

Simultaneous smoothing and detection of topological units of genome organization from sparse chromatin contact count matrices with matrix factorization

Da-Inn Lee¹ and Sushmita Roy^{1,2*}

¹Department of Biostatistics and Medical Informatics, University of Wisconsin-Madison,
Madison, WI 53715, USA

²Wisconsin Institute for Discovery, 330 N. Orchard Street, Madison, WI 53715, USA

*To whom correspondence should be addressed.

Abstract

The three-dimensional (3D) organization of the genome plays a critical role in gene regulation for diverse normal and disease processes. High-throughput chromosome conformation capture (3C) assays, such as Hi-C, SPRITE, GAM, and HiChIP, have revealed higher-order organizational units such as topologically associating domains (TADs), which can shape the regulatory landscape governing downstream phenotypes. Analysis of high-throughput 3C data depends on the sequencing depth, which directly affects the resolution and the sparsity of the generated 3D contact count map. Identification of TADs remains a significant challenge due to the sensitivity of existing methods to resolution and sparsity. Here we present GRiNCH, a novel matrix-factorization-based approach for simultaneous TAD discovery and smoothing of contact count matrices from high-throughput 3C data. GRiNCH TADs are enriched in known architectural proteins and chromatin modification signals and are stable to the resolution, and sparsity of the input data. GRiNCH smoothing improves the recovery of structure and significant interactions from low-depth datasets. Furthermore, enrichment analysis of 746 transcription factor motifs in GRiNCH TADs from developmental time-course and cell-line Hi-C datasets predicted transcription factors with potentially novel genome organization roles. GRiNCH is a broadly applicable tool for the analysis of high throughput 3C datasets from a variety of platforms including SPRITE and HiChIP to understand 3D genome organization in diverse biological contexts.

1 Introduction

2 The three-dimensional (3D) organization of the genome has emerged as an important layer of gene
3 regulation in developmental processes, disease progression, and evolution [1–6]. High-throughput chro-
4 mosome conformation capture (3C) assays such as Hi-C [7, 8], SPRITE [9], and GAM [6] provide a
5 comprehensive view of 3D organization by measuring interactions among chromosomal regions on a
6 genome-wide scale. High-throughput 3C data captured from diverse biological contexts and processes
7 has led to an improved understanding of DNA packaging in the nucleus, and the dynamics of 3D confor-
8 mation across developmental stages [10], and between normal and disease cellular states [4, 11]. Analysis
9 of such datasets has shown that chromosomal regions preferentially interact with one another, giving rise
10 to higher-order structural units such as chromosomal territories, compartments, and topologically asso-
11 ciating domains (TADs) which differ in the size of the structural unit and molecular features associated
12 with the constituent regions. Although the relationship between TADs and changes in gene expression
13 is debated [12–14], these units have been shown to be conserved across species [5, 15] and also associ-
14 ated with developmental [16] and disease processes [11, 17–19]. Accurate identification of TADs is an
15 important goal for linking 3D genome organization to cellular function.

16 Recently a large number of methods to identify TADs have been developed, utilizing different com-
17 putational frameworks, such as dynamic programming, [20, 21], community and subgraph detection
18 within networks [20, 22], Gaussian mixture modeling [23, 24], and signal processing approaches [25].
19 However, comparison of TAD-finding methods [26–28] have found large variability in the definition of
20 TADs and high sensitivity to the resolution (size of the genomic region), sequencing depth, and sparsity
21 of the input data. A lack of a clear definition for a TAD leads to difficulty in downstream interpretation
22 of these structures [29]. To address the sparsity of datasets, different smoothing based approaches have
23 been proposed, e.g. mean filter [30] and Gaussian filter [31]; however, it is unclear whether and to what
24 extent TAD identification can benefit from pre-smoothing the matrices.

25 Here, we present Graph Regularized Non-negative matrix factorization and Clustering for Hi-C
26 (GRiNCH), a novel matrix-factorization-based method for the analysis of high-throughput 3C datasets.
27 GRiNCH is based on non-negative matrix factorization (NMF), a powerful dimensionality reduction
28 method used to recover interpretable low-dimensional structure from high-dimensional datasets [32–34].
29 However, a standard application of NMF is not sufficient because of the strong distance dependence
30 found in Hi-C data, that is, regions that are close to each other on the linear genome tend to have more
31 interactions. We employ a graph regularized NMF approach, where the graph captures the distance de-

32 pendence of contact counts such that the learned lower-dimensional representation is smooth over the
33 graph structure [35]. Furthermore, by exploiting NMF's matrix completion property, which imputes
34 missing entries of a matrix from the product of the low-dimensional factors, GRiNCH can smooth a
35 sparse input matrix.

36 We perform a comprehensive comparison of GRiNCH and existing TAD-finding methods using a
37 number of metrics: similarity of interaction profiles of regions belonging to the same TAD, stability
38 to different resolutions and depth of input data, and enrichment of architectural proteins and histone
39 modification known to facilitate or correlate with 3D genome organization. Despite the general trend of
40 trade-off in performance among different criteria, e.g., a high performing method based on enrichment
41 of architectural proteins is not as stable to resolution and depth, GRiNCH consistently ranks among
42 the top across different measures. Furthermore, compared to existing smoothing approaches, GRiNCH-
43 based smoothing of downsampled data leads to the recovery of TADs and significant interactions best
44 in agreement with those from the original high-depth dataset. We apply GRiNCH to Hi-C data from
45 two different developmental time courses; we successfully recapitulate previously identified topological
46 changes around key genes, and predict novel boundary factors that could interact with known architec-
47 tural proteins to form topological domains. Taken together, GRiNCH is a robust and broadly applicable
48 approach to discover structural units and smooth sparse high-throughput 3C datasets from diverse plat-
49 forms including Hi-C, SPRITE and HiChIP.

50 Results

51 GRiNCH, a non-negative matrix factorization-based method for analyzing high-throughput 52 chromosome conformation capture datasets

53 GRiNCH uses graph-regularized Non-negative Matrix Factorization (NMF) to identify topologically
54 associating domains (TADs) from a high-dimensional 3C count matrix (**Figure 1, Methods**). GRiNCH
55 has several properties that make it attractive for analyzing these count matrices: (1) matrix factorization
56 methods including NMF have a “matrix completion” capability, which can be used to smooth noisy,
57 sparse matrices, (2) the low-dimensional factors provide a clustering of the row and column entities
58 that can be used to define chromosomal structural units, (3) the non-negativity constraint of the factors
59 provide a parts-based representation of the data and is well suited for count datasets (such as Hi-C
60 matrices), and (4) GRiNCH can be applied to any symmetric count matrix measuring chromosomal
61 interactions between genomic loci such as Hi-C, [36], SPRITE [9], and HiChIP [37] datasets.

62 For the ease of description, we will consider a Hi-C matrix as the input to GRiNCH. In GRiNCH,
63 the count matrix is approximated by the product of two lower dimensional matrices, U and V , both with
64 dimension $n \times k$, where n is the number of genomic regions in the given chromosome, and k is the
65 rank of the lower-dimensional space. Because Hi-C matrices have a strong distance dependence, we use
66 a constrained formulation of NMF, where the columns of the U and V matrices are smooth on a graph
67 of genomic regions (**Figure 1**), such that regions that are connected in the graph have similar sets of
68 values in the lower-dimensional space. The graph in turn captures the distance dependence using a local
69 neighborhood, where two regions i and j have an edge between them if they are within a particular radius
70 r of each other in linear distance along the chromosome. GRiNCH has three parameters, k , the rank of
71 the lower dimensional space, r to control the size of the neighborhood, and λ to control the strength
72 of graph regularization. After factorization, GRiNCH uses chain-constrained k -medoids clustering to
73 define clusters of contiguous regions, which we consider as TADs. We probed the impact of the three
74 parameters, k , r , and λ , on the resulting GRiNCH TADs. We determined that setting k to identify
75 TADs of size $\approx 1\text{Mb}$, with a neighborhood size of $r = 250\text{kb}$ and a small amount of regularization
76 ($\lambda = 1$), yields the best results (**Figure S1**). Notably, the regularization yields TADs with higher CTCF
77 enrichment than vanilla matrix factorization without any regularization (i.e. $\lambda = 0$).

78 **GRiNCH TADs are high quality and stable to varying resolution and depth of input Hi-C**
79 **data.**

80 To assess the quality of GRiNCH TADs, we considered seven existing TAD identification methods
81 (**Table 1**) and applied them along with GRiNCH to Hi-C data of five different cell lines from Rao et
82 al. [36] for comparison. The quality of a TAD was measured with two internal validation metrics used
83 for cluster evaluation, Davies-Bouldin index (DBI) and Delta Contact Count (DCC), both assessing the
84 similarity of interaction profiles of regions within defined TADs. DBI of a cluster measures how well
85 separated the given cluster is from other clusters; in our case, how distinct each TAD's interaction count
86 profile is from other TADs (**Methods**); a lower value for DBI indicates a more distinct, better-separated
87 cluster. DCC measures the difference between intra-TAD interaction counts and inter-TAD interaction
88 counts, with higher difference associated with better TADs. For each TAD-finding algorithm, we esti-
89 mated the proportion of TADs with significantly better DBI or DCC value than randomly shuffled TADs.
90 GRiNCH and HiCseg yield the highest proportion of TADs with significantly better DBI or DCC values
91 compared to randomly shuffled TADs (**Figure 2A**), suggesting these methods provide the most coherent
92 set of TADs.

93 Many TAD-calling methods are sensitive to the input data resolution (size of genomic region), with
94 the resulting TAD lengths varying greatly as a function of resolution [28]. A robust method is expected
95 to yield consistent length distribution of TADs when given the same user-specified parameter settings,
96 regardless of the change in resolution. Therefore, we next assessed the ability of GRiNCH and the
97 seven TAD calling methods for their ability to recover stable TADs across different resolutions, 10kb,
98 25kb, and 50kb. When comparing the median length of TADs across different resolutions (**Figure 2B**),
99 GRiNCH and Directionality Index are the most stable, with the exception of NHEK where Directionality
100 index learns longer TADs at 10k resolution. This suggests that GRiNCH is robust to different resolutions,
101 recovering consistently-sized TADs across different resolutions.

102 TAD-calling methods can be sensitive to the sparsity of the Hi-C matrices due to low sequencing
103 depth [28]. To assess the robustness of each method to low-depth, sparse datasets with many zero entries,
104 we first took the highest-depth dataset (GM12878, 86 million reads total) and downsampled to the depth
105 and sparsity level of lower-depth data from other cell lines (e.g. K562, the second “deepest” cell line with
106 16 million reads). We then compared the similarity of the TADs from the original high-depth data and
107 those from the downsampled counterpart (**Figure 3A, Methods**). We utilized metrics that can quantify
108 the similarity of pairs of clustering results: Rand Index and Mutual Information (**Methods**). Intuitively,

109 Rand Index is a measure of cluster membership consistency; it measures whether two data points (in
110 our case, two genomic regions) that belonged to the same cluster (TAD) in one clustering result also
111 stayed together in the other result, and whether two data points that belonged to different clusters stayed
112 separate. Rand Index ranges from 0 to 1, with 1 being perfect concordance. Mutual Information is an
113 informational-theoretic metric measuring the dependency between two random variables, where each
114 variable indicates a clustering result. A Mutual Information of 0 indicates complete disagreement and
115 the higher the Mutual Information value the better the agreement between the corresponding clustering
116 results. Based on Rand Index, TopDom, HiCseg, and GRiNCH yield the most reproducible TADs across
117 different depths, particularly at the lower depths of HMEC, HUVEC, and NHEK cell lines. Based on
118 Mutual Information, TopDom is the most consistent followed by GRiNCH and HiCseg. Other methods
119 were generally less consistent based on the Mutual Information metric.

120 A third hindrance in the interpretation of results from TAD finding methods is the disagreement
121 on the TAD definitions [28,29]. Hence, we further evaluated whether different TAD-calling methods
122 yielded relatively similar TADs, and which sets of methods yielded the most similar TADs to one another.
123 Here again, we used Rand Index and Mutual Information as metrics to compare the sets of TADs from
124 different methods. All pairwise comparisons of TAD-calling methods yielded high values of Rand Index
125 (>0.8) and high Mutual Information (**Figure 3B,C**). Furthermore, GRiNCH and TopDom yield the most
126 similar sets of TADs, followed by rGMAP across all cell lines. This pattern is fairly consistent even
127 when analyzed for each cell line individually (**Figure S2**).

128 To summarize, our internal validation and stability analysis showed that the top performing methods
129 depends upon the evaluation criteria. However, GRiNCH is among the top performing methods for all
130 the criteria we examined (**Table 1**), producing TADs that are as good or better than existing methods and
131 are stable to varying resolution and depth.

132 **GRiNCH TADs are enriched in architectural proteins and histone modification signals.**

133 We next characterized GRiNCH TADs as well as TADs from other methods for their ability to capture
134 well-known one-dimensional signal enrichment patterns (**Table 1**). In particular, one hallmark of TADs
135 is the enrichment of architectural proteins such as CTCF and cohesin elements (RAD21, SMC3) on the
136 boundaries of TADs [29,38]. We tested the TAD boundaries from each method for the enrichment of
137 peaks of CTCF, RAD21, and SMC3 in the five Rao et al. cell lines with Hi-C data (**Figure 4A, Methods**).
138 All methods identified boundaries enriched for peaks of these proteins; however, the methods varied
139 in their relative performance across cell lines. GRiNCH TAD boundaries have comparable or better

140 enrichment as the other top performing methods, namely, Directionality Index and Insulation Score in
141 most cell lines, and HiCseg in K562 and Huvec. All these methods including GRiNCH have significantly
142 higher enrichment than 3DNetMod, rGMAP, Armatus across different cell lines.

143 As histone modifications have been shown to be associated with three-dimensional organization [39],
144 we next measured the proportion of TADs with significant levels of mean histone modification signals
145 (**Figure 4B**) compared to randomly shuffled TADs (**Methods**). The histone modification signals in-
146 clude promoter- (H3K4me3, H3k4me2), elongation- (H3K79me2, H3k36me3), and enhancer-associated
147 marks (H3K27ac), and repressive chromatin marks (H3K27me3). A larger proportion of GRiNCH
148 TADs, along with Armatus and HiCseg TADs, are consistently enriched for the activating histone marks
149 such as H3K27ac, and the elongation marks, H3K36me3 and H3K79me2. Interestingly, with the ex-
150 ception of GM12878, the enrichment of histone marks in the TADs from Insulation and Directionality
151 index was much lower than the other methods suggesting these methods tend to find TADs defined by
152 CTCF and might miss other types of TADs [38]. These enrichment metrics show that when considering
153 existing methods, there is a tradeoff in the ability to recover TADs that are associated with CTCF and
154 TADs that are associated with significant histone modifications. However, GRiNCH ranks among the
155 top methods for both criteria suggesting that GRiNCH TADs capture a diverse types of TADs.

156 **GRiNCH smoothing of low-depth datasets help recover structure and significant** 157 **interactions.**

158 Our analysis so far compared different TAD finding methods for their ability to recover stable and bio-
159 logically meaningful topological units. However, most Hi-C datasets are sparse, which can influence the
160 TAD predictions significantly. Smoothing the input Hi-C matrix to impute missing values can enhance
161 the visualization of topological units on the matrix and improve the agreement among biological repli-
162 cates [30, 31]. Unlike existing TAD-calling methods, the matrix factorization framework of GRiNCH
163 provides a natural matrix completion solution that can generate a smoothed version of the sparse input
164 Hi-C matrix. We next compared GRiNCH's smoothing functionality to common smoothing techniques
165 such as mean filter and Gaussian filter, which are used in imaging domains and also for Hi-C data [30].
166 We used two metrics to assess the quality of smoothing: (a) recovery of TADs and (b) recovery of sig-
167 nificant interaction after smoothing downsampled data. To perform these comparisons, we again used
168 the downsampled GM12878 datasets.

169 To assess TAD recovery, we identified TADs on the original high-depth GM12878 dataset and com-
170 pared them to the TADs identified in the downsampled and smoothed data matrices using Rand Index

171 and Mutual Information. Here, to avoid any bias in our interpretation, we used the Directionality Index
172 method to call TADs. We find that using both Rand Index and Mutual Information, TADs recovered on
173 GRiNCH smoothed matrices are the most similar to the TADs from the high-depth dataset across differ-
174 ent parameter settings of the mean filter and Gaussian filters (**Figure 5A**). The usefulness of GRiNCH
175 is more apparent for lower-depth datasets such as NHEK. To assess the recovery of significant interac-
176 tions, we applied Fit-Hi-C [40] on the original GM12878 dataset and on the downsampled and smoothed
177 datasets to identify significant interactions (q -value < 0.05). Treating the significant interactions in
178 the original high-depth dataset as the ground truth, we measured precision and recall as a function of
179 the statistical significance of interactions from the smoothed datasets and computed the Area Under
180 Precision-Recall curve (AUPR). The higher the AUPR, the better the recovery of significant interactions
181 after smoothing. GRiNCH has the highest AUPR compared to mean filter and Gaussian filter (**Figure**
182 **5B**) across multiple parameter configurations. Overall, our experiments suggest that GRiNCH offers
183 superior smoothing functionality compared to standard smoothing techniques enabling better recovery
184 of TAD structures and long-range interactions.

185 **GRiNCH application to chromosomal organization during development.**

186 To assess the value of GRiNCH in primary cells and to examine dynamics in chromosomal organiza-
187 tion, we applied GRiNCH to two time-course Hi-C datasets profiling 3D genome organization during (a)
188 mouse neural development [41] and (b) pluripotency reprogramming in mouse [42]. Bonev et al. [41]
189 used high-resolution Hi-C experiments to measure 3D genome organization during neuronal differenti-
190 ation from the embryonic stem cell state (mESC) to neural progenitor cells (NPCs) and cortical neurons
191 (CNs). We applied GRiNCH on all chromosomes for all three cell types and compared them based on
192 the overall similarity of TADs between the cell lines. Based on the two metrics of Mutual Information
193 and Rand Index, the overall TAD similarity captured the temporal ordering of the cells, with CNs being
194 closer to NPCs and ESCs the most distinct (**Figure S3A**). We next focused on a specific 4Mbp region
195 around the *Zfp608* gene, which was found by Bonev et al. as a neural-specific gene associated with
196 a changing TAD boundary. In both NPCs and CNs, GRiNCH predicts a TAD near the *Zfp608* gene,
197 which is not present in the mESC state. *Zfp608* was also associated with increased expression, and ac-
198 tivating marks, H3K27ac and H3K4me3 at these time points, which is consistent with *Zfp608* being a
199 neural-specific gene (**Figure 6A**).

200 We next examined another time-course dataset which studied the 3D genome organization during
201 reprogramming of mouse pre-B cells to pluripotent stem cells (PSC), with four intermediate time points

202 (Day 2, 4, 6, and 8, see **Methods**). As in the neural developmental time course, we applied GRiNCH
203 to all chromosomes from each time point and compared the overall 3D genome configuration over time.
204 Here too we observed that time points closer to each other generally had greater similarity in their TAD
205 structure, as well as two different replicates within the same time point displaying even greater TAD
206 similarity (**Figure S4B**). We examined the interaction profile in the 1.3 Mbp around the Sox2 gene, a
207 known pluripotency gene (**Figure 6B**). We see a gradual formation of a boundary around Sox2, which is
208 also associated with concordant increase in expression, accessibility and the presence of H3K4me2, an
209 active promoter mark.

210 As chromatin accessibility data was also measured at each timepoint during reprogramming, we
211 asked if we could identify additional regulatory proteins that could play a role in establishing TADs
212 (**Methods**). Briefly, we tested the GRiNCH TAD boundaries from each mouse cell type, from pre-
213 B cell to pluripotent cells, for enrichment of accessible motif instances of 746 transcription factors in
214 the JASPAR 2020 core vertebrate motif database [43]. We ranked the TFs based on their significant
215 enrichment in each cell type (**Figure 6C**, **Table S2**). The top-ranking TF across the cell types was
216 CTCF, which is consistent with its role as an architectural protein in establishing TADs (**Figure 6C**).
217 We also found other factors in the same zinc finger protein family as CTCF [44], such as ZBTB14,
218 Plagl2/1, ZIC1/3/4/5, CTCFL, YY1/2 that were enriched across the cell types. YY1 and YY2 (which are
219 65 and 56% identical in their DNA and protein sequence respectively in humans [45]), are of interest,
220 as YY1 has been identified as an enforcer of long-range enhancer-promoter loops [46]. Interestingly, we
221 found several hematopoietic lineage factors, such as STAT3 and FOXP3, ranked highly in the pre-B cell
222 TADs compared to other time points. STAT3 is needed for B cell development [47]. FOXP3 is a master
223 regulator of T cells [48], but could be involved in the suppression of B cells. We also found a number of
224 HOX transcription factors, HOXA4, HOXA5, HOXB2, HOXB5, HOXB7, and the transcription factor
225 MEIS3 to be ranked highly in the B cells. The HOX genes depend upon MEIS3 [49] to bind to their
226 targets, supporting the simultaneous enrichment of these factors.

227 We repeated this analysis for the Rao et al. cell lines (**Table S3**). Here too we found CTCF and YY1/2
228 proteins highly enriched across cell lines. However, there was lesser degree of cell-line specificity for this
229 dataset. Taken together, this analysis suggests that GRiNCH captures high-quality TADs, which can be
230 used to define global similarities and difference between cell types. Furthermore, the GRiNCH boundary
231 enrichment analysis identified novel transcription factors that could be followed up with downstream
232 functional studies to examine their role in 3D genome organization.

GRiNCH can be used for a variety of 3D conformation capture technologies

Although Hi-C is still the most widely used technology to map 3D genome structure, recently several new methods have been developed to measure chromosomal contacts on a genome-wide scale [6]. To assess the applicability of GRiNCH to these technologies, we considered two complementary techniques to measure 3D genome organization: Split-Pool Recognition of Interactions by Tag Extension (SPRITE) [9] and HiChIP [37]. SPRITE measures multi-way chromatin interactions, and captures interactions across larger spatial distances than Hi-C. In HiChIP, long-range chromatin contacts are first established *in situ* in the nucleus before lysis; then chromatin immunoprecipitation (ChIP) is performed with respect to a specific protein or histone mark, directly capturing interactions associated with a protein or histone mark of interest [37]. A common property of both technologies is that they generate a contact count matrix, which is suitable for GRiNCH.

We applied GRiNCH to GM12878 contact matrices measured with SPRITE [9], cohesin HiChIP [37], and H3k27ac HiChIP [50]. A visual comparison between these datasets for an 8Mb region of chr8 shows regions of good concordance between datasets (**Figure 7A-D**). We quantified the global similarity of GRiNCH TADs from the four different datasets, for all chromosomes, with Rand Index (**Figure 7E**) and Mutual Information (**Figure 7F**). Interestingly, the GRiNCH TADs from Hi-C are the most similar to those from cohesin HiChIP and this similarity measure is higher than between the two HiChIP datasets. This is consistent with cohesin being a major determinant for the formation of loops detected in HiC datasets. The H3K27ac HiChIP data is as close to Hi-C as it is to cohesin HiChIP. Finally the most distinct set of TADs are identified by SPRITE, which is consistent with SPRITE capturing multi-way interactions and longer-distance interactions. Despite the differences in the cluster, overall the datasets look similar across different platforms (Rand Index >0.97). Taken together, this shows that GRiNCH is broadly apply to different experimental platforms for measuring genome-wide chromosome conformation.

Discussion

We present GRiNCH, a graph-regularized matrix factorization framework that enables reliable identification of high-quality genome organizational units, such as TADs, from high-throughput chromosome conformation capture datasets. GRiNCH is based on a novel constrained matrix factorization and clustering approach that enables recovery of contiguous blocks of genomic regions sharing similar interaction patterns as well as smoothing sparse input datasets.

A lack of gold standards for TADs emphasizes the need to probe both the statistical and biological nature of inferred TADs. Through extensive comparison of GRiNCH to existing methods with good performance in other benchmarking studies, we identified strengths and weaknesses of existing approaches. In particular, methods like Insulation Score identify TADs that are generally more enriched for signals such as CTCF and cohesin; however, when comparing statistical properties such as stability across resolutions and cluster coherence, this method does not necessarily rank the best. GRiNCH was among the top methods for both criteria, identifying clusters of genomic regions with high degree of similarity in their interaction profiles, stable to low-depth, sparse datasets, and enriched in architectural proteins and histone modification signals with known roles in chromatin organization.

A unique advantage of GRiNCH lies in its smoothing capability via matrix completion. Smoothing has been an independent task from TAD-calling and a key processing step in downstream analysis of Hi-C data (e.g. measuring reproducibility or concordance between Hi-C replicates [31]). We find that GRiNCH smoothing outperforms existing smoothing methods (mean filter and Gaussian filter) in its ability to retain TAD-level and interaction-level features of the input Hi-C data. Furthermore, GRiNCH is applicable to datasets from a wide variety of platforms, including SPRITE and HiChIP. Application of GRiNCH shows that Hi-C and HiChIP datasets capture more similar topological units than SPRITE. Interestingly, TADs from Hi-C and cohesin HiChIP are much closer than the two HiChIP datasets we compared. This shows that GRiNCH is capturing TADs that are reproducible across platforms. To study the ability of GRiNCH to identify dynamic topological changes along a time course, we applied GRiNCH to published developmental time-course datasets. GRiNCH recapitulated global temporal relationships in 3D organization and also transitions in topological units around key developmental genes. Thus, GRiNCH should be broadly applicable for analysis of chromosome conformation capture datasets with different experimental design, sequencing depths, and platforms.

The 3D organization of the genome is determined through a complex interplay of architectural proteins such as CTCF, cohesin elements, and other transcription factors such as WAPL [51]. Application

288 of GRiNCH to Hi-C datasets representing cell lines and temporally related conditions identified known
289 and novel transcription factors that could be important for establishing these boundaries in a cell-type-
290 specific or generic manner. In particular, we recovered YY1/2 proteins that have been shown to interact
291 with CTCF to establish long-range regulatory programs during lineage commitment [52]. Among the
292 novel factors that were present in both the cell lines as well as the mouse reprogramming dataset, were
293 several zinc finger proteins, e.g. PLAGL1, ZIC1, ZIC4/5, ZBTB14; such proteins can be investigated for
294 their role in establishing organizational units in mammalian genomes. We also found several factors that
295 were specific to cell lines and time points. For example, FOXI1, a forkhead protein, was ranked highly
296 in K562. Forkhead proteins are involved in genome organization and replication timing in yeast [53] and
297 zebra fish [54], but their role in mammalian genome organization is not well known. The time course data
298 identified additional unique TFs that are likely involved in determining specific lineages, e.g. STAT3,
299 MEIS3, FOXP3 and HOX genes in pre-B cells. HOX genes [55], FOXP3 [56], and STAT3 [47] in partic-
300 ular have been shown to play critical roles in B cell and T cell development. While MEIS1 and MEIS2
301 are involved in the hematopoietic lineage, MEIS3 specifically is involved in the binding of HOX TFs to
302 target genes in the brain [49]. Therefore the simultaneous enrichment of MEIS3 and HOX sites is con-
303 sistent with HOX proteins requiring MEIS3 for binding; however, its specific role in the hematopoietic
304 lineage is yet unknown. Investigating the interactions of these proteins with well-known architectural
305 proteins such as CTCF and cohesin could provide mechanistic insight into the factors governing 3D
306 genome organization [29, 57].

307 There are several directions of future work that are natural extensions to our framework. Although
308 our current approach of analyzing temporal organization in time-course data extracted interesting bio-
309 logical insights, TADs are identified independently for each time point, making it difficult to study the
310 conservation and specificity of individual TADs. One area of future work is to allow joint identification
311 of TADs or similar structural units across multiple conditions [58, 59]. Another direction is to leverage
312 one-dimensional features to potentially inform the TAD-finding algorithm. The GRiNCH framework
313 makes use of a distance dependence graph of regions; however, one could use the similarity of epige-
314 nomic profiles to construct an additional graph to constrain the NMF solution.

315 In conclusion, GRiNCH offers a unified solution, applicable to diverse platforms, to discover reliable
316 and biologically meaningful topological units, while handling sparse high-throughput chromosome con-
317 formation capture datasets. The outputs from GRiNCH can be used to predict novel boundary elements,
318 enabling us to test possible hypotheses of other mechanisms for TAD boundary formation. We have
319 made GRiNCH available at roy-lab.github.io/grinch, with a comprehensive installation and usage man-

320 ual. As efforts to map the three-dimensional genome organization expand to more conditions, platforms,
321 and species, a method such as GRiNCH will serve as a powerful analytical tool for understanding the
322 role of genome 3D organization in diverse complex processes.

Materials and Methods

Graph-regularized Non-negative Matrix Factorization (NMF) and Clustering for Hi-C data (GRiNCH) framework

GRiNCH is based on a regularized version of non-negative matrix factorization (NMF) [35] that is applicable to high-dimensional chromosome conformation capture data such as Hi-C (**Figure 1**). Below we describe the components of GRiNCH: NMF, graph regularization, and clustering for TAD identification.

Non-negative matrix factorization (NMF) and graph regularization

Non-negative matrix factorization is a popular dimensionality reduction method that aims to decompose a non-negative matrix, $X \in \mathbb{R}^{(n \times m)}$ into two lower dimensional non-negative matrices, $U \in \mathbb{R}^{(n \times k)}$ and $V \in \mathbb{R}^{(m \times k)}$, such that the product $X^* = UV^T$, well approximates the original X . We refer to the U and V matrices as factors. Here $k \ll n, m$ is the rank of the factors and is user-specified.

In application of NMF to Hi-C data, we represent the Hi-C data for each chromosome as a symmetric matrix $X = [x_{ij}] \in \mathbb{R}^{(n \times n)}$ where x_{ij} represents the contact count between region i and region j . We note that in the case of a symmetric matrix, U and V are the same or related by a scaling constant.

The goal of NMF is to minimize the following objective: $\|X - UV^T\|^2$, s.t. $U \geq 0, V \geq 0$ [32]. A number of algorithms to optimize this objective have been proposed; here we used the multiplicative update algorithm, where the entries of U and V are updated in an alternating manner each iteration:

$$u_{ik} \leftarrow u_{ik} \frac{(XV)_{ik}}{(UV^T V)_{ik}}, v_{jk} \leftarrow v_{jk} \frac{(X^T U)_{jk}}{(VU^T U)_{jk}} \quad (1)$$

Here u_{ik} corresponds to the i^{th} row of column $U(:, k)$ and v_{jk} corresponds to the j^{th} row of column $V(:, k)$.

Standard application of NMF to Hi-C data is ignorant of the strong distance dependence of the count matrix, that is, genomic regions that are close to each other tend to interact more with each other. To address this issue we apply an constrained version of NMF with graph regularization, where the graph represents additional constraints on the row (and/or column) entities [35]. Graph regularization enables the learned columns of U and V to be smooth over the input graph. In our application of NMF to Hi-C data, we define a graph composed of genomic regions as nodes, with edges connecting neighboring regions in the linear chromosome, where the size of the neighborhood is an input parameter. Specifically, we define a symmetric nearest-neighbor graph, W :

$$W_{ij} = \begin{cases} 1, & \text{if } x_i \in N_r(x_j) \text{ and } x_j \in N_r(x_i) \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

350 where $N_r(x_i)$ denotes r nearest neighbors in linear distance to region x_i .

351 Graph regularized NMF has the following objective:

$$\|X - UV^T\|^2 + \lambda \text{Tr}(V^T L V) + \lambda \text{Tr}(U^T L U), \quad (3)$$

352 where D is a diagonal matrix whose entries are column (or row, since W is symmetric) sums of
 353 W , i.e., $D_{ii} = \sum_j W_{ij}$. $L = D - W$ denotes the graph Laplacian and encodes the graph topology.
 354 The second and third terms are the regularization term and measures the smoothness of U and V with
 355 respect to the graph. Here λ is the regularization hyperparameter. This new objective has the effect of
 356 encouraging the factors to be smooth on the local neighborhood defined by the graph. Accordingly, the
 357 multiplicative update rule from (1) gains regularization terms [35]:

$$u_{ik} \leftarrow u_{ik} \frac{(XV + \lambda WU)_{ik}}{(UV^T V + \lambda DU)_{ik}}, \quad v_{jk} \leftarrow v_{jk} \frac{(X^T U + \lambda WV)_{jk}}{(VU^T U + \lambda DV)_{jk}} \quad (4)$$

358 Both r (neighborhood radius) and λ are parameters that can be specified, with λ setting the strength
 359 of regularization ($\lambda = 0$ makes this equivalent to basic NMF). See section on “Estimating GRINCH
 360 hyper-parameters” below.

361 **Chain-constrained k -medoids clustering for clustering assignment and TAD calling**

362 The factors U (or V) can be used to extract clusters of the row (or column) entities of the input matrix.
 363 Because X is symmetric in our application, either U or V can be used to define the clusters (the factors
 364 are equivalent up to a scaling constant). Assuming we use U , there are two common approaches for
 365 finding clusters from NMF factors: (1) assign each row entity i to its most dominant factor, i.e., assign
 366 it to cluster $c_i = \text{argmax}_{j \in \{1, \dots, k\}} u_{ij}$, or (2) apply k -means clustering on the rows of U . However, both
 367 approaches fall short in our application. The first approach is sensitive to extreme values which can still
 368 be present in the smoother factors, yielding non-informative clusters. Furthermore, neither approach
 369 reinforces contiguity of genomic regions in each cluster along their chromosomal position. As a result,
 370 a single cluster could potentially contain genomic regions from two opposite ends of the chromosomes
 371 instead of being a contiguous local structural unit. To address this problem, we apply chain-constrained
 372 k -medoids clustering. k -medoids clustering is similar to k -means clustering, except that the “center”

373 of each cluster is always an actual data point, rather than the mean of the datapoints in the cluster.
374 In its chain-constrained version (Algorithm 1), adopted from spatially connected k -medoids clustering
375 [60]: each cluster grows outwards from initial medoids along the linear chromosomal coordinates. The
376 algorithm assigns a genomic region to a valid medoid region either upstream or downstream along the
377 chromosome, ensuring the contiguity of the clusters and resilience to noise or extreme outliers provided
378 by using a robust ‘median’-like cluster center rather than a ‘mean’-like center used in k -means clustering.

Algorithm 1: Chain-constrained k -medoids clustering

Input: $U \in \mathbb{R}^{n \times k}$, one of the factors from NMF, and $maxIter$, the maximum number of iterations

Output: The cluster assignments, $\mathcal{C} \in \{c_1, c_2, \dots, c_n\}$, for each of the chromosomal bins

```

1 Initialize  $k$  medoids to be the rows with the largest value from each column of  $U$ 
2 Initialize an empty priority queue  $Q$ 
3 while  $numIter < maxIter$  do
4     Add current medoids to priority queue  $Q$ , with priority value of 0
                                     //  $Q$  orders bins by ascending priority values.
5     while  $Q$  is not empty do
6         Pop bin  $b$  from  $Q$ 
7         if  $b$  is not assigned to a cluster yet then
8             /* First, assign bin to cluster */
9             if  $b$  is a medoid then
10                Assign  $b$  to its own cluster
11            else
12                Assign  $b$  to either: the same cluster as its nearest upstream neighbor along the chromosome
13                already assigned to a cluster,  $u$ , or the same cluster as its nearest downstream neighbor
14                along the chromosome already assigned to a cluster,  $d$ , based on the similarity between the
15                latent feature vectors of  $b$  and the cluster medoids, i.e.,
16                 $\min_{c \in \{u, d\}} \|U[b, :] - U[\text{medoid of } c, :]\|$ 
17                /* Next, add any unassigned neighbor to priority queue: */
18                for each immediate upstream or downstream neighbor  $i$  of  $b$  not assigned to a cluster do
19                    Add  $i$  to  $Q$  with priority = priority of  $b$  +  $\|U[b, :] - U[i, :]\|$ 
20            end if
21        end while
22    Update medoids
23    if sum of distances between each bin and its cluster medoid didn't change from last iteration then
24        Break

```

379

380 **Selecting GRiNCH hyperparameters**

381 GRiNCH has three hyper-parameters: (a) k , the dimension of the lower-dimensional factors which can
382 alternately be viewed as the number of latent features or clusters, (b) r , the radius of the neighborhood
383 in the graph used for regularization, and (c) λ controlling the strength of regularization.

384 The parameter k determines the number of latent features to recover and the resulting number of
385 GRiNCH TADs. We can yield subTAD-, TAD-, or metaTAD-scale clusters (**Figure S5A**) by setting k
386 such that the expected size of a cluster is 500kb, 1Mb, or 2Mb, i.e., k equals the given chromosome's
387 length divided by the expected size. We find that a larger portion of subTAD-scale clusters (i.e. expected
388 TAD size = 500kb) have significant internal validation metric values (**Figure S5B**). SubTAD-scale clus-
389 ters tend to be more stable to depth and sparsity (**Figure S5C**), and are also more enriched in boundary
390 elements like CTCF (**Figure S6A**). As a tradeoff, higher proportion of metaTAD-scale clusters (i.e. ex-
391 pected cluster size = 2Mb) are enriched in histone modification marks (**Figure S6B**). Based on the use
392 case of GRiNCH, k can be set dynamically by the user; by default, GRiNCH sets k such that the expected
393 size of a cluster is 1Mb, or at TAD-scale.

394 For regularization strength, $\lambda \in \{0, 1, 10, 100, 100\}$ were considered, with $\lambda = 0$ equivalent to stan-
395 dard NMF without regularization. For neighborhood radius, $r \in \{25K, 50K, 100K, 250K, 500K, 1M\}$
396 were considered, where $r = 100K$ in a Hi-C dataset of 25Kb resolution will use 4 bins on either side
397 of a given region as its neighbors. We find that some regularization, with $\lambda = 1$, yields better CTCF
398 enrichment than other λ values (**Figure S1A**). With regularization, a neighborhood radius of 100Kb or
399 larger yields higher CTCF enrichment (**Figure S1B**). Based on these results, the default regularization
400 parameters for GRiNCH are set at $\lambda = 1$ and $r = 250kb$.

401 **Stability and initialization of NMF**

402 The NMF algorithm is commonly initialized with random non-negative values for the entries of U and
403 V . The initial values can significantly impact the final values of U and V [61]. This leads to instability of
404 the final factors hinging on the randomization schemes or changing seeds. To address the instability, we
405 used Non-Negative Double Singular Value Decomposition (NNDSVD), which initializes U and V with
406 a sparse SVD approximation of the input matrix X [62]. Since the derivation of exact singular values
407 can considerably slow down the initialization step, we use a randomized SVD algorithm which derives
408 approximate singular vectors [63]. NNDSVD initialization with randomized SVD results in lower loss,
409 i.e. factors that can better approximate the original Hi-C matrix, in fewer iterations (**Figure S7A,B**), and

410 more stable results than direct random initialization (**Figure S7C,D**).

411 **Datasets used in experiments and analysis**

412 **High-throughput chromosome conformation capture data**

413 We applied GRiNCH to SQRTVC-normalized Hi-C matrices from five cell lines, GM12878, NHEK,
414 HMEC, HUVEC, and K562 at 10kb, 25kb, and 50kb resolution from Rao et al. [36] (GEO accession:
415 GSE63525). We also applied GRiNCH to datasets from other technologies that capture the 3D genome
416 structure and chromatin interactions: Split-Pool Recognition of Interactions by Tag Extension (SPRITE)
417 [9] and HiChIP [37]. We used the SPRITE data for GM12878 cell line (GEO accession: GSE114242).
418 For HiChIP, we applied GRiNCH to the contact matrices from HiChIP with cohesin (GEO accession:
419 GSE80820) [37] and HiChIP with H3k27ac (GEO accession: GSE101498) [50].

420 We applied GRiNCH to two different mouse developmental time course data: (a) neural differ-
421 entiation Hi-C data from embryonic stem cells (mESC), neural progenitors (NPC), and cortical neu-
422 rons (CN) [41] and (b) Hi-C data from reprogramming pre-B cells to induced pluripotent state [42]
423 (GEO accession: GSE96553). For (a) neural differentiation dataset, Juicer Straw tool [64] was used to
424 obtain 25kb Hi-C matrices with vanilla-coverage square-root normalization (original GEO accession:
425 GSE96107). For (b) pluripotency reprogramming, we applied GRiNCH to published normalized Hi-C
426 data from pre-B cells, B α cells, day 2 of reprogramming, day 4, day 6, day 8, and finally, pluripotent
427 cells.

428 **ChIP-seq, DNaseq, ATACseq, and motif datasets**

429 To interpret the GRiNCH results and for comparison to other methods, we obtained a number of ChIP-
430 seq datasets. For CTCF, ChIP-seq narrow-peak datasets available as ENCODE Uniform TFBS com-
431 posite track [65] were downloaded from the UCSC genome browser (wgEncodeEH000029, wgEn-
432 codeEH000075, wgEncodeEH000054, wgEncodeEH000042, wgEncodeEH000063).

433 As ChIP-seq data for SMC3 and RAD21 is not available in the five cell lines from Rao et al [36],
434 we generated a list of cell-line specific accessible motif sites. Accessible motif sites are defined as the
435 intersection of motif-match regions and DNase-accessible regions in the given cell line. The SMC3
436 and RAD21 motif matches to the human genome (hg19) was obtained from [66]. To create a union
437 of DNase hotspot regions from replicates within a cell line, BEDtools [67] merge program was used.
438 Finally, the intersection of DNase hotspot regions and motif match regions was calculated for each cell

439 line using BEDtools intersect program. DNase hotspot data was obtained from the ENCODE consor-
440 tium [68, 69]: ENCFF856MFN, ENCFF235KUD, ENCFF491BOT, ENCFF946QPV, ENCFF968KGT,
441 ENCFF541JWD, ENCFF978UNU, ENCFF297CKS, ENCFF569UYX.

442 We obtained ChIP-seq datasets for histone modification marks from the ENCODE consortium [68,
443 69]. To generate genome-wide histone modification levels for each mark, fastq reads were aligned to the
444 human genome (hg19) with bowtie2 [70], and aggregated into a base-pair signal coverage profile using
445 SAMtools [71], and BEDtools [67]. The base-pair signal coverage was averaged within each 25kb bin to
446 match the resolution of Hi-C dataset. The aggregated signal was normalized by sequencing depth within
447 each replicate; the replicates were collapsed into a single value by taking the median.

448 In order to identify additional novel transcription factors that could play a role in 3D genome orga-
449 nization, we obtained motifs of 746 different transcription factors from JASPAR core vertebrate collec-
450 tion [43]. Next, we obtained their accessible motif match sites to hg19 and mm10 for the five
451 cell lines from [36] using the same process that was used for SMC3 and RAD21 motifs. To identify
452 the accessible motif sites for mouse cells during pluripotency reprogramming [42], we aligned ATACseq
453 fastq reads to the mouse genome (mm10) with bowtie2 [70] and deduplicated with SAMtools [71]. Ac-
454 cessible peaks were called with MACS2 [72]. The ATACseq peaks were then used in place of DNaseq
455 hotspots to find the accessible motif sites as was done for SMC3 and RAD21 motifs.

456 **TAD calling methods**

457 GRiNCH was benchmarked against 7 other TAD-calling methods: Directionality Index method [23],
458 Armatus [20], Insulation Score method [25], rGMAP [24], 3DNetMod [22], HiCseg [73] and TopDom
459 [74]. For all methods, default or recommended parameters values were used when available. Execution
460 scripts containing the parameter values used for these methods are available to download.

461 **Directionality index**

462 Directionality index uses a hidden Markov model (HMM) on estimated Directionality Index (DI) scores.
463 The DI score for a genomic region is determined by whether the region preferentially interacts with up-
464 stream or with downstream regions. A bin can take on one of three states: upstream-biased, downstream-
465 biased, or not biased, with directionally biased bins becoming TAD boundaries. TADs were called using
466 the directionality index method implementation in TADtool [75], version as of April 23, 2018.

467 **Armatus**

468 Armatus uses dynamic programming to find subgraphs in a network where the nodes are the genomic
469 regions, and the edge weights are the interaction counts. The objective is to find the set of dense sub-
470 graphs; subgraph density is defined as the ratio of the sum of edge weights to the number of nodes within
471 the subgraph. Armatus version 2.3 was used for comparison.

472 **Insulation score**

473 In the insulation score method, each bin is assigned an insulation score, calculated as the mean of the
474 interaction counts in the window (of a predefined size) centered on the given bin. Bins corresponding
475 to the local minima in the vector formed by these insulation scores are treated as TAD boundaries.
476 TADtool [75] implementation of insulation score method, version as of April 23, 2018, was used in our
477 experiments.

478 **3DNetMod**

479 3DNetMod employs a Louvain-like algorithm to partition a network of genomic regions into communi-
480 ties where the edge weights in the network are the interaction counts. It uses greedy dynamic program-
481 ming to maximize modularity, a metric of network structure measuring the density of intra-community
482 edges compared to random distribution of links between nodes. Version 1.0 (10/06/17) was used in our
483 comparison.

484 **rGMAP**

485 rGMAP trains a two-component Gaussian mixture model to group interactions into intra-domain or
486 inter-domain contacts. Putative TAD boundary bins are identified by those with significantly higher
487 intra-domain counts in its upstream window or downstream window of predefined size. The chromosome
488 is then segmented into TADs flanked by these boundaries. Version as of April 23, 2018 was used for
489 comparison.

490 **HiCseg**

491 HiCseg treats the Hi-C matrix as a 2D image to be segmented, with each block-diagonal segment corre-
492 sponding to a TAD. The counts within each block are modeled to be drawn from a certain distribution
493 (e.g. Gaussian distribution for normalized Hi-C data). Using dynamic programming, HiCseg finds a set

494 of block boundaries that would maximize the log likelihood of counts in each block being drawn from
495 an estimated distribution. Version 1.1 was used in our experiments.

496 **TopDom**

497 TopDom generates a score for each bin along the chromosome, where the score is the mean interaction
498 count between the given bin and a set of upstream and downstream neighbors (neighborhood size is a
499 user-specified parameter). Putative TAD boundaries are picked from a set of bins whose score forms a
500 local minimum; false positive boundaries are filtered out with a significance test. Version 0.0.2 was used
501 in our analysis.

502 **TAD evaluation criteria**

503 We evaluated the quality of TADs using different enrichment metrics as well as internal validation met-
504 rics used for comparing clustering algorithms.

505 **Enrichment analysis**

506 **Enrichment of known architectural proteins** We estimated the enrichment of three known ar-
507 chitectural proteins (CTCF, RAD21 and SMC3) in the TAD boundaries of five cell lines from Rao et
508 al [36]. TAD boundaries are defined by the starting bin and the ending bin of each predicted TAD, along
509 with one preceding the starting bin and one following the ending bin. Let N be the total number of bins
510 in a chromosome, n_{BIND} be the number of bins with one or more ChIP-seq peaks or accessible motif
511 sites, n_{TAD} be the number of TAD boundary bins, and $n_{\text{TAD-BIND}}$ be the number of TAD-boundary bins
512 with a binding event (ChIP-seq peak or accessible motif match site). The fold enrichment for a particular
513 protein is calculated as: $\frac{n_{\text{TAD-BIND}}/n_{\text{TAD}}}{n_{\text{BIND}}/N}$. Within each cell line, the fold enrichment across all chromosomes
514 was averaged; then the mean across cell lines was used to rank the TAD-calling methods (**Table S1F**,
515 Supplementary Data).

516 **Histone modification enrichment** We used the percentage of TADs enriched in histone modifi-
517 cation signals as a validation metric to assess the quality of TADs, similar to Zufferey et al., [28]. For
518 each TAD, the mean histone modification ChIP-seq signal was calculated for the regions within the
519 TAD. Next, for each TAD-calling method, TADs and non-TAD stretches were shuffled within each chro-
520 mosome 10 times to yield randomized TADs. The empirical p-value of a TAD was calculated as the
521 proportion of randomized TADs with higher mean ChIP-seq signal than that of the given TAD. A TAD

522 was considered significantly enriched if its p-value was less than 0.05. The mean proportion of TADs
523 with significant enrichment across cell lines was used to rank the TAD-calling methods (**Table S1G**,
524 Supplementary Data).

525 **Internal validation metrics**

526 Since a TAD represents a cluster of contiguous regions that tend to interact more among each other than
527 with regions from another TAD or cluster, we used two internal validation or cluster quality metrics,
528 Davies-Bouldin Index and Delta count, to evaluate the similarity of interaction profiles among regions
529 within a TAD.

530 **Davies-Bouldin Index (DBI)** The DBI for a single cluster C_i is defined as its similarity to its closest
531 cluster C_j , where $i, j \in \{1, \dots, k\}, i \neq j$: $DBI_i = \max_{i \neq j} S_{ij}$. The similarity metric, S_{ij} , between C_i
532 and C_j is defined as:

$$S_{ij} = \frac{d_i + d_j}{\text{distance}_{ij}} \quad (5)$$

533 where d_i is the average distance between each data point in cluster C_i and the cluster centroid and
534 distance_{ij} is the distance between the cluster centroids of C_i and C_j . In applying DBI to Hi-C data, a data
535 point consists of a vector of a genomic region's interaction counts with other regions in the chromosome
536 (e.g. an entire row or column in the Hi-C matrix); a cluster corresponds to a group of regions within
537 the same TAD; the cluster centroid is a mean vector of rows that belong to the same cluster/TAD. The
538 smaller the DBI, the more distinct the clusters are from one another.

539 For each TAD-calling method, we first computed the DBI for each TAD. Next, TADs and non-TAD
540 stretches were shuffled within each chromosome 10 times to yield randomized TADs. The empirical
541 p-value of a TAD was calculated as the proportion of randomized TADs with lower DBI (recall a lower
542 DBI means better clustering) than that of the given TAD.

543 **Delta Contact Count (DCC)** DCC for cluster C_i is defined as follows: let in_i denote the mean
544 interaction counts between pairs of regions that are both in C_i , and out_i denote the mean interaction
545 counts between pairs of regions where one region is in cluster C_i and the other region is not. Then
546 $DCC_i = \text{in}_i - \text{out}_i$.

547 We expect that for a good cluster, the pairs of regions within the cluster should have higher contact
548 counts. Therefore, the higher the value of DCC, the higher the quality of the cluster. Again, a cluster
549 corresponds to a group of regions within the same TAD. Given the DCC values for each TAD, the

550 empirical p-value of a TAD was calculated as the proportion of randomized TADs with higher delta
551 count than that of the given TAD.

552 A TAD was considered to have significant DBI or DCC if its p-value was less than 0.05. The
553 mean proportion of TADs with significant DBI/DCC across cell lines was used to rank the TAD-calling
554 methods (Table S1A,B, Supplementary Data).

555 TAD similarity and stability metrics

556 When assessing the similarity or stability of TADs, we used cluster comparison metrics, Rand Index and
557 Mutual Information. First, TADs were converted to clusters so that regions in the same TAD were all
558 assigned to the same cluster; all non-TAD regions, if a TAD-calling algorithm should have them, were
559 assigned to a single cluster together.

560 For Rand Index, each genomic region is treated as a node in a graph; two nodes are connected by an
561 edge if they are in the same cluster. Then, the number of edges that were preserved between clustering
562 result A and clustering result B is divided by the total number of pairs of nodes, i.e. number of edges
563 in a fully connected graph. Rand Index of 1 corresponds to perfect concordance between two clustering
564 results; Rand Index of 0 means no agreement.

565 Mutual Information (MI) is an information-theoretic metric measuring the dependency between two
566 random variables, where each variable can be a clustering result. Specifically,

$$567 \text{MI}(A; B) = \sum_{a \in A} \sum_{b \in B} p_{(A,B)}(a, b) \log \left(\frac{p_{(A,B)}(a, b)}{p_A(a)p_B(b)} \right) \quad (6)$$

568 where A, B are random variables derived from clustering results, e.g. A is the cluster assignment cor-
569 responding to TADs from high-depth data and B is the cluster assignment based on TADs from down-
570 sampled data. Mutual Information is 0 if the joint distribution of A and B equals the product of each
571 marginal distribution, i.e. A and B are independent, or in an information-theoretic sense, knowing A
572 does not provide any information about B . The higher the Mutual Information value, the greater the in-
573 formation conveyed by the variables about each other; in the context of measuring clustering agreement,
574 one clustering result is similar to the other.

575 Both metrics were used to evaluate the stability of TADs across depth, the similarity of TADs from
576 different TAD-calling methods, the recovery of TADs from smoothed Hi-C data, the similarity of TADs
577 along time-course data, and the consistency of GRiNCH TADs from different 3D genome capturing
578 technologies (e.g. SPRITE, HiChIP). In ranking TAD-calling methods for stability across depth, the

578 mean Rand Index or Mutual Information across cell lines was used (Table S1D,E, Supplementary Data).

579 **Robustness to low-depth data**

580 To assess the robustness or stability of TADs to low-depth input data, the TADs from a high-depth dataset
581 (GM12878) [36] were compared to the TADs from a downsampled, low-depth dataset. If the original set
582 of TADs are similar to the set of TADs from downsampled data, they are considered to be stable to low
583 depth. The similarity metrics used are described in the “TAD similarity and stability metrics” section.

584 In order to downsample a high-depth Hi-C matrix (e.g. from GM12878) to similar levels as a lower
585 depth one (e.g. from HMEC), a distance-stratified approach was used to match both the mean of non-
586 zero counts and sparsity level between the two datasets. First, for each distance threshold d , let μ_d^h denote
587 the mean of the non-zero counts in the high-depth dataset and μ_d^l denote the mean on non-zero counts in
588 the low-depth dataset. The scaled down value for each non-zero entry of the original high-depth dataset
589 is: $\tilde{x}_{ij} = \frac{x_{ij}^h}{\mu_d^h/\mu_d^l}$. where x_{ij}^h is the value for the i,j bin pair in the high-depth dataset. Then, to increase
590 the sparsity of the high-depth dataset, z_d of the non-zero counts in the high-depth dataset at distance d is
591 randomly set to zero, where z_d is the number of additional entries in the low-depth dataset that are zero
592 compared to the high-depth dataset.

593 **Identification of novel factor enrichment at GRiNCH TAD boundaries**

594 A similar procedure to CTCF boundary enrichment was used to identify novel boundary elements, by
595 assessing whether the accessible motif sites of 746 transcription factors from the JASPAR core vertebrate
596 collection [43] are enriched in GRiNCH TAD boundaries. This procedure was applied to the five cell
597 lines from Rao et al [36] and the cell types or time points from mouse reprogramming data [42]. One
598 change to the procedure was that instead of calculating fold enrichment per chromosome, all counts were
599 aggregated across all chromosomes within the given cell line, cell type, or time point. The hypergeo-
600 metric test was used to calculate the significance of the number of TF sites in the boundaries and were
601 ranked based on their p-value.

602 **Smoothing methods**

603 **Smoothing with GRiNCH via matrix completion** GRiNCH smooths a noisy input Hi-C matrix
604 by using the matrix completion aspect of NMF. Specifically, the reconstructed matrix $X^s = UV^T$ is
605 the smoothed matrix. The effectiveness of GRiNCH matrix completion as a smoothing method was

606 compared to that of mean filter and Gaussian filter, two methods used in image blurring [76] and Hi-C
607 datasets [30].

608 **Mean filter** Mean filtering is used in HiCRep [30] as a preprocessing step to measure reproducibility
609 of Hi-C datasets. To create a smoothed matrix X^s from a given input matrix X with a mean filter,
610 each element in x_{ij}^s is estimated from the mean of its neighboring elements within radius r : $x_{ij}^s =$
611 $\frac{1}{(2r+1)^2} \sum_{a=i-r}^{i+r} \sum_{b=j-r}^{j+r} x_{ab}$. Three different values for the radius r were considered: $r \in \{3, 6, 11\}$.

612 **Gaussian filter** A Gaussian filter uses a weighted mean of the neighborhood of a particular contact
613 count entry, x_{ij} , where the weight is determined by the distance of the neighbor from the given position:

$$x_{ij}^s = \frac{1}{2\pi\sigma^2} \sum_{a=i-n}^{i+n} \sum_{b=j-n}^{j+n} e^{-\frac{(i-a)^2+(j-b)^2}{2\sigma^2}} x_{ab} \quad (7)$$

614 Three different values of (σ) were considered, $\sigma \in \{1, 2, 3\}$ and n was set to $4 * \sigma$.

615 **Assessment of benefits from smoothing**

616 **Recovery of TADs from smoothed downsampled data** To assess whether smoothing helps pre-
617 serve or recover structure from low-depth data, downsampled datasets (see “Robustness to low-depth
618 data”) were smoothed with methods described above (see “Smoothing methods”). The Directionality
619 Index (DI) TAD finding method was applied to the high and low depth datasets. Then the similarity of
620 the TADs from the original high depth and the TADs from the smoothed data were measured (see “TAD
621 similarity and stability metrics”). Higher similarity metric values imply better recovery of structure from
622 smoothing.

623 **Recovery of significant interactions using Fit-Hi-C** Fit-Hi-C [40] was used to call significant
624 interactions in the original and the smoothed Hi-C datasets. Interactions from the original high-depth
625 Hi-C dataset with Fit-Hi-C q-value < 0.05 was defined as the set of “true” significant interactions. From
626 the downsampled then smoothed matrices, each smoothed interaction count was assigned a “prediction
627 score” of $1 -$ its Fit-Hi-C q-value. Precision and recall curves were then computed using the “true”
628 interactions and the “prediction scores.” The recovery of significant interactions was measured with the
629 Area under the Precision-Recall curve (AUPR).

630 **Implementation and availability**

631 Source code (implemented in C++), installation instructions (supported in Linux distributions), docu-
632 mentation, and tutorial for visualization (scripts implemented in Python) can be found at roy-lab.github.io/grinch.
633 Scripts used to analyze the results and generate the figures are available to download.

634 **Funding**

635 This work is supported by the National Institutes of Health (NIH) through the grant NHGRI R01-
636 HG010045-01.

637 **Author contributions**

638 Lee and Roy conceptualized the overall framework and algorithm. Lee implemented the algorithm,
639 designed and performed experiments, and wrote the manuscript. Roy designed the experiments and
640 wrote the manuscript.

641 **Acknowledgements**

642 We thank Shilu Zhang and Alireza Fotuhi Siahipirani for providing scripts for data processing and in-
643 terpretation of results. We also thank the Center for High Throughput Computing at UW Madison for
644 computational resources.

References

- 645
- 646 [1] Boyan Bonev and Giacomo Cavalli. Organization and function of the 3d genome. *Nature Reviews*
647 *Genetics*, 17(11):661–678, November 2016.
- 648 [2] Clemens B. Hug and Juan M. Vaquerizas. The Birth of the 3d Genome during Early Embryonic
649 Development. *Trends in Genetics*, 0(0), October 2018.
- 650 [3] Jordan Rowley, Michael Nichols, Xiaowen Lyu, Masami Ando-Kuri, Rivera, Karen Hermetz, Ping
651 Wang, Yijun Ruan, and Victor Corces. Evolutionarily Conserved Principles Predict 3d Chromatin
652 Organization. *Molecular Cell*, 67(5):837–852.e7, September 2017.
- 653 [4] Peter Hugo Lodewijk Krijger and Wouter de Laat. Regulation of disease-associated gene expres-
654 sion in the 3d genome. *Nature Reviews Molecular Cell Biology*, 17(12):771–782, December 2016.
- 655 [5] Quentin Szabo, Frédéric Bantignies, and Giacomo Cavalli. Principles of genome folding into
656 topologically associating domains. *Science Advances*, 5(4):eaaw1668, April 2019.
- 657 [6] Rieke Kempfer and Ana Pombo. Methods for mapping 3D chromosome architecture. *Nature*
658 *Reviews Genetics*, December 2019.
- 659 [7] E. Lieberman-Aiden, N. L. van Berkum, L. Williams, M. Imakaev, T. Ragozy, A. Telling, I. Amit,
660 B. R. Lajoie, P. J. Sabo, M. O. Dorschner, R. Sandstrom, B. Bernstein, M. A. Bender, M. Groudine,
661 A. Gnirke, J. Stamatoyannopoulos, L. A. Mirny, E. S. Lander, and J. Dekker. Comprehensive
662 Mapping of Long-Range Interactions Reveals Folding Principles of the Human Genome. *Science*,
663 326(5950):289–293, October 2009.
- 664 [8] M. Jordan Rowley and Victor G. Corces. Organizational principles of 3d genome architecture.
665 *Nature Reviews Genetics*, page 1, October 2018.
- 666 [9] Sofia Quinodoz, Noah Ollikainen, Barbara Tabak, Ali Palla, Jan Schmidt, Elizabeth Detmar, Mason
667 Lai, Alexander Shishkin, Prashant Bhat, Yodai Takei, Vickie Trinh, Erik Aznauryan, Pamela Rus-
668 sell, Christine Cheng, Marko Jovanovic, Amy Chow, Long Cai, Patrick McDonel, Manuel Garber,
669 and Mitchell Guttman. Higher-Order Inter-chromosomal Hubs Shape 3D Genome Organization in
670 the Nucleus. *Cell*, June 2018.
- 671 [10] Hui Zheng and Wei Xie. The role of 3D genome organization in development and cell differen-
672 tiation. *Nature Reviews Molecular Cell Biology*, 20(9):535–550, September 2019. Number: 9
673 Publisher: Nature Publishing Group.

- 674 [11] Abhijit Chakraborty and Ferhat Ay. The role of 3d genome organization in disease: From compart-
675 ments to single nucleotides. *Seminars in Cell & Developmental Biology*, July 2018.
- 676 [12] Somi Kim, Nam-Kyung Yu, and Bong-Kiun Kaang. CTCF as a multifunctional protein in genome
677 regulation and gene expression. *Experimental & Molecular Medicine*, 47(6):e166–e166, June
678 2015. Number: 6 Publisher: Nature Publishing Group.
- 679 [13] Yad Ghavi-Helm, Aleksander Jankowski, Sascha Meiers, Rebecca R. Viales, Jan O. Korbelt, and
680 Eileen E. M. Furlong. Highly rearranged chromosomes reveal uncoupling between genome topol-
681 ogy and gene expression. *Nature Genetics*, 51(8):1272–1282, August 2019. Number: 8 Publisher:
682 Nature Publishing Group.
- 683 [14] Bas van Steensel and Eileen E. M. Furlong. The role of transcription in shaping the spatial or-
684 ganization of the genome. *Nature Reviews Molecular Cell Biology*, 20(6):327–337, June 2019.
685 Number: 6 Publisher: Nature Publishing Group.
- 686 [15] Ittai E. Eres, Kaixuan Luo, Chiaowen Joyce Hsiao, Lauren E. Blake, and Yoav Gilad. Reorganiza-
687 tion of 3D Genome Structure May Contribute to Gene Regulatory Evolution in Primates. *bioRxiv*,
688 page 474841, November 2018.
- 689 [16] Ralph Stadhouders, Guillaume J. Filion, and Thomas Graf. Transcription factors and 3D genome
690 conformation in cell-fate decisions. *Nature*, 569(7756):345–354, May 2019. Number: 7756 Pub-
691 lisher: Nature Publishing Group.
- 692 [17] William A. Flavahan, Yotam Drier, Brian B. Liao, Shawn M. Gillespie, Andrew S. Venteicher,
693 Anat O. Stemmer-Rachamimov, Mario L. Suvà, and Bradley E. Bernstein. Insulator dysfunction
694 and oncogene activation in IDH mutant gliomas. *Nature*, 529(7584):110–114, January 2016.
- 695 [18] Dirk A. Kleinjan and Laura A. Lettice. Long-range gene control and genetic disease. *Advances in*
696 *genetics*, 61:339–388, 2008.
- 697 [19] Anne-Laure Valton and Job Dekker. TAD disruption as oncogenic driver. *Current opinion in*
698 *genetics & development*, 36:34–40, February 2016.
- 699 [20] Darya Filippova, Rob Patro, Geet Duggal, and Carl Kingsford. Identification of alternative topo-
700 logical domains in chromatin. *Algorithms for Molecular Biology*, 9(1):14+, May 2014.
- 701 [21] Caleb Weinreb and Benjamin J. Raphael. Identification of hierarchical chromatin domains. *Bioin-*
702 *formatics*, pages btv485+, August 2015.

- 703 [22] Heidi K. Norton, Daniel J. Emerson, Harvey Huang, Jesi Kim, Katelyn R. Titus, Shi Gu, Danielle S.
704 Bassett, and Jennifer E. Phillips-Cremins. Detecting hierarchical genome folding with network
705 modularity. *Nature Methods*, 15(2):119–122, January 2018.
- 706 [23] Jesse R. Dixon, Siddarth Selvaraj, Feng Yue, Audrey Kim, Yan Li, Yin Shen, Ming Hu, Jun S. Liu,
707 and Bing Ren. Topological domains in mammalian genomes identified by analysis of chromatin
708 interactions. *Nature*, 485(7398):376–380, April 2012.
- 709 [24] Wenbao Yu, Bing He, and Kai Tan. Identifying topologically associating domains and subdomains
710 by Gaussian Mixture model And Proportion test. *Nature Communications*, 8(1), September 2017.
- 711 [25] Emily Crane, Qian Bian, Rachel P. McCord, Bryan R. Lajoie, Bayly S. Wheeler, Edward J. Ralston,
712 Satoru Uzawa, Job Dekker, and Barbara J. Meyer. Condensin-driven remodelling of X chromosome
713 topology during dosage compensation. *Nature*, 523(7559):240–244, June 2015.
- 714 [26] Mattia Forcato, Chiara Nicoletti, Koustav Pal, Carmen M. Livi, Francesco Ferrari, and Silvio Bic-
715 ciato. Comparison of computational methods for Hi-C data analysis. *Nature Methods*, 14(7):679–
716 685, June 2017.
- 717 [27] Rola Dali and Mathieu Blanchette. A critical assessment of topologically associating domain pre-
718 diction tools. *Nucleic acids research*, 45(6):2994–3005, April 2017.
- 719 [28] 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):217, December 2018.
- 722 [29] Elzo de Wit. TADs as the Caller Calls Them. *Journal of Molecular Biology*, page
723 S0022283619305923, October 2019.
- 724 [30] Tao Yang, Feipeng Zhang, Galip Gurkan Yardimci, Fan Song, Ross C. Hardison, William Stafford
725 Noble, Feng Yue, and Qunhua Li. HiCRep: assessing the reproducibility of Hi-C data using a
726 stratum- adjusted correlation coefficient. *Genome Research*, page gr.220640.117, August 2017.
- 727 [31] Oana Ursu, Nathan Boley, Maryna Taranova, Y. X. Rachel Wang, Galip Gurkan Yardimci, William
728 Stafford Noble, and Anshul Kundaje. GenomeDISCO: a concordance score for chromosome
729 conformation capture experiments using random walks on contact map graphs. *Bioinformatics*,
730 34(16):2701–2707, August 2018.
- 731 [32] Daniel D. Lee and H. Sebastian Seung. Algorithms for Non-negative Matrix Factorization. In *In*
732 *NIPS*, volume 13, pages 556–562, 2000.

- 733 [33] Yan Wu, Pablo Tamayo, and Kun Zhang. Visualizing and Interpreting Single-Cell Gene Expression
734 Datasets with Similarity Weighted Nonnegative Embedding. *Cell Systems*, December 2018.
- 735 [34] Genevieve L. Stein-O'Brien, Raman Arora, Aedin C. Culhane, Alexander V. Favorov, Lana X.
736 Garmire, Casey S. Greene, Loyal A. Goff, Yifeng Li, Aloune Ngom, Michael F. Ochs, Yanxun Xu,
737 and Elana J. Fertig. Enter the Matrix: Factorization Uncovers Knowledge from Omics. *Trends in*
738 *Genetics*, 0(0), August 2018.
- 739 [35] Deng Cai, Xiaofei He, Jiawei Han, and Thomas S. Huang. Graph Regularized Nonnegative Matrix
740 Factorization for Data Representation. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(8):1548–1560,
741 August 2011.
- 742 [36] Suhas S. P. Rao, Miriam H. Huntley, Neva C. Durand, Elena K. Stamenova, Ivan D. Bochkov,
743 James T. Robinson, Adrian L. Sanborn, Ido Machol, Arina D. Omer, Eric S. Lander, and Erez L.
744 Aiden. A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin
745 Looping. *Cell*, 159(7):1665–1680, December 2014.
- 746 [37] Maxwell R. Mumbach, Adam J. Rubin, Ryan A. Flynn, Chao Dai, Paul A. Khavari, William J.
747 Greenleaf, and Howard Y. Chang. HiChIP: efficient and sensitive analysis of protein-directed
748 genome architecture. *Nature Methods*, 13(11):919–922, November 2016.
- 749 [38] Li-Hsin Chang, Sourav Ghosh, and Daan Noordermeer. TADs and their borders: free movement
750 or building a wall? *Journal of Molecular Biology*, page S0022283619307429, December 2019.
- 751 [39] Guillaume Andrey, Robert Schöpflin, Ivana Jerković, Verena Heinrich, Daniel M. Ibrahim,
752 Christina Paliou, Myriam Hochradel, Bernd Timmermann, Stefan Haas, Martin Vingron, and Ste-
753 fan Mundlos. Characterization of hundreds of regulatory landscapes in developing limbs reveals
754 two regimes of chromatin folding. *Genome Research*, 27(2):223–233, February 2017. Company:
755 Cold Spring Harbor Laboratory Press Distributor: Cold Spring Harbor Laboratory Press Institu-
756 tion: Cold Spring Harbor Laboratory Press Label: Cold Spring Harbor Laboratory Press Publisher:
757 Cold Spring Harbor Lab.
- 758 [40] Ferhat Ay, Timothy L. Bailey, and William Stafford Noble. Statistical confidence estimation for
759 Hi-C data reveals regulatory chromatin contacts. *Genome Research*, 24(6):999–1011, June 2014.
- 760 [41] Boyan Bonev, Netta Mendelson Cohen, Quentin Szabo, Lauriane Fritsch, Giorgio L. Papadopoulos,
761 Yaniv Lubling, Xiaole Xu, Xiaodan Lv, Jean-Philippe Hugnot, Amos Tanay, and Giacomo Cavalli.
762 Multiscale 3D Genome Rewiring during Mouse Neural Development. *Cell*, 171(3):557–572.e24,
763 October 2017.

- 764 [42] Ralph Stadhouders, Enrique Vidal, François Serra, Bruno Di Stefano, François Le Dily, Javier
765 Quilez, Antonio Gomez, Samuel Collombet, Clara Berenguer, Yasmina Cuartero, Jochen Hecht,
766 Guillaume Filion, Miguel Beato, Marc Marti-Renom, and Thomas Graf. Transcription factors
767 orchestrate dynamic interplay between genome topology and gene regulation during cell repro-
768 gramming. *Nature Genetics*, 50(2):238–249, January 2018.
- 769 [43] Oriol Fornes, Jaime A. Castro-Mondragon, Aziz Khan, Robin van der Lee, Xi Zhang, Phillip A.
770 Richmond, Bhavi P. Modi, Solenne Correard, Marius Gheorghe, Damir Baranašić, Walter Santana-
771 Garcia, Ge Tan, Jeanne Chèneby, Benoit Ballester, François Parcy, Albin Sandelin, Boris Lenhard,
772 Wyeth W. Wasserman, and Anthony Mathelier. JASPAR 2020: update of the open-access database
773 of transcription factor binding profiles. *Nucleic Acids Research*, 48(D1):D87–D92, January 2020.
774 Publisher: Oxford Academic.
- 775 [44] Matteo Cassandri, Artem Smirnov, Flavia Novelli, Consuelo Pitolli, Massimiliano Agostini,
776 Michal Malewicz, Gerry Melino, and Giuseppe Raschellà. Zinc-finger proteins in health and dis-
777 ease. *Cell Death Discovery*, 3(1):1–12, November 2017. Number: 1 Publisher: Nature Publishing
778 Group.
- 779 [45] Xiao-nan Wu, Tao-tao Shi, Yao-hui He, Fei-fei Wang, Rui Sang, Jian-cheng Ding, Wen-juan Zhang,
780 Xing-yi Shu, Hai-feng Shen, Jia Yi, Xiang Gao, and Wen Liu. Methylation of transcription factor
781 YY2 regulates its transcriptional activity and cell proliferation. *Cell Discovery*, 3(1):1–22, October
782 2017. Number: 1 Publisher: Nature Publishing Group.
- 783 [46] Abraham S. Weintraub, Charles H. Li, Alicia V. Zamudio, Alla A. Sigova, Nancy M. Hannett,
784 Daniel S. Day, Brian J. Abraham, Malkiel A. Cohen, Behnam Nabet, Dennis L. Buckley, Yang Eric
785 Guo, Denes Hnisz, Rudolf Jaenisch, James E. Bradner, Nathanael S. Gray, and Richard A. Young.
786 YY1 Is a Structural Regulator of Enhancer-Promoter Loops. *Cell*, 171(7):1573–1588.e28, Decem-
787 ber 2017.
- 788 [47] Wei-Chun Chou, David E. Levy, and Chien-Kuo Lee. STAT3 positively regulates an early step in
789 B-cell development. *Blood*, 108(9):3005–3011, November 2006.
- 790 [48] Ling Lu, Joseph Barbi, and Fan Pan. The regulation of immune tolerance by FOXP3. *Nature*
791 *Reviews Immunology*, 17(11):703–717, November 2017. Number: 11 Publisher: Nature Publishing
792 Group.
- 793 [49] Rosa A. Uribe and Marianne E. Bronner. Meis3 is required for neural crest invasion of the gut
794 during zebrafish enteric nervous system development. *Molecular Biology of the Cell*, 26(21):3728–

795 3740, November 2015.

- 796 [50] Maxwell R. Mumbach, Ansuman T. Satpathy, Evan A. Boyle, Chao Dai, Benjamin G. Gowen,
797 Seung Woo Cho, Michelle L. Nguyen, Adam J. Rubin, Jeffrey M. Granja, Katelynn R. Kazane,
798 Yuning Wei, Trieu Nguyen, Peyton G. Greenside, M. Ryan Corces, Josh Tycko, Dimitre R. Sime-
799 onov, Nabeela Suliman, Rui Li, Jin Xu, Ryan A. Flynn, Anshul Kundaje, Paul A. Khavari, Alexan-
800 der Marson, Jacob E. Corn, Thomas Quertermous, William J. Greenleaf, and Howard Y. Chang.
801 Enhancer connectome in primary human cells identifies target genes of disease-associated DNA
802 elements. *Nature Genetics*, 49(11):1602–1612, November 2017. Number: 11 Publisher: Nature
803 Publishing Group.
- 804 [51] Judith H. I. Haarhuis, Robin H. van der Weide, Vincent A. Blomen, J. Omar Yáñez-Cuna, Mario
805 Amendola, Marjon S. van Ruiten, Peter H. L. Krijger, Hans Teunissen, René H. Medema, Bas van
806 Steensel, Thijn R. Brummelkamp, Elzo de Wit, and Benjamin D. Rowland. The Cohesin Release
807 Factor WAPL Restricts Chromatin Loop Extension. *Cell*, 169(4):693–707.e14, May 2017.
- 808 [52] Jonathan Beagan, Michael Duong, Katelyn Titus, Linda Zhou, Zhendong Cao, Jingjing Ma, Car-
809 oline Lachanski, Daniel Gillis, and Jennifer Phillips-Cremins. YY1 and CTCF orchestrate a 3D
810 chromatin looping switch during early neural lineage commitment. *Genome Research*, 27(7):1139–
811 1152, May 2017.
- 812 [53] Simon R. V. Knott, Jared M. Peace, A. Zachary Ostrow, Yan Gan, Alexandra E. Rex, Christopher J.
813 Viggiani, Simon Tavaré, and Oscar M. Aparicio. Forkhead Transcription Factors Establish Origin
814 Timing and Long-Range Clustering in *S. cerevisiae*. *Cell*, 148(1):99–111, January 2012. Publisher:
815 Elsevier.
- 816 [54] Jizhou Yan, Lisha Xu, Gregory Crawford, Zenfeng Wang, and Shawn M. Burgess. The Fork-
817 head Transcription Factor FoxI1 Remains Bound to Condensed Mitotic Chromosomes and Stably
818 Remodels Chromatin Structure. *Molecular and Cellular Biology*, 26(1):155–168, January 2006.
819 Publisher: American Society for Microbiology Journals Section: ARTICLES.
- 820 [55] R. A. Alharbi, R. Pettengell, H. S. Pandha, and R. Morgan. The role of HOX genes in normal
821 hematopoiesis and acute leukemia. *Leukemia*, 27(5):1000–1008, April 2013.
- 822 [56] Zhiyuan Li, Dan Li, Andy Tsun, and Bin Li. FOXP3 + regulatory T cells and their functional
823 regulation. *Cellular & Molecular Immunology*, 12(5):558–565, September 2015. Number: 5
824 Publisher: Nature Publishing Group.

- 825 [57] Caelin Cubeñas-Potts and Victor G. Corces. Architectural Proteins, Transcription, and the Three-
826 dimensional Organization of the Genome. *FEBS letters*, 589(20 0 0):2923–2930, October 2015.
- 827 [58] Alireza Fotuhi Siahpirani, Ferhat Ay, and Sushmita Roy. A multi-task graph-clustering approach
828 for chromosome conformation capture data sets identifies conserved modules of chromosomal in-
829 teractions. *Genome Biology*, 17(1):114, May 2016.
- 830 [59] Yang Yang, Yang Zhang, Bing Ren, Jesse R. Dixon, and Jian Ma. Comparing 3D Genome Organi-
831 zation in Multiple Species Using Phylo-HMRF. *Cell Systems*, 8(6):494–505.e14, 2019.
- 832 [60] S. Soor, A. Challa, S. Danda, B. S. Daya Sagar, and L. Najman. Extending K-Means to Preserve
833 Spatial Connectivity. In *IGARSS 2018 - 2018 IEEE International Geoscience and Remote Sensing*
834 *Symposium*, pages 6959–6962, July 2018.
- 835 [61] Mark Belford, Brian Mac Namee, and Derek Greene. Stability of topic modeling via matrix fac-
836 torization. *Expert Systems with Applications*, 91:159–169, January 2018.
- 837 [62] C. Boutsidis and E. Gallopoulos. SVD based initialization: A head start for nonnegative matrix
838 factorization. *Pattern Recognition*, 41(4):1350–1362, April 2008.
- 839 [63] Sergey Voronin and Per-Gunnar Martinsson. Rsvdpack: An implementation of randomized al-
840 gorithms for computing the singular value, interpolative, and cur decompositions of matrices on
841 multi-core and gpu architectures. 2015.
- 842 [64] Neva C. Durand, Muhammad S. Shamim, Ido Machol, Suhas S. P. Rao, Miriam H. Huntley, Eric S.
843 Lander, and Erez Lieberman Aiden. Juicer Provides a One-Click System for Analyzing Loop-
844 Resolution Hi-C Experiments. *Cell Systems*, 3(1):95–98, July 2016. Publisher: Elsevier.
- 845 [65] Kate R. Rosenbloom, Cricket A. Sloan, Venkat S. Malladi, Timothy R. Dreszer, Katrina Learned,
846 Vanessa M. Kirkup, Matthew C. Wong, Morgan Maddren, Ruihua Fang, Steven G. Heitner, Brian T.
847 Lee, Galt P. Barber, Rachel A. Harte, Mark Diekhans, Jeffrey C. Long, Steven P. Wilder, Ann S.
848 Zweig, Donna Karolchik, Robert M. Kuhn, David Haussler, and W. James Kent. ENCODE Data
849 in the UCSC Genome Browser: year 5 update. *Nucleic Acids Research*, 41(Database issue):D56–
850 D63, January 2013.
- 851 [66] Pouya Kheradpour and Manolis Kellis. Systematic discovery and characterization of regulatory
852 motifs in ENCODE TF binding experiments. *Nucleic Acids Research*, 42(5):2976–2987, March
853 2014. Publisher: Oxford Academic.
- 854 [67] Aaron R. Quinlan and Ira M. Hall. BEDTools: a flexible suite of utilities for comparing genomic
855 features. *Bioinformatics*, 26(6):841–842, March 2010.

- 856 [68] ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human
857 genome. *Nature*, 489(7414):57–74, September 2012.
- 858 [69] Cricket A. Sloan, Esther T. Chan, Jean M. Davidson, Venkat S. Malladi, J. Seth Strattan, Ben-
859 jamin C. Hitz, Idan Gabdank, Aditi K. Narayanan, Marcus Ho, Brian T. Lee, Laurence D. Rowe,
860 Timothy R. Dreszer, Greg Roe, Nikhil R. Podduturi, Forrest Tanaka, Eurie L. Hong, and J. Michael
861 Cherry. ENCODE data at the ENCODE portal. *Nucleic acids research*, 44(D1):gkv1160+, January
862 2016.
- 863 [70] Ben Langmead and Steven L. Salzberg. Fast gapped-read alignment with Bowtie 2. *Nature Meth-*
864 *ods*, 9(4):357–359, April 2012.
- 865 [71] Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth,
866 Goncalo Abecasis, Richard Durbin, and 1000 Genome Project Data Processing Subgroup. The
867 Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*, 25(16):2078–
868 2079, August 2009.
- 869 [72] Yong Zhang, Tao Liu, Clifford Meyer, Jérôme Eeckhoutte, David Johnson, Bradley Bernstein, Chad
870 Nussbaum, Richard Myers, Myles Brown, Wei Li, and Shirley Liu. Model-based Analysis of ChIP-
871 Seq (MACS). *Genome Biology*, 9(9):R137, September 2008.
- 872 [73] Celine Lévy-Leduc, M. Delattre, T. Mary-Huard, and S. Robin. Two-dimensional segmentation for
873 analyzing Hi-C data. *Bioinformatics (Oxford, England)*, 30(17):i386–392, September 2014.
- 874 [74] Hanjun Shin, Yi Shi, Chao Dai, Harianto Tjong, Ke Gong, Frank Alber, and Xianghong Jasmine
875 Zhou. TopDom: an efficient and deterministic method for identifying topological domains in
876 genomes. *Nucleic Acids Research*, 44(7):e70, April 2016.
- 877 [75] Kai Kruse, Clemens B. Hug, Benjamín Hernández-Rodríguez, and Juan M. Vaquerizas. TADtool:
878 visual parameter identification for TAD-calling algorithms. *Bioinformatics*, 32(20):3190–3192,
879 October 2016.
- 880 [76] E. R. Davies. *Machine Vision: Theory, Algorithms, Practicalities*. Elsevier, December 2004.

	Validation		Stability			Consistency	Enrichment	
	DBI	DCC	Resolution	RI	MI		CTCF	Histone
GRiNCH								
3DNetMod								
rGMAP								
Armatus								
Directionality								
Insulation								
HiCseg								
TopDom								

Table 1: Shown are different criteria of evaluation. A medal denotes whether the given TAD-calling method is among the top 3 methods for a particular criteria (gold/yellow: 1st place; silver/grey: 2nd place; bronze/brown: 3rd place). DBI: proportion of TADs with significant Davies-Bouldin Index; DCC: proportion of TADs with significant Delta Contact Counts; Resolution: stability of median TAD size to Hi-C resolution; RI, MI: stability to depth and sparsity of input data, measured by Rand Index (RI) and Mutual Information (MI); Consistency: a group of methods yielding TADs with highest similarity, with gold for the pair of methods with highest similarity according to hierarchical clustering; CTCF: fold enrichment of CTCF and cohesin elements in TAD boundaries; Histone: proportion of TADs with significant mean histone signal. See **Table S1** and Supplementary Data for more details.

List of Figures

1	Overview of GRiNCH	39
2	Characterizing TADs with internal validation metrics and TAD size.	40
3	Evaluating the stability and similarity of different TAD-calling methods.	41
4	Evaluating the quality of TADs from different TAD-calling methods using boundary element and chromatin mark enrichment.	42
5	Evaluating the benefits of smoothing in GRiNCH.	43
6	GRiNCH applied to Hi-C datasets along developmental time courses.	44
7	Applying GRiNCH to datasets from other 3D genome conformation capture technologies.	45

Figure 1

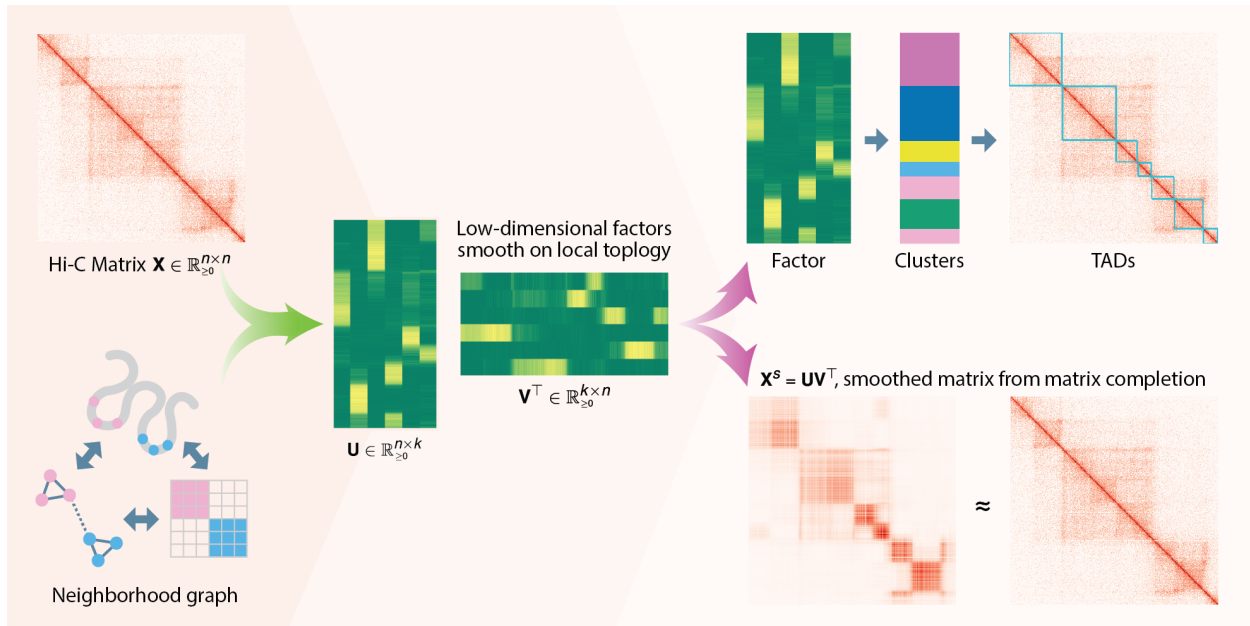


Figure 1: Overview of GRiNCH. GRiNCH applies Non-negative Matrix Factorization (NMF) to a Hi-C or a similar high-throughput 3C matrix to find clusters of densely interacting genomic regions. NMF recovers low-dimensional factors \mathbf{U} and \mathbf{V} of the input matrix \mathbf{X} that can be used to reconstruct the input matrix. As nearby genomic regions tend to interact more with each other, we regularize the factor matrices with a neighborhood graph to encourage neighboring regions to have a similar lower-dimensional representation, and subsequently belong to the same cluster. We cluster the regions by treating one of the factor matrices as a set of latent features and applying k -medoids clustering. The clusters represent topological units such as TADs. The factor matrices can be multiplied together to yield a smoothed version of the input matrix which is often sparse and noisy.

Figure 2

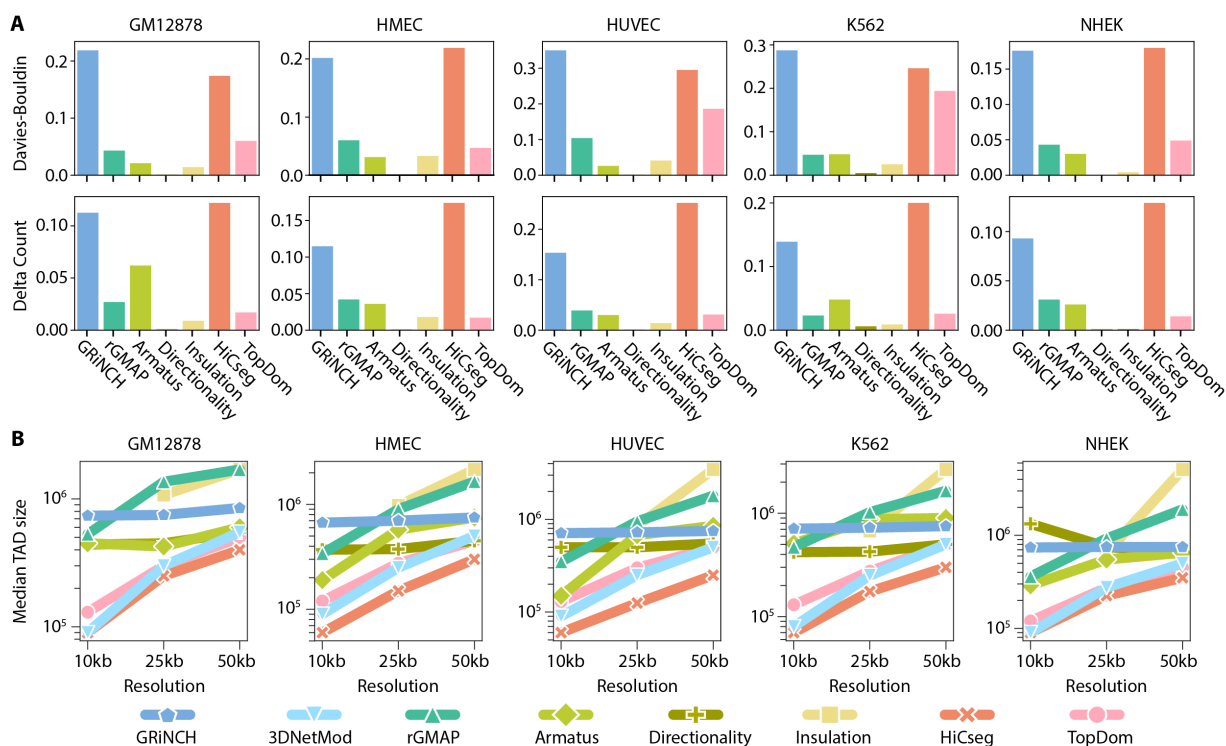


Figure 2: Characterizing TADs with internal validation metrics and median TAD size. **A.** Proportion of TADs with empirical p-value < 0.05 for internal validation metrics Davies-Bouldin Index and Delta Count, shown for GRiNCH and six other methods. Note: 3DNetMod outputted overlapping TADs, even when run under non-hierarchical settings and was excluded from this analysis which involves TAD randomization/shuffling. Treating each TAD as a cluster of genomic regions, we evaluate how distinct each cluster is to other clusters (Davies-Bouldin Index) and how much higher the intra-cluster interactions are compared to inter-cluster interactions (Delta Count). The p-value of each cluster's metric value is derived against the empirical distribution of the metric values in randomly shuffled TADs. The higher the bar, the better a method. **B.** The median size of TADs identified by different methods from different Hi-C resolutions. A method is considered stable to the resolution (size of the region) of the data if the median TAD size does not change substantially with the resolution, given the same user-defined parameter settings.

Figure 3

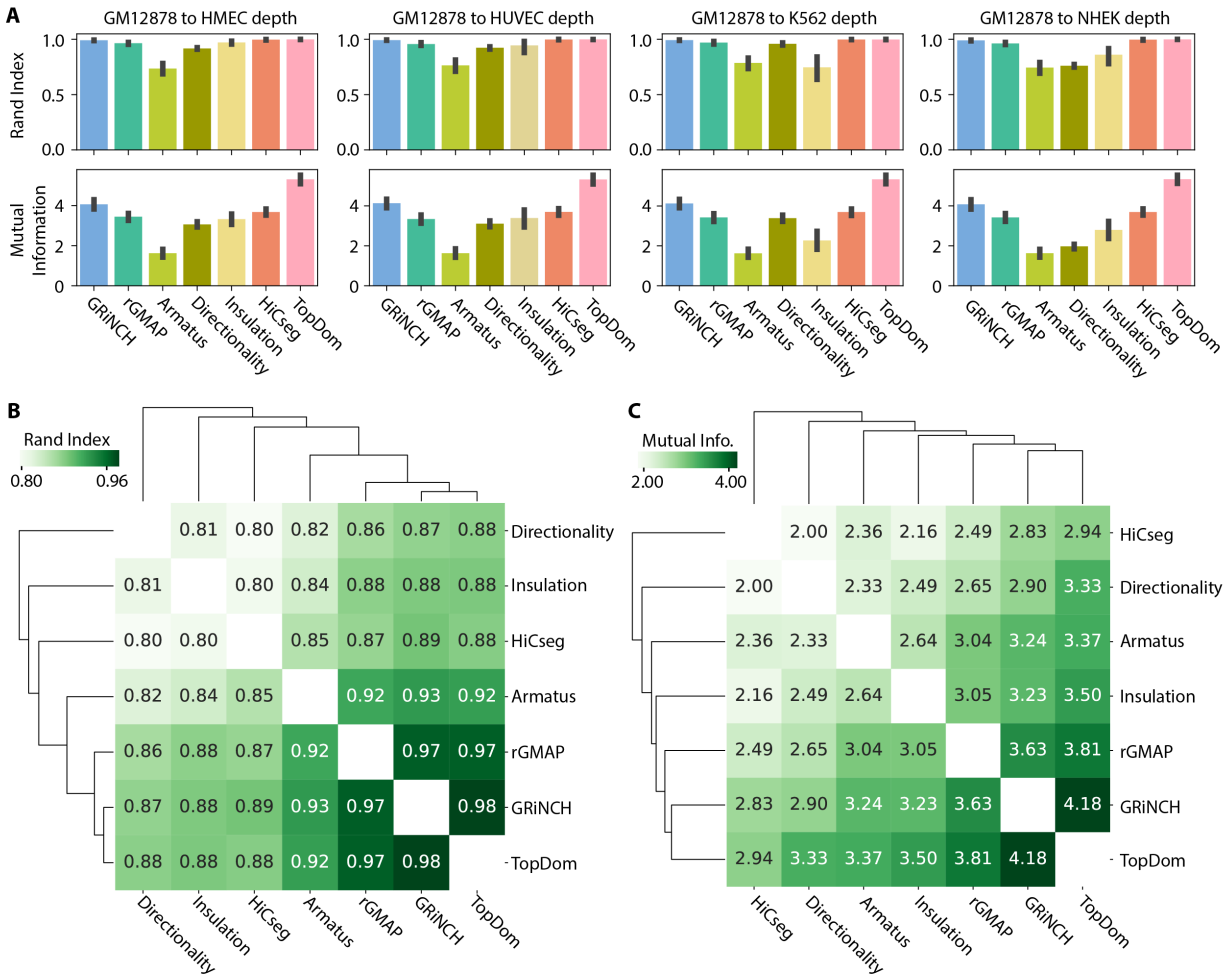


Figure 3: Evaluating the stability and similarity of different TAD-calling methods. **A.** The mean similarity, across chromosomes, between TADs from high-depth GM12878 dataset and TADs from low-depth GM12878 datasets obtained by downsampling the GM12878 dataset to different depths observed in our five cell-line dataset. The similarity of the TADs is measured by Rand Index and Mutual Information. The error bar denotes the standard deviation from the mean. **B.** Similarity of TADs from pairs of TAD-calling methods (e.g. GRiNCH vs. TopDom), measured by Rand index. The higher the number, the higher the similarity. **C.** Similarity of TADs from pairs of TAD-calling methods measured by Mutual information. Note: 3DNetMod outputted overlapping TADs, even when run under non-hierarchical settings; it was excluded from this analysis because of the requirement of distinct within-TAD measurements.

Figure 4

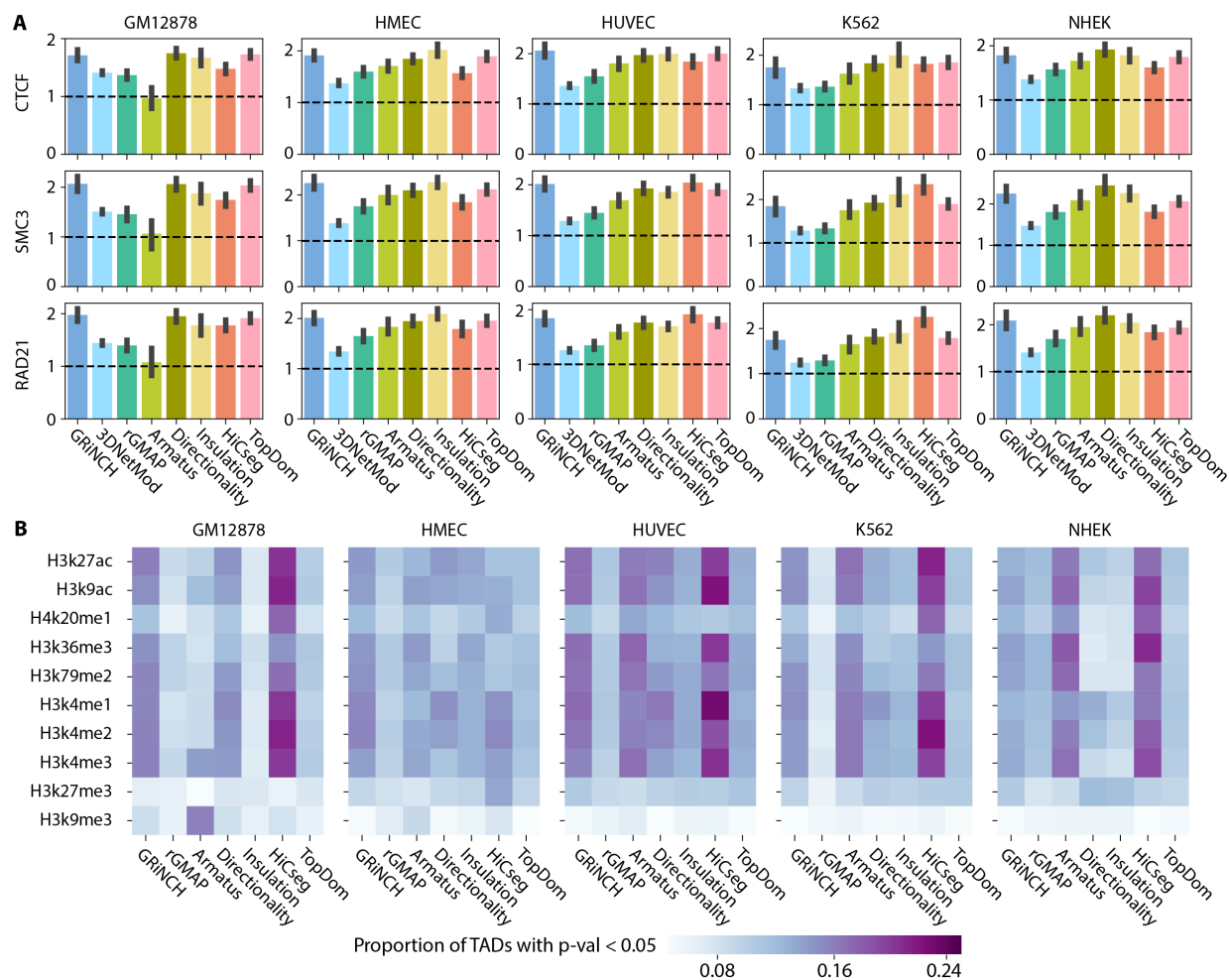


Figure 4: Evaluating the quality of TADs from different TAD-calling methods using enrichment of boundary elements and regulatory signals. **A.** Fold enrichment of binding signals of architectural protein in TAD boundaries. Shown are the mean fold enrichment of CTCF ChIP-seq peaks and accessible motif instances of cohesin proteins, RAD21 and SMC3, estimated across multiple chromosomes. The error bar denotes the standard deviation from the mean. **B.** Proportion of TADs with significant mean histone modification signal (i.e. empirical p-value < 0.05). The darker the entry the higher the proportion of TADs with significant histone enrichment. The average ChIP-seq signal for each histone modification mark was taken from within each TAD; the p-value of each TAD is derived from an empirical null distribution of mean signals in randomly shuffled TADs. Note: 3DNetMod outputted overlapping TADs, even when run under non-hierarchical settings; and it was excluded from this analysis as it involves TAD randomization/shuffling.

Figure 5

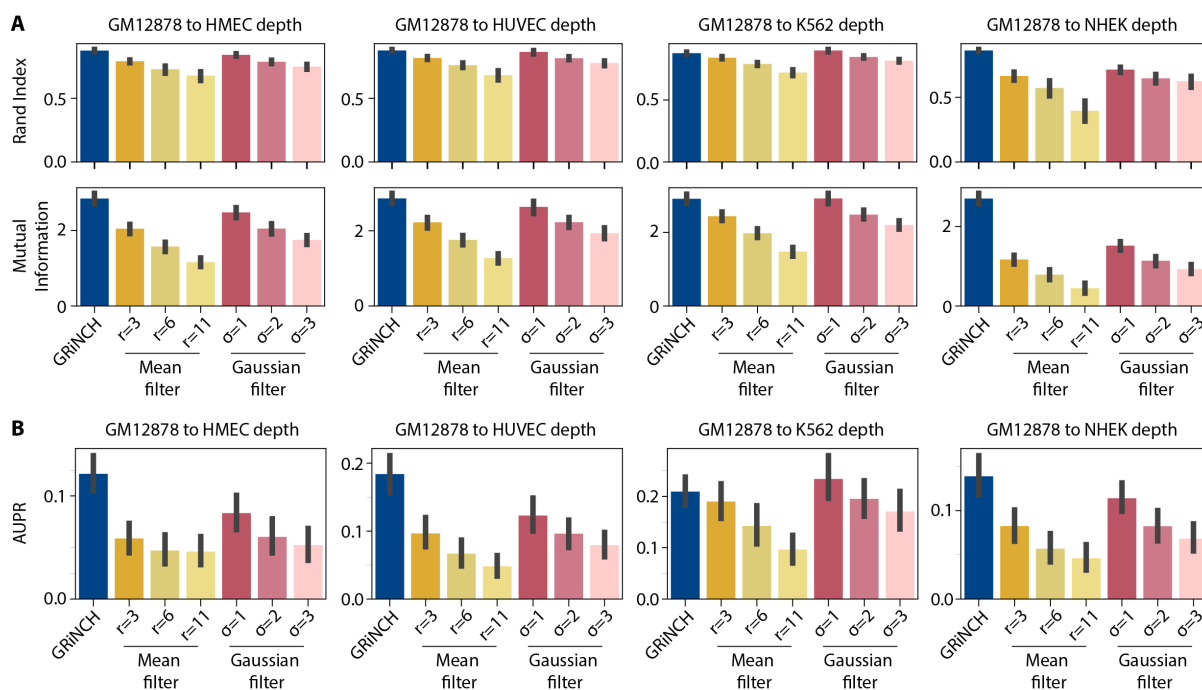


Figure 5: Evaluating the benefits of smoothing in GRiNCH. **A.** Recovery of Directionality Index TADs in downsampled then smoothed data. Shown is the mean similarity (measured by Rand Index and Mutual Information) between Directionality Index TADs from high-depth GM12878 dataset and Directionality TADs from downsampled datasets smoothed by different methods (GRiNCH, Mean Filter, Gaussian Filter). The mean is computed across chromosomes and the error bar denotes deviation from the mean. Directionality was used as a TAD-calling method independent of any of the smoothing methods, i.e., GRiNCH. **B.** Recovery of Fit-Hi-C significant interactions, as measured by the Area Under Precision-Recall curve (AUPR), with precision and recall measured for significant interactions from downsampled and smoothed datasets against the “true” interactions defined as the significant Fit-Hi-C interactions from the high-depth GM12878 dataset.

Figure 6

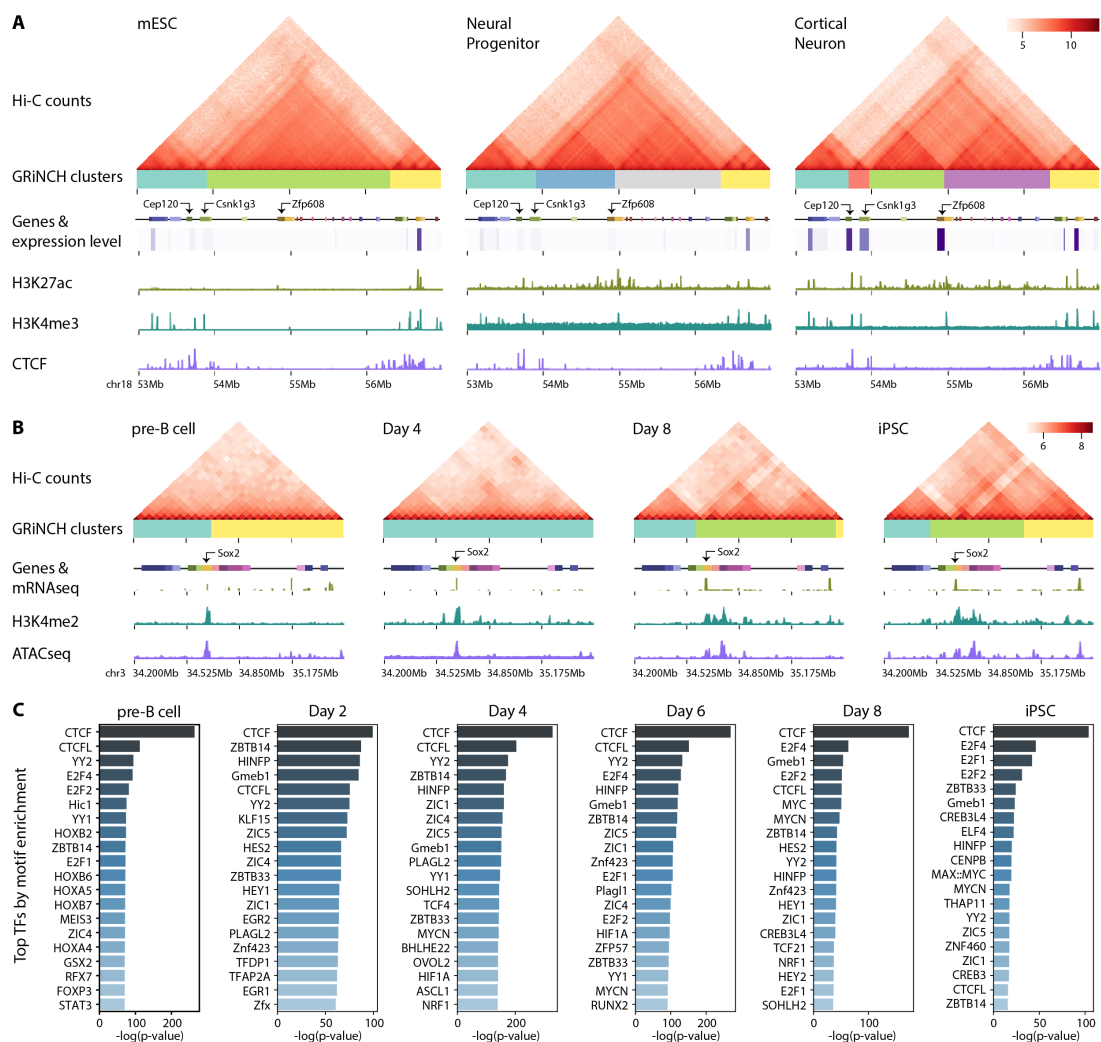


Figure 6: GRiNCH applied to Hi-C datasets along developmental time courses. A. Interaction profile near the Zfp608 gene in mouse embryonic stem cells (mESC), neural progenitors (NPC), and differentiated cortical neurons (CN). Heatmaps are of Hi-C matrices after log₂-transformation of interaction counts for better visualization. GRiNCH clusters are visualized as blocks of different colors under the heatmap of interaction counts. Genes in the nearby regions are marked by small boxes, and a heatmap of their corresponding RNA-seq levels (in TPM) is shown underneath each gene. ChIP-seq signals from H3K27ac, H3K4me3, and CTCF are shown as separate tracks. **B.** Interaction profile near the Sox2 gene in mouse pre-B cells, in day 4 and day 8 of reprogramming, and in induced pluripotent stem cell (iPSC). Heatmaps are of Hi-C matrices after log₂-transformation of interaction counts for better visualization. GRiNCH clusters are visualized as blocks of different colors under the heatmap of interaction counts. Genes in the nearby regions are marked by small boxes, and peaks of their corresponding RNA-seq levels are shown underneath each gene. ChIP-seq signals from H3K4me2 and ATAC-seq signals are shown as separate tracks. **C.** Top 20 TFs from a collection of 746 TFs ranked based on their motif enrichment in GRiNCH TAD boundaries from the mouse reprogramming time course data. The significance of their fold enrichment was calculated with the hypergeometric test and TFs were ranked by descending negative log p-value.

Figure 7

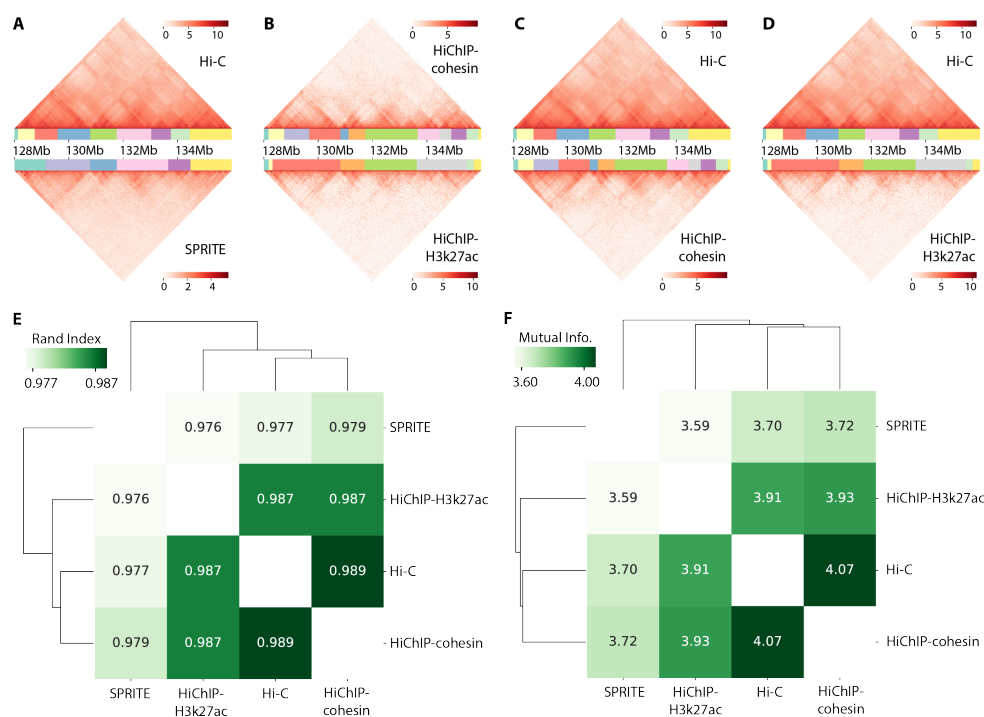


Figure 7: Applying GRiNCH to datasets from different 3D genome conformation capture technologies. Visual comparison of the interaction profile and GRiNCH TADs from a 8Mb region in chr8, GM12878 cell line. GRiNCH TADs are visualized as blocks of different colors under the heatmap of interaction counts. **A.** Hi-C vs SPRITE. The top heatmap and clusters are from Hi-C; bottom from SPRITE. **B.** HiChIP with cohesin (top) vs HiChIP with H3k27ac (bottom). **C.** Hi-C (top) vs HiChIP with cohesin (bottom). **D.** Hi-C (top) vs HiChIP with H3K37ac (bottom). For visualization purposes all interaction counts were \log_2 -transformed. **E.** Measuring the similarity of GRiNCH TADs from Hi-C and other 3D genome conformation capture platform (e.g. SPRITE, HiChIP with cohesin, or HiChIP with H3k27ac) in the same GM12878 cell line, with Rand Index. The dendrogram depicts the relative similarity between samples. **F.** Mutual-Information-based similarity of GRiNCH TADs from Hi-C and other technologies.