

1 **Dr.seq2: a quality control and analysis pipeline for**
2 **parallel single cell transcriptome and epigenome data**

3 Chengchen Zhao¹, Sheng'en Hu¹, Xiao Huo¹, Yong Zhang^{1*}

4 ¹Translational Medical Center for Stem Cell Therapy & Institute for Regenerative Medicine,
5 Shanghai East Hospital, School of Life Science and Technology, Shanghai Key Laboratory of
6 Signaling and Disease Research, Tongji University, Shanghai 20092, China

7 * Corresponding author

8 E-mail: yzhang@tongji.edu.cn

9 **Abstract**

10 An increasing number of single cell transcriptome and epigenome technologies,
11 including single cell ATAC-seq (scATAC-seq), have been recently developed as
12 powerful tools to analyze the features of many individual cells simultaneously.
13 However, the methods and software were designed for one certain data type and only
14 for single cell transcriptome data. A systematic approach for epigenome data and
15 multiple types of transcriptome data is needed to control data quality and to perform
16 cell-to-cell heterogeneity analysis on these ultra-high-dimensional transcriptome and
17 epigenome datasets. Here we developed Dr.seq2, a Quality Control (QC) and analysis
18 pipeline for multiple types of single cell transcriptome and epigenome data, including
19 scATAC-seq and Drop-ChIP data. Application of this pipeline provides four groups
20 of QC measurements and different analyses, including cell heterogeneity analysis.
21 Dr.seq2 produced reliable results on published single cell transcriptome and
22 epigenome datasets. Overall, Dr.seq2 is a systematic and comprehensive QC and
23 analysis pipeline designed for parallel single cell transcriptome and epigenome data.
24 Dr.seq2 is freely available at: <http://www.tongji.edu.cn/~zhanglab/drseq2/> and
25 <https://github.com/ChengchenZhao/DrSeq2>.

26 **Keywords:** Single cell ATAC-seq; Quality control and analysis pipeline; Single cell
27 transcriptome data analysis; Single cell epigenome data analysis; Drop-seq; MARS-
28 seq; 10x genomics; Drop-ChIP;

29 **Introduction**

30 To better understand cell-to-cell variability, an increasing number of transcriptome
31 technologies, such as Drop-seq [1, 2], Cyto-seq [3], 10x genomics [4], MARS-seq [5],
32 and epigenome technologies, such as Drop-ChIP [6], single cell ATAC-seq

33 (scATAC-seq) [7], have been developed in recent years. These technologies can
34 easily provide a large amount of single cell transcriptome information or epigenome
35 information at minimal cost, which makes it possible to perform analysis of cell
36 heterogeneity on the transcriptome and epigenome levels, deconstruction of a cell
37 population, and detection of rare cell populations. However, different single cell
38 transcriptome technologies have their own features given their specific experimental
39 design, such as cell sorting methods, RNA capture rates, and sequencing depths. But
40 the methods and software such as Dr.seq [8] were developed for one single cell data
41 type with certain functions (S1 File). Furthermore, the quality control step of single
42 cell epigenome data is more challenging than for transcriptome data given the
43 amplification noise caused by the limit number of DNA copy in single cell epigenome
44 experiments. But few quality control and analysis method was developed specific for
45 single cell epigenome data. Thus a comprehensive QC pipeline suitable for multiple
46 types of single cell transcriptome data and epigenome data is urgently needed. Here,
47 we provide Dr.seq2, a QC and analysis pipeline for multiple types of parallel single
48 cell transcriptome and epigenome data, including recently published scATAC-seq
49 data. Dr.seq2 can systematically generate specific QC, analyze, and visualize
50 unsupervised cell clustering for multiple types of single cell data. For single cell
51 transcriptome data, the QC steps of Dr.seq2 are primarily derived from Dr.seq [8] and
52 the output of Dr.seq2 on these data will not be described in details in this paper.

53 **Materials and methods**

54 **Drop-seq data**

55 The Drop-seq samples were obtained from NCBI Gene Expression Omnibus (GEO)
56 database under accession GSM1626793.

57

58 **MARS-seq data**

59 The MARS-seq samples were obtained from NCBI Gene Expression Omnibus (GEO)
60 database under accession GSE54006. These samples were combined as a MARS-seq
61 dataset and analyzed by Dr.seq2 using three different dimension reduction methods.

62 **10x genomics data**

63 The 10x genomics datasets were obtained from 10x genomic data support
64 (<https://support.10xgenomics.com/single-cell/datasets>). The sample named “50%: 50%
65 Jurkat: 293T Cell Mixture” was analyzed by Dr.seq2 using three different dimension
66 reduction methods.

67 **scATAC-seq data**

68 The scATAC-seq datasets were obtained from NCBI Gene Expression Omnibus
69 (GEO) database under accession GSE65360. We combined 288 scATAC datasets
70 (GSM1596255 ~ GSM1596350, GSM1596735 ~ GSM1596830, GSM1597119 ~
71 GSM1597214) from three cell types and analyzed by Dr.seq2. Cell clustering was
72 conducted for the combined scATAC-seq dataset. We also plotted the cell type labels
73 using different colors on the clustering plot and found consistent classifications with
74 the clustering results.

75 **Drop-ChIP data**

76 The Drop-ChIP datasets were obtained from NCBI Gene Expression Omnibus (GEO)
77 database under accession GSE70253.

78 **Implementation of Dr.seq2**

79 Dr.seq2 was implemented using Python and R. Linux or MacOS environment with
80 Python (version = 2.7) and R (version >= 2.14.1) was suitable for Dr.seq2. It was
81 distributed under the GNU General Public License version 3 (GPLv3). A detailed
82 tutorial was provided on the Dr.seq2 webpage
83 (<http://www.tongji.edu.cn/~zhanglab/drseq2>) and source code of Dr.seq2 was
84 available on github (<https://github.com/ChengchenZhao/DrSeq2>).

85 **Quality control components**

86 Dr.seq2 conducted four groups of QC measurements on single cell epigenome data: (i)
87 reads level QC; (ii) bulk-cell level QC; (iii) individual-cell level QC; and (iv) cell-
88 clustering level QC.

89 **Reads level QC and bulk-cell level QC**

90 We used a published package called RseQC [9] for reads level QC of Drop-ChIP data
91 and scATAC-seq data to measure the general sequence quality. In bulk-cell level QC,
92 a Drop-ChIP dataset (or scATAC-seq datasets combined from several scATAC-seq
93 samples) was regarded as a bulk-cell ChIP-seq (or bulk-cell ATAC-seq) data. Next,
94 “combined peaks” were detected with total reads from the “bulk-cell” data using
95 MACS[10] for output and the following steps. Different MACS parameters were
96 applied to Drop-ChIP and scATAC-seq data. We used the published package CEAS
97 to measure the performance of ChIP for ChIP-seq data (or Tn5 digestion for scATAC-
98 seq data) [11].

99 **Individual-cell level QC**

100 The reads number distribution was calculated by counting the number of reads
101 assigned to each single cell. A single cell referred to a unique cell barcode in Drop-
102 ChIP data. For scATAC-seq data, the peak number in each cell was defined as the
103 number of “combined peaks” occupied by the reads in the cell. The distribution of

104 different peak numbers in each cell indicated the different amount of information the
105 cell contains.

106 **Cell-clustering level QC**

107 Cells were first clustered based on their occupancy of “combined peaks” using
108 hierarchical clustering. Next, cells in each cluster were regarded as the same cell type
109 (or same cell sub-type), and reads from the same cell type were merged. For each cell
110 type, unique peaks from other cell types were defined as specific peaks in this cell
111 type. Specific peaks in different cell types were displayed with different colors
112 according to genomic locations. Silhouette method is used to interpret and validate the
113 consistency within clusters defined in previous steps.

114 Note that reads with no overlap with “combined peaks” were discarded in this step
115 and the following steps. Clusters containing less than 3 single cells were also
116 discarded.

117 **Simulation of scATAC-seq datasets**

118 To measure the tolerance of Dr.seq2 for low sequencing depth and small numbers of
119 cells of a certain cell type, we simulated datasets from 3 cell types with different cell
120 proportions and sequencing depths using scATAC-seq data (Table 1). To test the
121 effect of low sequencing depth, we sampled the reads count from 10,000 reads to
122 100,000 reads for each cell and compared these results with the Goodman-Kruskal’s
123 lambda index [12] of clustering results using cells with a certain number of reads.

124

125

126

127

128 **Table 1. Meta data and accession ID for the scATAC-seq data used in simulation**
129 **for pipeline tolerance evaluation.**

Accession ID	Cell line	Cell type	Target/regular cell
GSM1596255-GSM1596350	H1	human embryonic stem cell line	Target
GSM1596735-GSM1596830	GM12878	lymphoblastoid cells	Regular
GSM1597119-GSM1597214	K562	chronic myeloid leukemia cells	Regular

130 We defined the 1 out of 3 cell types as “target cell type”, while the other cell types were defined as
131 “regular cell type”.

132

133 To test the effect of low cell numbers of a certain cell type (defined as a target cell
134 type) on cell clustering, we defined 1 of the 3 cell types as the “target cell type”,
135 whereas the other cell types were defined as the “regular cell type”, and sampled cells
136 with following compositions: 10:70:70 (10 for target cell type, 70 for the two regular
137 cell types), 15:67:67, 20:65:65, 25:62:62, 30:60:60, 35:57:57, 40:55:55, 45:52:52 and
138 50:50:50. Then, we called “combined peaks” and clustered cells on the simulated
139 dataset. The Goodman-Kruskal’s lambda index [12] was calculated to evaluate the
140 cell clustering performance. The average Goodman-Kruskal’s lambda index and 95%
141 confidence intervals were calculated from 20 simulations.

142 **Results and discussion**

143 **Dr.seq2 overview**

144 The Dr.seq2 QC and analysis pipeline is suitable for both single cell transcriptome
145 data and epigenome data. Multiple types of single cell transcriptome data (including
146 scRNA-seq, Drop-seq, inDrop, MARS-seq and 10x genomics data) and epigenome
147 data (including scATAC-seq and Drop-ChIP) are acceptable for Dr.seq2 with relevant
148 functions (S1 Fig).

149 Recently many methods and software were developed for single cell RNA-seq data.
150 However most of them were suitable for certain data types with limited functions. We
151 compared the major function of Dr.seq2 to existing state-of-the-art methods (Table 2).
152 Dr.seq2 provides two advantages: 1) Dr.seq2 supports different types of single cell
153 transcriptome data and single cell epigenome data. 2) Dr.seq2 provides both
154 multifaceted QC reports and cell clustering results. Then We used the simulated
155 single cell RNA-seq data from seven RNA-seq datasets from ENCODE (S2 File) to
156 estimate the performance of our Dr.seq2 pipeline (using different dimensional
157 reduction methods: SIMLR and t-SNE) in cell clustering comparing to three existing
158 methods (SINCERA, SNN-Cliq, BackSPIN). We applied these five methods on ten
159 datasets with different numbers of reads per cell range from 100 to 10,000 to measure
160 the accuracy and time cost of each method on different sequencing depth. SIMLR
161 shows more accurate clustering results than t-SNE on the datasets with small number
162 of reads per cell and comparable clustering results on the datasets with large number
163 of reads per cell. And Dr.seq2 (using either SIMLR or t-SNE) shows better clustering
164 accuracy than SNN-Cliq, and comparable clustering accuracy with BackSPIN and
165 SINCERA on the datasets with large number of reads per cell. On the datasets with
166 small number of reads per cell, SINCERA clustering result shows better accuracy
167 than Dr.seq2 (using either SIMLR or t-SNE) and SNN-Cliq. However SINCERA
168 takes a great amount of time on all these datasets comparing with Dr.seq2. As for
169 BackSPIN, it does not support for these datasets with small number of reads per cell.
170 Overall, Dr.seq2 (using either SIMLR or t-SNE) provides reliable cell clustering
171 results with acceptable time cost (S2 Fig).
172

173 **Table 2. Comparison of functions between Dr.seq2 and other software developed**
 174 **for single cell transcriptome data.**

Name	Supporting parallel single cell data (e.g. Drop-seq, MARS-seq)	Supporting single cell epigenome data (e.g. scATAC and Drop-ChIP)	Reads level QC	Individual cell level QC	Highly variable gene detection	Noise reduction	Informative cell selection	Cell clustering (sub cell type identification)	Differential expressed gene detection	Pseudo-temporal ordering	Reference
Dr.seq2	Y	Y	Y	Y	Y	Y	Y	Y	Y		-
Dr.seq	Y		Y	Y	Y	Y	Y	Y	Y		[8]
BASICS					Y	Y		Y			[13]
scLVM					Y	Y			Y		[14]
SINCERA			Y	Y	Y	Y		Y	Y		[15]
OEFinder					Y	Y					[16]
ZIFA						Y		Y			[17]
Destiny						Y		Y			[18]
SNN-Cliq						Y		Y			[19]
RaceID						Y		Y	Y		[20]
SCUBA						Y		Y	Y	Y	[21]
BackSPIN								Y	Y		[22]
PAGODA					Y	Y		Y			[23]
MAST									Y		[24]
SCDE								Y	Y		[25]
scDD									Y		[26]
Monocle										Y	[27]
Waterfall								Y		Y	[28]
Sincell								Y		Y	[29]
Oscope								Y		Y	[30]
Wanderlust								Y		Y	[31]
CellTree								Y		Y	[32]
SinQC			Y	Y		Y					[33]
ASAP						Y		Y	Y		[34]

175 We compare the major function of Dr.seq2 to existing state-of-the-art methods. Each column shows
 176 different functions of these methods and software.

177

178 **QC and analysis workflow**

179 Dr.seq2 uses raw sequencing files in FASTQ format or alignment results in
180 SAM/BAM format as input with relevant commands and generates four steps of QC
181 measurements and analysis results (Fig 1).

182

183 **Fig 1. Flowchart illustrating the Dr.seq2 pipeline with default parameters.** The
184 workflow of the Dr.seq2 pipeline includes QC and analysis components for parallel
185 single cell transcriptome and epigenome data. The QC component contains reads
186 level, bulk-cell level, individual-cell level and cell-clustering level QC.

187

188 For transcriptome data, the QC steps of Dr.seq2 are primarily derived from Dr.seq [8].
189 However, almost all data types are now supported, and more dimension reduction
190 methods, including PCA, t-SNE and SIMLR[35], are supported. For single cell
191 epigenome data, technologies like scATAC-seq and Drop-ChIP are increasingly
192 common. However few quality control and analysis approaches have been developed
193 for these data. Dr.seq2 conducts QC measurements on single cell epigenome data
194 from four aspects: (i) reads level QC, including sequence quality, nucleotide
195 composition and GC content of reads inherited from previous work; (ii) bulk-cell
196 level QC, including genomic distribution of “combined peaks” and average profile on
197 regulatory regions; (iii) individual-cell level QC, including the distribution of the
198 number of reads and the peak number distribution; and (iv) cell-clustering level QC,
199 including Silhouette score[36] and cell type-specific peak detection.

200 **Cell clustering for different single cell transcriptome data** 201 **types using different dimension reduction methods**

202 We applied our pipeline to three different types of single cell transcriptome data
203 (Drop-seq, MARS-seq and 10x genomics data) using three different dimension

204 reduction methods (PCA, t-SNE and SIMLR[35]) to evaluate the performance of
205 Dr.seq2 on different types of single cell transcriptome data (Fig 2). Due to the
206 different distance calculation method and kernel function the method used, Dr.seq2
207 represented cluster results from different dimensions.

208

209 **Fig 2. Dimensional reduction results for different single cell transcriptome data**

210 **types.**

211 (A-I) Cell clustering results using dimensional reduction methods (PCA, t-SNE and
212 SIMLR) on different types of single cell transcriptome data (Drop-seq, 10x genomics
213 and MARS-seq).

214

215 **Bulk-cell level QC of scATAC-seq data to measure the**
216 **performance of Tn5 digestion**

217 To evaluate the performance of Dr.seq2 on single cell epigenome data, we combined
218 288 scATAC datasets (GSM1596255 ~ GSM1596350, GSM1596735 ~ GSM1596830,
219 GSM1597119 ~ GSM1597214) from three cell types and applied Dr.seq2 to it.
220 “Combined peaks” were detected with total reads from the combined dataset using
221 MACS for output and the following steps. We measured the scATAC data quality in
222 bulk-cell level from 4 aspects (Fig 3): 1) Peak distribution on each chromosome; 2)
223 Open regions distributed over the genome along with their scores; 3) Average
224 profiling on different genomic features; 4) Fragment length distribution. The peak
225 distribution on each chromosome and the open region distributed over the genome
226 showed the general quality of Tn5 digestion. The average profiling on different
227 genomic features represented the quality of Tn5 digestion around specific regions.

228 And the periodicity fragment length distribution indicated factor occupancy and
229 nucleosome positions due to different Tn5 digestion degrees.

230

231 **Fig 3. Bulk-cell level QC for scATAC-seq datasets.** A) Peak region number
232 distribution on each chromosome. The blue bars represent the percentages of the
233 whole tiled or mappable regions in the chromosomes (genome background) and the
234 red bars showed the percentages of the whole open region. These percentages are also
235 marked right next to the bars. P-values for the significance of the relative enrichment
236 of open regions with respect to the genome background are shown in parentheses next
237 to the percentages of the red bars. B) Open region distribution over the genome along
238 with their scores or peak heights. The line graph on the top left corner illustrates the
239 distribution of peak score. The x-axis of the main plot represents the actual
240 chromosome sizes. C) Average profiling on different genomic features. The panels on
241 the first row display the average enrichment signals around TSS and TTS of genes,
242 respectively. The bottom panel represents the average signals on the meta-gene of 5
243 kb. D) Red line shows number distribution of different fragment length.

244

245 **Cell clustering for scATAC-seq datasets with three clusters** 246 **that were consistent with the cell type labels**

247 To measure the cell clustering performance of Dr.seq2 on epigenome data, cells from
248 the combined scATAC-seq dataset were firstly clustered based on their occupancy of
249 “combined peaks” using hierarchical clustering. Then cell type labels were marked by
250 different colors according to the original cell type information (red stand for H1 cells,
251 yellow stand for GM12878 cells and blue stand for K562 cells). Cells were clearly

252 separated into different groups that were consistent with the cell type labels by
253 Dr.seq2 (Fig 4A).

254

255 **Fig 4. Cell-clustering level QC and single-cell level QC for scATAC-seq data. A)**

256 Upper panel shows cell-clustering results for combined scATAC samples generated

257 from 3 different cell types. Bottom panel shows corresponding cell type labels of each

258 cell marked by different colors (red stand for H1 cells, yellow stand for GM12878

259 cells and blue stand for K562 cells). The clustering step of Dr.seq2 clearly separated

260 the scATAC-seq samples from three different cell types into different groups that

261 were consistent with the cell type labels. B) Distribution of peak number for each

262 single cell. C) Cell Clustering tree and peak region in each cell. The upper panel

263 represents the hieratical clustering results based on each single cell. The second panel

264 with different colors represents decision of cell clustering. The bottom two panels

265 (heatmap and color bar) represent the “combined peaks” occupancy of each single

266 cell. D) Barplot shows Silhouette score of each cluster. Silhouette method is used to

267 interpret and validate the consistency within clusters defined in previous steps. E)

268 Cluster specific regions in each chromosome. Specific regions for different cell

269 clusters are marked by different colors and ordered according to genomic loci.

270

271 **Single-cell level QC and post analysis of scATAC-seq data**

272 In the single-cell level QC of Dr.seq2 on scATAC-seq data, the peak number of in

273 each cell was defined as the number of “combined peaks” occupied by the reads in the

274 cell. The distribution of different peak numbers in each cell indicated the different

275 amount of information the cell contains (Fig 4B). Cell clustering was conducted based

276 on the peak information in each cell using hierarchical clustering and open region was

277 shown in the order of genomic location (Fig 4C). And Silhouette score [36] validated
278 the consistency of each cluster (Fig 4D). Then cells in the same clusters were
279 considered as cells in the same cell type and combined for the detection of cell type
280 specific regions, which were defined as the peak regions that only covered in this cell
281 type. Specific regions for different cell clusters were marked by different colors and
282 ordered according to genomic loci (Fig 4E).

283 **Cell clustering stability on simulated scATAC-seq data**

284 To measure the tolerance of Dr.seq2 for low sequencing depth and small numbers of
285 cells of a certain cell type, we simulated datasets with different cell proportions and
286 sequencing depths by using scATAC-seq datasets from three cell types (Table 1).

287 We selected cells in different proportion with 100,000 reads per cell and then
288 performed cell clustering using Dr.seq2. The performance of cell clustering methods
289 was evaluated by Goodman-Kruskal's lambda index. And the average Goodman-
290 Kruskal's lambda index calculated from 20 simulations indicated that Dr.seq2 was
291 suitable for cell clustering with different cell proportions (Fig 5A). We also selected
292 fifty cells from each cell type with the reads count range from 10,000 reads to
293 100,000 reads for each cell to measure the tolerance of Dr.seq2 on low sequence
294 depth. Dr.seq2 produced stable clustering results with greater than 40,000 reads per
295 cell (Fig 5B).

296

297 **Fig 5. Cell clustering stability on simulated scATAC-seq data.** A) Clustering
298 stability of Dr.seq2 on simulated data with different numbers of reads per cell. The
299 lambda index (y-axis) is plotted as a function of the number of reads per cell (x-axis).
300 Error bars represent 95% confidence intervals calculated from 20 simulations. B)
301 Clustering stability of Dr.seq2 on simulated data with different cell proportion depths.

302 The lambda index (y-axis) is plotted as a function of the target cell number (x-axis).

303 Error bars represent 95% confidence intervals calculated from 20 simulations.

304

305 **Computational cost of Dr.seq2**

306 We also measured the computational time cost of Dr.seq2 by applied Dr.seq2 on

307 combined scATAC-seq datasets (Table 3). The running time of each step was

308 calculated using a single CPU (Intel® Xeon® CPU E5-2640 v2 @ 2.00 GHz).

309

310 **Table 3. Running time of each QC and analysis step for scATAC datasets.**

Steps	Time (s/CPU)	Percentage (%)
Merge Cells	1507	39.72
Bulk-cell level QC	1654	43.60
Individual-cell level QC and cell-clustering QC	626	16.50
Post-analysis	5	0.13
Summary Report	2	0.05

311 288 scATAC datasets from three cell types were used to evaluate the runtime of Dr.seq2. The

312 running time for each step was calculated using a single CPU (Intel® Xeon® CPU E5-2640

313 v2 @ 2.00 GHz).

314 **Conclusions**

315 In summary, Dr.seq2 is designed for QC and analysis components of parallel single

316 cell transcriptome and epigenome data. Parallel single cell transcriptome data

317 generated by different technologies can be transformed to the standard input for

318 Dr.seq2 with contained functions. Using relevant commands, Dr.seq2 can also be

319 used to report quality measurements based on four aspects and generate detailed

320 analysis results for scATAC-seq and Drop-ChIP datasets.

321 **Acknowledgments**

322 We thank Shiyang Zeng and Yiyang Lang for their suggestions.

323 **Funding**

324 This work was supported by National Natural Science Foundation of China
325 (31571365, 31322031 and 31371288), National Key Research and Development
326 Program of China (2016YFA0100400), Specialized Research Fund for the Doctoral
327 Program of Higher Education (20130072110032), and Program of Shanghai
328 Academic Research Leader (17XD1403600).

329 *Conflict of Interest:* none declared.

330 **References**

- 331 1. Macosko EZ, Basu A, Satija R, Nemesh J, Shekhar K, Goldman M, et al.
332 Highly Parallel Genome-wide Expression Profiling of Individual Cells Using
333 Nanoliter Droplets. *Cell*. 2015;161(5):1202-14. doi: 10.1016/j.cell.2015.05.002.
334 PubMed PMID: 26000488; PubMed Central PMCID: PMC4481139.
- 335 2. Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, et al.
336 Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells.
337 *Cell*. 2015;161(5):1187-201. doi: 10.1016/j.cell.2015.04.044. PubMed PMID:
338 26000487; PubMed Central PMCID: PMC4441768.
- 339 3. Fan HC, Fu GK, Fodor SP. Expression profiling. Combinatorial labeling of
340 single cells for gene expression cytometry. *Science*. 2015;347(6222):1258367. doi:
341 10.1126/science.1258367. PubMed PMID: 25657253.

- 342 4. Zheng GX, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al.
343 Massively parallel digital transcriptional profiling of single cells. *Nat Commun.*
344 2017;8:14049. doi: 10.1038/ncomms14049. PubMed PMID: 28091601; PubMed
345 Central PMCID: PMCPMC5241818.
- 346 5. Jaitin DA, Kenigsberg E, Keren-Shaul H, Elefant N, Paul F, Zaretsky I, et al.
347 Massively parallel single-cell RNA-seq for marker-free decomposition of tissues into
348 cell types. *Science.* 2014;343(6172):776-9. doi: 10.1126/science.1247651. PubMed
349 PMID: 24531970; PubMed Central PMCID: PMCPMC4412462.
- 350 6. Rotem A, Ram O, Shoresh N, Sperling RA, Goren A, Weitz DA, et al. Single-
351 cell ChIP-seq reveals cell subpopulations defined by chromatin state. *Nat Biotechnol.*
352 2015;33(11):1165-72. doi: 10.1038/nbt.3383. PubMed PMID: 26458175; PubMed
353 Central PMCID: PMCPMC4636926.
- 354 7. Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP,
355 et al. Single-cell chromatin accessibility reveals principles of regulatory variation.
356 *Nature.* 2015;523(7561):486-90. doi: 10.1038/nature14590. PubMed PMID:
357 26083756; PubMed Central PMCID: PMCPMC4685948.
- 358 8. Huo X, Hu S, Zhao C, Zhang Y. Dr.seq: a quality control and analysis pipeline
359 for droplet sequencing. *Bioinformatics.* 2016;32(14):2221-3. doi:
360 10.1093/bioinformatics/btw174. PubMed PMID: 27153611.
- 361 9. Wang L, Wang S, Li W. RSeQC: quality control of RNA-seq experiments.
362 *Bioinformatics.* 2012;28(16):2184-5. doi: 10.1093/bioinformatics/bts356. PubMed
363 PMID: 22743226.
- 364 10. Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, et al.
365 Model-based analysis of ChIP-Seq (MACS). *Genome Biol.* 2008;9(9):R137. doi:

- 366 10.1186/gb-2008-9-9-r137. PubMed PMID: 18798982; PubMed Central PMCID:
367 PMCPMC2592715.
- 368 11. Shin H, Liu T, Manrai AK, Liu XS. CEAS: cis-regulatory element annotation
369 system. *Bioinformatics*. 2009;25(19):2605-6. doi: 10.1093/bioinformatics/btp479.
370 PubMed PMID: 19689956.
- 371 12. Goodman LA, Kruskal WH. Measures of association for cross-classification. *J*
372 *Am Stat Assoc*. 1954;49:732-64.
- 373 13. Vallejos CA, Marioni JC, Richardson S. BASiCS: Bayesian Analysis of
374 Single-Cell Sequencing Data. *PLoS Comput Biol*. 2015;11(6):e1004333. doi:
375 10.1371/journal.pcbi.1004333. PubMed PMID: 26107944; PubMed Central PMCID:
376 PMCPMC4480965.
- 377 14. Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, et
378 al. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-
379 sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol*.
380 2015;33(2):155-60. doi: 10.1038/nbt.3102. PubMed PMID: 25599176.
- 381 15. Guo M, Wang H, Potter SS, Whitsett JA, Xu Y. SINCERA: A Pipeline for
382 Single-Cell RNA-Seq Profiling Analysis. *PLoS Comput Biol*. 2015;11(11):e1004575.
383 doi: 10.1371/journal.pcbi.1004575. PubMed PMID: 26600239; PubMed Central
384 PMCID: PMCPMC4658017.
- 385 16. Leng N, Choi J, Chu LF, Thomson JA, Kendziorski C, Stewart R. OEFinder: a
386 user interface to identify and visualize ordering effects in single-cell RNA-seq data.
387 *Bioinformatics*. 2016;32(9):1408-10. doi: 10.1093/bioinformatics/btw004. PubMed
388 PMID: 26743507; PubMed Central PMCID: PMCPMC4848403.

- 389 17. Pierson E, Yau C. ZIFA: Dimensionality reduction for zero-inflated single-cell
390 gene expression analysis. *Genome Biol.* 2015;16:241. doi: 10.1186/s13059-015-0805-
391 z. PubMed PMID: 26527291; PubMed Central PMCID: PMC4630968.
- 392 18. Angerer P, Haghverdi L, Buttner M, Theis FJ, Marr C, Buettner F. destiny:
393 diffusion maps for large-scale single-cell data in R. *Bioinformatics.* 2016;32(8):1241-
394 3. doi: 10.1093/bioinformatics/btv715. PubMed PMID: 26668002.
- 395 19. Xu C, Su Z. Identification of cell types from single-cell transcriptomes using a
396 novel clustering method. *Bioinformatics.* 2015;31(12):1974-80. doi:
397 10.1093/bioinformatics/btv088. PubMed PMID: 25805722.
- 398 20. Grun D, Lyubimova A, Kester L, Wiebrands K, Basak O, Sasaki N, et al.
399 Single-cell messenger RNA sequencing reveals rare intestinal cell types. *Nature.*
400 2015;525(7568):251-5. doi: 10.1038/nature14966. PubMed PMID: 26287467.
- 401 21. Marco E, Karp RL, Guo G, Robson P, Hart AH, Trippa L, et al. Bifurcation
402 analysis of single-cell gene expression data reveals epigenetic landscape. *Proc Natl*
403 *Acad Sci U S A.* 2014;111(52):E5643-50. doi: 10.1073/pnas.1408993111. PubMed
404 PMID: 25512504; PubMed Central PMCID: PMC4284553.
- 405 22. Zeisel A, Munoz-Manchado AB, Codeluppi S, Lonnerberg P, La Manno G,
406 Jureus A, et al. Brain structure. Cell types in the mouse cortex and hippocampus
407 revealed by single-cell RNA-seq. *Science.* 2015;347(6226):1138-42. doi:
408 10.1126/science.aaa1934. PubMed PMID: 25700174.
- 409 23. Fan J, Salathia N, Liu R, Kaeser GE, Yung YC, Herman JL, et al.
410 Characterizing transcriptional heterogeneity through pathway and gene set
411 overdispersion analysis. *Nat Methods.* 2016;13(3):241-4. doi: 10.1038/nmeth.3734.
412 PubMed PMID: 26780092; PubMed Central PMCID: PMC4772672.

- 413 24. Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, et al. MAST:
414 a flexible statistical framework for assessing transcriptional changes and
415 characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.*
416 2015;16:278. doi: 10.1186/s13059-015-0844-5. PubMed PMID: 26653891; PubMed
417 Central PMCID: PMCPMC4676162.
- 418 25. Kharchenko PV, Silberstein L, Scadden DT. Bayesian approach to single-cell
419 differential expression analysis. *Nat Methods.* 2014;11(7):740-2. doi:
420 10.1038/nmeth.2967. PubMed PMID: 24836921; PubMed Central PMCID:
421 PMCPMC4112276.
- 422 26. Korthauer KD, Chu LF, Newton MA, Li Y, Thomson J, Stewart R, et al. A
423 statistical approach for identifying differential distributions in single-cell RNA-seq
424 experiments. *Genome Biol.* 2016;17(1):222. doi: 10.1186/s13059-016-1077-y.
425 PubMed PMID: 27782827; PubMed Central PMCID: PMCPMC5080738.
- 426 27. Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, et al. The
427 dynamics and regulators of cell fate decisions are revealed by pseudotemporal
428 ordering of single cells. *Nat Biotechnol.* 2014;32(4):381-6. doi: 10.1038/nbt.2859.
429 PubMed PMID: 24658644; PubMed Central PMCID: PMCPMC4122333.
- 430 28. Shin J, Berg DA, Zhu Y, Shin JY, Song J, Bonaguidi MA, et al. Single-Cell
431 RNA-Seq with Waterfall Reveals Molecular Cascades underlying Adult
432 Neurogenesis. *Cell Stem Cell.* 2015;17(3):360-72. doi: 10.1016/j.stem.2015.07.013.
433 PubMed PMID: 26299571.
- 434 29. Julia M, Telenti A, Rausell A. Sincell: an R/Bioconductor package for
435 statistical assessment of cell-state hierarchies from single-cell RNA-seq.
436 *Bioinformatics.* 2015;31(20):3380-2. doi: 10.1093/bioinformatics/btv368. PubMed
437 PMID: 26099264; PubMed Central PMCID: PMCPMC4595899.

- 438 30. Leng N, Chu LF, Barry C, Li Y, Choi J, Li X, et al. Oscope identifies
439 oscillatory genes in unsynchronized single-cell RNA-seq experiments. *Nat Methods*.
440 2015;12(10):947-50. doi: 10.1038/nmeth.3549. PubMed PMID: 26301841; PubMed
441 Central PMCID: PMC4589503.
- 442 31. Bendall SC, Davis KL, Amir el AD, Tadmor MD, Simonds EF, Chen TJ, et al.
443 Single-cell trajectory detection uncovers progression and regulatory coordination in
444 human B cell development. *Cell*. 2014;157(3):714-25. doi:
445 10.1016/j.cell.2014.04.005. PubMed PMID: 24766814; PubMed Central PMCID:
446 PMC4045247.
- 447 32. duVerle DA, Yotsukura S, Nomura S, Aburatani H, Tsuda K. CellTree: an
448 R/bioconductor package to infer the hierarchical structure of cell populations from
449 single-cell RNA-seq data. *BMC Bioinformatics*. 2016;17(1):363. doi:
450 10.1186/s12859-016-1175-6. PubMed PMID: 27620863; PubMed Central PMCID:
451 PMC45020541.
- 452 33. Jiang P, Thomson JA, Stewart R. Quality control of single-cell RNA-seq by
453 SinQC. *Bioinformatics*. 2016;32(16):2514-6. doi: 10.1093/bioinformatics/btw176.
454 PubMed PMID: 27153613; PubMed Central PMCID: PMC4978927.
- 455 34. Gardeux V, David F, Shajkofci A, Schwalie PC, Deplancke B. ASAP: a Web-
456 based platform for the analysis and inter-active visualization of single-cell RNA-seq
457 data. *bioRxiv*. 2016. doi: 10.1101/096222.
- 458 35. Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S. Visualization and
459 analysis of single-cell RNA-seq data by kernel-based similarity learning. *Nat*
460 *Methods*. 2017;14(4):414-6. doi: 10.1038/nmeth.4207. PubMed PMID: 28263960.
- 461 36. Rousseeuw PJ. Silhouettes: A graphical aid to the interpretation and validation
462 of cluster analysis. *Comput Appl Math*. 1987;20:53-65.

463

464 **Supporting information**

465 **S1 Fig. Workflow displays the software structure and detailed QC steps of**

466 **Dr.seq2. A)** Dr.seq2 provides QC and analysis for three major data types: single cell

467 transcriptome data (DrSeq part), Drop-ChIP data (DrChIP part) and scATAC-seq data

468 (ATAC part). For single cell RNA-seq data, two additional step-by-step functions are

469 included: 1. Expression matrix generation for amounts of single cell RNA-seq

470 datasets (GeMa step) and 2. Cell clustering and analysis for the single cell expression

471 matrix (comCluster step). For different parallel single cell RNA-seq technologies,

472 input data are standardized for DrSeq part. **B)** Four groups of QC measurements are

473 conducted on single cell transcriptome data and epigenome data: 1. Reads level QC

474 including reads quality, reads nucleotide composition and reads GC content 2. Bulk-

475 cell level QC including reads alignment summary and gene body coverage for

476 transcriptome data; peak distribution; average profile on regulatory region and the

477 distribution of different numbers of fragment length for epigenome data. 3.

478 Individual-cell level QC including duplicate rate distribution, covered gene number

479 and intron rate distribution and intron rate distribution for transcriptome data; peak

480 number distribution and fragment length distribution for epigenome data. 4. Cell-

481 clustering level QC including Gap statistics score and Silhouette score for

482 transcriptome data, h-clustering and cluster specific peaks for epigenome data.

483 **S2 Fig. Comparing the performance of Dr.seq2 and three existing state-of-the art**

484 **methods on cell clustering. A)** Clustering accuracy measured by the Goodman-

485 Kruskal's lambda index of Dr.seq2 t-SNE, Dr.seq2 SIMLR methods and three

486 published methods on simulated data with different numbers of reads per cell. The

487 lambda index (y-axis) is plotted as a function of the number of reads per cell (x-axis).

488 **B)** Running time of Dr.seq2 t-SNE, Dr.seq2 SIMLR methods and three published

489 methods on simulated data with different numbers of reads per cell. The running time

490 (y-axis) is plotted as a function of the number of reads per cell (x-axis). The running

491 time for each method was calculated using a single CPU (Intel® Xeon® CPU E5-

492 2640 v2 @ 2.00 GHz).

493 **S1 File. Comparison of functions between Dr.seq2 and other software developed**

494 **for single cell transcriptome data.**

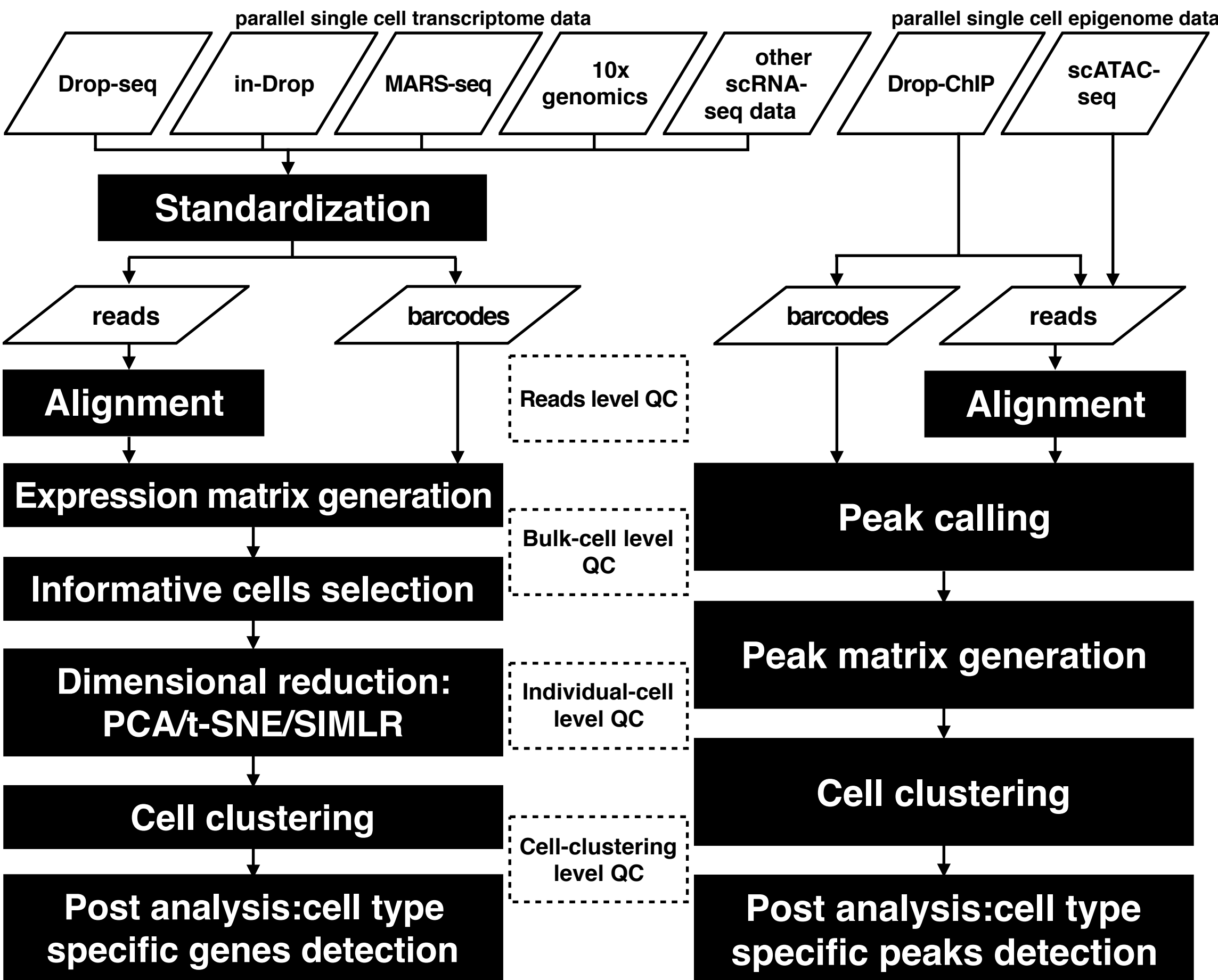
495 **S2 File. Meta data and accession ID for the bulk-cell RNA-seq data used in**

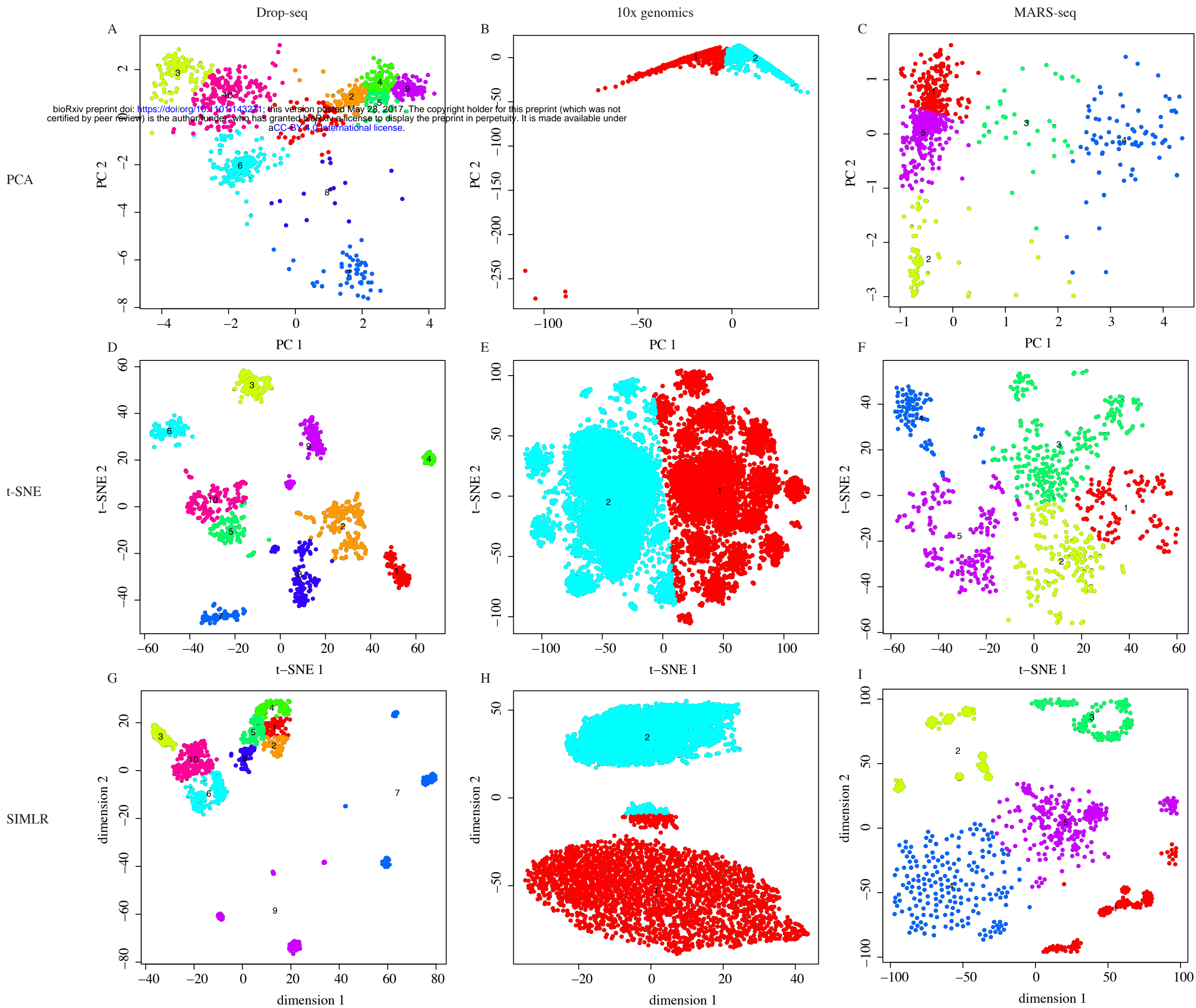
496 **simulation.**

497 **S3 File. Dr.seq2 QC and analysis output report for the scATAC-seq dataset.**

498 **S4 File. Dr.seq2 QC and analysis output report for the Drop-ChIP dataset.**

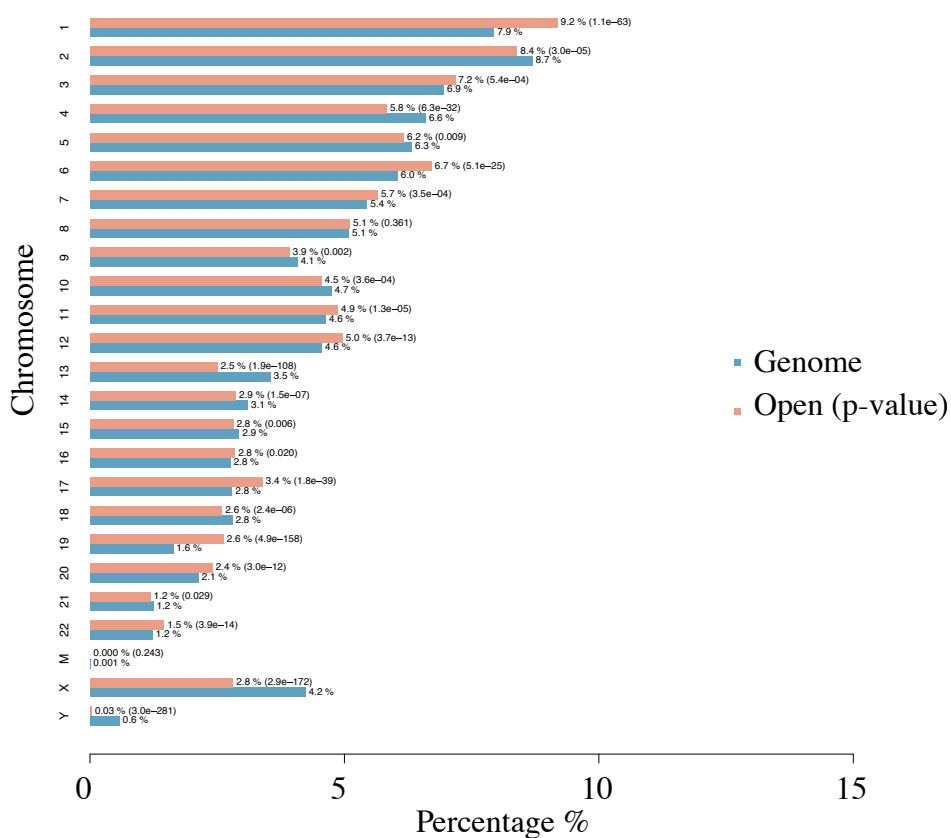
499 **S5 File. Dr.seq2 QC and analysis output report for the 10x genomics dataset.**





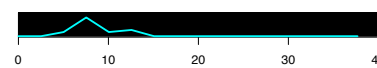
A

Chromosomal Distribution of Open Regions

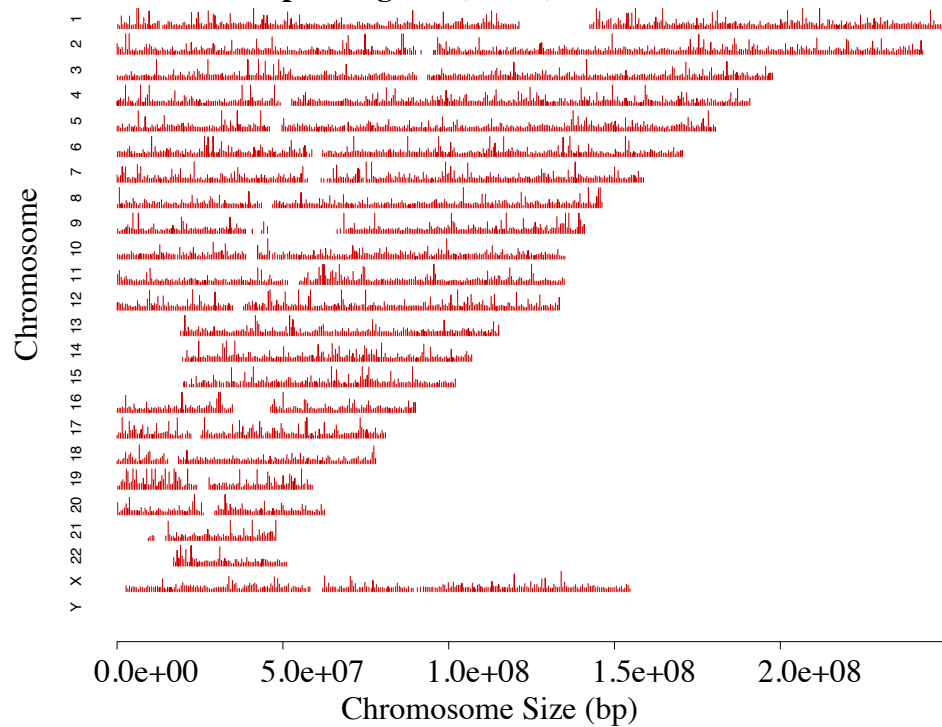


B

Distribution of Peak Heights

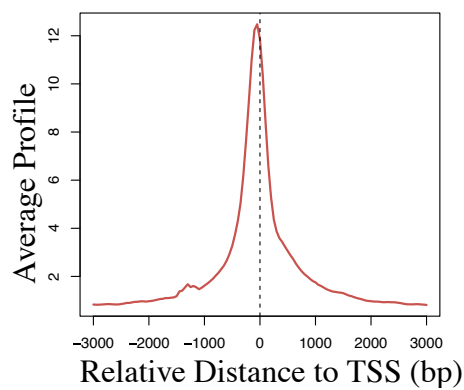


Open Regions (Peaks) over Chromosomes

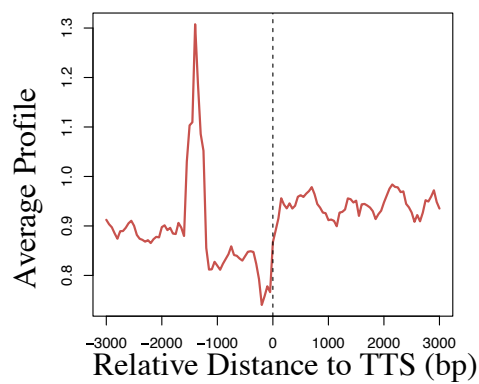


C

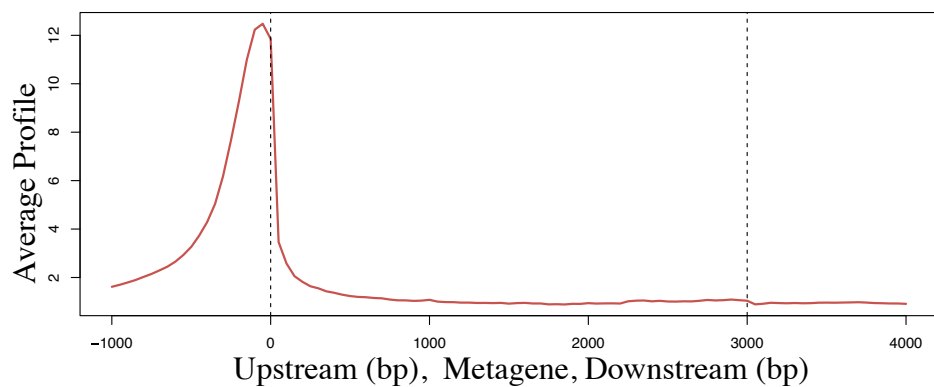
Average Profile near TSS



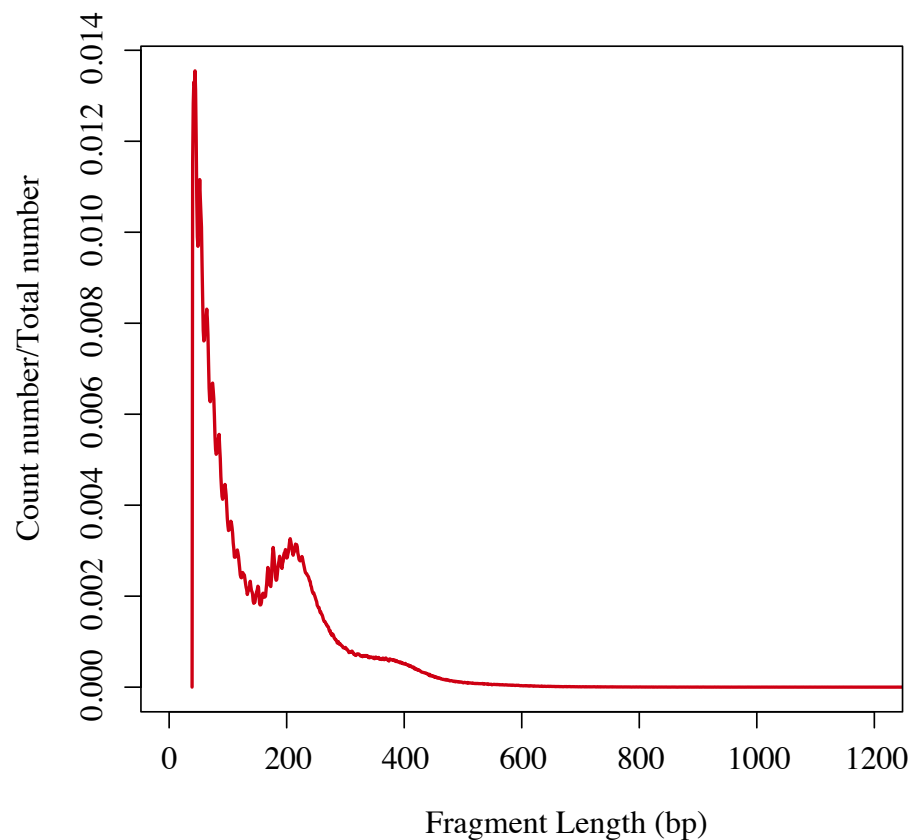
Average Profile near TTS



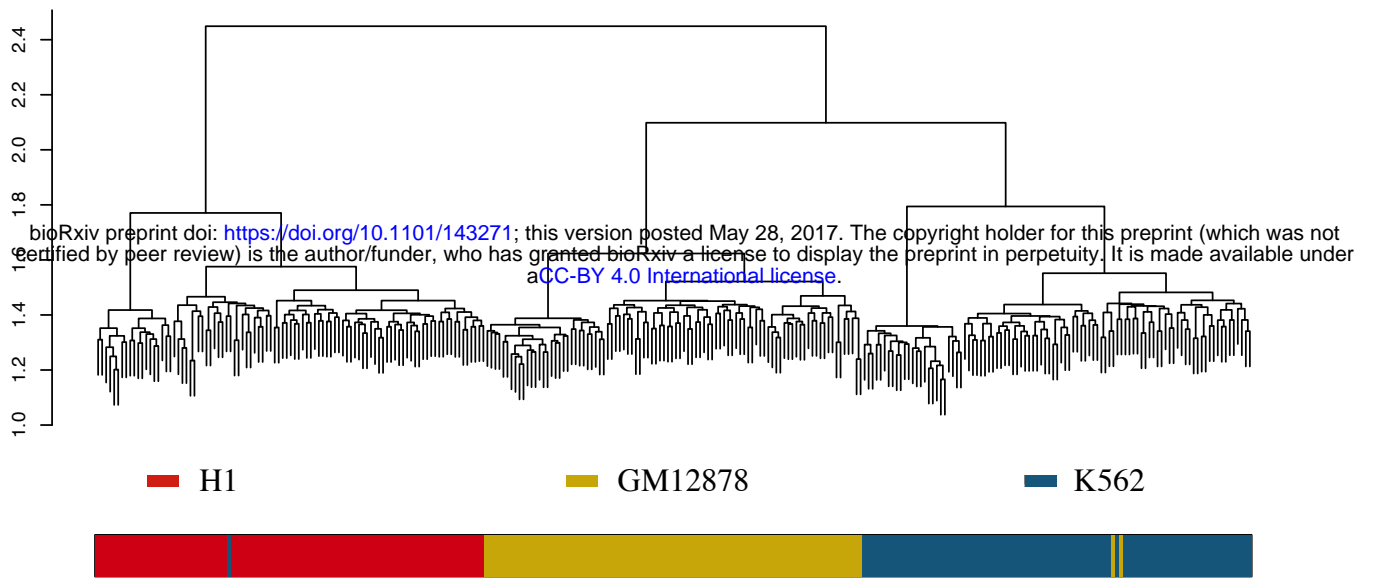
Average Gene Profile



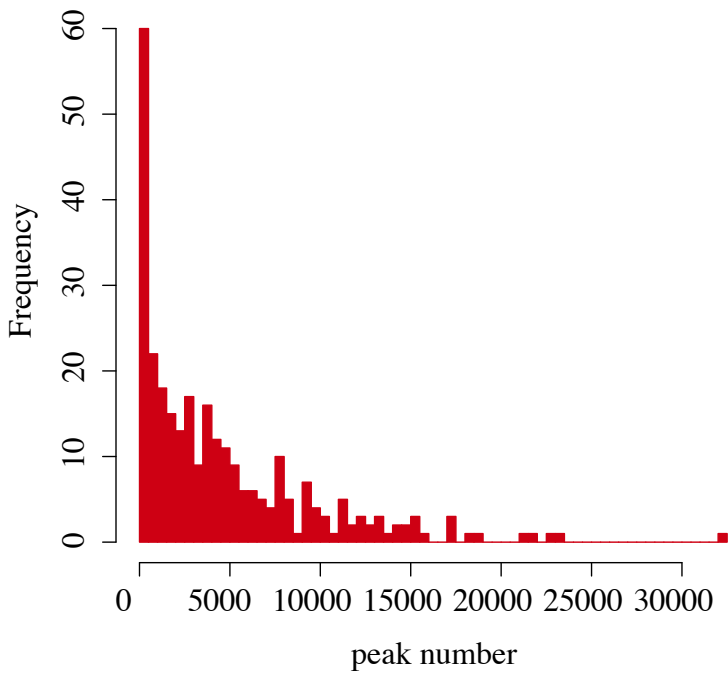
D



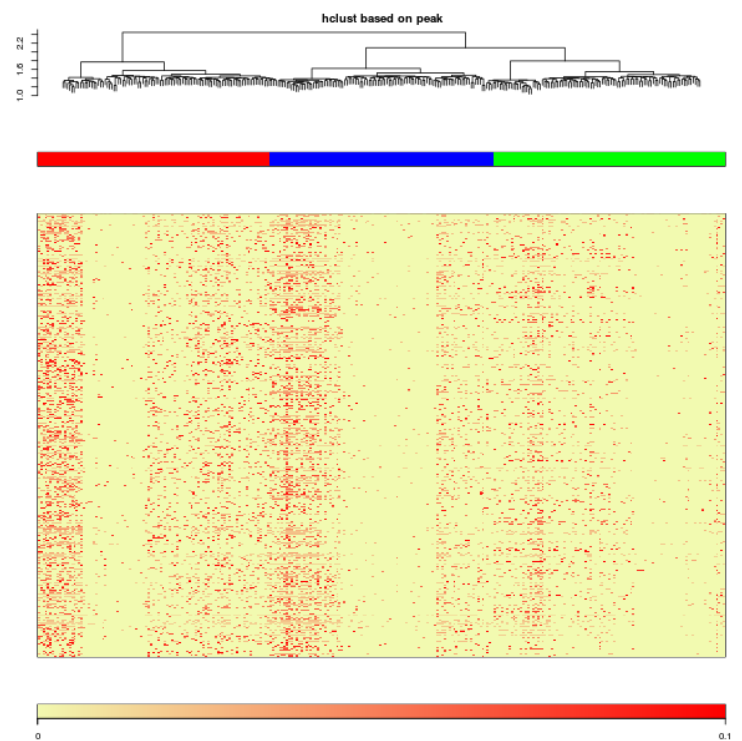
A



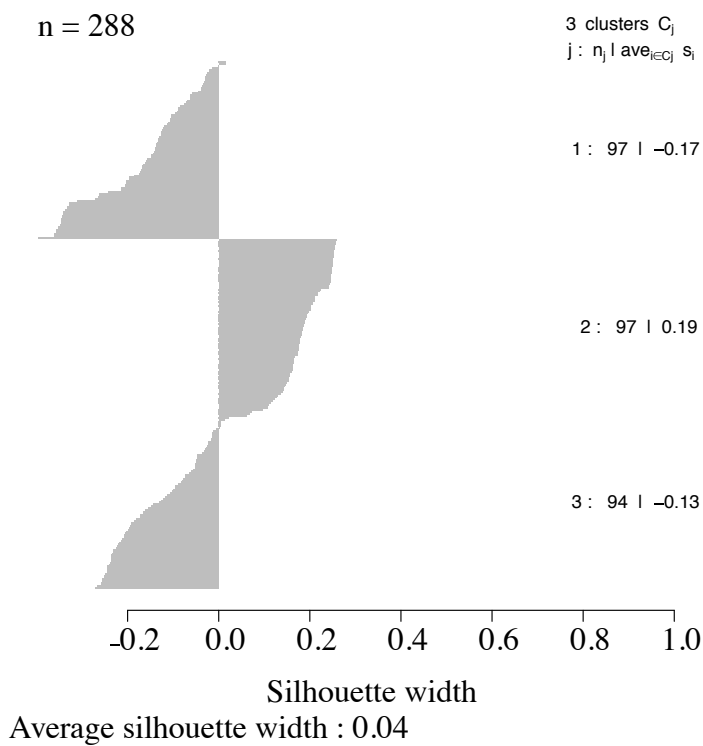
B



C



D



E

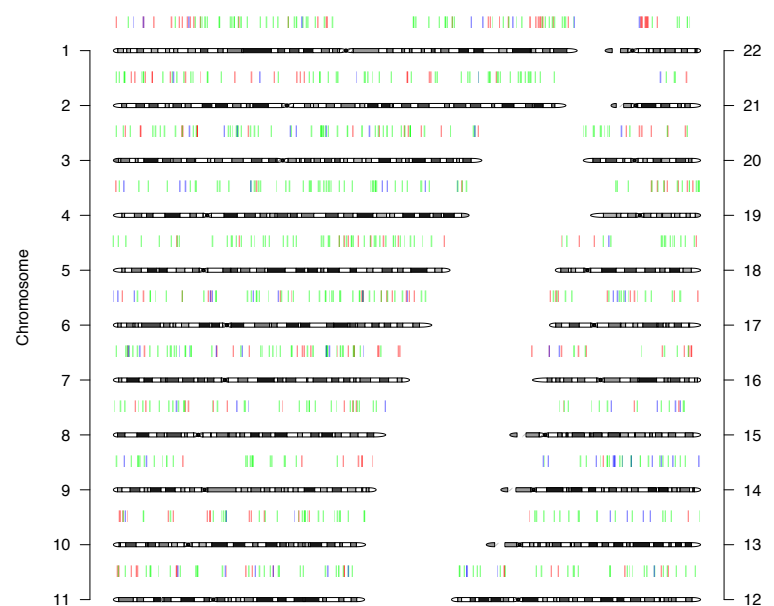
3 clusters C_j
 $j: n_j | \text{ave}_{i \in C_j} s_i$

1: 97 | -0.17

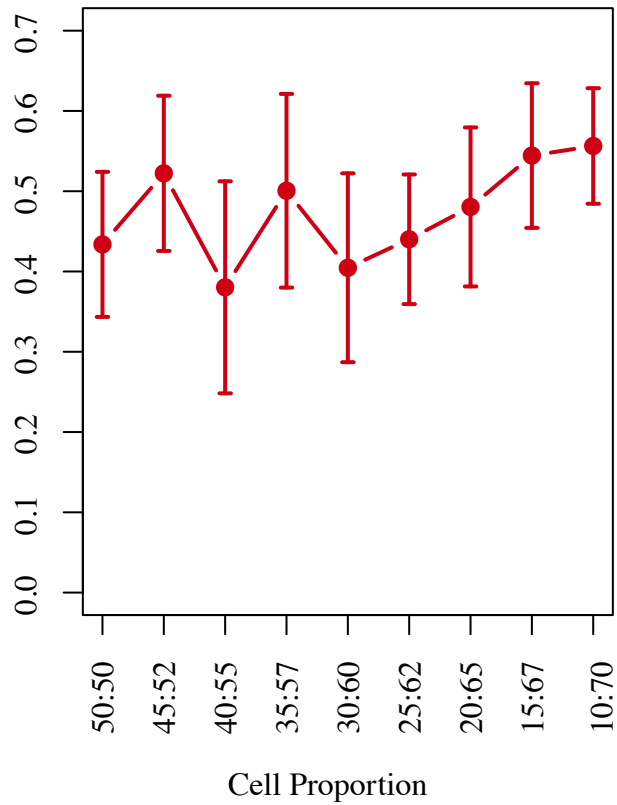
2: 97 | 0.19

3: 94 | -0.13

Cluster specific region on genomic ideogram



A Goodman-Kruskals lambda index



B Goodman-Kruskals lambda index

