixon, Inkyung Jung, Siddarth Selvaraj, Yin Shen, Jessica E Antosiewicz-Bourget,
678 Ah Young Lee, Zhen Ye, Audrey Kim, Nisha Rajagopal, Wei Xie, et al. Chromatin architecture
679 reorganization during stem cell differentiation. *Nature*, 518(7539):331, 2015.
- 680 [21] Sven Heinz, Christopher Benner, Nathanael Spann, Eric Bertolino, Yin C Lin, Peter Laslo,
681 Jason X Cheng, Cornelis Murre, Harinder Singh, and Christopher K Glass. Simple combi-
682 nations of lineage-determining transcription factors prime cis-regulatory elements required for
683 macrophage and b cell identities. *Molecular cell*, 38(4):576–589, 2010.
- 684 [22] Mark D Robinson, Davis J McCarthy, and Gordon K Smyth. edgeR: a bioconductor package for
685 differential expression analysis of digital gene expression data. *Bioinformatics*, 26(1):139–140,
686 2010.
- 687 [23] Michael I Love, Wolfgang Huber, and Simon Anders. Moderated estimation of fold change
688 and dispersion for rna-seq data with deseq2. *Genome biology*, 15(12):550, 2014.
- 689 [24] Aaron TL Lun and Gordon K Smyth. diffhic: a bioconductor package to detect differential
690 genomic interactions in hi-c data. *BMC bioinformatics*, 16(1):258, 2015.
- 691 [25] John C Stansfield, Kellen G Cresswell, Vladimir I Vladimirov, and Mikhail G Dozmorov.
692 Hiccompare: an r-package for joint normalization and comparison of hi-c datasets. *BMC*
693 *bioinformatics*, 19(1):1–10, 2018.

- 694 [26] Mohamed Nadhir Djekidel, Yang Chen, and Michael Q Zhang. Find: differential chromatin
695 interactions detection using a spatial poisson process. *Genome research*, 28(3):412–422, 2018.
- 696 [27] Kate B Cook, Borislav H Hristov, Karine G Le Roch, Jean Philippe Vert, and William Stafford
697 Noble. Measuring significant changes in chromatin conformation with accost. *Nucleic acids
698 research*, 48(5):2303–2311, 2020.
- 699 [28] Phillippa C Taberlay, Joanna Achinger-Kawecka, Aaron TL Lun, Fabian A Buske, Kenneth
700 Sabir, Cathryn M Gould, Elena Zotenko, Saul A Bert, Katherine A Giles, Denis C Bauer,
701 et al. Three-dimensional disorganization of the cancer genome occurs coincident with long-
702 range genetic and epigenetic alterations. *Genome research*, 26(6):719–731, 2016.
- 703 [29] Kellen G Cresswell and Mikhail G Dozmorov. Tadcompare: An r package for differential and
704 temporal analysis of topologically associated domains. *Frontiers in Genetics*, 11:158, 2020.
- 705 [30] Raphaël Mourad. Tadreg: a versatile regression framework for tad identification, differential
706 analysis and rearranged 3d genome prediction. *BMC bioinformatics*, 23(1):1–14, 2022.
- 707 [31] Joachim Wolff, Leily Rabbani, Ralf Gilsbach, Gautier Richard, Thomas Manke, Rolf Backofen,
708 and Björn A Grüning. Galaxy hicexplorer 3: a web server for reproducible hi-c, capture hi-c
709 and single-cell hi-c data analysis, quality control and visualization. *Nucleic Acids Research*,
710 48(W1):W177–W184, 2020.
- 711 [32] Joachim Wolff, Vivek Bhardwaj, Stephan Nothjunge, Gautier Richard, Gina Renschler, Ralf
712 Gilsbach, Thomas Manke, Rolf Backofen, Fidel Ramírez, and Björn A Grüning. Galaxy hic-
713 explorer: a web server for reproducible hi-c data analysis, quality control and visualization.
714 *Nucleic acids research*, 46(W1):W11–W16, 2018.
- 715 [33] Fidel Ramírez, Vivek Bhardwaj, Laura Arrigoni, Kin Chung Lam, Björn A Grüning, José
716 Villaveces, Bianca Habermann, Asifa Akhtar, and Thomas Manke. High-resolution tads reveal
717 dna sequences underlying genome organization in flies. *Nature communications*, 9(1):1–15,
718 2018.
- 719 [34] Marie Zufferey, Daniele Tavernari, Elisa Oricchio, and Giovanni Ciriello. Comparison of com-
720 putational methods for the identification of topologically associating domains. *Genome biology*,
721 19(1):1–18, 2018.
- 722 [35] Tao Yang, Feipeng Zhang, Galip Gürkan Yardımcı, Fan Song, Ross C Hardison,
723 William Stafford Noble, Feng Yue, and Qunhua Li. Hicrep: assessing the reproducibility of
724 hi-c data using a stratum-adjusted correlation coefficient. *Genome research*, 27(11):1939–1949,
725 2017.
- 726 [36] Alan Agresti and Maria Kateri. *Categorical data analysis*. Springer, 2011.
- 727 [37] Philip A Knight and Daniel Ruiz. A fast algorithm for matrix balancing. *IMA Journal of
728 Numerical Analysis*, 33(3):1029–1047, 2013.
- 729 [38] Yoav Benjamini and Yosef Hochberg. Controlling the false discovery rate: a practical and pow-
730 erful approach to multiple testing. *Journal of the Royal statistical society: series B (Method-
731 ological)*, 57(1):289–300, 1995.

- 732 [39] Galip Gürkan Yardımcı, Hakan Ozadam, Michael EG Sauria, Oana Ursu, Koon-Kiu Yan, Tao
733 Yang, Abhijit Chakraborty, Arya Kaul, Bryan R Lajoie, Fan Song, et al. Measuring the
734 reproducibility and quality of hi-c data. *Genome biology*, 20(1):1–19, 2019.
- 735 [40] ENCODE Project Consortium et al. The encode (encyclopedia of dna elements) project.
736 *Science*, 306(5696):636–640, 2004.
- 737 [41] ENCODE Project Consortium et al. An integrated encyclopedia of dna elements in the human
738 genome. *Nature*, 489(7414):57, 2012.
- 739 [42] Mihaela Pertea, Daehwan Kim, Geo M Pertea, Jeffrey T Leek, and Steven L Salzberg.
740 Transcript-level expression analysis of rna-seq experiments with hisat, stringtie and ballgown.
741 *Nature protocols*, 11(9):1650–1667, 2016.
- 742 [43] Brad T Sherman, Richard A Lempicki, et al. Systematic and integrative analysis of large gene
743 lists using david bioinformatics resources. *Nature protocols*, 4(1):44, 2009.
- 744 [44] Natalie Sauerwald, Akshat Singhal, and Carl Kingsford. Analysis of the structural variability
745 of topologically associated domains as revealed by hi-c. *NAR genomics and bioinformatics*, 2
746 (1):lqz008, 2020.
- 747 [45] Evonne McArthur and John A Capra. Topologically associating domain boundaries that are
748 stable across diverse cell types are evolutionarily constrained and enriched for heritability. *The*
749 *American Journal of Human Genetics*, 108(2):269–283, 2021.