

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

SnapHiC: a computational pipeline to map chromatin contacts from single cell Hi-C data

Miao Yu^{1,2,*}, Armen Abnousi^{3,*}, Yanxiao Zhang², Guoqiang Li², Lindsay Lee³, Ziyin Chen¹, Rongxin Fang^{2,4}, Jia Wen⁵, Quan Sun⁵, Yun Li⁵, Bing Ren^{2,6,#} and Ming Hu^{3,#}

1. School of Life Sciences, Fudan University, Shanghai, China.

2. Ludwig Institute for Cancer Research, La Jolla, CA, USA.

3. Department of Quantitative Health Sciences, Lerner Research Institute, Cleveland Clinic Foundation, Cleveland, OH, USA.

4. Howard Hughes Medical Institute, Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA, USA.

5. Department of Genetics, University of North Carolina, Chapel Hill, NC, USA.

6. Center for Epigenomics, Department of Cellular and Molecular Medicine, University of California, San Diego, La Jolla, CA, USA.

*Contributed equally.

#Correspondence to: biren@health.ucsd.edu and hum@ccf.org.

24 **Abstract**

25 Single cell Hi-C (scHi-C) analysis has been increasingly used to map the chromatin architecture
26 in diverse tissue contexts, but computational tools to define chromatin contacts at high resolution
27 from scHi-C data are still lacking. Here, we describe SnapHiC, a method that can identify
28 chromatin loops at high resolution and accuracy from scHi-C data. We benchmark SnapHiC
29 against HiCCUPS, a common tool for mapping chromatin contacts in bulk Hi-C data, using scHi-
30 C data from 742 mouse embryonic stem cells. We further demonstrate its utility by analyzing
31 single-nucleus methyl-3C-seq data from 2,869 human prefrontal cortical cells. We uncover cell-
32 type-specific chromatin loops and predict putative target genes for non-coding sequence variants
33 associated with neuropsychiatric disorders. Our results suggest that SnapHiC could facilitate the
34 analysis of cell-type-specific chromatin architecture and gene regulatory programs in complex
35 tissues.

36

37 **Main text**

38 Transcriptional regulatory elements communicate with each other dynamically in the 3D nuclear
39 space to direct cell-type-specific gene expression during development¹⁻³. Understanding the
40 transcriptional regulatory programs requires a high resolution view of the 3D chromatin
41 architecture in the cell. Technologies have been developed to map chromatin architecture in
42 single cells to explore the heterogeneity of chromatin organization in complex tissues⁴⁻¹³. However,
43 it is still challenging to identify chromatin loops at the necessary resolution to delineate spatial
44 proximity between transcriptional regulatory elements due to the extreme sparsity of the single
45 cell chromatin contact matrix. The current strategy to identify chromatin loops from aggregated
46 single cell Hi-C data from the same cell type with existing loop calling methods¹⁴⁻¹⁸ requires a
47 large number of cells (>500-1,000), which is both cost prohibitive and impractical for the rare cell
48 types in a complex tissue. Simulation studies¹⁹ showed that the sensitivity of existing loop calling
49 methods decays exponentially with the decrease in the number of contacts. Here, we report single
50 nucleus analysis pipeline for Hi-C (SnapHiC), a new computational framework that fully exploits
51 the power of single cell Hi-C (scHi-C) data to identify chromatin loops at high resolution and
52 accuracy.

53
54 SnapHiC identifies chromatin loops at 10-kilobase (Kb) resolution from scHi-C data by maximizing
55 the usage of information from each single cell (**Fig. 1a** and **Methods**). Specifically, SnapHiC first
56 imputes chromatin contact probability between all intra-chromosomal bin pairs with the random
57 walk with restart (RWR) algorithm²⁰ in each individual cell. Next, it converts the imputed contact
58 probability into the normalized contact probability stratified based on linear genomic distances.
59 SnapHiC then applies the paired *t*-test using all cells to identify loop candidates (see details in
60 **Methods**). To remove false positives, SnapHiC considers a bin pair as a loop candidate only
61 when it has significantly higher normalized contact probability than expected by chance based on
62 both the global background and the local background. Finally, SnapHiC groups the loop
63 candidates into discrete clusters using the Rodriguez and Laio's algorithm²¹, and identifies the
64 summit(s) within each cluster.

65
66 To benchmark the performance of SnapHiC against a commonly used method, HiCCUPS¹⁴
67 designed for bulk Hi-C data analysis, we applied it to the published scHi-C data⁵ generated from
68 mouse embryonic stem (mES) cells. We sub-sampled 10, 25, 50, 75, 100, 200, 300, 400, 500,
69 600, 700 and 742 cells from this dataset, and determined the intra-chromosomal loops at 10Kb
70 resolution from 100Kb to 1Mb range. For each sub-sampling, we also pooled the scHi-C data and

71 identified chromatin loops at 10Kb resolution within the same distance range using HiCCUPS.
72 For each sub-sampling dataset, SnapHiC found more chromatin loops than HiCCUPS, suggesting
73 that SnapHiC has a much higher sensitivity than HiCCUPS (**Fig. 1b** and **Supplementary Table**
74 **1-3**). Even from 75 cells, SnapHiC identified 1,219 loops, whereas HiCCUPS found only 2 loops.
75 Additionally, HiCCUPS-identified loops tended to be a subset of SnapHiC-identified loops. For
76 example, SnapHiC and HiCCUPS identified 15,896 and 559 loops from 742 cells, respectively,
77 and 511 (91.4%) of HiCCUPS-identified loops are re-captured by SnapHiC (**Supplementary**
78 **Table 1**). Moreover, SnapHiC achieves higher reproducibility than HiCCUPS for loop calling
79 between replicates (for each replicate with 371 cells, 50.8% vs. 38.7%, paired *t*-test *p*-value =
80 7.86e-8, see details in **Methods**).

81
82 We used the F1 score, the harmonic mean of the precision and recall, to evaluate the overall
83 performance of each method (see details in **Methods**). To calculate the F1 score, we combined
84 long-range chromatin interactions identified by HiCCUPS from bulk *in situ* Hi-C data²², with
85 interactions identified by MAPS from H3K4me3 PLAC-seq data²³, cohesin²⁴ and H3K27ac HiChIP
86 data²⁵, all from mES cells as a reference loop list (**Supplementary Table 4**). At each sub-
87 sampling of scHi-C data, SnapHiC consistently attained a greater F1 score than HiCCUPS (**Fig.**
88 **1c**, **Supplementary Fig. 1**). The reliability of SnapHiC-identified loops can be further supported
89 by two additional lines of evidence: 1) Significantly focal enrichment can be observed from
90 aggregate peak analysis (APA) plots of SnapHiC-identified loops from the different number of
91 cells (except for 10 cells) on aggregated scHi-C contact matrix from 742 cells (**Supplementary**
92 **Fig. 2**); 2) For the SnapHiC-identified loops that have CTCF binding on both ends, there is a clear
93 preference in convergent orientation – ranging from 63.6% to 78.7% when at least 50 cells are
94 used for loop calling (**Supplementary Table 5**), as predicted by the loop extrusion model^{14,26}.

95
96 The advantage of SnapHiC is more obvious when the number of cells profiled is limited. As
97 illustrated in **Fig. 1d** (see also **Supplementary Fig. 3**), SnapHiC detected previously verified long-
98 range interactions at *Sox2*, *Wnt6*, and *Mtnr1a* loci^{27,28} with as few as 75 or 100 cells, whereas
99 HiCCUPS required at least 300-600 cells to detect the same loops. Taken together, the above
100 results suggest that SnapHiC allows for the identification of chromatin loops from a small number
101 of cells with high sensitivity and accuracy, underlining its potential utility in scHi-C data generated
102 from complex tissues.

103

104 To demonstrate the utility of SnapHiC for analysis of scHi-C data from complex tissues, we applied
105 SnapHiC to the published single-nucleus methyl-3C-seq (sn-m3C-seq) data¹³ from human
106 prefrontal cortex, which simultaneously profiled DNA methylome and chromatin organization from
107 the same cells. In this study, 14 major cell types were identified using CG and non-CG methylation.
108 We applied SnapHiC to each of the 14 cell clusters and identified 817 ~ 27,379 loops at 10Kb
109 resolution (**Fig. 2a** and **Supplementary Table 6**). Consistent with our observation on mES cells,
110 SnapHiC identified more loops than HiCCUPS for all cell clusters, and more than 78% of
111 HiCCUPS-identified loops are captured by SnapHiC (**Supplementary Table 7-8**). Except for
112 oligodendrocytes, which have >1,000 cells, SnapHiC found ~4-70 folds more loops than
113 HiCCUPS in other 13 cell types. We also calculated the F1 scores of SnapHiC- and HiCCUPS-
114 identified chromatin loops in oligodendrocytes, microglia, and eight neuronal subtypes, and
115 benchmarked against promoter-centered chromatin contacts previously identified from H3K4me3
116 PLAC-seq analysis of purified oligodendrocytes, microglia, astrocytes and neurons
117 (**Supplementary Table 9**)²⁹. Again, SnapHiC achieved much greater F1 scores than HiCCUPS
118 in each cell cluster (**Fig. 2b** and **Supplementary Fig. 4**).

119
120 The accuracy and sensitivity of SnapHiC are further supported by several lines of evidence. First,
121 APA analysis confirms that SnapHiC-identified loops show significant enrichment of contacts
122 compared to their local background on the aggregated contact matrix from cells in the
123 corresponding cluster (**Supplementary Fig. 5**). Next, SnapHiC-identified loops correlate with cell-
124 type-specific chromatin accessibility, histone acetylation, and gene expression. For this analysis,
125 we focused on four distinct cell types, astrocytes, L2/3 excitatory neurons, oligodendrocytes and
126 microglia, in which ATAC-seq, H3K27ac ChIP-seq and RNA-seq data are available^{29,30}. To
127 minimize the effect of cell number variation between different cell types, we randomly selected
128 the same number of cells (N=261) from astrocytes, oligodendrocytes and microglia to match the
129 number of cells available from L2/3 excitatory neurons, and applied SnapHiC to identify loops
130 from these sub-sampled data (**Supplementary Table 10**). We found that most chromatin loops
131 are cell-type-specific (**Supplementary Table 11**, see details in **Methods**). Further analysis
132 showed that the anchors of cell-type-specific loops show significantly higher ATAC-seq and
133 H3K27ac ChIP-seq signals in the matched cell type compared to those in the other three cell
134 types (**Fig. 2c**). In addition, we found 407, 616, 860 and 1,002 genes whose promoters link to
135 astrocyte-, microglia-, oligodendrocyte- and L2/3 excitatory neurons-specific loops, respectively
136 (**Supplementary Table 12**). These genes show significantly higher expression levels in the
137 matched cell type than those in the other three cell types (**Fig. 2c**) and are associated with gene

138 ontology terms³¹ related to cell-type-specific biological processes (**Fig. 2d**). Taken together, our
139 results suggest that SnapHiC can detect chromatin contacts reliably from single cell Hi-C data in
140 complex tissues.

141
142 How sequence variations determine the phenotypic traits and propensity to human diseases is
143 one of the fundamental questions in biology³². It is generally believed that many disease-
144 associated non-coding variants contribute to disease etiology by perturbing the transcriptional
145 regulatory sequences and affecting target gene expression³³⁻³⁵. The current catalogs of genes
146 and candidate regulatory sequences in the human genome³³⁻³⁷ still lack the information about the
147 target genes of annotated candidate *cis*-regulatory elements, making it a challenge to interpret
148 the biological roles of non-coding risk variants. We used SnapHiC-identified loops in the four brain
149 cell types (astrocytes, microglia, oligodendrocytes and L2/3 excitatory neurons) to assign
150 candidate target genes to non-coding GWAS SNPs. We first collected 30,262 genome-wide
151 significant (p -value $<5e-8$) non-coding GWAS SNP-trait associations from seven neuropsychiatric
152 disorders and traits, including Alzheimer's diseases³⁸ (AD), attention deficit hyperactivity
153 disorder³⁹ (ADHD), autism spectrum disorder⁴⁰ (ASD), bipolar disorder⁴¹ (BIP), intelligence
154 quotient⁴² (IQ), major depressive disorder⁴³ (MDD) and schizophrenia⁴⁴ (SCZ), resulting in a total
155 of 28,099 unique GWAS SNPs (**Supplementary Table 13**). We then focused on 3,639 SNP-
156 disease associations (3,471 unique GWAS SNPs), where the corresponding SNPs reside within
157 active enhancers of astrocytes, neurons, microglia or oligodendrocytes defined in the previous
158 study²⁹ (**Supplementary Table 13**). Using SnapHiC loops from the matching cell types (L2/3
159 excitatory neurons to represent neurons, all four cell types with 261 cells), we found 788 SNP-
160 disease-loop-gene linkages, connecting 445 SNP-disease associations (416 unique GWAS SNPs)
161 to 189 genes via 175 loops (**Supplementary Table 14**). Notably, such a list of GWAS SNP-
162 interacting genes includes several known disease risk genes, including *APOE* (AD), *GRIN2A* (IQ),
163 *INPP5D* (AD), *RAB27B* (MDD), *SORL1* (AD), *THRB* (IQ), and *ZNF184* (SCZ and MDD). **Fig. 2e**
164 shows an illustrative example of gene *APOE*, which is specifically expressed in astrocyte. We
165 found two astrocyte-specific chromatin loops, connecting the TSS of *APOE* to two active
166 enhancers in astrocyte, ~150Kb and ~200Kb downstream, respectively. These two enhancers
167 also contain two AD-associated GWAS SNPs, rs112481437 and rs138137383. Our data suggest
168 that *APOE* is the putative target gene of these two GWAS SNPs only in astrocytes.

169
170 In summary, we describe SnapHiC, a novel method customized for sparse single cell Hi-C
171 datasets to identify chromatin loops at high resolution and accuracy. Re-analysis of published

172 single cell Hi-C data from mES cells demonstrate that SnapHiC greatly boosts the statistical
173 power in loop detection. Application of SnapHiC to sn-m3C-seq data from human prefrontal
174 cortical cells reveals cell-type-specific loops, which can be used to predict putative target genes
175 of non-coding GWAS SNPs. SnapHiC has the potential to facilitate the study of cell-type-specific
176 chromatin spatial organization in complex tissues.

177

178 **Code availability**

179 SnapHiC software package with a detailed user tutorial and sample input and output files can be
180 found at: <https://github.com/HuMingLab/SnapHiC>.

181

182 **Acknowledgements**

183 We thank 4D Nucleome consortium investigators for comments and suggestions on the early
184 version of this work. This study was funded by U54DK107977, UM1HG011585 (to B.R. and M.H.),
185 and U01DA052713, R01GM105785 and P50HD103573 (to Y.L.).

186

187 **Author Contributions**

188 This study was conceived and designed by M.H. and B.R.; Data analysis was performed by M.H.,
189 M.Y., A.A., Y.Z., G.L., L.L., Z.C., R.F., J.W., Q.S. and Y.L.; SnapHiC software package was
190 developed by A.A. and M.H.; Manuscript was written by M.H., M.Y. and B.R. with input from all
191 authors.

192

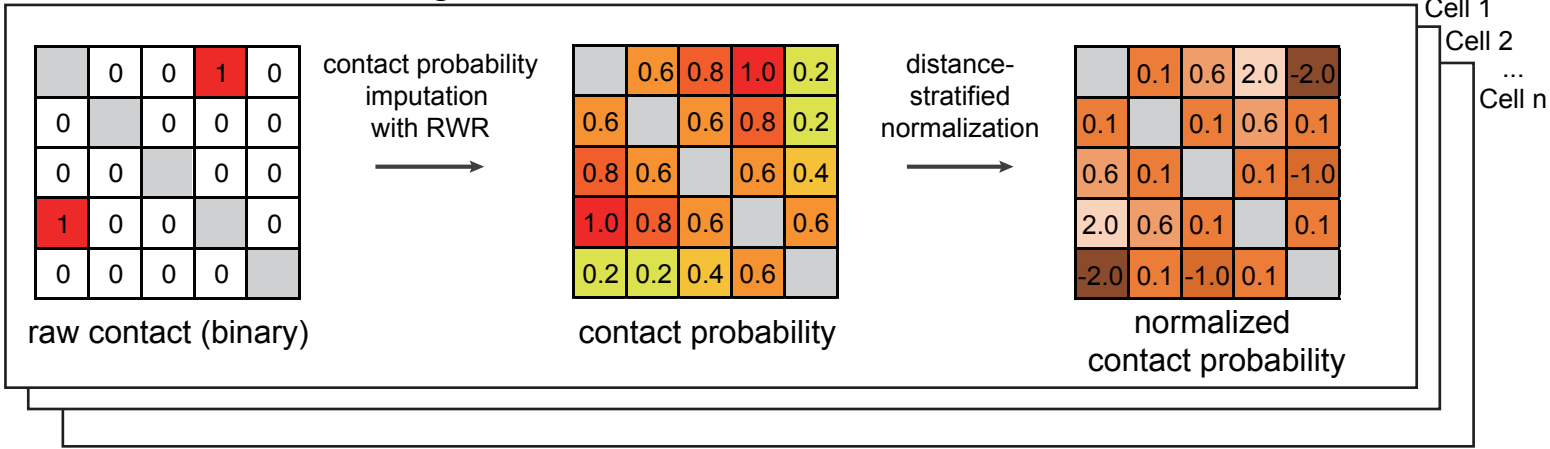
193 **Competing interests**

194 B.R. is co-founder and shareholder of Arima Genomics and Epigenome Technologies. The other
195 authors declare that they have no competing interests.

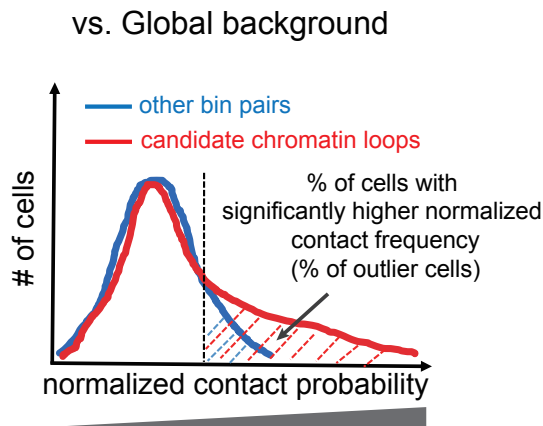
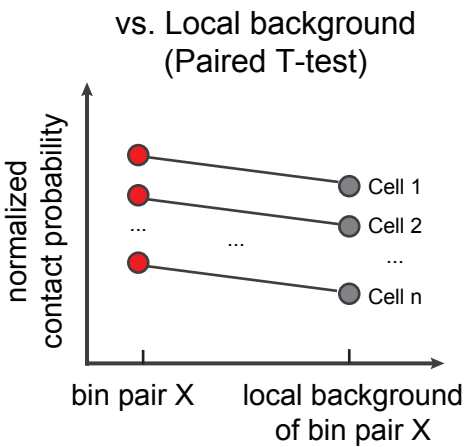
196

Figure 1

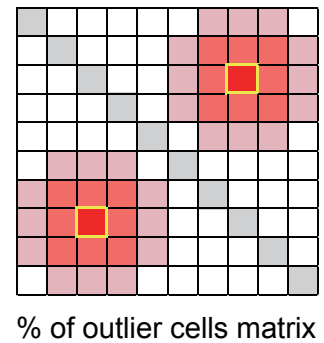
a 1. Convert matrix on each single cell



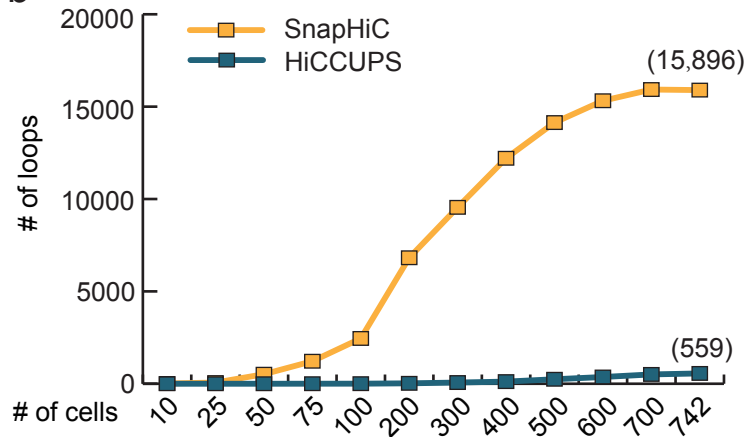
2. Identify candidate chromatin loops



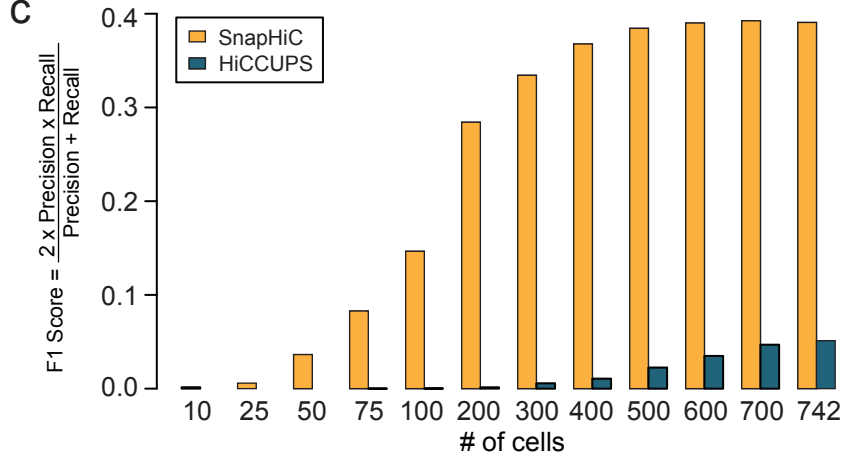
3. Merge loops, find summit(s) and visualization



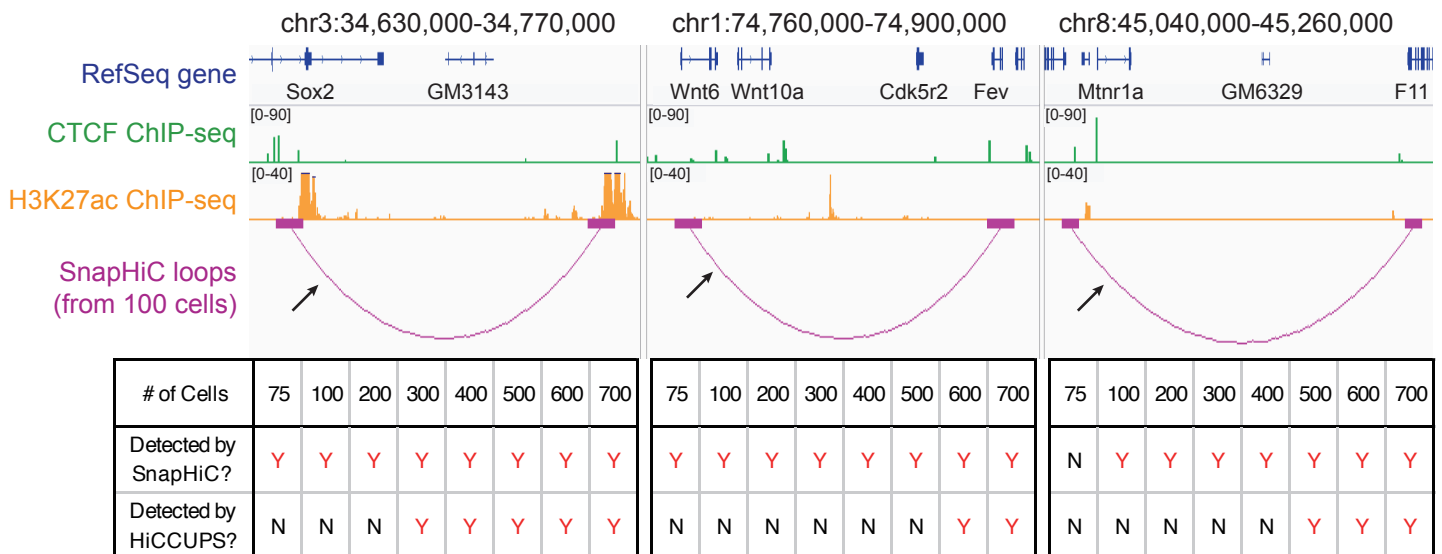
b



c



d



197 **Figure 1. SnapHiC reveals chromatin loops at high resolution and accuracy. (a)** Overview
198 of SnapHiC workflow. The first step of SnapHiC is to convert the binary contact matrix to
199 normalized contact frequency for each individual cell. Next, SnapHiC applies the paired *t*-test to
200 identify candidate chromatin loops by comparing the normalized contact frequency of any given
201 bin pair with its local and global background. Finally, SnapHiC merges nearby candidate loops
202 into clusters and identifies the summit(s). Due to the sparsity of the raw count matrix of scHi-C
203 data, the SnapHiC-identified loops can be visualized by the percentage of the outlier cells matrix.
204 **(b)** The number of chromatin loops at 10Kb resolution identified by SnapHiC and HiCCUPS from
205 different numbers of mES cells. **(c)** F1 score (the harmonic mean of the precision and recall) of
206 SnapHiC- and HiCCUPS-identified loops from different numbers of mES cells. **(d)** (Top)
207 Chromatin loops around *Sox2* (left), *Wnt6* (middle), and *Mtnr1a* (right) gene identified from 100
208 mES cells using SnapHiC at 10Kb resolution. The black arrow points to the interaction verified in
209 the previous publications^{27,28} with CRISPR/Cas9 deletion or 3C-qPCR. (Bottom) Comparison of
210 the performance of SnapHiC and HiCCUPS (applied on aggregated scHi-C data) from the
211 different number of mES cells at these three regions. If the previously verified interaction (black
212 arrow) is recaptured, it is labeled as “Y”; otherwise, it is labeled as “N”.

213

214

215

216

217

218

219

220

221

222

223

224

225

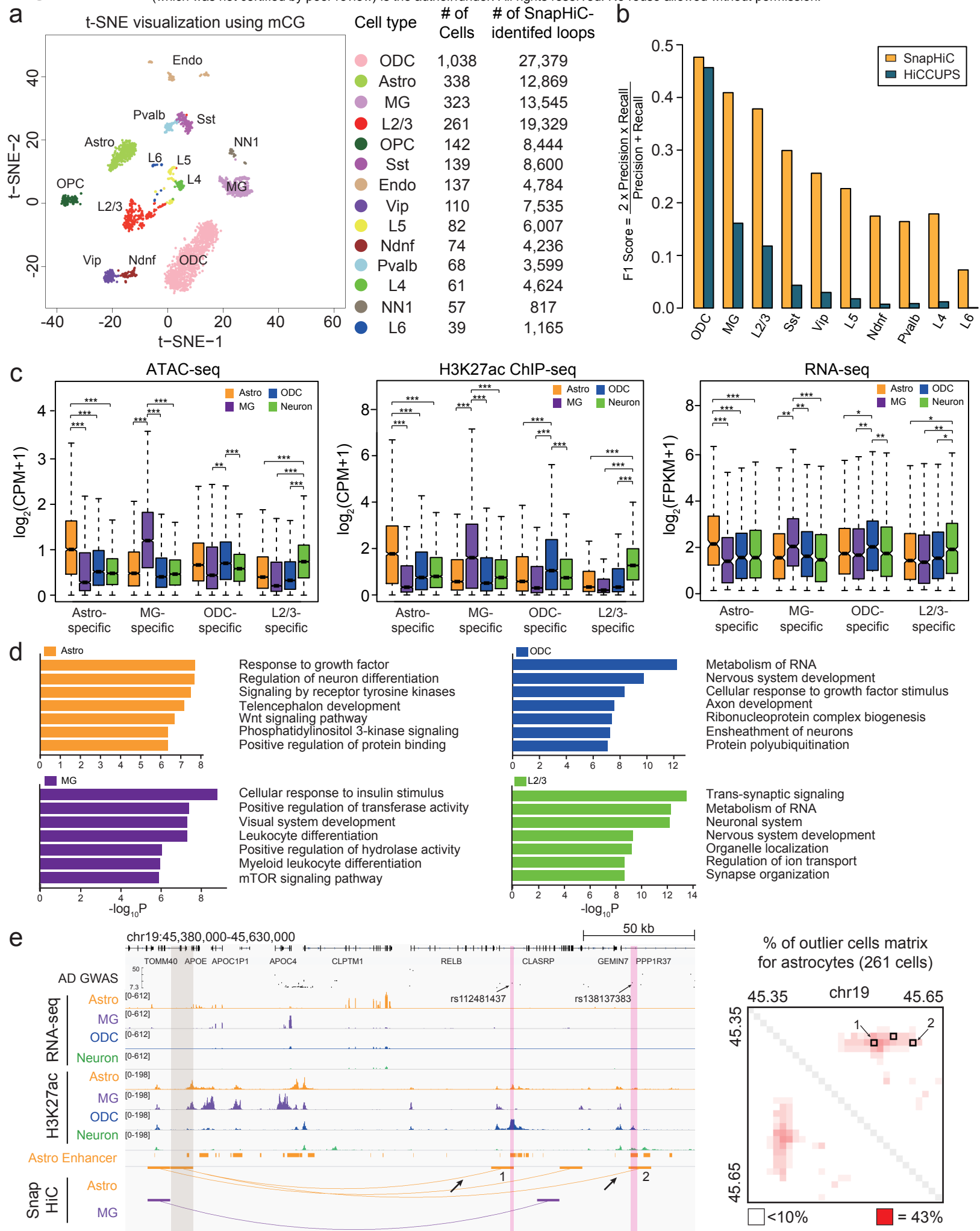
226

227

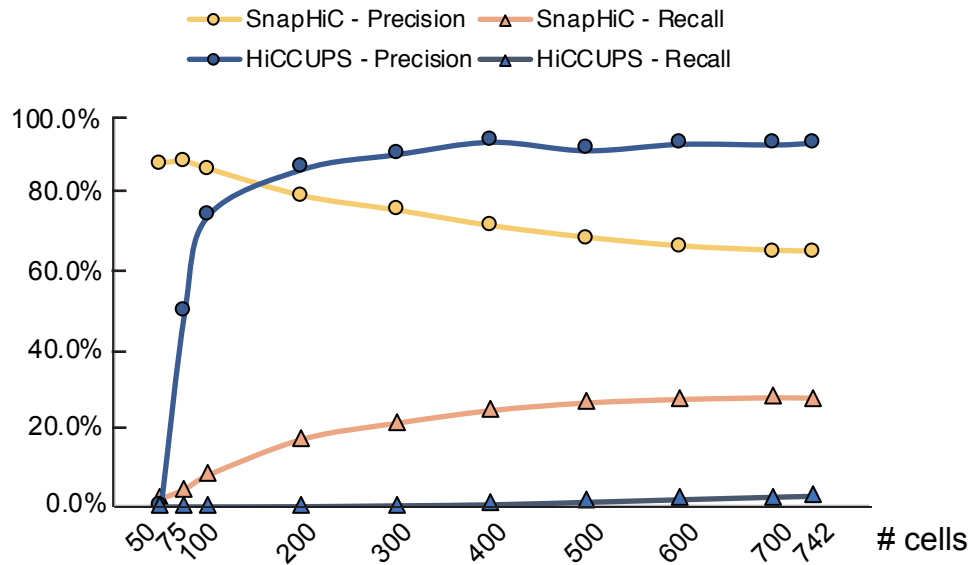
228

229

Figure 2

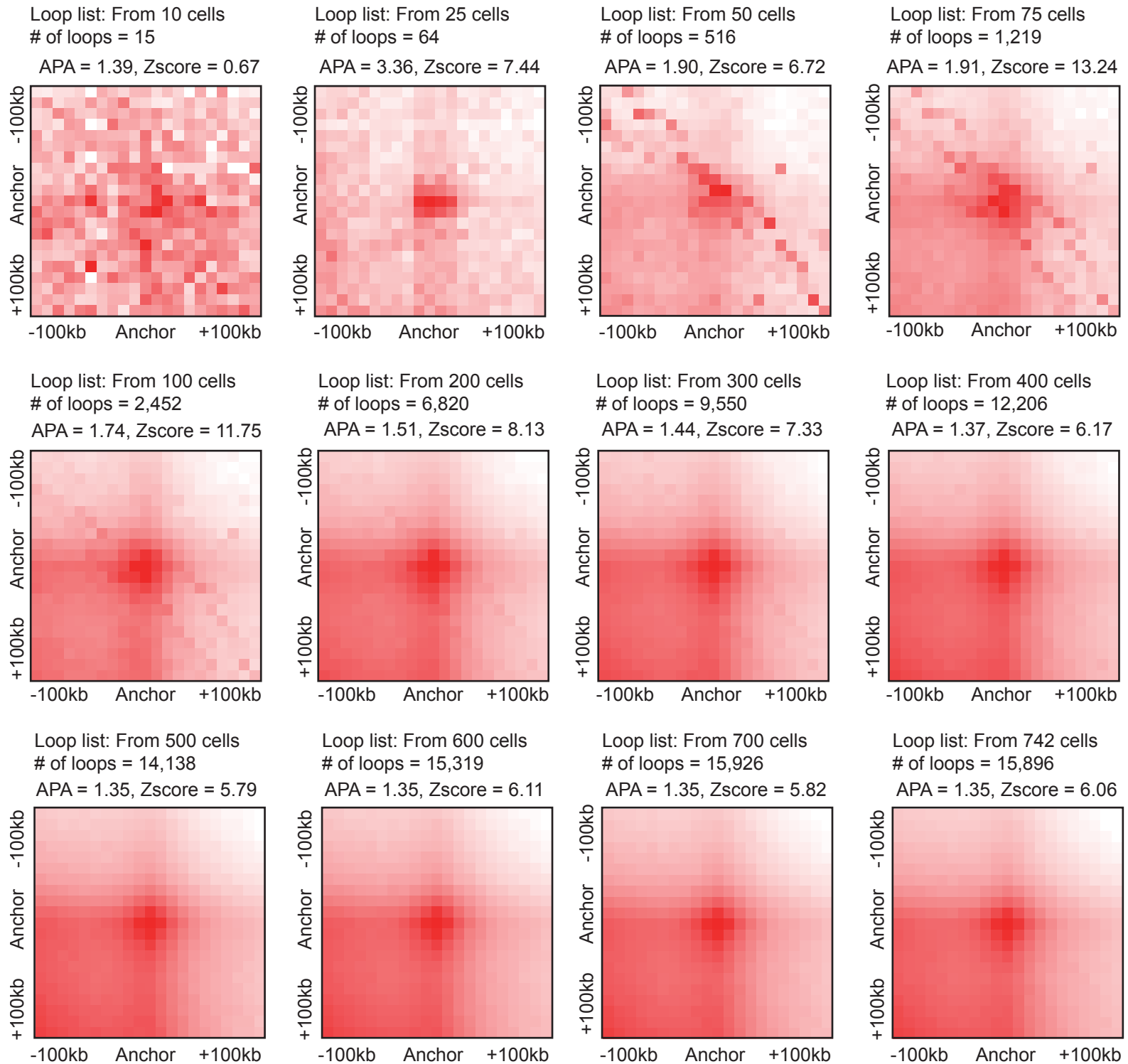


230 **Figure 2. Application of SnapHiC to sn-m3C-seq data from human prefrontal cortex**
231 **uncovered chromatin loops in diverse brain cell types. (a)** (Left) t-SNE visualization of 14
232 major cell types identified in human prefrontal cortex in Lee et al. study¹³ using CG methylation of
233 non-overlapping 100Kb genomic bins. ODC: oligodendrocyte. Astro: astrocyte. MG: microglia.
234 OPC: oligodendrocyte progenitor cell. Endo: endothelial cell. L2/3, L4, L5 and L6: excitatory
235 neuron subtypes located in different cortical layers. Pvalb and Sst: medial ganglionic eminence-
236 derived inhibitory subtypes. Ndnf and Vip: CGE-derived inhibitory subtypes. NN1: non-neuronal
237 cell type 1. (Right) The number of cells and SnapHiC-identified loops in each of the 14 cell types.
238 **(b)** F1 score (the harmonic mean of the precision and recall) of SnapHiC- and HiCCUPS-identified
239 loops for oligodendrocytes (ODC), microglia (MG) and eight neuronal subtypes. **(c)** Boxplot of
240 ATAC-seq $\log_2(\text{CPM}+1)$ value (left), H3K27ac ChIP-seq $\log_2(\text{CPM}+1)$ value (middle) and RNA-
241 seq $\log_2(\text{FPKM}+1)$ value (right) in astrocyte, microglia, oligodendrocytes and neurons at the
242 anchors of Astro-specific, MG-specific, ODC-specific, L2/3-specific SnapHiC loops summarized
243 in **Supplementary Table 11**. *** $p < 2.2e-16$; ** $p < 1e-10$; * $p < 1e-7$ by the paired Wilcoxon signed-
244 rank test. **(d)** Top seven enriched gene ontology (GO) terms of genes associated with cell-type-
245 specific SnapHiC loops. **(e)** (Left) SnapHiC-identified loops from astrocyte and microglia around
246 gene *APOE*. There is no loop identified in this genomic region from oligodendrocytes or L2/3
247 excitatory neurons, so no corresponding tracks are shown. Two astrocyte-specific loops linking
248 the *APOE* promoter (highlighted in grey) and the active enhancers in astrocyte (highlighted in
249 pink) containing two AD-associated GWAS SNPs are marked by black arrows. Only *APOE* TSS-
250 distal AD-associated GWAS SNPs are shown in the figures (residing in the region chr19:
251 45,440,000-45,630,000). (Right) Matrix of the percentage of cells with significantly higher
252 normalized contact frequency (percentage of outlier cells with normalized contact frequency > 1.96)
253 for 261 astrocytes. The SnapHiC-identified loops from astrocyte are marked by black squares.
254
255
256
257
258
259
260
261
262
263

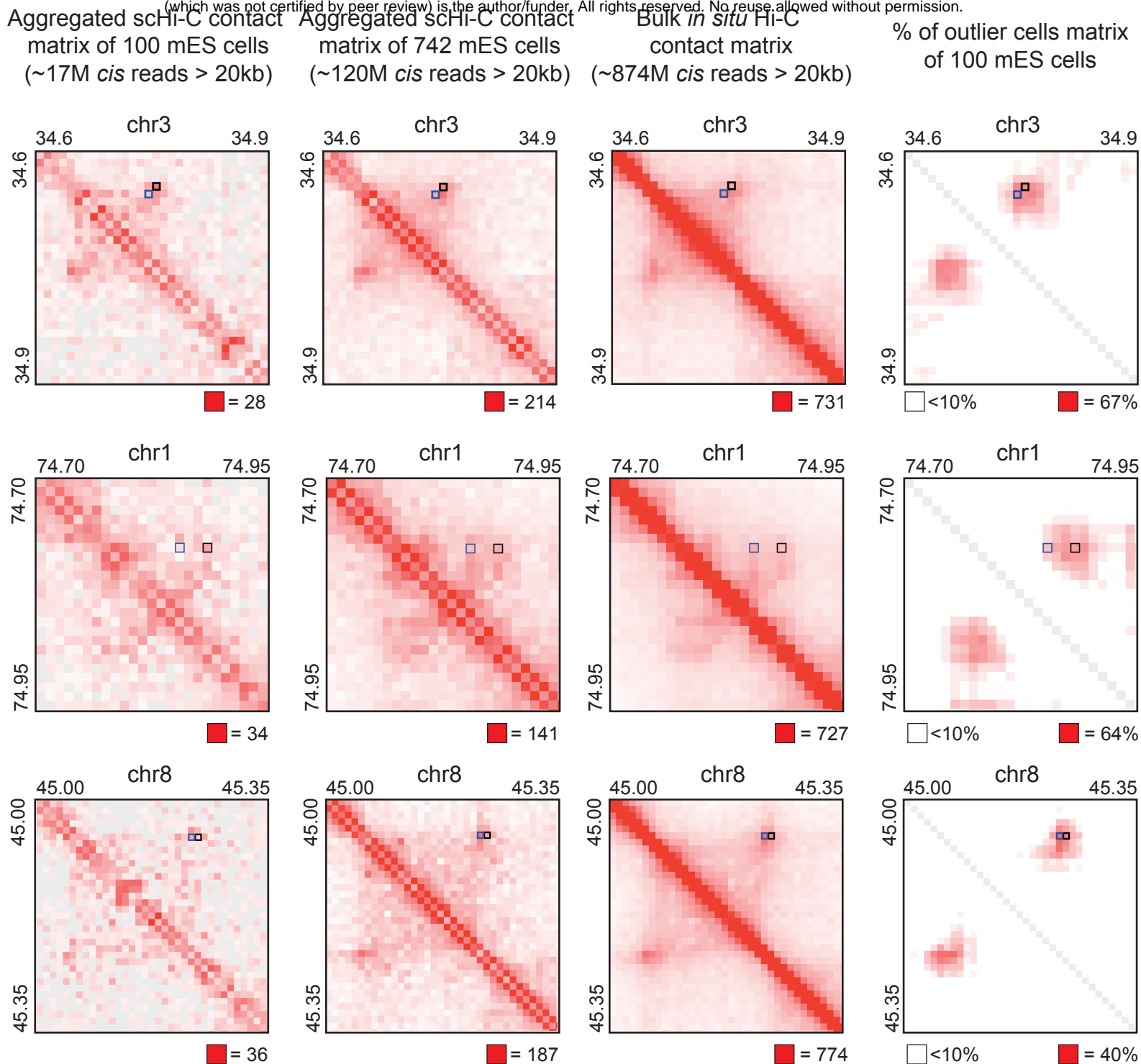


Supplementary Figure 1. Comparison of the precision and recall values of SnapHiC- and HiCCUPS-identified loops from mES cells. The precision and recall values are calculated for the loops identified by SnapHiC and HiCCUPS from different numbers of mES cells. These values are also used to calculate the F1 score in **Fig. 1c**.

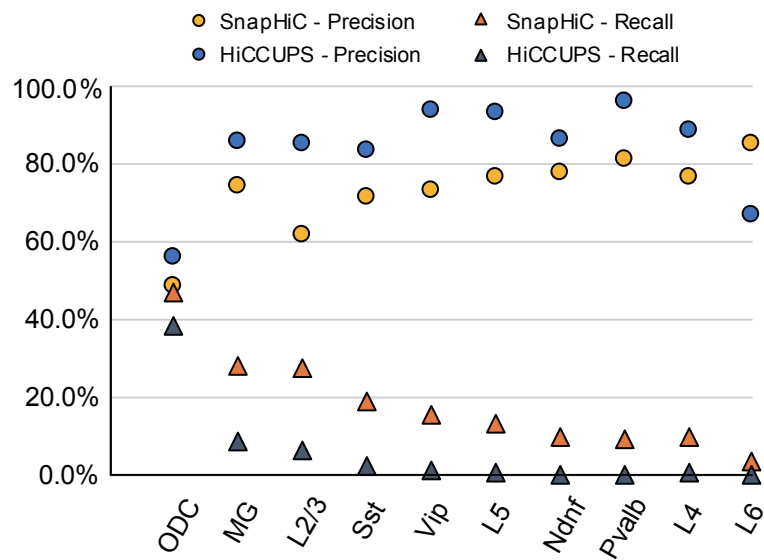
Map: Aggregated scHi-C contact matrix from 742 mES cells



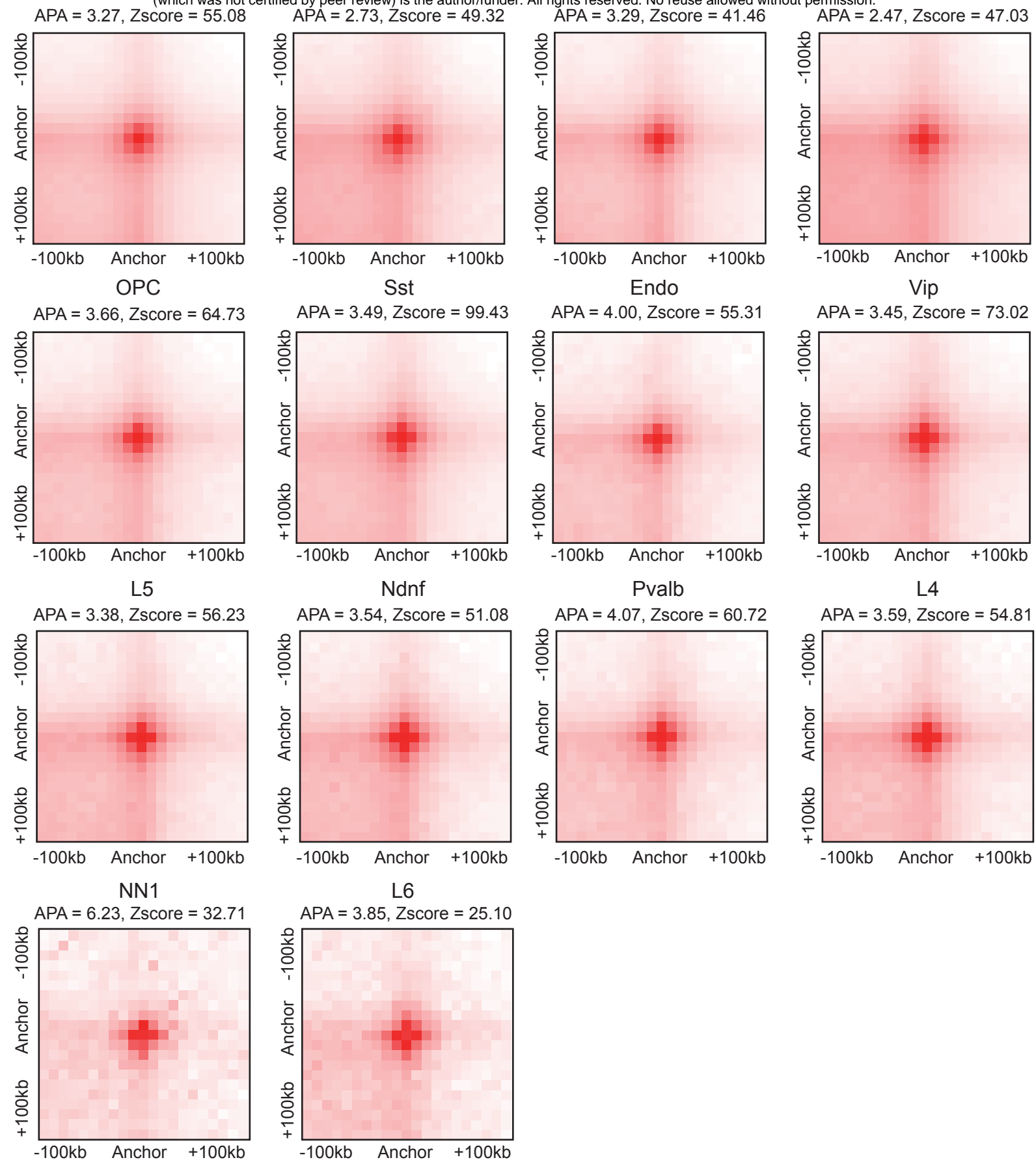
Supplementary Figure 2. SnapHiC-identified loops from different sub-sampling of mES cells show significant enrichment over their local background. Aggregate peak analysis (APA) of SnapHiC-identified loops from different sub-sampling of mES cells examined on aggregated scHi-C contact matrix of 742 cells.



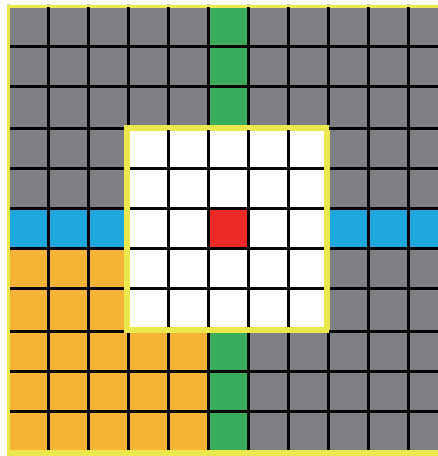
Supplementary Figure 3. Visualization of selected SnapHiC-identified loops. From left to right: aggregated scHi-C contact matrix of 100 mES cells, aggregated scHi-C contact matrix of 742 mES cells, bulk *in situ* Hi-C contact matrix from mES cells (replicate 1 from Bonev et al. study²²) and % of outlier cells matrix of 100 mES cells at 10Kb resolution; from top to bottom: *Sox2* locus, *Wnt6* locus, and *Mtnr1a* locus. Black squares represent the SnapHiC-identified loops from 100 mES cells, which are shown in Fig. 1d as purple arcs. For comparison, the HiCCUPS-identified loops from the deepest available bulk *in situ* Hi-C data of mES cells (combining all four replicates from Bonev et al. study²²) are marked as blue squares.



Supplementary Figure 4. Comparison of the precision and recall values of SnapHiC- and HiCCUPS-identified loops for ten cell clusters from human prefrontal cortex. The precision and recall values are calculated for the loops identified by SnapHiC and HiCCUPS for oligodendrocytes (ODC), microglia (MG), and eight neuronal subtypes. These values are also used to calculate the F1 score in **Fig. 2b**.

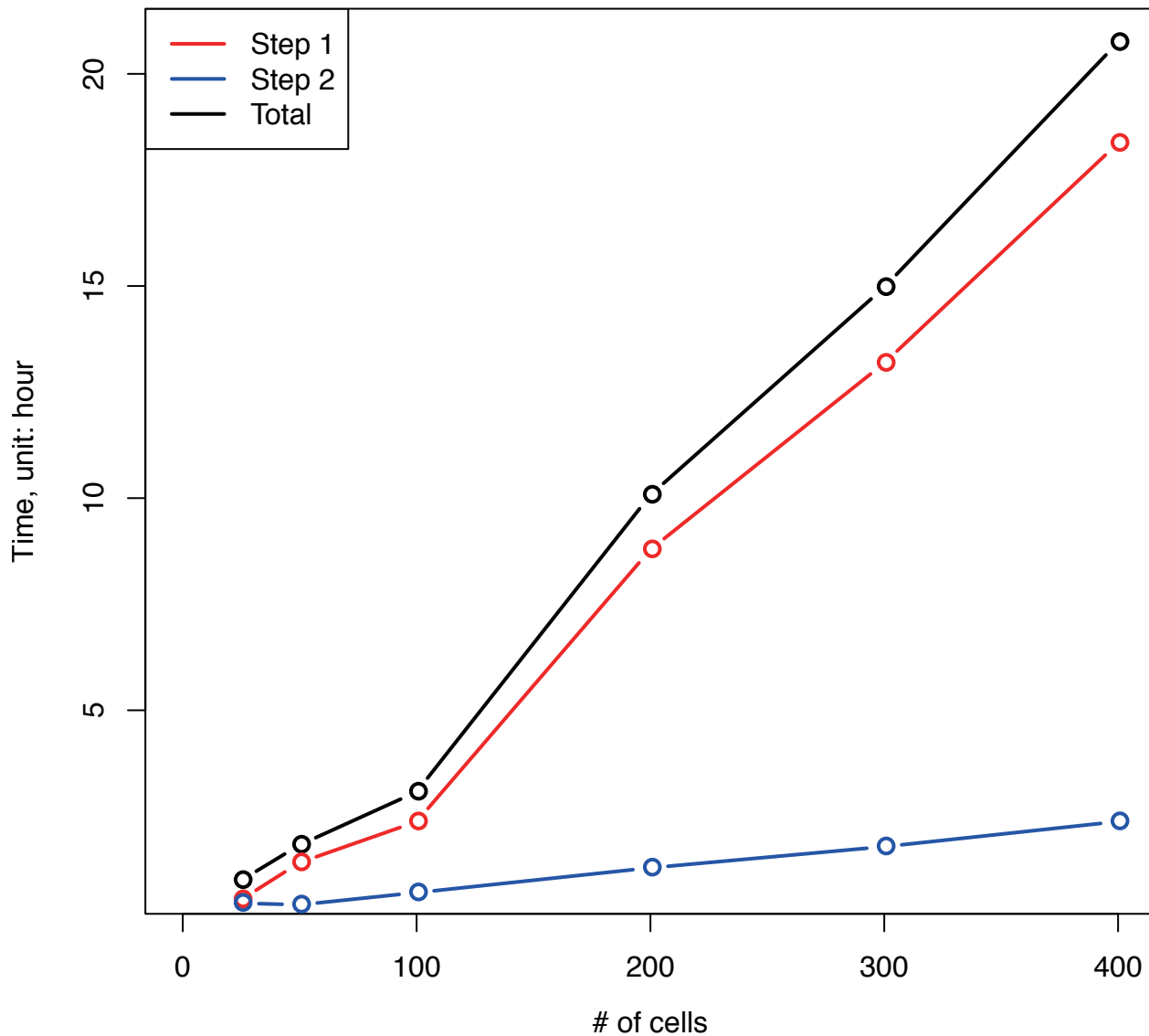


Supplementary Figure 5. SnapHiC-identified loops from each of the 14 cell clusters identified from sn-m3C-seq data of the human prefrontal cortex show significant enrichment over their local background. Aggregate peak analysis (APA) of SnapHiC-identified loops for each of the 14 cell clusters demonstrated in Fig. 2a examined on the aggregated contact matrix from the matching cell clusters.



Supplementary Figure 6. Illustration of different types of the local background used for SnapHiC loop calling. For each 10Kb bin pair of interest (red), its horizontal background, vertical background, lower left background and donut background are the blue, green, yellow and grey areas, respectively. The circle background, which is also the local neighborhood, is the union of the blue, green, yellow and grey areas.

SnapHiC running time



Supplementary Figure 7. The relationship between the number of cells and the running time of SnapHiC analysis. We tested the running time of SnapHiC on scHi-C data from 25, 50, 100, 200, 300 and 400 mES cells (10Kb resolution, searching for loops between 100Kb to 1Mb genomic distance). SnapHiC consists of two steps: (1) applying the random walk with restart (RWR) algorithm to impute contact probability within every single cell, and (2) integrating imputed contact probability matrices from all single cells to identify chromatin loops. The running time of each step and the sum of both steps against the number of cells is plotted.

264 **Supplementary Table Legends**

265

266 **Supplementary Table 1.** Summary of SnapHiC- and HiCCUPS-identified loops from mES scHi-
267 C data. Related to **Fig. 1b**.

268

269 **Supplementary Table 2.** SnapHiC-identified loops from 10, 25, 50, 75, 100, 200, 300, 400, 500,
270 600, 700 and 742 mES cells (Column D in **Supplementary Table 1**).

271

272 **Supplementary Table 3.** HiCCUPS-identified loops from 10, 25, 50, 75, 100, 200, 300, 400, 500,
273 600, 700 and 742 mES cells (after filtering, Column F in **Supplementary Table 1**).

274

275 **Supplementary Table 4.** HiCCUPS-identified loops from bulk *in situ* Hi-C data, MAPS-identified
276 significant interactions from H3K4me3 PLAC-seq, cohesin HiChIP, and H3K27ac HiChIP data,
277 which are used as the reference loop list after pooling to calculate precision, recall values and the
278 F1 score in **Fig. 1c** and **Supplementary Fig. 1**.

279

280 **Supplementary Table 5.** CTCF motif orientation analysis for SnapHiC-identified loops from
281 different numbers of mES cells.

282

283 **Supplementary Table 6.** SnapHiC-identified loops from 14 different cell clusters demonstrated
284 in **Fig. 2a** (Column D in **Supplementary Table 8**).

285

286 **Supplementary Table 7.** HiCCUPS-identified loops from 14 different cell clusters demonstrated
287 in **Fig. 2a** (after filtering, Column F in **Supplementary Table 8**).

288

289 **Supplementary Table 8.** Summary of SnapHiC- and HiCCUPS-identified loops from 14 different
290 cell clusters demonstrated in **Fig. 2a**.

291

292 **Supplementary Table 9.** MAPS-identified interaction lists of human microglia, oligodendrocytes,
293 and neurons based on the Nott. et al. study²⁹, which are used to calculate precision, recall values
294 and the F1 score in **Fig. 2b** and **Supplementary Fig. 4**.

295

296 **Supplementary Table 10.** SnapHiC-identified loops from astrocytes, microglia and
297 oligodendrocytes after sub-sampling (261 cells for each cell type).

298 **Supplementary Table 11.** Cell-type-specific SnapHiC loops identified from astrocytes, microglia,
299 oligodendrocytes, and L2/3 excitatory neurons after sub-sampling (261 cells for each cell type).

300

301 **Supplementary Table 12.** Genes whose promoter overlaps with cell-type-specific SnapHiC loops
302 identified from astrocytes, microglia, oligodendrocytes, and L2/3 excitatory neurons after sub-
303 sampling (261 cells for each cell type). Related to [Fig. 2c](#), [2d](#) and [Supplementary Table 11](#).

304

305 **Supplementary Table 13.** Non-coding GWAS SNPs associated with seven neuropsychiatric
306 disorders with p -value $< 5 \times 10^{-8}$ and the SNPs residing in the active enhancers of astrocytes,
307 microglia, oligodendrocytes or neurons defined in the previous publication²⁹.

308

309 **Supplementary Table 14.** Predicted 788 SNP-disease-loop-gene quadruplets using SnapHiC-
310 identified loops in astrocytes, microglia, oligodendrocytes and L2/3 excitatory neurons (261 cells
311 for each cell type).

312

313 **Methods**

314 **Single-cell Hi-C (scHi-C) data processing**

315 For scHi-C data from mES cells⁵, we downloaded the raw fastq files of all diploid serum cells (in
316 total 1,175 cells). We first aligned scHi-C read pairs for each single cell to mm10 genome with
317 BWA-MEM with the “-5” option to report the most 5’ end alignment as the primary alignment, and
318 the “-P” option to perform Smith-Waterman algorithm to rescue chimeric reads. We only used
319 primary alignments in the next steps. We then de-duplicated read pairs with the Picard tool to
320 keep only one read pair at the exact same position. We further applied two filtering steps to
321 remove read duplications: (1) we split each chromosome into consecutive non-overlapping 1Kb
322 bins, and only kept one contact for each 1Kb bin pair, (2) we removed 1Kb bins which contact
323 with more than 10 other 1Kb bins, since they are likely mapping artifacts. We found that the
324 number of contacts per cell for these 1,175 cells has a bimodal distribution, therefore we selected
325 the top 742 cells with >150,000 contacts per cell for downstream analysis.

326

327 **Single-nucleus methyl-3C-seq (sn-m3C-seq) data processing**

328 For sn-m3C-seq data from human prefrontal cortex, we performed data processing using
329 reference genome hg19 as described in the previous study¹³. After this processing, we also
330 applied two additional filtering steps to remove read duplications as described in the “**Single-cell**
331 **Hi-C (scHi-C) data processing**” section. Similar to scHi-C data from mES cells, we also observed
332 a bimodal distribution in the number of contacts per cell for all 4,238 cells. Again, we selected the
333 top 2,869 cells with >150,000 contacts per cell for downstream analysis. The method for clustering
334 and cell type annotation for these 2,869 cells was the same as previously described¹³.

335

336 **SnapHiC algorithm**

337 **Step A. Contact probability imputation using the random walk with restart (RWR) algorithm.**

338 We first partitioned each autosomal chromosome into consecutive non-overlapping bins at a pre-
339 specified resolution (10Kb in this study) and dichotomized contact for each 10Kb bin pair (binary
340 contact matrix with 1 indicating non-zero contact and 0 otherwise). Next, we modeled each
341 autosomal chromosome as an unweighted graph, where each 10Kb bin is one node, and each
342 non-zero contact between any two 10Kb bins is one edge. We also added edges to all adjacent
343 10Kb bins. We then implemented the random walk with restart (RWR) algorithm²⁰ with the restart
344 probability 0.05 to impute the contact probability between all intra-chromosomal 10Kb bin pairs.
345 We used the Python “NetworkX” package to construct the graph, and adopted the “linalg.solve”
346 function in the Python “SciPy” package to solve the linear equation in the RWR algorithm. In

347 addition, we distributed the analysis for different chromosomes in different cells between different
348 processors using the Python “mip4py” package to speed up the computation.

349
350 We further evaluated whether the contact probability imputed by the RWR algorithm in each single
351 cell contains systematic biases, including effective fragment size, GC content and mappability,
352 which are known systematic biases in bulk Hi-C data⁴⁵. Specifically, for each of the 742 mES
353 scHi-C profiles, we used the RWR algorithm to impute the contact probability between all intra-
354 chromosomal 10Kb bin pairs (i, j) within 1Mb genomic distance, denoted as x_{ij} . Let F_i , GC_i and
355 M_i represent the effective fragment size, GC content and mappability of the 10Kb bin i , which are
356 calculated according to our previous work⁴⁵. We define $f_{ij} = F_i * F_j$, $gc_{ij} = GC_i * GC_j$, and $m_{ij} =$
357 $M_i * M_j$, as the measure of three types of bias for each 10Kb bin pair. We then calculated the
358 Pearson Correlation Coefficient between the contact probability x_{ij} and f_{ij} , gc_{ij} and m_{ij} ,
359 respectively, for each of the 19 autosomal chromosomes in one cell. Next, we used the average
360 Pearson Correlation Coefficient (aPCC) across all chromosomes as the measurement of bias in
361 each cell. Among all 742 cells, the mean of aPCC is 0.0110, 0.0085 and -0.0016 for effective
362 fragment size, GC content and mappability, respectively. The standard deviation of aPCC is
363 0.0068, 0.0113 and 0.0029 for effective fragment size, GC content and mappability, respectively.
364 These results suggest that the systematic biases in imputed contact probabilities in scHi-C data
365 are negligible, thus normalization against effective fragment size, GC content or mappability is
366 not needed.

367
368 **Step B. Contact probability normalization based on 1D genomic distance.**

369 Since the contact probability between any two genomic loci is strongly dependent on their 1D
370 genomic distance, normalization of the imputed contact probability against 1D genomic distance
371 is needed before loop calling. To achieve this, we first removed the bin pairs residing in the first
372 50Kb or the last 50Kb of each chromosome, which often have unusually high imputed contact
373 probability due to the edge effect of the RWR algorithm. We then stratified all intra-chromosomal
374 10Kb bin pairs by their 1D genomic distance. Specifically, let x_{ij} represent the contact probability
375 between bin i and bin j . Define the set A_d as all bin pairs (i, j) with the 1D genomic distance d .
376 For simplicity, we only considered bin pairs (i, j) in the upper triangle of the contact matrix where
377 $i < j$. We removed the top 1% bin pairs in A_d with the highest contact probability, and then
378 computed the mean μ_d and the standard deviation σ_d of the contact probability using the
379 remaining bin pairs in A_d . We further calculated the normalized contact probability (i.e., Z-score),

380 defined as $z_{ij} = (x_{ij} - \mu_d)/\sigma_d$, for all bin pairs in A_d . For single cells with very few contacts, the
381 imputed contact probabilities x_{ij} at specific 1D genomic distance d are close to zero, leading to
382 very small standard deviation σ_d and numerical errors in the Z-score transformation. To avoid this
383 issue, when σ_d is less than 1e-6, we defined $z_{ij} = 0$ for all bin pairs in A_d . After the calculation
384 described above, bin pair (i, j) with higher normalized contact probability z_{ij} suggests that bin i
385 and bin j are more likely to interact with each other than the other genomic loci pairs.

386

387 **Step C. Identification of loop candidates.**

388 To minimize false positives in loop calling results, we defined a bin pair as a loop candidate only
389 if it shows higher contact probability compared to both its global and local background. Specifically,
390 we required the loop candidate to satisfy the following criteria:

391

392 (1) Its average normalized contact probability of all single cells is greater than 0 (i.e., with respect
393 to global background).

394

395 (2) More than 10% of all single cells have normalized contact probability above 1.96 at the loop
396 candidate (i.e., Z-score > 1.96, corresponding to p -value < 0.05, with respect to global background).

397

398 (3) For each 10Kb bin pair (i, j) , we defined its local neighborhood as all 10Kb bin pairs (m, n)
399 such that $30\text{Kb} \leq \max\{d(i, m), d(j, n)\} \leq 50\text{Kb}$ (**Supplementary Fig. 6**), where $d(i, m)$ is the
400 genomic distance between the center of bin i and the center of bin m . Here we did not consider
401 the bin pairs within 20Kb of bin pair (i, j) as part of its local neighborhood because they can be
402 part of the same loop cluster centered at bin pair (i, j) . We then compared the normalized contact
403 probability at bin pair (i, j) with the mean of the normalized contact probability of all 96 10Kb bin
404 pairs within its local neighborhood region, and applied the paired t -test across all single cells to
405 obtain a p -value. We further converted p -values into false discovery rates (FDRs) using the
406 Benjamin-Hochberg procedure, again stratified by 1D genomic distance. The loop candidates
407 must have FDR < 10% and t -statistics greater than 3 in the paired t -test (i.e., with respect to local
408 background).

409

410 (4) Motivated by the HiCCUPS algorithm¹⁴, we also required the loop candidate to have at least
411 33% higher average normalized contact frequency than its circle, donut and lower left background

412 and 20% higher average normalized contact frequency than its horizontal and vertical background
413 (**Supplementary Fig. 6**) (i.e., with respect to local background).

414
415 (5) Finally, we removed the loop candidates with either end having low mappability score (≤ 0.8),
416 or overlapping with the ENCODE blacklist regions
417 ([http://mitra.stanford.edu/kundaje/akundaje/release/blacklists/mm10-](http://mitra.stanford.edu/kundaje/akundaje/release/blacklists/mm10-mouse/mm10.blacklist.bed.gz)
418 [mouse/mm10.blacklist.bed.gz](http://mitra.stanford.edu/kundaje/akundaje/release/blacklists/mm10-mouse/mm10.blacklist.bed.gz) for mm10 and
419 <https://www.encodeproject.org/files/ENCFF001TDO/> for hg19). The sequence mappability for
420 each 10Kb bin is calculated based on our previous study⁴⁵, and it can be downloaded from
421 http://enhancer.sdsc.edu/yunjiang/resources/genomic_features/.

422
423 **Step D. Clustering of loop candidates and identifying the summit(s) as final outputs.**

424 For each loop candidate (i, j) , we defined its surrounding area as all 10Kb bin pairs (m, n) such
425 that $\max\{d(i, m), d(j, n)\} \leq 20\text{Kb}$, where $d(i, m)$ is the genomic distance between the center of
426 bin i and the center of bin m . We defined a loop candidate as a singleton if there is no other loop
427 candidate within its surrounding area, and removed all singletons from downstream analysis since
428 the singletons are likely to be false positives.

429
430 To group the remaining non-singleton loop candidates into clusters, we adopted the Rodriguez
431 and Laio's algorithm²¹. Specifically, for each loop candidate (i, j) , we first counted the number of
432 loop candidates in its adjacent neighborhood regions: $(m, n): \max\{d(i, m), d(j, n)\} \leq 10\text{Kb}$, and
433 defined this number as its local density $\rho(i, j)$. Next, we calculated the minimum Euclidean
434 distance between the loop candidate (i, j) and any other loop candidate with higher local density
435 on the same chromosome, defined as $\delta(i, j)$:

$$436 \quad \delta(i, j) = \min_{(m, n): \rho(m, n) > \rho(i, j)} \sqrt{(i - m)^2 + (j - n)^2}.$$

437 If the loop candidate (i, j) has the highest local density (i.e., $\rho(i, j) = \rho$), $\delta(i, j)$ is defined as:

$$438 \quad \delta(i, j) = \max_{(m, n)} \sqrt{(i - m)^2 + (j - n)^2}.$$

439 We then selected the loop candidates which have high local density ρ , and are relatively far away
440 from the other loop candidates with higher local density, i.e., high δ , as loop cluster centers. To
441 determine the cutoff values of ρ and δ for such centers, we implemented an algorithm similar to
442 the ROSE algorithm⁴⁶, which is used to identify super-enhancers. Specifically, let ρ_{max} and δ_{max}
443 represent the maximal value of ρ and δ of all loop candidates on each chromosome, respectively.
444 We defined $\rho'(i, j) = \rho(i, j)/\rho_{max}$ and $\delta'(i, j) = \delta(i, j)/\delta_{max}$ such that both $\rho'(i, j)$ and $\delta'(i, j)$ are

445 within range $[0,1]$. We then defined $\eta(i,j) = \rho'(i,j) * \delta'(i,j)$, ordered all loop candidates by their
446 η in the descending order, and plotted the rank of η against the value of η . In this plot, we selected
447 the reflection point such that the slope at the reflection point is one. All loop candidates with η
448 larger than η at the reflection point were chosen to be the loop cluster centers. After finding the
449 loop cluster centers, we assigned each remaining loop candidate to the same loop cluster as its
450 nearest neighbor with higher local density ρ .

451
452 Within each loop cluster, we defined the loop candidate with the lowest FDR as the first summit
453 of the cluster. For the first summit (i,j) , we defined its surrounding area as all 10Kb bin pairs
454 (m,n) such that $\max\{d(i,m),d(j,n)\} \leq 20\text{Kb}$, and removed all loop candidates within its
455 surrounding area. Next, we selected the loop candidate with the lowest FDR among the remaining
456 ones (if there is any) as the second summit of this cluster. We then removed all loop candidates
457 within the surrounding area of the second summit in the same way as we did for the first summit,
458 and searched for the third summit (if there is any) with the lowest FDR among the remaining loop
459 candidates. Such procedure was iterated until there are no loop candidates left in this cluster.
460 Notably, one loop cluster may contain multiple summits. SnapHiC algorithm outputs a file
461 containing the summit(s) of each loop cluster as its final chromatin loop list.

462
463 **Identification of chromatin loops with SnapHiC.**
464 We applied SnapHiC to scHi-C data from 10, 25, 50, 75, 100, 200, 300, 400, 500, 600, 700 and
465 742 mES cells and each of the 14 cell clusters from sn-m3C-seq data of human prefrontal cortex
466 to call chromatin loops at 10Kb resolution between 100Kb and 1Mb region on autosomal
467 chromosomes.

468
469 We did not take bin pairs within 100Kb into consideration because they do not have complete
470 information in their local neighborhood (refer to “**SnapHiC algorithm**”). We also evaluated the bin
471 pairs beyond 1Mb distance. When we extended the maximal genomic distance from 1Mb to 2Mb
472 for loop calling using scHi-C data from 742 mES cells, only 4.6% SnapHiC-identified loops (758
473 out of 16,654) are between 1Mb and 2Mb. Therefore, we restricted our loop calling from 100Kb
474 to 1Mb genomic distance for all the datasets mentioned in this study. In practice, we also suggest
475 using 1Mb as the maximal 1D genomic distance for loop calling to save computational cost.

476
477 **Visualization of scHi-C and sn-m3C-seq data using percentage (%) of outlier cells matrix.**

478 We first computed the % of outlier cells (i.e., the proportion of cells with normalized contact
479 probability > 1.96), and then took the integer ceiling of $100 * (\% \text{ of outlier cells})$ to create a count
480 matrix. We then used the Juicer⁴⁷ software to convert the count matrix into a .hic file and visualize
481 it in Juicebox⁴⁸.

482

483 **Computational cost (memory, time) of SnapHiC.**

484 To assess the relationship between the number of cells and running time, we tested the running
485 time of SnapHiC on 25, 50, 100, 200, 300 and 400 mES cells (10Kb resolution, searching for
486 loops between 100Kb to 1Mb genomic distance) and found its running time increases linearly with
487 the increase of cell numbers (**Supplementary Fig. 7**).

488

489 As described in our GitHub website (<https://github.com/HuMingLab/SnapHiC>), SnapHiC consists
490 of two steps: (1) applying the random walk with restart (RWR) algorithm to impute contact
491 probability within each single cell, and (2) integrating imputed contact probability matrices from all
492 single cells to identify significant chromatin loops. Since the RWR algorithm can be applied to
493 each chromosome in each single cell in parallel, in step 1, using as many processors as possible
494 (e.g., maximal $N = \# \text{ of cells} * \# \text{ of chromosomes}$) can speed up the computation. Resolution and
495 chromosome size are two important factors to determine the required memory per processor in
496 step 1. For human or mouse genome at 10Kb resolution, we recommend allocating at least 30GB
497 of memory for each processor. In the benchmarking experiments shown in **Supplementary Fig.**
498 **7**, we used 45 processors (15 nodes, 3 processors per node) for step 1, where each node has
499 96GB of memory, and it takes around 2.4 hours to process 100 cells.

500

501 In step 2, since the computation is performed jointly for all cells and separately for each
502 chromosome, we recommend using the same number of processors as the number of
503 chromosomes. Using more processors than that will be a waste of computing resources. It is also
504 important to ensure that each processor has access to sufficient memory for the computation over
505 all cells, and the amount of memory needed is correlated with the range of 1D genomic distance,
506 the bin resolution, and to a less extent to the number of cells. Increasing the number of cells,
507 slightly adds to the memory usage, however, since we only load the indices in the matrix that are
508 used in each step of the computation, this increase in memory usage is sublinear in regard to the
509 increase in the number of cells. In the benchmarking experiments shown in **Supplementary Fig.**
510 **7**, we used 20 processors (5 nodes, 4 processors per node) for step 2, where each node has
511 96GB of memory, and it takes around 0.7 hours to process 100 cells in step 2.

512

513 **Generation of aggregated contact matrix for scHi-C and sn-m3C-seq data.**

514 We pooled contacts from single cells of interest to create the aggregated contact matrix in .hic
515 format using Juicer with KR normalization⁴⁷. Only intra-chromosomal contacts >2Kb away are
516 used.

517

518 **Identification of HiCCUPS loops from aggregated contact matrix.**

519 We applied the HiCCUPS¹⁴ to the aggregated contact matrix after pooling the contacts from single
520 cells of interest and calling loops at 10Kb resolution with the following parameters: "--
521 ignore_sparsity -r 10000 -k KR -f.1 -p 2 -i 5 -t 0.02,1.5,1.75,2 -d 20000". Due to the sparsity of the
522 aggregated contact matrix generated using single cell data, KR normalization may not always
523 converge. Therefore, for some datasets, no HiCCUPS loops can be identified on specific
524 chromosomes where KR-normalized matrices are not available.

525

526 To ensure a fair comparison of HiCCUPS-identified loops with SnapHiC-identified loops, we
527 further filtered the HiCCUPS-identified loops by selecting the intra-chromosomal ones within
528 genomic distance 100Kb~1Mb and removing the loops whose anchor bins have low mappability
529 (≤ 0.8) or overlap with the ENCODE blacklist regions (refer to **Step C** in "**SnapHiC algorithm**").

530

531 **Definition of loop overlap.**

532 Let bin pair (i, j) represent a loop in set A . We define it overlaps with a loop in set B , if and only if
533 there exists a loop (m, n) in set B such that $\max(d_{im}, d_{jn}) \leq 20\text{Kb}$, where d_{im} is the 1D genomic
534 distance between the middle base pair of bin i and the middle base pair of bin m . We allow up to
535 20Kb gap in the definition of loop overlap, since SnapHiC outputs summits, and bin pairs within
536 20Kb of the summit can be part of the same loop cluster.

537

538 **Sub-sampling of scHi-C and sn-m3C-seq data.**

539 For scHi-C data from mES cells, we randomly permuted the order of all 742 cells, and selected
540 the first 10, 25, 50, 75, 100, 200, 300, 400, 500, 600, 700 cells from all 742 cells to create a series
541 of sub-sampled datasets. Notably, the dataset with fewer cells is always a subset of the dataset
542 with more cells.

543

544 For sc-m3C-seq data from human prefrontal cortex, we randomly permuted the order of all 338
545 astrocytes, 323 microglia and 1,038 oligodendrocytes and selected the first 261 astrocytes,

546 microglia and oligodendrocytes to create the sub-sampled datasets for astrocytes, microglia and
547 oligodendrocytes, respectively.

548

549 **Reproducibility of SnapHiC- and HiCCUPS-identified loops.**

550 Suppose we have two sets of loop list A and B . Let P_A represent the proportion of loops in set A
551 overlapped with loops in set B (up to 20Kb gap, see **Definition of loop overlap**) and let P_B
552 represent the proportion of loops in set B overlapped with loops in set A . We used $(P_A + P_B)/2$ to
553 measure the reproducibility of loops in the two sets.

554

555 To test the reproducibility of SnapHiC and HiCCUPS, we first randomly split all 742 mES cells
556 into two groups where each group consists of 371 cells, and then applied SnapHiC and HiCCUPS
557 to identify loops for each group. The reproducibility of SnapHiC- and HiCCUPS-identified loops
558 between two sets of 371 cells are calculated as described above. We repeated such random
559 splitting and loop calling analysis ten times, and reported the mean of reproducibility of SnapHiC-
560 and HiCCUPS-identified loops. We further used the paired t -test to evaluate the statistical
561 significance of the difference in reproducibility between these two methods.

562

563 **Generation of the reference loop lists for calculation of precision, recall and F1 score.**

564 For mES cells, the HiCCUPS loops at 10Kb resolution from bulk *in situ* Hi-C data were called as
565 previously described²³ using the pooled datasets of all 4 biological replicates from Bonev *et al.*
566 study²². MAPS pipeline was applied to H3K4me3 PLAC-seq data²³, cohesin HiChIP data²⁴ and
567 H3K27ac HiChIP data²⁵ to call significant interactions at 10Kb resolution within 1Mb genomic
568 distance. We combined the above four loop lists and further filtered by selecting the intra-
569 chromosomal loops within genomic distance 100Kb~1Mb and removing loops where anchor bins
570 have low mappability (≤ 0.8) or overlap with the ENCODE blacklist regions to create the final
571 reference loop list (**Supplementary Table 4**).

572

573 For oligodendrocytes, microglia and eight neuronal subtypes from human prefrontal cortex, we
574 used MAPS-identified interactions from H3K4me3 PLAC-seq data of purified oligodendrocytes,
575 microglia and neurons as their reference loop list, respectively (provided in Supplementary Table
576 5 in Nott *et al.* study²⁹). We first filtered the list by selecting the intra-chromosomal loops with
577 genomic distance 100Kb~1Mb and removing loops where anchor bins have low mappability (≤ 0.8)
578 or overlap with the ENCODE blacklist regions. We further selected the loops in which at least one

579 end contains active promoters of the corresponding cell type to create the final reference loop list
580 (**Supplementary Table 9**).

581

582 **Calculation of precision, recall and F1 score.**

583 Let N represent the number of loops in the reference loop list for the cell type of interest. Suppose
584 SnapHiC (or HiCCUPS) identifies M loops from the same cell type, and m of them overlapped
585 with loops in the reference loop list (see **Definition of loop overlap**). The precision is calculated
586 as m/M . Suppose among all N loops in the reference loop list, n loops overlapped with SnapHiC-
587 (or HiCCUPS-) identified loops. The recall is calculated as n/N . Notably, m and n may not be
588 equal since we allow up to a 20Kb gap between two overlapped loops. The F1 score is the
589 harmonic mean of the precision and recall and is calculated as below:

590

$$591 \quad F1 \text{ score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}} = 2 * \frac{m/M * n/N}{m/M + n/N}$$

592

593 For mES cells, we used all SnapHiC- or HiCCUPS-identified loops for the above calculation. For
594 oligodendrocytes, microglia and eight neuronal subtypes, we only selected the SnapHiC- or
595 HiCCUPS-identified loops in which at least one anchor contains active promoters of the
596 corresponding cell type for this calculation, since the available reference loop lists are called from
597 H3K4me3 PLAC-seq data, which can only detect interactions centered at promoter regions.

598

599 **Aggregate peak analysis (APA).**

600 We used the Juicer⁴⁷ software with the command “java -jar juicer_tools_1.19.02.jar apa -r 10000
601 -k KR -u input.hic loops.txt APA” to perform the aggregate peak analysis. We reported “P2LL”
602 (also known as the APA score) and “ZscoreLL” to evaluate the enrichment of SnapHiC-identified
603 loops with respect to the lower left background.

604

605 **CTCF motif orientation analysis.**

606 We obtained the CTCF ChIP-seq peaks of mES cells from a previous study⁴⁹, and used FIMO⁵⁰
607 with default parameters and the CTCF motif (MA0139.1) from the JASPAR⁵¹ database to search
608 for CTCF sequence motifs among those CTCF ChIP-seq peaks. Based on this CTCF motif list,
609 we then selected a subset of testable SnapHiC-identified loops in which both ends contain either
610 a single CTCF motif or multiple CTCF motifs in the same direction. Finally, we calculated the
611 proportion of convergent, tandem and divergent CTCF motif pairs among all testable loops.

612

613 **Visualization of CTCF and H3K27ac ChIP-seq data from mES cells.**

614 We downloaded the signal tracks from the ENCODE portal^{33,52} (<https://www.encodeproject.org/>)
615 with the following identifiers: ENCFF230RNU (for H3K27ac) and ENCFF069PTO (for CTCF) for
616 **Fig. 1d**.

617

618 **Definition of cell-type-specific SnapHiC loops.**

619 We used the SnapHiC loops identified from sub-sampled astrocytes, microglia, oligodendrocytes
620 datasets, and L2/3 excitatory neurons (all with 261 cells) to define cell-type-specific loops.
621 Specifically, we defined a loop identified from one cell type as cell-type-specific, if it did not overlap
622 (up to 20Kb gap, see **Definition of loop overlap**) with loops identified from any of the other three
623 cell types.

624

625 **Selection of genes associated with cell-type-specific SnapHiC loops.**

626 We first used the Gencode v34 (GRCh37) to obtain the location of transcription start site (TSS)
627 for 19,079 protein-coding genes in human autosomal chromosomes, and then selected genes
628 where TSS overlaps cell-type-specific loops for astrocytes, L2/3 excitatory neurons, microglia and
629 oligodendrocytes, respectively.

630

631 **Processing of ATAC-seq and H3K27ac ChIP-seq data from four brain cell types.**

632 The ATAC-seq and H3K27ac ChIP-seq data from human astrocytes, oligodendrocytes, microglia
633 and neurons are from the previous study²⁹ and are processed with ENCODE ATAC-seq and ChIP-
634 seq pipelines as previously described²⁹. The normalized bigwig tracks with RPKM as the Y-axis
635 are generated for visualization in **Fig. 2e**.

636

637 **Processing of RNA-seq from four brain cell types.**

638 The RNA-seq data from human astrocytes, oligodendrocytes, microglia and neurons are acquired
639 from the previous study³⁰. The alignment and quantification are performed with pipeline:
640 <https://github.com/ren-lab/rnaseq-pipeline>. Briefly, we first aligned RNA-seq raw reads to hg19.
641 Next, we used Gencode GTF `gencode.v19.annotation.gtf` for hg19 with STAR⁵³ following the
642 'ENCODE' options outlined in the STAR manual
643 ([http://labshare.cshl.edu/shares/gingeraslab/www-](http://labshare.cshl.edu/shares/gingeraslab/www-data/dobin/STAR/STAR_posix/doc/STARmanual.pdf)
644 [data/dobin/STAR/STAR_posix/doc/STARmanual.pdf](http://labshare.cshl.edu/shares/gingeraslab/www-data/dobin/STAR/STAR_posix/doc/STARmanual.pdf)). We then used Picard
645 (<http://broadinstitute.github.io/picard/>) to remove PCR duplicates. We also generated the

646 normalized bigwig tracks with RPKM (reads per kilobase of a transcript, per million mapped reads)
647 as the Y-axis for visualization in **Fig. 2e**.

648

649 **Enrichment analysis of ATAC-seq or H3K27ac ChIP-seq signals at cell-type-specific loops.**

650 To quantify the intensity of ATAC-seq or H3K27ac ChIP-seq signals at cell-type-specific loops in
651 the cell type of interest, we first calculated reads per million (CPM) values in each 10Kb anchor
652 of the cell-type-specific loops using ATAC-seq or H3K27ac ChIP-seq data from the cell type of
653 interest. To minimize the background noise, we only considered the reads falling into the ATAC-
654 seq or H3K27ac ChIP-seq peak regions defined in the cell type of interest but not all the reads in
655 the entire 10Kb bin. If there are multiple ATAC-seq or H3K27ac ChIP-seq peaks in the same 10Kb
656 bin, we then added up the CPM values and took the sum as the value for that 10Kb bin. Since
657 each loop has two anchors, we took their average CPM to represent the intensity of ATAC-seq or
658 H3K27ac ChIP-seq signal for that loop in the cell type of interest. Lastly, we applied the paired
659 Wilcoxon signed-rank test on $\log_2(\text{CPM}+1)$ values from different combinations of cell types of
660 interest and the cell-type-specific loop sets to test whether there is a significantly difference (**Fig.**
661 **2c**).

662

663 **Gene expression analysis at cell-type-specific loops.**

664 We obtained the FPKM values of each protein-coding genes in human astrocytes, neurons,
665 microglia and oligodendrocytes from Supplementary Table 4 provided in the previous study (Col
666 P-U for astrocytes, Col AB for neurons, Col AC-AG for oligodendrocytes, and Col AH-AJ for
667 microglia in the “Human data only” tab)³⁰. For each gene, we took the average of FPKM across
668 biological replicates of the same cell type. For the selected genes where promoters are
669 overlapped with cell-type-specific loops, we applied the Wilcoxon signed-rank test to evaluate
670 whether they are highly expressed in the matched cell type.

671

672 **Gene ontology enrichment analysis.**

673 We used Metascape³¹ to perform gene ontology enrichment analysis for selected genes where
674 promoters overlapped with cell-type-specific loops, and reported the top seven enriched biological
675 processes.

676

677

678 **Reference**

- 679 1 Zheng, H. & Xie, W. The role of 3D genome organization in development and cell
680 differentiation. *Nature reviews. Molecular cell biology* 20, 535-550, doi:10.1038/s41580-
681 019-0132-4 (2019).
- 682 2 Schmitt, A. D., Hu, M. & Ren, B. Genome-wide mapping and analysis of chromosome
683 architecture. *Nature reviews. Molecular cell biology* 17, 743-755,
684 doi:10.1038/nrm.2016.104 (2016).
- 685 3 Yu, M. & Ren, B. The Three-Dimensional Organization of Mammalian Genomes. *Annu*
686 *Rev Cell Dev Biol* 33, 265-289, doi:10.1146/annurev-cellbio-100616-060531 (2017).
- 687 4 Nagano, T. et al. Single-cell Hi-C reveals cell-to-cell variability in chromosome structure.
688 *Nature* 502, 59-64, doi:10.1038/nature12593 (2013).
- 689 5 Nagano, T. et al. Cell-cycle dynamics of chromosomal organization at single-cell
690 resolution. *Nature* 547, 61-67, doi:10.1038/nature23001 (2017).
- 691 6 Stevens, T. J. et al. 3D structures of individual mammalian genomes studied by single-cell
692 Hi-C. *Nature* 544, 59-64, doi:10.1038/nature21429 (2017).
- 693 7 Ramani, V. et al. Massively multiplex single-cell Hi-C. *Nature methods* 14, 263-266,
694 doi:10.1038/nmeth.4155 (2017).
- 695 8 Flyamer, I. M. et al. Single-nucleus Hi-C reveals unique chromatin reorganization at
696 oocyte-to-zygote transition. *Nature* 544, 110-114, doi:10.1038/nature21711 (2017).
- 697 9 Collombet, S. et al. Parental-to-embryo switch of chromosome organization in early
698 embryogenesis. *Nature* 580, 142-146, doi:10.1038/s41586-020-2125-z (2020).
- 699 10 Tan, L., Xing, D., Chang, C. H., Li, H. & Xie, X. S. Three-dimensional genome structures
700 of single diploid human cells. *Science (New York, N.Y.)* 361, 924-928,
701 doi:10.1126/science.aat5641 (2018).
- 702 11 Tan, L., Xing, D., Daley, N. & Xie, X. S. Three-dimensional genome structures of single
703 sensory neurons in mouse visual and olfactory systems. *Nature structural & molecular*
704 *biology* 26, 297-307, doi:10.1038/s41594-019-0205-2 (2019).
- 705 12 Li, G. et al. Joint profiling of DNA methylation and chromatin architecture in single cells.
706 *Nature methods* 16, 991-993, doi:10.1038/s41592-019-0502-z (2019).
- 707 13 Lee, D. S. et al. Simultaneous profiling of 3D genome structure and DNA methylation in
708 single human cells. *Nature methods* 16, 999-1006, doi:10.1038/s41592-019-0547-z
709 (2019).
- 710 14 Rao, Suhas S. P. et al. A 3D Map of the Human Genome at Kilobase Resolution Reveals
711 Principles of Chromatin Looping. *Cell* 159, 1665-1680 (2014).
- 712 15 Ay, F., Bailey, T. L. & Noble, W. S. Statistical confidence estimation for Hi-C data reveals
713 regulatory chromatin contacts. *Genome Res* 24, 999-1011, doi:10.1101/gr.160374.113
714 (2014).
- 715 16 Kaul, A., Bhattacharyya, S. & Ay, F. Identifying statistically significant chromatin contacts
716 from Hi-C data with FitHiC2. *Nature protocols* 15, 991-1012, doi:10.1038/s41596-019-
717 0273-0 (2020).
- 718 17 Xu, Z. et al. A hidden Markov random field-based Bayesian method for the detection of
719 long-range chromosomal interactions in Hi-C data. *Bioinformatics (Oxford, England)* 32,
720 650-656, doi:10.1093/bioinformatics/btv650 (2016).
- 721 18 Xu, Z., Zhang, G., Wu, C., Li, Y. & Hu, M. FastHiC: a fast and accurate algorithm to detect
722 long-range chromosomal interactions from Hi-C data. *Bioinformatics (Oxford, England)* 32,
723 2692-2695 (2016).
- 724 19 Li, X., An, Z. & Zhang, Z. Comparison of computational methods for 3D genome analysis
725 at single-cell Hi-C level. *Methods*, doi:10.1016/j.ymeth.2019.08.005 (2019).

- 726 20 Zhou, J. et al. Robust single-cell Hi-C clustering by convolution- and random-walk-based
727 imputation. *Proceedings of the National Academy of Sciences of the United States of*
728 *America* 116, 14011-14018, doi:10.1073/pnas.1901423116 (2019).
- 729 21 Rodriguez, A. & Laio, A. Clustering by fast search and find of density peaks. *Science (New*
730 *York, N.Y.)* 344, 1492-1496, doi:10.1126/science.1242072 (2014).
- 731 22 Bonev, B. et al. Multiscale 3D Genome Rewiring during Mouse Neural Development. *Cell*
732 171, 557-572.e524, doi:10.1016/j.cell.2017.09.043 (2017).
- 733 23 Juric, I. et al. MAPS: model-based analysis of long-range chromatin interactions from
734 PLAC-seq and HiChIP experiments. *PLoS computational biology* In press.,
735 doi:10.1101/411835 (2019).
- 736 24 Mumbach, M. R. et al. HiChIP: efficient and sensitive analysis of protein-directed genome
737 architecture. *Nature methods* 13, 919-922, doi:10.1038/nmeth.3999 (2016).
- 738 25 Mumbach, M. R. et al. Enhancer connectome in primary human cells identifies target
739 genes of disease-associated DNA elements. *Nature genetics* 49, 1602-1612,
740 doi:10.1038/ng.3963 (2017).
- 741 26 Fudenberg, G. et al. Formation of Chromosomal Domains by Loop Extrusion. *Cell Rep* 15,
742 2038-2049, doi:10.1016/j.celrep.2016.04.085 (2016).
- 743 27 Li, Y. et al. CRISPR reveals a distal super-enhancer required for Sox2 expression in
744 mouse embryonic stem cells. *PLoS One* 9, e114485, doi:10.1371/journal.pone.0114485
745 (2014).
- 746 28 Schoenfelder, S. et al. The pluripotent regulatory circuitry connecting promoters to their
747 long-range interacting elements. *Genome Res* 25, 582-597, doi:10.1101/gr.185272.114
748 (2015).
- 749 29 Nott, A. et al. Brain cell type-specific enhancer-promoter interactome maps and disease-
750 risk association. *Science (New York, N.Y.)* 366, 1134-1139, doi:10.1126/science.aay0793
751 (2019).
- 752 30 Zhang, Y. et al. Purification and Characterization of Progenitor and Mature Human
753 Astrocytes Reveals Transcriptional and Functional Differences with Mouse. *Neuron* 89,
754 37-53, doi:10.1016/j.neuron.2015.11.013 (2016).
- 755 31 Zhou, Y. et al. Metascape provides a biologist-oriented resource for the analysis of
756 systems-level datasets. *Nature communications* 10, 1523, doi:10.1038/s41467-019-
757 09234-6 (2019).
- 758 32 Lander, E. S. et al. Initial sequencing and analysis of the human genome. *Nature* 409,
759 860-921, doi:10.1038/35057062 (2001).
- 760 33 Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome.
761 *Nature* 489, 57-74, doi:10.1038/nature11247 (2012).
- 762 34 Roadmap Epigenomics, C. et al. Integrative analysis of 111 reference human epigenomes.
763 *Nature* 518, 317-330, doi:10.1038/nature14248 (2015).
- 764 35 Maurano, M. T. et al. Systematic localization of common disease-associated variation in
765 regulatory DNA. *Science (New York, N.Y.)* 337, 1190-1195, doi:10.1126/science.1222794
766 (2012).
- 767 36 Shen, Y. et al. A map of the cis-regulatory sequences in the mouse genome. *Nature* 488,
768 116-120, doi:10.1038/nature11243 (2012).
- 769 37 Andersson, R. et al. An atlas of active enhancers across human cell types and tissues.
770 *Nature* 507, 455-461, doi:10.1038/nature12787 (2014).
- 771 38 Jansen, I. E. et al. Genome-wide meta-analysis identifies new loci and functional pathways
772 influencing Alzheimer's disease risk. *Nature genetics* 51, 404-413, doi:10.1038/s41588-
773 018-0311-9 (2019).
- 774 39 Demontis, D. et al. Discovery of the first genome-wide significant risk loci for attention
775 deficit/hyperactivity disorder. *Nature genetics* 51, 63-75, doi:10.1038/s41588-018-0269-7
776 (2019).

- 777 40 Grove, J. et al. Identification of common genetic risk variants for autism spectrum disorder.
778 Nature genetics 51, 431-444, doi:10.1038/s41588-019-0344-8 (2019).
- 779 41 Stahl, E. A. et al. Genome-wide association study identifies 30 loci associated with bipolar
780 disorder. Nature genetics 51, 793-803, doi:10.1038/s41588-019-0397-8 (2019).
- 781 42 Savage, J. E. et al. Genome-wide association meta-analysis in 269,867 individuals
782 identifies new genetic and functional links to intelligence. Nature genetics 50, 912-919,
783 doi:10.1038/s41588-018-0152-6 (2018).
- 784 43 Howard, D. M. et al. Genome-wide meta-analysis of depression identifies 102 independent
785 variants and highlights the importance of the prefrontal brain regions. Nature neuroscience
786 22, 343-352, doi:10.1038/s41593-018-0326-7 (2019).
- 787 44 Pardinás, A. F. et al. Common schizophrenia alleles are enriched in mutation-intolerant
788 genes and in regions under strong background selection. Nature genetics 50, 381-389,
789 doi:10.1038/s41588-018-0059-2 (2018).
- 790 45 Hu, M. et al. HiCNorm: removing biases in Hi-C data via Poisson regression.
791 Bioinformatics (Oxford, England) 28, 3131-3133, doi:10.1093/bioinformatics/bts570
792 (2012).
- 793 46 Whyte, W. A. et al. Master transcription factors and mediator establish super-enhancers
794 at key cell identity genes. Cell 153, 307-319, doi:10.1016/j.cell.2013.03.035 (2013).
- 795 47 Durand, N. C. et al. Juicer Provides a One-Click System for Analyzing Loop-Resolution
796 Hi-C Experiments. Cell systems 3, 95-98, doi:10.1016/j.cels.2016.07.002 (2016).
- 797 48 Durand, N. C. et al. Juicebox Provides a Visualization System for Hi-C Contact Maps with
798 Unlimited Zoom. Cell systems 3, 99-101, doi:10.1016/j.cels.2015.07.012 (2016).
- 799 49 Kubo, N. et al. CTCF Promotes Long-range Enhancer-promoter Interactions and Lineage-
800 specific Gene Expression in Mammalian Cells. 2020.2003.2021.001693,
801 doi:10.1101/2020.03.21.001693 %J bioRxiv (2020).
- 802 50 Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.
803 Bioinformatics (Oxford, England) 27, 1017-1018, doi:10.1093/bioinformatics/btr064 (2011).
- 804 51 Khan, A. et al. JASPAR 2018: update of the open-access database of transcription factor
805 binding profiles and its web framework. Nucleic acids research 46, D260-d266,
806 doi:10.1093/nar/gkx1126 (2018).
- 807 52 Davis, C. A. et al. The Encyclopedia of DNA elements (ENCODE): data portal update.
808 Nucleic acids research 46, D794-d801, doi:10.1093/nar/gkx1081 (2018).
- 809 53 Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics (Oxford,
810 England) 29, 15-21, doi:10.1093/bioinformatics/bts635 (2013).

811