

1 **Assessment of computational methods for the analysis of single-cell ATAC-seq data**

2 Huidong Chen^{1,2,3,4}, Caleb Lareau^{1,4,5,8}, Tommaso Andreani^{1,2,3,6,8}, Michael E.
3 Vinyard^{1,2,3,4,7,8}, Sara P. Garcia¹, Kendell Clement^{1,2,3,4}, Miguel A Andrade-Navarro⁶, Jason
4 D. Buenrostro^{4,5}, Luca Pinello^{1,2,3,4}

5

6 ¹Molecular Pathology Unit, Massachusetts General Hospital Research Institute,
7 Charlestown, MA 02129, USA

8 ²Center for Cancer Research, Massachusetts General Hospital, Charlestown, MA 02129,
9 USA

10 ³Department of Pathology, Harvard Medical School, Boston, MA 02115, USA.

11 ⁴Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA.

12 ⁵Department of Stem Cell and Regenerative Biology, Harvard University, Cambridge,
13 MA 02138, USA

14 ⁶Faculty of Biology, Computational Biology and Data Mining Lab. Johannes Gutenberg
15 University of Mainz, 55128 Mainz, Germany.

16 ⁷Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA
17 02142, USA.

18 ⁸These authors contributed equally.

19 Correspondence should be addressed to L.P. (email: lpinello@mgh.harvard.edu)

20

21 **Abstract**

22

23 **Background**

24 Recent innovations in single-cell Assay for Transposase Accessible Chromatin using
25 sequencing (scATAC-seq) enable profiling of the epigenetic landscape of thousands of
26 individual cells. scATAC-seq data analysis presents unique methodological challenges.
27 scATAC-seq experiments sample DNA, which, due to low copy numbers (diploid in
28 humans) lead to inherent data sparsity (1-10% of peaks detected per cell) compared to
29 transcriptomic (scRNA-seq) data (20-50% of expressed genes detected per cell). Such
30 challenges in data generation emphasize the need for informative features to assess cell
31 heterogeneity at the chromatin level.

32

33 **Results**

34 We present a benchmarking framework that was applied to 10 computational methods
35 for scATAC-seq on 13 synthetic and real datasets from different assays, profiling cell
36 types from diverse tissues and organisms. Methods for processing and featurizing
37 scATAC-seq data were evaluated by their ability to discriminate cell types when
38 combined with common unsupervised clustering approaches. We rank evaluated
39 methods and discuss computational challenges associated with scATAC-seq analysis
40 including inherently sparse data, determination of features, peak calling, the effects of
41 sequencing coverage and noise, and clustering performance. Running times and
42 memory requirements are also discussed.

43

44 **Conclusions**

45 This reference summary of scATAC-seq methods offers recommendations for best
46 practices with consideration for both the non-expert user and the methods developer.
47 Despite variation across methods and datasets, SnapATAC, *Cusanovich2018*, and
48 cisTopic outperform other methods in separating cell populations of different coverages
49 and noise levels in both synthetic and real datasets. Notably, SnapATAC was the only
50 method able to analyze a large dataset (> 80,000 cells).

51

52 **Keywords:** scATAC-seq, feature matrix, benchmarking, regulatory genomics,
53 clustering, visualization, featurization, dimensionality reduction

54

55 **Background**

56 Individual cell types within heterogenous tissues coordinate to perform complex
57 biological functions, many of which are not fully understood. Recent technological
58 advances in single-cell methodologies have resulted in an increased capacity to study
59 cell-to-cell heterogeneity and the underlying molecular regulatory programs that drive
60 such variation.

61

62 To date, most single-cell profiling efforts have been performed via quantification of
63 RNA by sequencing (scRNA-seq). While this provides snapshots of inter- and intra-
64 cellular variability in gene expression, investigation of the *epigenomic* landscape in

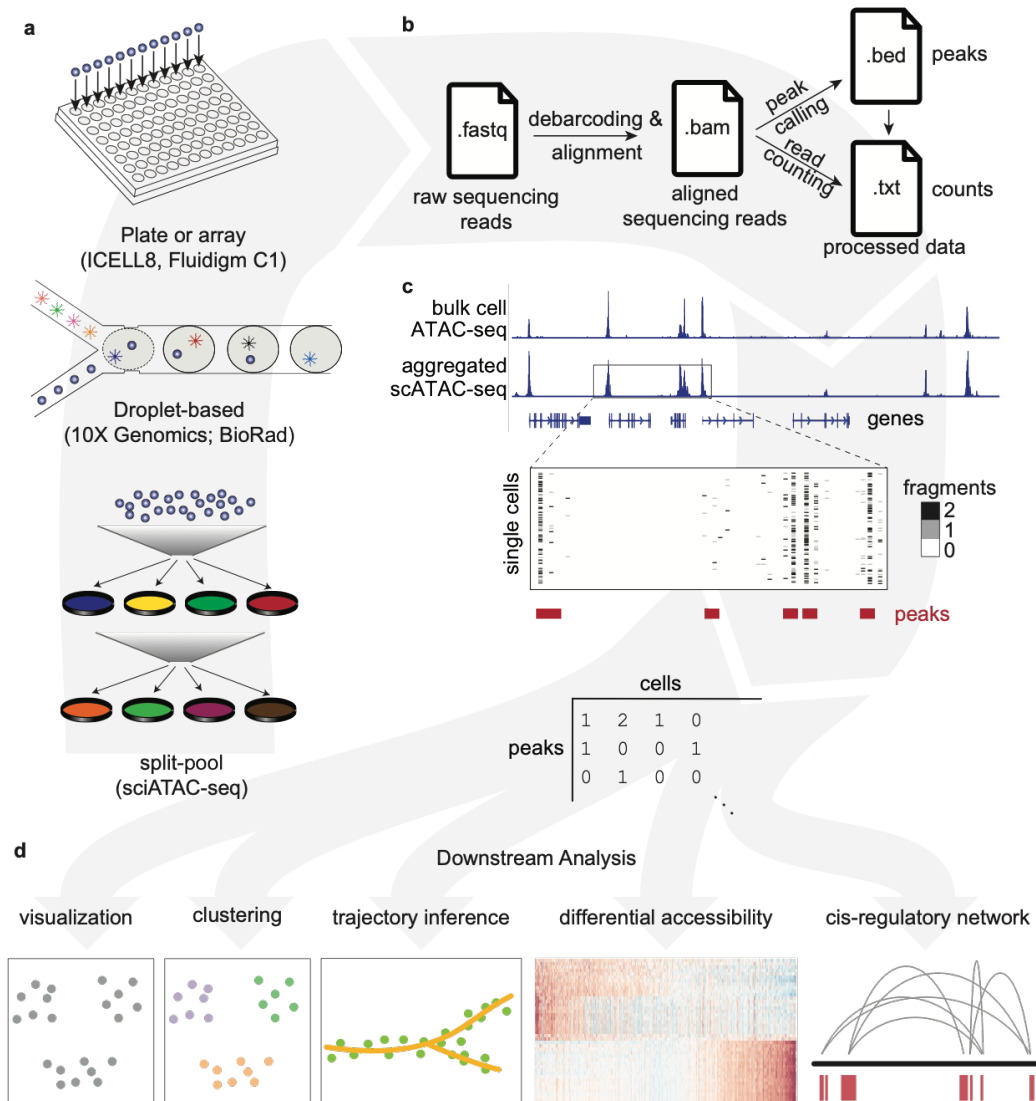
65 single cells holds great promise for uncovering an important component of the
66 regulatory logic of gene expression programs. Enabled by advances in array-based
67 technologies, droplet microfluidics and combinatorial indexing through split-pooling[1]
68 (**Fig. 1a**), single-cell Assay for Transposase Accessible Chromatin using sequencing
69 (scATAC-seq) has recently overcome previous limitations of technology and scale to
70 generate chromatin accessibility data for thousands of single cells in a relatively easy
71 and cost-effective manner.

72

73 However, the analysis of scATAC-seq data presents methodological challenges distinct
74 from those of single-cell transcriptomic (scRNA-seq) data. The primary difficulty arises
75 from a difference in the number of RNA vs DNA molecules available for profiling in
76 single cells. While for an expressed gene several RNA molecules are present in a single
77 cell, scATAC-seq assays profile DNA, a molecule which is present in only few copies
78 per cell (two in a diploid organism). The low copy number results in an inherent per-
79 cell data sparsity, where only 1-10% of expected accessible peaks are detected in single
80 cells from scATAC-seq data, compared to 20-50% of expressed genes detected in single
81 cells from scRNA-seq data. This emphasizes the need to recover informative features
82 from sparse data to assess variability between cells in scATAC-seq analyses. Further,
83 determination of which features best define cell state is currently unclear.

84

85 The difference in readout (gene expression versus chromatin accessibility) has also
86 motivated a variety of approaches to selecting informative features in scATAC-seq
87 methods. While most processing pipelines share common upstream processing steps
88 (i.e. alignment, peak calling, and counting; **Fig. 1b**), existing computational approaches
89 differ in the way they obtain a feature matrix for downstream analyses. For example,
90 some methods select features based on the sequence content of accessible regions (e.g. k -
91 mer frequencies[2, 3] or transcription factor (TF) motifs [3]), whereas other methods
92 select features based on the genomic coordinates of the accessible regions (e.g. extended
93 promoter regions to determine chromatin activity surrounding genes [2, 4]). Finally, the
94 potential feature set in scATAC-seq, which includes genome-wide regions of accessible
95 chromatin (**Fig. 1c**), is typically 10-20x the size of the feature set in scRNA-seq
96 experiments (which is defined and limited by the number of genes expressed). This
97 larger feature set could be valuable in distinguishing a wider variety of cell populations
98 and inferring the dynamics underlying cell organization into complex tissues[5].



99
100

101 **Figure 1.** Schematic overview of single cell ATAC-seq assays and analysis steps. (a)
102 Single cell ATAC libraries are created from single cells that have been exposed to the
103 Tn5 transposase using one of three protocols: 1) Single cells are individually barcoded
104 by a split-and-pool approach where unique barcodes added at each step can be used to
105 identify reads originating from each cell 2) microfluidic droplet-based technologies
106 provided by 10x Genomics and BioRad are used to extract and label DNA from each
107 cell or 3) each single cell is deposited into a multi-well plate or array from ICELL8 or
108 Fluidigm C1 for library preparation. (b) After sequencing, the raw reads obtained in
109 .fastq format for each single cell are mapped to a reference genome, producing aligned
110 reads in .bam format. Finally, peak calling and read counting return the genomic

111 position and the read count files in .bed and .txt format, respectively. Data in these file
112 formats is then used for downstream analysis. (c) ATAC-seq peaks in bulk samples can
113 generally be recapitulated in aggregated single cell samples, but not every single cell
114 has a fragment at every peak. A feature matrix can be constructed from single cells (e.g.,
115 by counting the number of reads at each peak for every cell). (d) Following construction
116 of the feature matrix, common downstream analyses including visualization, clustering,
117 trajectory inference, determination of differential accessibility, and the prediction of cis-
118 regulatory networks can be performed using the methods benchmarked in this
119 manuscript.

120

121 However, the novelty and assay-specific challenges associated with these large-scale
122 scATAC-seq datasets and the lack of analysis guidelines have resulted in diverging
123 computational strategies to aggregate data across such an immense feature space with
124 no clear indication as to which strategy or strategies are most advantageous.

125

126 Here, we provide the first benchmark assessment of computational methods for the
127 analysis of scATAC-seq data. We discuss the impact of feature matrix construction
128 strategies (e.g. sequence content-based vs. genomic coordinates) on common
129 downstream analysis, with a focus on clustering and visualization. This comprehensive
130 survey of current available methods provides user-specific recommendations for best
131 practices that aim to maximize inference-capability for current and future scATAC-seq
132 workflows. Importantly, we provide more than 100 well-documented Jupyter
133 Notebooks (<https://github.com/pinellolab/scATAC-benchmarking/>) to easily reproduce
134 our analyses. We anticipate that this will be a valuable resource for future scATAC-seq
135 benchmark studies.

136

137 **Results**

138 **Benchmark Framework**

139 For this benchmarking study we created an unbiased framework to qualitatively and
140 quantitatively survey the ability of available scATAC-seq methods to featurize
141 chromatin accessibility data. Evaluated using this framework were several datasets of
142 divergent size and profiling technologies. Using widely accepted quantitative metrics,

143 we explored how differences in feature matrix construction influence outcomes in
144 exploratory visualization and clustering, two common downstream analyses. The
145 general overview of our framework is presented in **Fig. 2**.

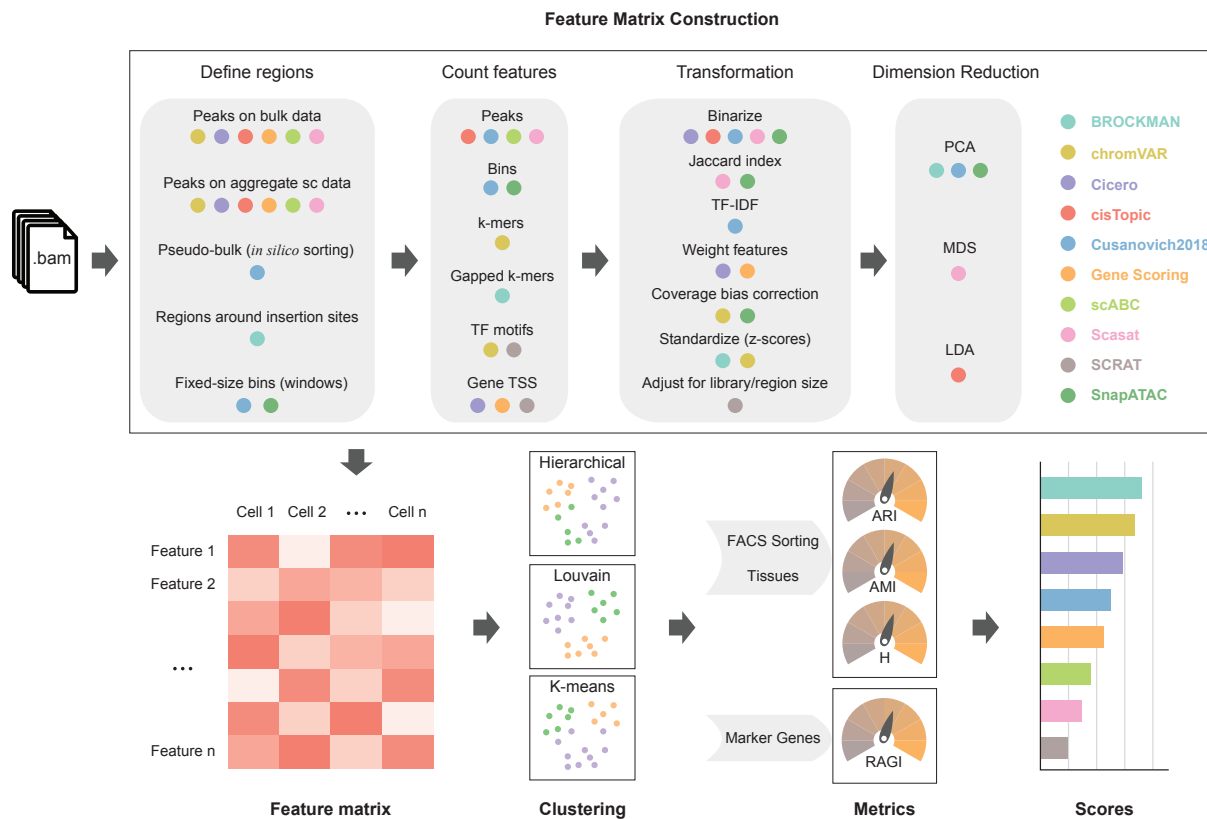
146 For this study we collected public data from three published studies (aligned files in
147 BAM format) and generated ten simulated datasets with various coverages and noise
148 levels (see **Methods**). To calculate feature matrices for downstream analysis, for each
149 method we followed the guidelines provided in the documentation in the original study
150 or as suggested by the respective authors. After feature matrix construction, we used
151 three commonly used clustering approaches (K-means, Louvain and Hierarchical
152 Clustering)[6] and UMAP[7] projection to find putative subpopulations and visualize
153 cell-to-cell similarities for each method. Next, the quality of the clustering solutions was
154 evaluated by adjusted random index (ARI), adjusted mutual information (AMI) and
155 homogeneity (H) when FACS-sorting labels or tissues were available (gold standard); or
156 by a proposed Gini-index-based metric called Residual Average Gini Index (RAGI)
157 when only known marker genes were available (silver standard). Finally, based on
158 these metrics, the methods were ranked by the quality of their clustering solutions
159 across datasets.

160

161 **Methods overview and featurization of chromatin accessibility data**

162 Several computational methods have been developed to address the inherent sparsity
163 and high dimensionality of single cell ATAC-seq data, including BROCKMAN[3],
164 chromVAR[2], Cicero[8], cisTopic[9], *Cusanovich2018*[1, 10, 11], Gene Scoring[12],
165 scABC[13], Scasat[14], SCRAT[4], and SnapATAC[15]. Based on the proposed workflow
166 of each method, we were able to compute different feature matrices defined as a
167 features-by-cells matrix (e.g. read counts for each cell (columns) in a given open
168 chromatin peak *feature* (rows)) that could then be readily used for downstream analyses
169 such as clustering. Starting from single cell BAM files, the feature matrix construction
170 can be roughly summarized into four different common modules: *define regions, count*
171 *features, transformation, and dimensionality reduction* as illustrated in **Fig. 2**. Not every
172 method uses all steps, therefore we provide below, a short summary of the strategies

173 adopted by each method and a *per module* discussion to highlight key similarities and
 174 differences (for a more detailed description of each strategy see **Methods**).



175

176 **Figure 2.** Benchmarking workflow. Starting from aligned read files in .bam format,
 177 feature matrices were constructed using each method. The feature matrix construction
 178 techniques used by each method were grouped into four broad categories: *Define regions*,
 179 *Count features*, *Transformation* and *Dimension Reduction*. A colored dot under a technique
 180 indicates that the method (signified by the respective color in the legend on the right)
 181 uses that technique. For each method, feature matrix files (defined as columns as cells
 182 and rows as features) are calculated and used to perform hierarchical, Louvain and k-
 183 means clustering analysis. For datasets with a ground truth such as FACS-sorting labels
 184 or known tissues, clustering evaluation was performed according to the Adjusted
 185 Random Index (ARI), Adjusted Mutual Information (AMI) and homogeneity (H) scores.
 186 For datasets without ground truth, the clustering solutions were evaluated according to
 187 a Residual Average Gini Index (RAGI), a metric that compares cluster separation based
 188 on known marker genes against housekeeping genes. Lastly, a final score is assigned to
 189 each method.

190 Briefly, BROCKMAN[3] represents genomic sequences by gapped k-mers (short DNA
191 sequences of length k) within transposon integration sites and infers the variation in k-
192 mer occupancy using principal component analysis (PCA). chromVAR[2] estimates the
193 dispersion of chromatin accessibility within peaks sharing the same feature, e.g. motifs
194 or k-mers. Cicero[8] calculates a gene activity score based on accessibility at a promoter
195 region and the regulatory potential of peaks nearby. cisTopic[9] applies Latent Dirichlet
196 Allocation (LDA) (a Bayesian topic modeling approach commonly used in natural
197 language processing) to identify cell states from topic-cell distribution and explore cis-
198 regulatory regions from region-topic distribution. Previous approaches that utilize
199 latent semantic indexing (LSI) (termed here as *Cusanovich2018*)[1, 10, 11] first partition
200 the genome into windows, normalize reads within windows using the term frequency-
201 inverse document frequency transformation (TF-IDF), reduce dimensionality using
202 singular value decomposition (SVD), and perform a first-round of clustering (referred
203 to as ‘in silico cell sorting’) to generate clades and call peaks within them. Finally, the
204 clusters are refined with a second-round of clustering after TF-IDF and SVD based on
205 read counts in peaks. The Gene Scoring method[12] assigns each gene an accessibility
206 score by summarizing peaks near its transcription start site (TSS) and weighting them
207 by an exponential decay function based on their distances to the TSS. scABC[13] first
208 calculates a global weight for each cell by taking into account the number of distinct
209 reads in the regions flanking peaks (to estimate the expected background). Based on
210 these weights, it then uses weighted k-medoids to cluster cells based on the reads in
211 peaks. Scasat[14] binarizes peak accessibility and uses multidimensional scaling (MDS)
212 based on the Jaccard distance to reduce dimensionality before clustering. SCRAT[4]
213 summarizes read counts on different regulatory features (e.g. transcription factor
214 binding motifs, gene TSS regions). SnapATAC[15] segments the genome into
215 uniformly-sized bins and adjusts for differences in library size between cells using a
216 regression-based normalization method; finally PCA is performed to select the most
217 significant components for clustering analysis.

218 Define Regions

219 An essential aspect of feature matrix construction is the selection of a set of regions to
220 describe the data (e.g. putative regulatory elements such as peaks, promoters etc.). Most
221 methods described above, including chromVAR, Cicero, cisTopic, Gene Scoring, scABC,
222 and Scasat, define regions based on peak calling from either a reference bulk ATAC-seq
223 profile or an aggregated single cell ATAC-seq profile. *Cusanovich2018*, as briefly

224 mentioned above, instead of aggregating single cell to call peaks, first creates pseudo-
225 bulk clades by performing hierarchical clustering on the TF-IDF and SVD transformed
226 matrix using the top frequently accessible windows. Then peaks are called by
227 aggregating cells within each pseudo-bulk clade. In addition to relying on peaks, some
228 methods have proposed different strategies. BROCKMAN uses the union of regions
229 around transposon integration sites. *Cusanovich2018* (before *in silico* sorting) and
230 SnapATAC segment the genomes into fixed-size bins (windows) and count features
231 within each bin.

232 Count Features

233 Once feature regions are defined, raw features within these regions are counted. Note
234 that some methods (e.g. chromVAR) may support the counting of multiple features. For
235 *cisTopic*, *Cusanovich2018*, *scABC*, and *Scasat*, reads overlapping peaks are counted. For
236 *Cusanovich2018* (before the *in silico* sorting step) and SnapATAC, reads overlapping bins
237 are counted. k-mers are counted under peaks for chromVAR while gapped k-mers are
238 counted for BROCKMAN around transposase cut sites. Similarly, transcription factor
239 motifs (e.g. from the JASPAR database[16]) can be used as features by counting reads
240 overlapping their binding sites in peaks (chromVAR) or genome-wide (SCRAT). If
241 predefined genomic annotations such as coding genes are given, Gene Scoring, Cicero,
242 and SCRAT use gene TSSs as anchor points to calculate gene enrichment scores based
243 on reads nearby or just within peaks nearby.

244 Transformation

245 After building the initial raw feature matrix using the counting step, different
246 transformation methods can be performed. Binarization of read counts is used by five
247 out of the ten evaluated methods: Cicero, *cisTopic*, *Cusanovich2018*, *Scasat*, and
248 SnapATAC. (**Fig. 2**). This step is based on the assumption that each site is present at
249 most twice (for diploid genomes) and that the count matrix is inherently sparse.
250 Binarization is advantageous in alleviating challenges arising from sequencing depth or
251 PCR amplification artifacts. SnapATAC and *Scasat* convert the binary count matrix into
252 a cell-pairwise Jaccard index similarity matrix. *Cusanovich2018* normalizes the binary
253 count matrix using the TF-IDF transformation. Cicero weights feature sites by their co-
254 accessibility, while Gene Scoring weights sites by a decaying function based on its
255 distance to a gene TSS. Both chromVAR and SnapATAC perform a read coverage bias
256 correction to account for the influence of sample depth. *scABC* also implements a

257 similar step but calculates a weight for each cell; even if these weights are not used to
258 transform the matrix, they are used later in the clustering procedure. SCRAT adjusts for
259 both library size and region length. chromVAR creates ‘background’ peaks consisting of
260 an equal number of peaks matched for both average accessibility and GC content to
261 calculate bias-corrected deviation. Both BROCKMAN and chromVAR compute z-scores
262 to measure the gain or loss of chromatin accessibility across cells.

263 Dimensionality Reduction

264 In the final step before downstream analysis, several methods apply different
265 dimensionality reduction techniques to project the cells into a space of fewer
266 dimensions. This step can refine the feature space mitigating redundant features and
267 potential artifacts, and potentially reducing the computation time of downstream
268 analysis (**Fig. 2**). PCA is the most commonly used method (used by BROCKMAN,
269 SnapATAC, and *Cusanovich2018*). cisTopic uses latent Dirichlet allocation (LDA) to
270 generate two distributions including topic-cell distribution and region-topic
271 distribution. Choosing the top topics based on the topic-cell distribution reduces the
272 dimensionality. Scasat uses multidimensional scaling (MDS). When reviewing the
273 different methods to include in our benchmark, we noticed that not all methods
274 perform a dimensionality reduction step, which could skew the relative performance
275 across methods. Therefore, for chromVAR, Cicero (gene activity score), Gene Scoring,
276 scABC, and SCRAT, we considered in addition to the original feature matrix, also a new
277 feature matrix after PCA transformation, since this is simple and commonly used
278 technique for dimensionality reduction.

279 To better evaluate the effects of different modules including *define regions*, *count features*,
280 *transformation*, and *dimensionality reduction*, we also considered a simple control method,
281 referred to as Control-Naïve, by combining the most common and simple steps for
282 building a feature matrix, i.e. counting reads within peaks to obtain a peaks-by-cells
283 raw count matrix and then performing PCA on it (the number of top principal
284 components was determined based on the elbow plot for all the methods). Since the
285 feature matrix of scABC is also a peaks-by-cells raw count matrix, this matrix after PCA
286 will correspond to the one obtained by the Control-Naïve method (to avoid
287 redundancies, in our assessment we refer to this matrix as Control-Naïve).

288 We also noticed that some methods might slightly diverge from the proposed four
289 modules common framework. For example, Cicero calculates gene activity scores by
290 first performing two transformations (binarize and weight features) and then
291 performing the counting step around the annotated TSS. We believe the proposed
292 modularization of the of the feature matrix construction can still serve as a useful
293 framework to represent the core components of the different methods and provides an
294 intuitive and informative summary of the diverse scATAC-seq methodologies.

295 Once dimensionality reduction is completed, the transformed feature matrix can be
296 used for unbiased clustering, visualization, or other downstream analyses. Here we
297 have used the final feature matrices generated by each scATAC-seq analysis method,
298 and evaluated their performance in uncovering different populations by unsupervised
299 clustering.

300 **Clustering approaches and metrics used for performance evaluation**

301 This study employed three diverse types of commonly used unsupervised clustering
302 methods for single cell analysis [6]: K-means clustering, Hierarchical Clustering, and the
303 Louvain community detection algorithm (see **Methods**).

304 Clustering results were evaluated by three commonly used metrics: adjusted random
305 index (ARI), adjusted mutual information (AMI) and homogeneity when a gold
306 standard solution was available (known labels for the simulation data and FACS-sorted
307 cell populations or known tissues for the real datasets). We propose a Gini-index-based
308 metric called Residual Average Gini Index (RAGI), which was used to evaluate the
309 clustering results when no ground truth was available and only a few marker genes
310 were known by which populations could be discriminated (see **Methods**). For each
311 metric, we defined the *clustering score* as the highest score amongst the three clustering
312 methods, i.e. the score which corresponded to the clustering solution that maximized
313 the metric.

314 This framework allowed for benchmarking the ability of each strategy to featurize
315 chromatin accessibility data and its impact on important downstream analyses such as
316 clustering and visualization. The following sections present the results of this
317 evaluation for all above-described synthetic and real scATAC-seq datasets.

318 Clustering performance on simulated datasets

319 We simulated 10 scATAC-seq datasets using available bulk ATAC-seq datasets with
320 clear annotations from bone marrow and erythropoiesis[5, 17] using varying noise
321 levels and read coverages. Briefly, to generate the peak by cell matrices, we defined a
322 noise parameter (between 0 and 1) as the proportion of reads occurring in a random
323 peak from one of the sorted populations. The remaining proportion of reads was
324 distributed as a function of the bulk sample (see Methods). A feature matrix with a
325 noise level of 0 preserved perfectly the underlying cell type specificity of the reads
326 within peaks. Conversely, a feature matrix with a noise level of 1, contained no
327 information to discriminate cell types based on the reads within peaks. In our study, we
328 considered three noise levels: no noise (0), moderate noise (0.2) and high noise (0.4). To
329 better and more fairly evaluate the contribution of the core steps of each method (i.e.
330 *count features, transformation and dimensionality reduction*) regardless of the preprocessing
331 steps usually excluded from these methods (reads filtering, alignment, peak calling,
332 etc.), we compared the performance of each method using a set of predefined peak
333 regions from bulk ATAC-seq datasets. We selected the top 80,000 peaks based on the
334 number of cells in which peaks were observed (each peak that was present in at least
335 one cell) for all methods and all synthetic datasets.

336 Using the bulk ATAC-seq bone marrow dataset, we simulated five additional datasets
337 to explore the effect of coverage on clustering performance (5,000 fragments, 2,500
338 fragments, 1,000 fragments, 500 fragments, 250 fragments respectively per cell).

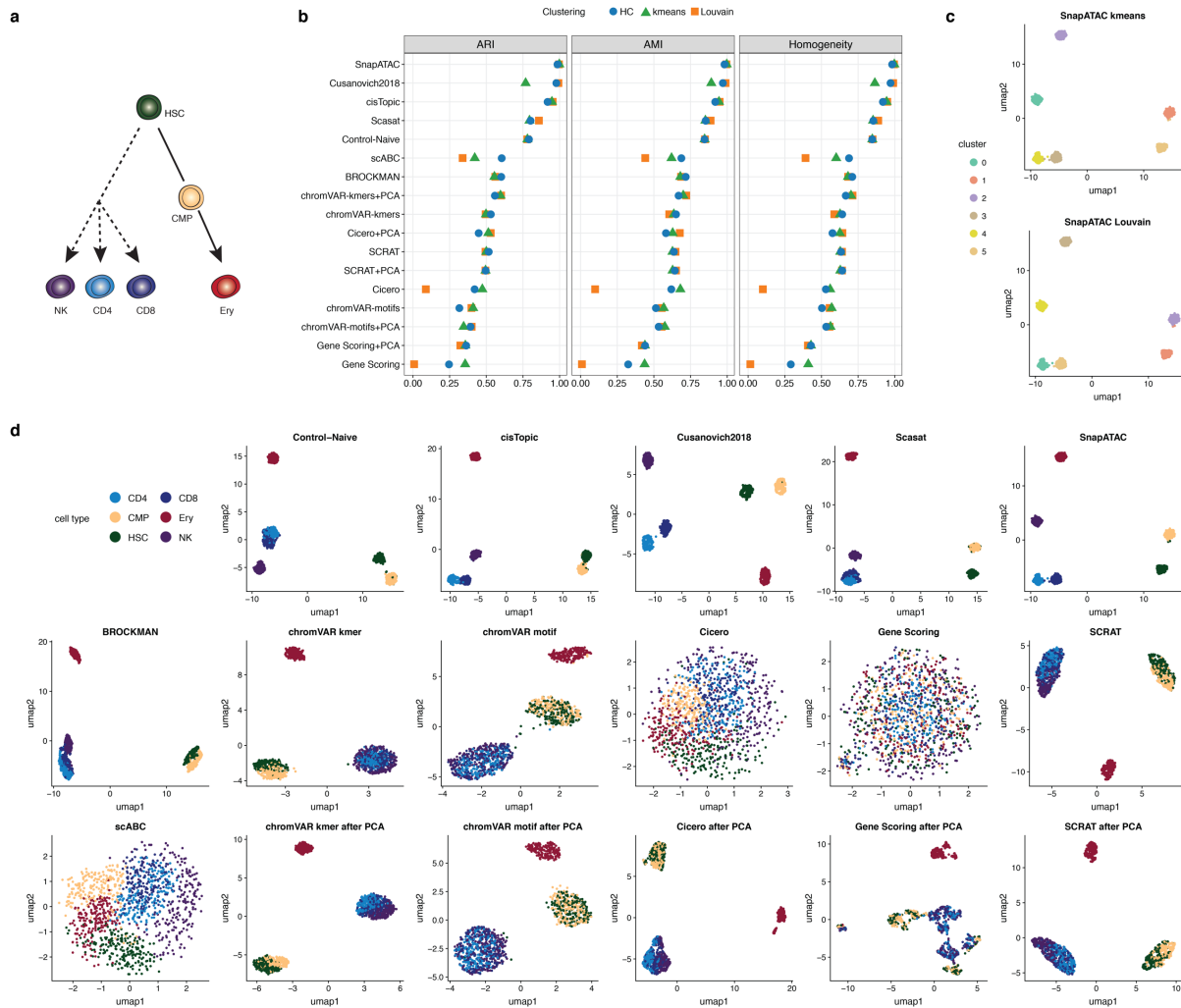
339 Each method was used to analyze all synthetic datasets as suggested in the method
340 documentation (see **Sup Note 1** and **Sup Fig. 1**).

341 Simulated bone marrow datasets

342

343 We generated chromatin accessibility profiles (2,500 fragments per cell) based on six
344 different FACS-sorted bulk cell populations: hematopoietic stem cells (HSCs), common
345 myeloid progenitor cells (CMPs), erythroid cells (Ery), and other three lymphoid cell
346 types: natural killer cells (NK), CD4 and CD8 T-cells (see **Fig. 3a**). We used ARI, AMI
347 and homogeneity metrics to compare the clustering solutions with the known cell type
348 labels (**Fig. 3b**, **Sup Fig. 2**, **Sup Table 1**). The top three methods based on these

349 simulation settings were cisTopic, *Cusanovich2018*, and SnapATAC. They performed
350 equally well with no noise and moderate noise (with clustering scores close to 1.0) (**Sup**
351 **Fig. 2, Sup Table 2**). At a noise level of 0.4, the methods showed more separation in
352 performance accordingly to the three metrics (**Fig. 3b, Sup Table 3**). SnapATAC,
353 *Cusanovich2018*, and cisTopic clearly outperformed the Control-Naïve method with
354 consistently higher clustering scores across all metrics. Scasat performed slightly better
355 than the Control-Naïve method, and the remaining methods under-performed relative
356 to the Control-Naïve method. For scABC (i.e. peaks-by-cells raw count matrix),
357 Hierarchical Clustering performs much better than the other two clustering methods.
358 chromVAR performance using k-mers as features was superior to the approach using
359 motifs. Another k-mer-based method, BROCKMAN demonstrated similar performance
360 to the k-mer-based chromVAR method. Motif-based SCRAT performed better than
361 motif-based chromVAR. Both Cicero gene activity scores and Gene Scoring (which
362 summarize the chromatin accessibility around coding annotations without a
363 dimensionality reduction step) generally performed poorly. PCA boosted performance
364 of scABC, Cicero, and Gene Scoring. This step improved clustering performance
365 regardless of the clustering method (also we noted again that scABC after PCA is
366 equivalent to the Control-Naïve method), especially for the Louvain approach. PCA
367 also slightly boosted performance of the k-mer-based chromVAR but did not markedly
368 improve the results of the motif-based chromVAR or SCRAT analyses.
369



370

371

372 **Figure 3.** Benchmarking results in simulated bone marrow datasets at a noise level of 0.4

373 and a coverage of 2,500 fragments. **(a)** Cell types used to create the simulated dataset. **(b)**

374 Dot plot of scores for each metric to quantitatively measure the clustering performance

375 of each method, sorted by maximum ARI score. **(c)** The two top-scoring pairings of

376 scATAC-seq analysis method and clustering technique. Cell cluster assignments from

377 each method are shown using the colors in the legend on the left. **(d)** UMAP visualization

378 of the feature matrix produced by each method for the simulated dataset. Individual cells

379 are colored indicating the cell type labels shown in (a).

380

381 We next investigated qualitatively the obtained clustering solutions, using the

382 respective feature matrices to project the cells onto a 2-D space using UMAP and

383 colored them based on the obtained clustering solutions (**Sup Fig. 3**) or based on the

384 true population labels used to generate the data (**Fig. 3d**). The top two clustering
385 solutions based on the ARI (SnapATAC with k-means and SnapATAC with Louvain)
386 are shown for ease of comparison (**Fig. 3c**).

387
388 *Cusanovich2018* and SnapATAC are the only two methods that clearly separated all six
389 populations. *cisTopic* slightly mixed CD4 and CD8 T-cells. *Scasat* and the Control-
390 Naïve method failed to separate CD4 and CD8 T-cell populations. BROCKMAN slightly
391 mixed NK with CD4 and CD8 T-cells and could not further separate CD4 and CD8 T-
392 cells. It also failed to clearly separate HSC and CMP. Both kmer-based and motif-based
393 chromVAR as well as SCRAT could only separate the Ery population while failing to
394 separate HSC and CMP as well as CD4, CD8 T-cells, and NK. The chromVAR k-mers-
395 based method mixed HSC and CMP to a lesser extent compared to the motifs-based
396 method. There was no clear separation of cells using *scABC* (the peaks-by-cells raw
397 count matrix), *Cicero*, or Gene Scoring. We observed that PCA clearly improved the
398 separation of cell populations for *Cicero* and Gene Scoring. It also slightly improved the
399 separation of CD4, CD8 T-cells, and NK populations by k-mer-based chromVAR. No
400 clear improvement was observed for the motif-based chromVAR, or SCRAT methods.
401 We further observed that a lack of visual separation of cell types in the UMAP plots
402 (*scABC*, *Cicero*, and Gene Scoring), corresponded with substantial variation between
403 the performances of the three clustering methods, showing better performance in the k-
404 means clustering (**Fig. 3b,d**).

405
406 All methods except for *Cusanovich2018* and SnapATAC demonstrated declining
407 performance with increased noise level (**Sup Fig. 2, 4a**). *Cusanovich2018* and SnapATAC
408 were more robust to noise, showing no noticeable changes at increasing noise levels,
409 while *cisTopic* was slightly more sensitive to noise; its performance dropped markedly
410 when the noise level was increased to 0.4.

411
412 Next, the effect of the coverage on clustering performance was investigated. We
413 progressively decreased the number of fragments per cell from a high coverage of 5,000
414 fragments, to a medium coverage of 2,500 fragments and 1,000 fragments, then to a low
415 coverage of 500 fragments and finally to 250 fragments. The performance of all methods
416 declined as coverage was decreased. (**Sup Fig. 4b, Sup Fig. 5, Sup Table 4-5-6-7-8**).
417 *Cusanovich2018*, SnapATAC, *Scasat*, and Control-Naïve are relatively robust to low

418 coverage and outperform other methods. cisTopic worked well with high coverage but
419 in contrast to the above listed methods, was more sensitive to lower coverages (**Sup Fig.**
420 **5e**).

421

422 Simulated erythropoiesis datasets

423 Following the simulation of discrete sorted cell populations, we simulated three
424 scATAC-seq datasets aimed at mimicking the continuous developmental erythropoiesis
425 process and encompassing the following twelve populations: hematopoietic stem cells
426 (HSC), common myeloid progenitors (CMP), megakaryocyte-erythroid progenitor
427 (MEP), multipotent progenitors (MPP), myeloid progenitors (MyP), colony forming
428 unit-erythroid (CFU-E), proerythroblasts (ProE1), proerythroblasts (ProE2), basophilic
429 erythroblasts (BasoE), polychromatic erythroblasts (PolyE), orthochromatic
430 erythroblasts (OrthoE) and OrthoE and reticulocytes (Orth/Ret). These datasets were
431 generated as before with three noise levels (0, 0.2 and 0.4) and with 2,500 fragments per
432 cell.

433

434 To first quantitatively evaluate the clustering solutions we used ARI, AMI and the
435 homogeneity metrics (**Sup Fig. 6 and Sup Table 9**). Without noise, SnapATAC, cisTopic
436 BROCKMAN, *Cusanovich2018*, and Scasat consistently outperform the Control-Naïve
437 across the three metrics (**Sup Fig.6a**). chromVAR as before, performs better using k-
438 mers as features than when using motifs. SCRAT and scABC work as well as k-mers-
439 based chromVAR. Again, methods such as Cicero and Gene Scoring that only
440 summarize chromatin accessibility around TSS perform poorly. For scABC, Cicero and
441 Gene Scoring, we also notice that there are significant discrepancies between the three
442 clustering methods, but their performances become similar after PCA (scABC after PCA
443 is equivalent to the Control-Naïve method). Again, we observe that PCA can
444 significantly improve the clustering performance of Louvain for scABC, Cicero and
445 Gene Scoring but not for chromVAR and SCRAT.

446

447 As before, to qualitatively assess population separation, we inspected UMAP
448 projections applied to the noise-free simulated dataset (**Sup Fig. 6a**). In accordance with
449 the quantitative comparison, cisTopic, *Cusanovich2018*, SnapATAC, and BROCKMAN
450 demonstrate better performance in separating cell types compared to the Control-Naïve
451 method and are able to further separate BasoE and PolyE. Moreover, SnapATAC can

452 clearly distinguish CFU-E, ProE1, ProE2 while cisTopic, *Cusanovich2018*, and
453 BROCKMAN are only able to separate ProE2 out of these three populations. Scasat
454 performs similarly to the Control-Naïve method. chromVAR with k-mers as features
455 and SCRAT are able to isolate six major groups including HSCs-MPPs, CMP, MEP,
456 Myp, CFU-E-ProE1-ProE2, and BasoE-PolyE-OrthoE-Orth/Ret. chromVAR with k-mers
457 performs well in preserving the order of CFU-E-ProE1-ProE2 and BasoE-PolyE-OrthoE-
458 Orth/Ret. SCRAT can further separate BasoE-PolyE from OrthoE-Orth/Ret while mixing
459 up CFU-E-ProE1-ProE2. As before, we noticed that chromVAR using k-mers as features
460 obtained a better separation of cell types than when using motifs. scABC is able to
461 preserve well the order of major groups in a continuous way but fails to separate CFU-
462 E-ProE1-ProE2 and OrthoE-Orth/Ret. Cicero gene activity score and Gene Scoring
463 mixed different cell types but after a simple PCA step they clearly separate cells into
464 three major groups. scABC did not perform well and produced small noisy clusters
465 with different cell types mixed together.

466

467 As expected, we observed that increasing the level of noise resulted in clustering
468 performance decrease and a decline of visual separation of cell types for all the methods
469 (**Sup Fig. 4c, Sup Fig. 6, Sup Table 10-11**). SnapATAC, cisTopic, and *Cusanovich2018*
470 performed reasonably well when increasing the noise level, with SnapATAC the most
471 robust among the three.

472

473 **Clustering performance on real datasets**

474 Following the benchmark of the synthetic datasets, we assessed the performance of the
475 methods on real datasets. These datasets were generated using different technologies:
476 the Fluidigm C1 array[18], the 10X Genomics droplet based scATAC platform, and a
477 recently-optimized split-pool protocol[1]. Each real dataset used was fundamentally
478 different in its cellular makeup as well as size and subpopulation organization. Notably,
479 as ‘true positive’ labels are not always available, in addition to the metrics used on the
480 simulated datasets, here we introduced the RAGI, a simple metric based on the Gini
481 Index that can be adopted when marker genes for the expected populations are known
482 (**see Methods**). In our assessment of *Cusanovich2018*, to make a fair comparison, we use
483 first the same set of peaks used for other methods instead of the peaks called from its
484 pseudo-bulk-based procedure. However, since this strategy may be important for the

485 final clustering performance, the pseudo-bulk based peak calling strategy is tested and
486 discussed in a subsequent section.

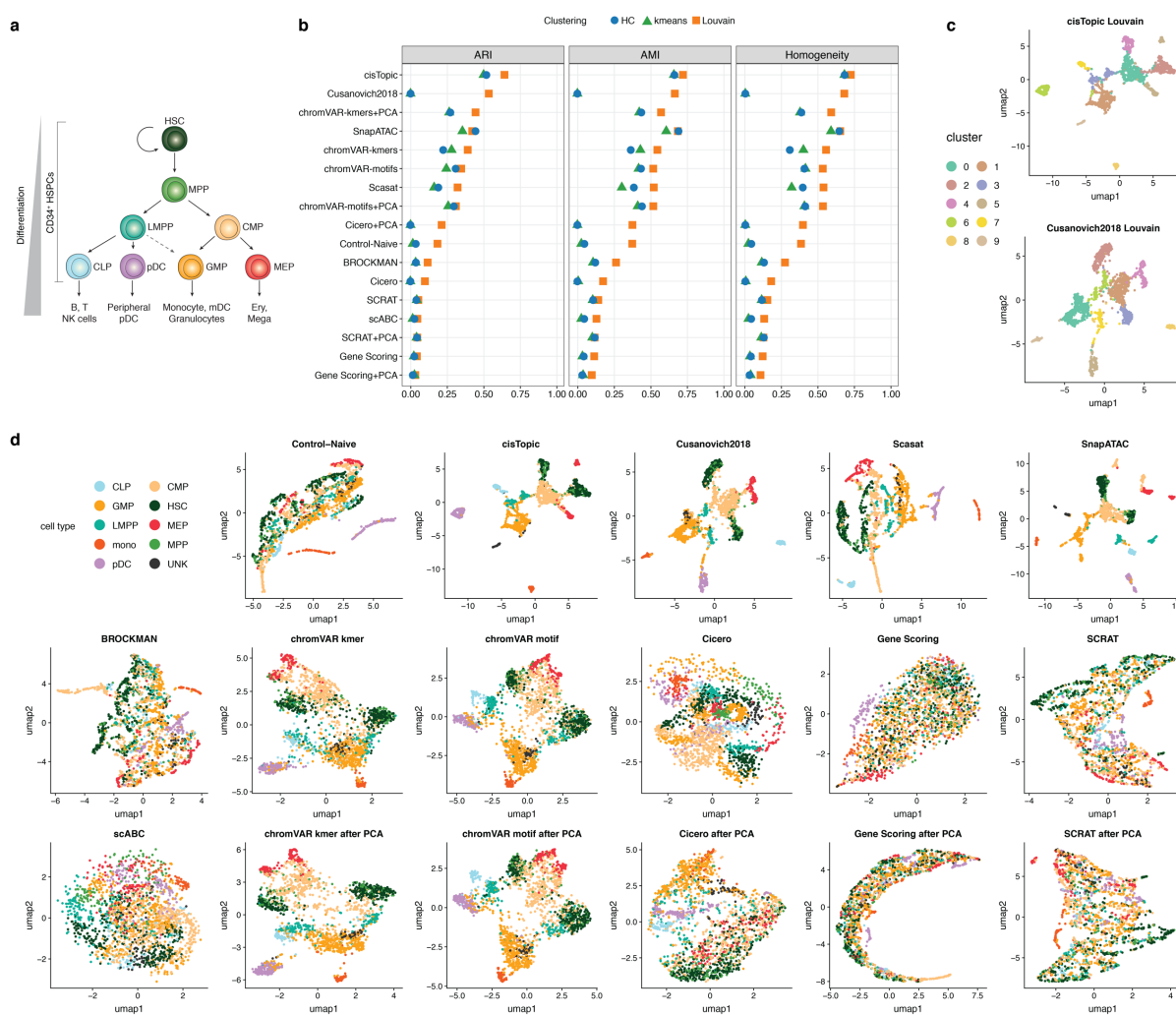
487 *Buenrostro2018 dataset*

488 The first and smallest dataset we used in our benchmarking contains single cell ATAC-
489 seq data from the human hematopoietic system (hereafter *Buenrostro2018*)[18]. This
490 dataset consists of 2034 hematopoietic cells that were profiled and FACS-sorted from 10
491 cell populations including hematopoietic stem cells (HSCs), multipotent progenitors
492 (MPPs), lymphoid-primed multipotent progenitors (LMPPs), common myeloid
493 progenitors (CMPs) and granulocyte-macrophage progenitors (GMPs), GMP-like cells,
494 megakaryocyte-erythroid progenitors (MEPs), common lymphoid progenitor (CLPs),
495 monocytes (mono) and plasmacytoid dendritic cells (pDCs). **Fig. 4a** illustrates the
496 roadmap of hematopoietic differentiation. For this dataset, the FACS-sorting labels are
497 used as gold standard. The analysis details for each method are documented in **Sup**
498 **Note 2.**

499
500 We started by evaluating the clustering solutions based on the feature matrices
501 generated by the different methods. We used the same metrics used for the synthetic
502 datasets: ARI, AMI and homogeneity (**Fig. 4b, Sup Table 12**). *cisTopic*, *Cusanovich2018*,
503 *chromVAR*, *SnapATAC*, and *Scasat* outperform the other methods across all three
504 metrics. We also observed that *chromVAR* with k-mers or TF motifs and with or
505 without PCA performs consistently well. As before, k-mers-based features work better
506 than motif-based features. This can be also observed when comparing *BROCKMAN*,
507 another k-mers-based method, with *SCRAT*, which is a motifs-based method. TSS based
508 methods including *Cicero* and *Gene Scoring* did not perform well. *Cicero* requires a
509 preprocessing step to assess cell similarity; poor performance might be due to the
510 internally incorrectly inferred coordinates (our assessment used the t-SNE procedure as
511 suggested in their documentation). Implementing PCA consistently improves the
512 performance of *scABC* (as mentioned before, *scABC* after PCA is equivalent to the
513 *Control-Naïve* method) and *Cicero* but does not impact the performance of *chromVAR*,
514 *SCRAT*, and *Gene Scoring*. We also observed that for this dataset, Louvain algorithm
515 works consistently well across different metrics and methods and performs better than
516 hierarchical clustering and k-means in almost all the cases.

517

518 We also qualitatively assessed the separation of different cell types by visualizing cells
 519 in UMAP projections based on the FACS-sorted labels (**Fig. 4d**) and clustering solutions
 520 (**Sup Fig. 7**). **Fig. 4c** shows the best two combinations based on ARI: cisTopic with
 521 Louvain and *Cusanovich2018* with Louvain (the complete ranking is presented in **Sup**
 522 **Table 12**).
 523



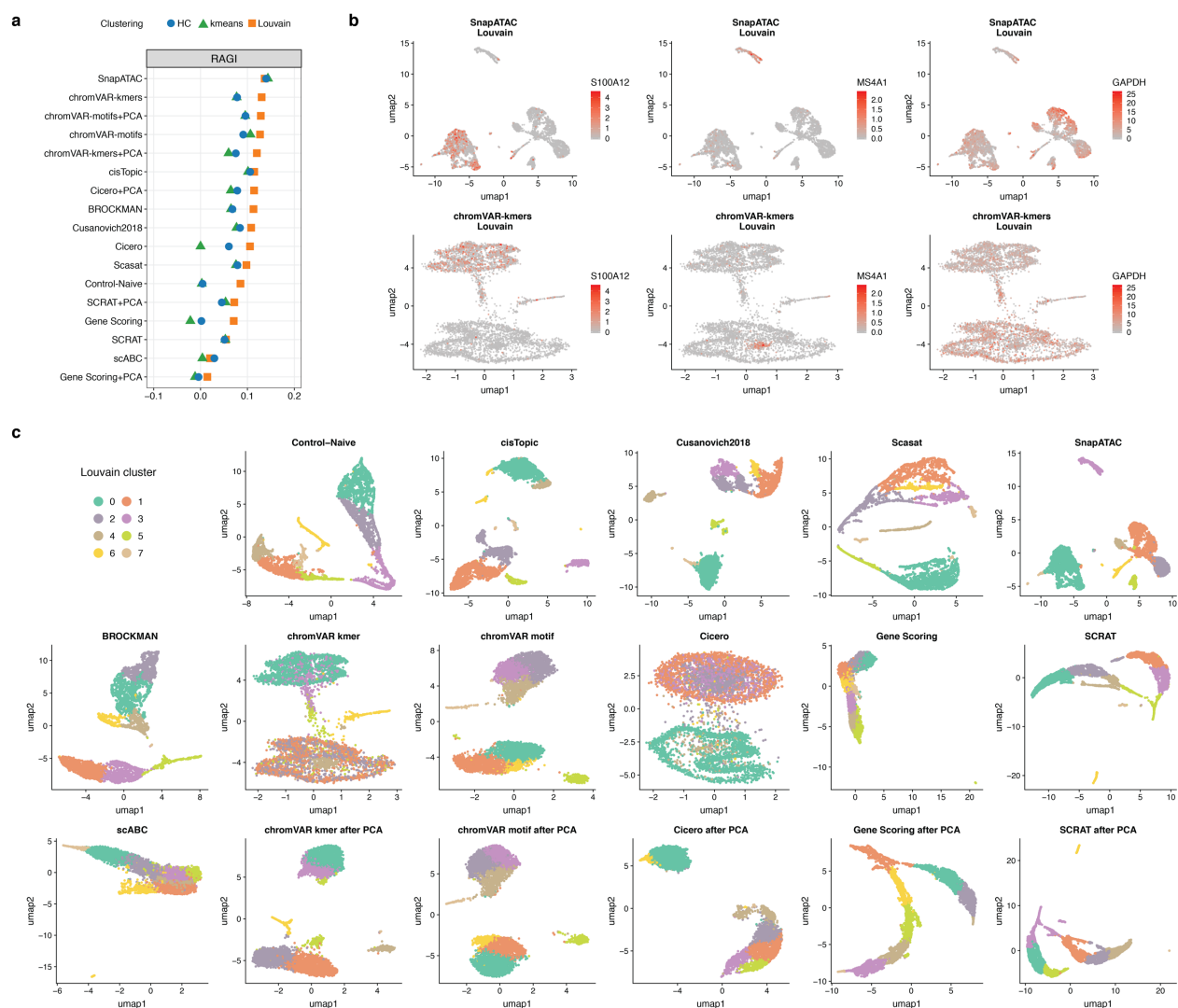
524
 525 **Figure 4.** Benchmarking results using the *Buenrostro2018* scATAC-seq dataset. **(a)**
 526 Developmental roadmap of cell types analyzed. **(b)** Dot plot of scores for each metric to
 527 quantitatively measure the clustering performance of each method, sorted by maximum
 528 ARI score. **(c)** The two top-scoring pairings of scATAC-seq analysis method and
 529 clustering technique. UMAP visualization of the feature matrix produced by each method
 530 for the *Buenrostro2018* dataset. Individual cells are colored indicating the cell type labels
 531 shown in (a).

532 As **Fig. 4d** shows, in accordance with the clustering analyses, *cisTopic*, *Cusanovich2018*,
533 *Scasat*, *SnapATAC*, and *chromVAR* can generally separate cell types, and reasonably
534 capture the expected hematopoietic hierarchy. *cisTopic* and *SnapATAC* show a clear
535 and compact separation among groups, with *SnapATAC* recovering finer structure
536 within each cell type cluster. *chromVAR* with k-mers or motifs corresponds to a more
537 continuous progression of the different cell types. *Control-Naïve* and *BROCKMAN*
538 perform comparably in distinguishing cell types and preserving the continuous
539 hematopoietic differentiation. *Cicero* gene activity scores, *SCRAT*, and *scABC* show
540 ambiguous patterns of distinct cell populations while *Gene Scoring* fails to separate
541 different cell types. For *Cicero* gene activity score, after performing PCA, the separation
542 of different cells is noticeably improved. For *SCRAT*, performing PCA does not show
543 clear improvement.

544 *Peripheral blood mono nuclear cells (PBMCs) 10X dataset*

545 Next, we investigated a recent dataset produced by 10X Genomics profiling peripheral
546 blood mononuclear cells (PBMCs) from a single healthy donor. In this dataset, 5335
547 single nuclei were profiled (~42k read pairs per cell); no cell annotations are provided.
548 Based on recent studies [9, 19], we expected ~8 populations: CD34+, Natural Killer and
549 Dendritic cells, Monocytes, lymphocyte B and lymphocyte T cells, together with
550 terminally differentiated CD4 and CD8 cells. Therefore, we used 8 as the number of
551 expected populations for the clustering procedures. The analysis details for each
552 method are documented in **Sup Note 3**.

553 Several marker genes have been proposed to label the different populations or to
554 annotate clustering solutions for PMBCs [9, 19]. To measure cluster relevance based on
555 these marker genes, we can annotate the clusters (or alternatively any group of cells)
556 according to the accessibility values at those marker genes. In addition, accessibility at
557 marker genes should be more variable between clusters than accessibility at
558 housekeeping genes (since they should be, by definition, more equally expressed across
559 different populations). Based on these ideas, we proposed and calculated the Residual
560 Average Gini Index (RAGI) score (see **Methods**) contrasting marker and housekeeping
561 genes (**Fig. 5a, Sup Table 13**). For reasonable clustering solutions, we expect that the
562 accessibility of marker genes defines clear populations corresponding to one or few
563 clusters, whereas accessibility of the housekeeping genes is broadly distributed across
564 all the clusters.



565
 566 **Figure 5.** Benchmarking results using scATAC-seq data for 5k Peripheral blood
 567 mononuclear cells (PBMCs) from 10x Genomics. **(a)** Dot plot of RAGI scores for each
 568 method, sorted by the maximum RAGI score. A positive RAGI value indicates that a
 569 method is able to produce a clustering of PBMCs in which chromatin accessibility of each
 570 marker gene is high in only a few clusters relative to the number of clusters with high
 571 accessibility of housekeeping genes. **(b)** UMAP visualization of the feature matrix
 572 produced by the top two methods (top row: SnapATAC, bottom row: chromVAR using
 573 kmers). Chromatin accessibility of S100A12 (left, Monocyte marker gene), MS4A1 (center,
 574 B-cell marker gene) and GPDH (right, housekeeping gene) are projected onto the
 575 visualization. **(c)** UMAP visualization of the feature matrix produced by each method for
 576 the 5k PBMCs dataset from 10x genomics. Individual cells are colored indicating cluster
 577 assignments using Louvain clustering.

578 As expected, methods with the highest performance such as SnapATAC and
579 chromVAR, showed a higher average accessibility for just one cluster for the same
580 marker gene, while lower performing methods such as SCRAT or Gene Scoring showed
581 higher average accessibility in multiple clusters for the same marker gene, further
582 motivating the use of the RAGI metric (**Sup Fig. 8**). **Fig. 5b** shows for the top two
583 performing methods based on RAGI (SnapATAC and chromVAR with k-mers) the gene
584 accessibility patterns for 3 genes (S100A12 - Monocytes-specific, MS4A1 - B cells specific
585 and GAPDH - housekeeping.) The same three genes are also shown in UMAP plots of
586 the other methods (**Sup Fig. 9**). Again, we observed that Louvain algorithm performed
587 better than k-means and hierarchical clustering for almost all scATAC-seq methods.
588 Importantly, negative RAGI score for a method (see for example the solutions obtained
589 by the Gene Scoring in **Fig. 5a**, **Sup Fig. 9**) may suggest that its clustering solutions are
590 defined by housekeeping genes rather than informative marker genes

591 We also qualitatively evaluated the clustering solutions of the different methods using
592 UMAP projections (**Fig. 5c**, **Sup Fig. 10**). We observed two major groups for all methods
593 except for scABC. Among these methods, the UMAP projections based on feature
594 matrices obtained by Control-Naïve, cisTopic, *Cusanovich2018*, Scasat SnapATAC,
595 BROCKMAN and chromVAR showed additional smaller groups and finer structures.
596 For Cicero gene activity scores, performing PCA helps to improve the separation of
597 more putative cell types. Instead for SCRAT and Gene Scoring, the PCA step did not
598 improve the separation.

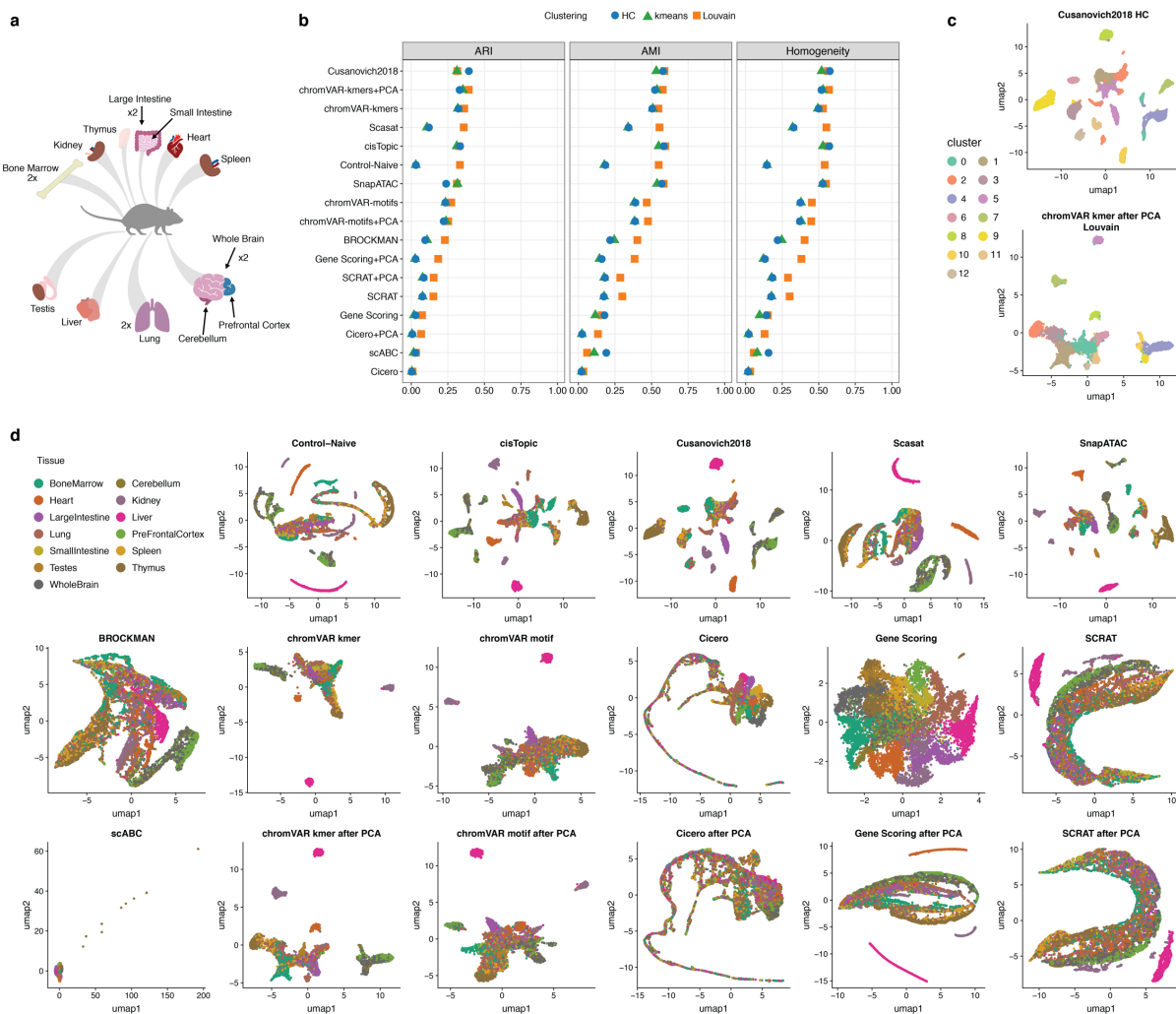
599 Given that the ranking of methods in datasets with ground truth is similar to the
600 ranking based on the RAGI metric, we believe this simple approach is a reasonable
601 surrogate metric that can be useful for evaluating unannotated datasets, a common
602 scenario in single cell omics studies.

603 *sci-ATAC-seq mouse dataset*

604

605 The last dataset analyzed in our benchmark consists of sciATAC-seq data from 13 adult
606 mouse tissues (bone marrow, cerebellum, heart, kidney, large intestine, liver, lung, pre-
607 frontal cortex, small intestine, spleen, testes, thymus and whole brain), of which 4 were
608 analyzed in duplicate for a total of 17 samples and 81,173 single cells[1]. Each tissue can

609 be interpreted as a coarse ground truth, used later to evaluate clustering solutions (Fig.
 610 6a). The analysis details for each method are documented in Sup Note 4.
 611



612
 613
 614 **Figure 6.** Benchmarking results using the downsampled sci-ATAC-seq mouse dataset
 615 from 13 adult mouse tissues. **(b)** Dot plot of scores for each metric to quantitatively
 616 measure the clustering performance of each method, sorted by maximum ARI score. **(c)**
 617 The two top-scoring pairings of scATAC-seq analysis method and clustering technique.
 618 Cell cluster assignments from each method are shown using the colors in the legend on
 619 the left. **(d)** UMAP visualization of the feature matrix produced by each method for the
 620 downsampled sci-ATAC-seq mouse dataset. Individual cells colors indicate the cell type.
 621

622 Despite using a machine with 1 TB of memory, almost all the methods failed to even
623 load this dataset, owing to its size. The only method capable of processing this dataset
624 in a reasonable time was SnapATAC (~700 minutes). The other methods failed to run
625 due to memory requirements. To understand the causes of this failure we did an in-
626 depth analysis of their scalability looking at their source code (**Sup Note 5**). Briefly, we
627 found that the majority of the methods try to load the entire dataset in the central
628 memory while SnapATAC uses a custom file format (.snap) based on HDF5
629 (<https://support.hdfgroup.org/HDF5/whatishdf5.html>), allowing out of core
630 computation by efficiently and progressively loading in the central memory only the
631 data chunks required at any given moment of the analysis.

632 On this dataset, SnapATAC was able to correctly cluster cells of the following tissues:
633 kidney, lung, heart, cerebellum, whole brain and thymus. However, for the other
634 tissues, including bone marrow and small intestine, cells are distributed in groups of
635 mixed cell types (**Sup Fig.11**), as reflected by the score of the three metrics used for the
636 other datasets evaluation (**Sup Table 14**), i.e. ARI= (HC=0.24, k-means=0.34,
637 Louvain=0.39), AMI=(HC=0.55, k-means=0.55, Louvain=0.62), Homogeneity=(HC=0.52,
638 k-means=0.54 , Louvain=0.60).

639

640 To gain insight on the performance of the other methods on this this dataset, we
641 randomly selected 15% of cells from each sample to construct a smaller sciATAC-seq
642 dataset consisting of 12,178 cells.

643

644 As **Fig. 6b** shows *Cusanovich2018*, k-mer-based chromVAR, cisTopic, SnapATAC, Scasat
645 and Control-Naïve perform comparably well and have noticeably better clustering
646 scores than the other methods (**Sup Table 15**). Consistent with what we observed
647 previously, peaks or bins level methods generally work better. In this dataset, k-mers-
648 based chromVAR and its combination with PCA transformation performs equally well
649 as peaks or bins-level methods and better than the motifs-based methods. Simply
650 counting reads within peaks (scABC) and gene-level-featurization-based methods
651 (Gene Scoring and Cicero) perform poorly overall. Adding a PCA step improves
652 noticeably scABC (scABC after PCA is the same as Control-Naïve) and Gene Scoring. It
653 also slightly improves Cicero but it does not affect chromVAR and SCRAT.

654

655 As before, all the clustering solutions of the different methods were visualized in
656 UMAP plots (**Sup Fig. 12**). The top two combinations, i.e. *Cusanovich2018* and
657 chromVAR k-mers with PCA, are visualized in **Fig.6c**. To visually compare the
658 separation of the different tissues across methods, we also inspected UMAP plots where
659 cells are colored based on the tissue of origin. Similar to what we observed using the
660 clustering analysis, cisTopic, *Cusanovich2018*, and SnapATAC are able to separate cells
661 into the major tissues and also to capture finer discrete groups. The Control-Naïve
662 method and Scasat are also able to distinguish the major tissues but show some mixing
663 within each discrete cell population. K-mer-based chromVAR can separate out liver,
664 kidney, and heart tissues and present the other tissues within a continuous bulk
665 population while preserving the structure of the distinct tissues. We observed that after
666 running PCA, k-mer-based chromVAR can recover an additional group of cells within
667 the lung tissue and also detect finer structure within the cells from the brain. Compared
668 with k-mer-based features, motif-based chromVAR and its combination with PCA
669 transformation distinguished fewer tissue groups while mixing more cells from
670 different tissues. BROCKMAN recovered a continuous structure with the different
671 tissues but does not distinguished them clearly. Similarly, Gene Scoring put cells from
672 different tissues into a big bulk population with limited separation. PCA improved its
673 ability to separate out a few tissues, including liver, heart, and kidney. SCRAT and
674 Cicero gene activity scores mixed most of the cells from different tissues and performed
675 poorly on this dataset with or without PCA.

676

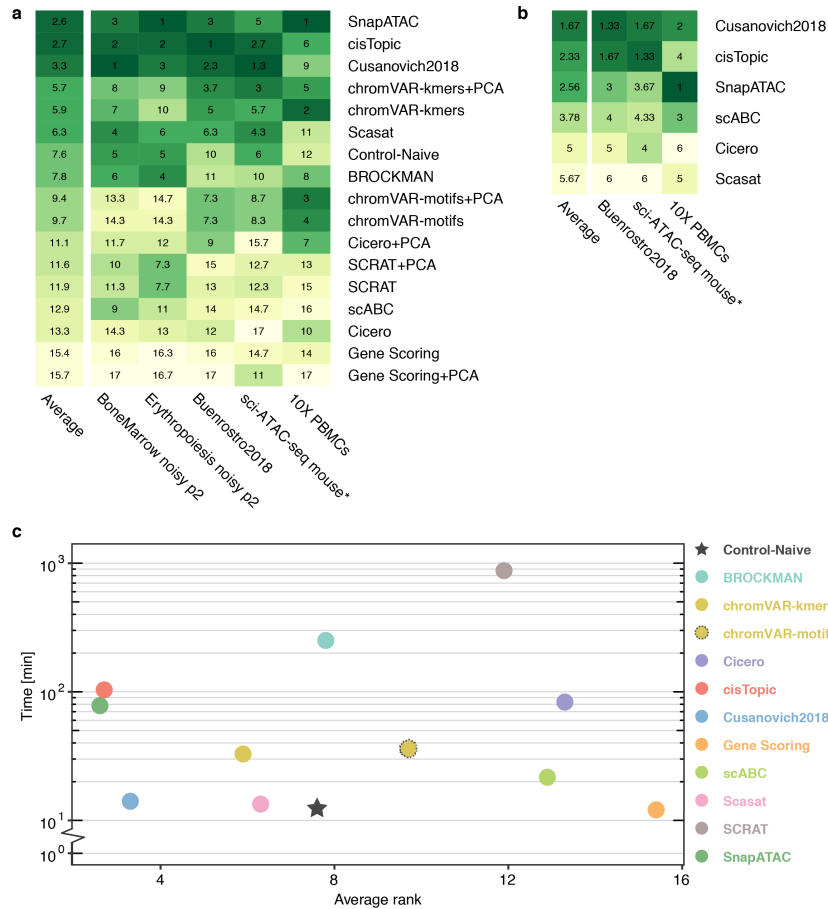
677

678 **Clustering performance summary**

679 To assess and compare the overall performance of scATAC-seq analysis methods, we
680 ranked the methods based on each metric (ARI, AMI, Homogeneity, RAGI) by taking
681 the best clustering solution for the three real datasets (*Buenrostro2018* dataset, *PBMCs*
682 *10X* dataset, and the down-sampled *sci-ATAC-seq mouse* dataset) and two synthetic
683 datasets (simulated bone marrow dataset and simulated erythropoiesis dataset with the
684 moderate noise level of 0.2 and a medium coverage of 2500 fragments per cell). Then for
685 each dataset except for the *PBMCs 10X* dataset, we calculated the average rank across
686 ARI, AMI, and Homogeneity. For the *PBMCs 10X* dataset, RAGI is calculated instead
687 (**Sup Fig.13a**). Lastly, we calculated the average rank across different datasets.

688 According to the average ranking, SnapATAC, cisTopic and *Cusanovich2018* are the top

689 three methods to create feature matrices that can be used to cluster single cells into
 690 biologically-relevant subpopulations (**Fig. 7a**). SnapATAC consistently performed well
 691 across all datasets. Both cisTopic and *Cusanovich2018* demonstrated satisfactory
 692 performance across all datasets except for the 10X PBMCs dataset.



693 **Figure 7.** Aggregate benchmark results. **(a)** For each method, the rank based on the best-
 694 performing clustering method is measured for each metric (e.g. ARI, AMI, H, or RAGI).
 695 The average metric ranks for each dataset were used to calculate a performance score for
 696 each method. Each method was then assigned a cumulative average score based on its
 697 performance across all datasets. * indicates a downsampled dataset of the indicated
 698 original dataset. **(b)** For methods that specify an end-to-end clustering pipeline, average
 699 rank and cumulative average scores for each method were calculated as in (a). **(c)** Plot of
 700 running time against performance for each method. Cumulative average scores, which
 701 were calculated in part (a) are shown on the x-axis, and the average running time across
 702 the three real datasets (*Buenrostro2018*, 10X PBMCs, and downsampled sci-ATAC-seq
 703 mouse) is shown on the y-axis.
 704

705 Generally, methods that implement a dimensionality reduction step work better
706 (SnapATAC, cisTopic, *Cusanovich2018*, Scasat, Control-Naïve, and BROCKMAN) than
707 those without it (SCRAT, scABC, Cicero, and Gene Scoring). We also observed that
708 chromVAR performs better in real datasets than in simulated datasets and that the
709 kmer-based version of chromVAR consistently outperforms motif-based chromVAR.
710 For the methods that do not implement dimensionality reduction, the PCA step does
711 not always improve the performance except for scABC and Cicero, in which the PCA
712 transformation consistently boosts the results. Interestingly, we observed that
713 regardless of the method, the PCA consistently improves the clustering solutions
714 obtained by the Louvain algorithm.

715

716 Keeping the first PC vs removing the first PC

717

718 In preparing this manuscript, we noticed that in some cases, the first principal
719 component (PC) may only capture variation in sequencing depth instead of biologically
720 meaningful variability. To make a thorough assessment of how the first PC affects the
721 clustering results, we compared the effect of keeping vs removing the first PC on the
722 three real datasets (for this comparison we consider both the methods that implemented
723 PCA and the combination of PCA and the methods that did not implement a
724 dimensionality reduction step) (**Sup Fig.14**). Across all three datasets, we observe that
725 for Control-Naïve, BROCKMAN, SCRAT-PCA, and Gene Scoring-PCA, removing the
726 first PC consistently helped in better separating the different populations in UMAP
727 projections and improved clustering performance. In contrast, the performance of
728 chromVAR-PCA with motifs as features consistently dropped after removing the first
729 PC. *Cusanovich2018* and SnapATAC performed similarly before and after removing the
730 first PC across all datasets. For Cicero-PCA, removing first PC did not clearly affect its
731 performance in *Buenrostro2018* and 10X PBMCs datasets but improved its performance
732 in the down-sampled sci-ATAC mouse dataset.

733

734 Generally, the methods that implement binarization (e.g. *Cusanovich2018*, SnapATAC)
735 or that implement cell coverage bias correction (e.g. chromVAR, SnapATAC), tend to be
736 less affected by the sample sequencing depths. Therefore, for these methods we believe
737 that the first PC does not capture the library size and removing it does not help to
738 improve the clustering results. On the contrary, for methods that do not implement any

739 specific step to correct for potential artifacts associated with sequencing depth, the first
740 PC is more likely to capture biologically irrelevant factors and therefore may reduce
741 biology-driven differences. However, this operation must be applied with caution, since
742 removing the first component could also in some cases remove some biological
743 variation (e.g. motif-based chromVAR).

744 Clustering performance when running methods as end to end pipelines

745 When designing this study, we reasoned that a benchmark procedure could be
746 approached from two very different perspectives. The first is the end user perspective,
747 i.e. a user that runs a method as a black box following the provided documentation with
748 the goal to obtain a reasonable clustering solution without worrying too much about the
749 internal design choices and procedures. In these settings, it is not trivial to
750 systematically compare the methods and understand which part related to the
751 featurization may influence the final clustering performance, especially if also the
752 clustering algorithms used are different. The second perspective that was used instead
753 in the rest of this benchmarking effort is the developer perspective, i.e. we tried to
754 understand what are the key steps of each method that can boost clustering
755 performance of common clustering approaches. Regardless, we reasoned that it is
756 important to provide some insights on the user perspective, since some readers will use
757 the tested methods as end-to-end pipelines. Therefore, we also compared the clustering
758 solutions produced by running the complete analysis pipelines as outlined in tutorials
759 for the methods that explicitly implement a clustering step (see **Sup Note 6**). We
760 evaluated the clustering results using ARI, AMI and Homogeneity for the
761 *Buenrostro2018* and sci-ATAC-seq mouse datasets, and RAGI for the PBMCs 10X dataset
762 (**Sup Table 16-17-18**). We observe the top three methods, i.e. Cusanovich2018, cisTopic
763 and SnapATAC, still outperform the other methods but with a slightly different
764 ranking. (Cusanovich2018 is ranked first followed by cisTopic and SnapATAC, Fig. 7b,
765 Sup Fig. 13b). Also, both scABC and Cicero performed better than Scasat in this
766 analysis. Interestingly, we observed that SnapATAC, cisTopic, Cusanovich2018, and
767 Scasat have even better clustering solutions in our benchmarking framework compared
768 to using their own clustering approach. On the other hand, scABC and Cicero had
769 better clustering results when running their own clustering procedure. scABC uses an
770 unsupervised clustering method tailored to single cell epigenomic data (including
771 scATAC-seq). Although it uses the naïve peaks-by-cells raw count as its feature matrix,

772 it calculates cells weights by considering their sequencing coverage and giving more
773 weight to cells with higher number of reads. Also, it performs two steps of clustering by
774 using weighted k-medoid algorithm based on Spearman rank correlation to find
775 landmarks first and then assigns cells to the landmarks. These specific steps help
776 improve its clustering performance. For the Cicero clustering workflow, we used the
777 gene activity scores and, as proposed in their tutorial, functions from Monocle2, to (i)
778 normalize the scores and (ii) reduce the dimensionality with tSNE by using the top PCs
779 before clustering cells. These extra steps helped in improving its clustering solutions.
780 This suggests that appropriate normalization steps need to be properly performed to
781 improve clustering analysis, in addition to simple transformations like binarizing
782 counts and/or performing a PCA.

783
784 Taken together, based on these analyses, we recommend using SnapATAC, cisTopic, or
785 *Cusanovich2018* to cluster cells in meaningful subpopulations. This step can be followed
786 by methods such as Cicero, Gene Scores or with TF motifs (e.g. chromVar) to annotate
787 clusters and to determine cell types in an integrative approach.

788
789 **Important considerations in defining informative regions for scATAC-seq analyses**

790 Feature sets of informative peaks for scATAC analyses may be computed from bulk
791 samples available through large scale consortia such as ENCODE[20] and
792 ROADMAP[21] or more precise tissue-specific cell types as in the
793 murine ImmGen Project[22]. However, scATAC-seq analyses often require *de*
794 *novo* inference of dataset-specific accessibility peaks in order to resolve cell types and
795 regulatory activity.

796 To date, there are three major methods for generating peak sets for scATAC
797 experiments. The first strategy (pseudo-bulk from all single cells, PB-All) for inferring
798 peaks is to call peaks on a pseudo-bulk sample composed of all the reads from all cells in
799 the library. The second (pseudo-bulk from FACS, PB-FACS) is to call peaks in *a priori*-
800 defined cell types isolated by FACS-sorting. A consensus peak set can be defined by
801 combining summits of individual peaks using an iterative algorithm [5, 18, 23]. Finally,
802 a third strategy (pseudo-bulk from clades, PB-Clades) uses a pre-clustering of cells to
803 define initial populations[1, 10]. Subsequent peak calling is performed in each initial

804 cluster. Aggregate peak sets can then be defined from synthesizing the summits of each
805 cluster-specific peak set as described above.

806 Bulk ATAC-seq peaks vs aggregated scATAC-seq peaks

807 To evaluate the effect of using peaks obtained from bulk ATAC-seq data versus peaks
808 obtained from aggregated single cell profiles, we reanalyzed the *Buenrostro2018* dataset
809 in which both are available (**Sup Fig. 15-16**). Here we considered only the methods that
810 use peaks as input (i.e. SnapATAC, SCRAT, BROCKMAN are excluded). For the
811 aggregated scATAC-seq peaks, we merged cells of the same cell type based on the
812 FACS sorting labels and performed peak calling within each cell type. Then peaks
813 defined within each cell type were merged. For most methods we did not observe clear
814 differences in performance between the two input peak strategies. For cisTopic,
815 *Cusanovich2018*, and Cicero, aggregated scATAC-seq peaks overall perform better
816 across all three metrics (**Sup Fig. 17a, Sup Table 19**).

817 We also tested the strategy of defining pseudo bulk samples from clades when no
818 sorting labels are provided. *Cusanovich2018* is the only method that provides a
819 workflow to identify initial clades and call peaks within each clade. It counts reads
820 within the fixed-size windows and pre-clusters cells using hierarchical clustering to
821 define initial clades from which peaks are called. We applied this strategy to all three
822 real datasets (**Sup Fig. 18**). We observed that in all three datasets, *Cusanovich2018*
823 performs well in identifying the isolated major groups and the identified clades match
824 well the labels provided, including FACS-sorted labels, cell-ranger clustering solutions,
825 and known tissues labels. Overall the *Cusanovich2018* ‘pseudo bulk’ strategy for
826 defining *de novo* peaks is able to capture the heterogeneity within single cell populations
827 and can serve as a promising unsupervised way to define pseudo bulk subpopulations
828 and to perform peak calling.

829 The effect of excluding regions using the ENCODE blacklist annotation

830
831 cisTopic, Scasat, SCRAT, and SnapATAC employ a blacklist filtering step to remove
832 features annotated by ENCODE as belonging to a subset of genomic regions, which
833 harbor the potential to produce artifacts in downstream analysis steps [24]. cisTopic and
834 Scasat perform a peak filtering in the pre-processing steps of their pipeline. Our
835 benchmarking pipeline makes use of the ENCODE ATAC-seq pre-processing pipeline,

836 which removes peaks overlapping with regions on the blacklist annotations list.
837 Therefore, we tested the remaining two methods, which do not use peaks as features,
838 SCRAT or SnapATAC. In particular, we wanted to test whether we would observe any
839 change in downstream clustering performance upon opting to perform a blacklist
840 removal step. Through a qualitative and quantitative comparison of clustering
841 performance, we determined that methods, which remove features according to
842 blacklist annotations show no considerable advantage over those that permitted such
843 features (**Sup Fig.19**).

844

845 *Rare cell type-specific peak detection*

846 As all cell identities may not be pre-defined in complex tissue types, we sought to
847 examine PB-All and PB-Clades strategies to infer a chromatin accessibility feature set
848 from the scATAC-seq libraries directly. To achieve this, we established a simulation
849 setting where we mixed bulk ATAC-seq data from three sorted populations (B-cells,
850 CD4+ T-cells, and monocytes from the PBMCs 10X dataset) that would be mixed in
851 complex tissue (i.e. peripheral blood mononuclear cells) (**Sup Fig. 17b**). After peak
852 calling on both the synthetic bulk and isolated reads from each cell type, we inferred the
853 proportion of cell type-specific peaks from the minor cell population that were captured
854 by the peak calling in the synthetic bulk mixture (see **Methods**).

855 Overall, the results indicate that cell type-specific peaks may be vastly underestimated
856 from performing peak calling on the mixture of single cells (PB-All) (**Sup Fig. 17b**).
857 Specifically, only ~18% of cell type-specific peaks from very rare (1% prevalence) or
858 ~40% from rare (5% prevalence) cell populations were detected when peaks were called
859 when treating the heterogenous source as a synthetic bulk experiment. Consequently, as
860 these peaks would be vastly under-represented in a consensus peak set, virtually all
861 computational algorithms will fail to identify rare populations. Moreover, as many
862 common quality-control measures for scATAC involve filtering based on the proportion
863 of reads in peaks, these cell populations may be under-represented in quality-controlled
864 datasets.

865 As observed in other studies [1, 25], these results suggest calling peaks on PB-All may
866 result in sub-optimal performance. Alternatively, when isolated populations have been
867 profiled (for example by FACS) peak sets can be defined by calling peaks using data
868 from cells in each pre-defined population separately as discussed in the previous

869 section since this enables the resolution of rare subpopulations (for example HSC in the
870 hematopoietic system).

871 Frequency-based peak selection vs intensity-based peak selection

872 *Cusanovich2018* selects peaks that are present in at least a specified percentage of cells
873 before performing TF-IDF transformation, while *scABC* selects peaks with the most
874 reads to cluster cells. To evaluate the effect of selecting peaks based on their
875 representation in the cell population or based on their intensity (defined as the sum of
876 reads in that peak in all samples), we focus on the two methods that implement the step
877 of peak selection, *Cusanovich2018* and Control-Naïve (equivalent to *scABC+PCA*).

878 To assess the two peak selection strategies, we ran both *Cusanovich2018* and Control-
879 Naïve on both simulated bone marrow dataset at noise level of 0.2 with a coverage of
880 2500 fragments and the *Buenrostro2018* dataset by varying the cutoffs for peak inclusion
881 (**Sup Fig. 20-21**). We calculated the intensity of peaks by counting the number of reads
882 across all cells and calculated the frequency of peaks by counting the number of cells in
883 which a peak is observed. For this analysis we selected the top peaks based on intensity
884 and frequency with the following cutoffs: top 100%, 80%, 60%, 40%, 20%, 10%, 8%, 6%,
885 4%, %2, 1%.

886 For both *Cusanovich2018* and Control-Naïve, the two peak selection strategies have
887 similar clustering result scores when varying the cutoff (**Sup Fig. 20a-b,21a-b**). We
888 observed reasonable and stable clustering performance using more than 20% of the
889 ranked peaks. As the number of peaks is reduced, the scores start to decline noticeably
890 and decrease almost monotonically. Below 1%, both methods perform poorly. In
891 addition, we observed that the Louvain method produces more stable results than
892 hierarchical clustering and k-means across the considered settings.

893 **Running time of different methods**

894 In our analysis, we also collected the running time of each method on both simulated
895 and real datasets (see **Sup Note 6**). For the simulated datasets, we only reported the
896 execution time necessary to build a feature matrix starting from a peaks-by-cells count
897 matrix. For real datasets, we considered the execution time to build a feature matrix
898 from bam files. The running times are shown in **Sup Fig. 22 (Sup Table 20)**. All the
899 tests were run on a machine with an Intel Xeon E5-2600 v4 X CPU with 44 cores and 1

900 TB of RAM with the CentOS 7 operating system. When analyzing real datasets with
901 methods that rely on peaks but do not provide an explicit function to construct a peaks-
902 by-cells matrix (*Cusanovich2018*, Cicero, Gene Scoring and Scasat), we ran the same
903 script on a Linux cluster to obtain the peaks-by-cells matrix such that the execution time
904 of this step is equivalent across these methods. It is worthwhile to mention that not all
905 the methods of this benchmark support parallel computing. For the methods that
906 support parallel computing, including SnapATAC, chromVAR, and cisTopic, the
907 execution time was reported using 10 cores. For the rest of methods, we run them using
908 a single core. We selected this number reasoning that a typical lab may not have access
909 to a machine with 44 cores and instead may use a mid-size computing node with 8-12
910 cores. Notably, SnapATAC was the only method capable of processing the full sci-
911 ATAC-seq mouse dataset (~80,000 single cells).

912 As shown in **Sup Fig. 22**, BROCKMAN and SCRAT have the largest greater execution
913 time in all the real datasets while the methods that use a custom script to obtain a
914 peaks-by-cells matrix tend to have shorter execution time (e.g. Scasat, *Cusanovich2018*,
915 Gene Scoring).

916 We also assessed the scalability of methods with respect to the increasing coverage (250,
917 500, 1000, 2500 and 5000 fragments per peaks). We observe that with the increase of
918 read coverage, for cisTopic there is an exponential increase of the running time whereas
919 for other methods, the running time stays stable or increases linearly (**Sup Fig. 22, Sup
920 table 21**).

921 Finally, we compared execution time vs clustering performance (**Fig. 7c**). Interestingly,
922 the most accurate methods (SnapATAC, cisTopic and *Cusanovich2018*) have a
923 reasonable running time while outperforming the other methods for clustering quality
924 across all the datasets. Considering the computational time as an important factor that
925 must be carefully evaluated before the implementation of any bioinformatics pipeline,
926 we believe that *Cusanovich2018* is the best in balancing clustering performance with
927 execution time.

928

929 **Discussion**

930 scATAC technologies enable the epigenetic profiling of thousands of single cells, and
931 many computational methods have been developed to analyze and interpret this data.
932 However, the sparsity of scATAC-seq datasets provides unique challenges that must be
933 addressed in order to perform essential analyses such as cluster identification,
934 visualization and trajectory inference [26, 27]. Moreover, the rapid technological
935 innovations that facilitate profiling accessible chromatin landscapes of 10^4 or 10^5 cells
936 provide additional computational challenges to efficiently store and analyze data.

937 In this study, we compared ten computational methods developed to construct
938 informative feature matrices for the downstream analysis of scATAC-seq data. We
939 developed a uniform processing framework that ranks methods based on their ability to
940 discriminate cell types when combined with three common unsupervised clustering
941 approaches, followed by evaluation of three well-accepted clustering metrics. We
942 evaluated these methods on thirteen datasets, three of those obtained using different
943 technologies (Fluidigm C1, 10X, and sci-ATAC), and five consisting of simulated data
944 with varying noise levels. These datasets comprise cells from different tissues in both
945 mouse and human.

946 In addition to identifying various methodologies that perform optimally on real and
947 simulated data, our benchmarking examination of scATAC-seq methodologies reveals
948 general principles that will inform the development of future algorithms. First, peak-
949 level or bin-level feature counting generally performs better in distinguishing different
950 cell types followed in turn by k-mer-level, TF motifs-level, and gene-centric level
951 summarization. We interpret this finding as an indication of the complexity of gene
952 regulatory circuits where precise enhancer elements may have distinct functions that
953 cannot be sufficiently approximated by sequence context or proximity to gene bodies
954 alone. Second, we note that the methods that implement a dimensionality reduction
955 step generally perform better in the separation of cell types, since this step may help to
956 remove the redundancy between a large number of raw features and to mitigate the
957 effect of noise. Third, for the methods that do not implement a dimensionality reduction
958 step, simply adding a PCA step could significantly improve the clustering results. In
959 fact, PCA generally boosts Louvain clustering results. For methods that do not account
960 for the differing sequencing coverage of cells, the first PC could be used to capture and
961 correct for sample depth differences. In this case, removing the first PC may improve
962 the performance of these methods. Fourth, we observe that the Louvain method overall

963 performs more consistently and accurately than k-means and hierarchical clustering. In
964 contrast, k-means and hierarchical clustering are more sensitive to outliers and may
965 result in suboptimal clustering solutions since some of clusters may correspond to
966 single or few outlier cells. Fifth, the robustness of different methods to noise and
967 coverage varies among different datasets. Among the top three methods, cisTopic is the
968 most penalized by low coverage. Sixth, it was also observed that inappropriate
969 transformations, such as log₂ transformation and normalization based on region size as
970 implemented in SCRAT may impact negatively clustering performance.

971 We observe that many methods fail to scale to larger datasets, which are now available
972 due to improvements in split-pool technology and droplet microfluidics. As
973 technologies improve and individual labs and international consortia lead efforts to
974 generate ever larger single-cell datasets, scalability will be an unavoidable goal of
975 method developments on a par with accuracy. As many of our evaluated methods were
976 designed in the context of data generated from the Fluidigm C1 platform (which
977 produces ~10² cells), such approaches were often incapable of analyzing large datasets.
978 In particular, the sci-ATAC-seq mouse dataset served as a useful resource to test the
979 scalability of the methods that were benchmarked (~80,000 cells). Notably, our
980 evaluation demonstrates that only SnapATAC was able to scale to process and analyze
981 this large dataset. Future methods must be capable of processing datasets of this size
982 especially adopting efficient data structures that allow out of core computing. Our
983 findings reinforce the need for methods that not only are accurate but highly scalable
984 for scATAC-seq data processing.

985 Defining regions is an important step in constructing feature matrices. Selecting
986 informative regions generally improves downstream analyses such as clustering to
987 capture heterogeneity within cell populations. Peak calling is a popular and
988 straightforward way to define regions of interest. We observe that clustering
989 performance is not generally impacted by using peaks defined from bulk ATAC-seq
990 data vs using peaks obtained from aggregating single cell data based on FACS-sorting
991 labels. However, performing peak calling by simply pooling reads from single cells may
992 obfuscate peaks specific to rare cell populations leading to failures in uncovering them.
993 In addition, the *Cusanovich2018* approach to identify pseudo-bulk clades is a promising
994 unsupervised way to perform *in silico*-sorting without relying on FACS-sorting labels.
995 This strategy potentially serves as a suitable way to preserve peaks specific to rare cell

996 types. Also choosing an appropriate number of peaks is important for improving the
997 downstream analysis (for example based on intensity/frequency-based given that they
998 perform similarly).

999 We are aware of current limitations in our benchmarking effort. We have compared
1000 single cell ATAC-seq methods based on their ability to separate discrete cell
1001 populations; however, this might not be ideal when dealing with a continuous cell
1002 lineage landscape. We observe that chromVAR generally works better in preserving a
1003 continuous space while SnapATAC tends to break a putative landscape into discrete
1004 populations. The choice of method is ultimately case-specific and may be driven by the
1005 downstream application. For example, the feature matrix obtained by chromVAR may
1006 be more suitable for trajectory inference [26] while the one obtained from SnapATAC
1007 may be more appropriate better identify discrete and well separated cell populations by
1008 clustering. We acknowledge also that not all tested methods were specifically designed
1009 to produce clustering results. For example, chromVAR, Cicero, and Gene Scoring were
1010 designed to determine important marker genes, their regulatory logic, or to infer
1011 enriched TF binding sites within accessible chromatin regions. However, because
1012 clustering is a critical part of single-cell analysis and researchers frequently use output
1013 from all methods to produce clustering results [1], we felt that evaluating the clustering
1014 abilities using feature matrices produced by each method was a useful measure. An
1015 additional limitation of our study is that it is impossible to create a simulation
1016 framework that models an experimental outcome with perfect accuracy. Several
1017 assumptions were made to enable our simulation of the data; these assumptions are
1018 described in the methods section of this manuscript, where we detail explicitly how the
1019 simulated data was generated.

1020 Interestingly, we learnt that some combinations of feature matrices with the simple
1021 clustering approaches included in our benchmarking framework perform even better
1022 than the original combination proposed by the respective authors. This highlights the
1023 value of this dual-characterization (*user vs designer perspective*) and provides a summary
1024 of both perspectives to the readers.

1025 We believe it is important to stress the distinction between biological realities and
1026 computational performance, especially in the context of unsupervised clustering. A big
1027 and critical assumption (or hope) of our field is that an unsupervised clustering

1028 procedure will provide clustering solutions that recapitulate different populations
1029 corresponding to different cell types/states. Given that for several real datasets the
1030 ground truth is not known, a current compromise during the exploratory clustering
1031 analysis is to use known marker genes, sorted populations or known tissues to validate
1032 the clustering solutions based on classic metrics. If we embrace this assumption,
1033 keeping in mind that additional validation is required to truly delineate the
1034 subpopulation structure of a population of cells, the two views, biological and
1035 computational can be reconciled. Our benchmark procedure is aimed to provide some
1036 guidelines based on explorative analyses that are currently adopted in several
1037 published papers.

1038 Looking forward, due to the wealth of data being produced by new scATAC
1039 technologies, we hypothesize that more powerful machine learning frameworks may be
1040 able to uncover complex *cis* and *trans* relationships that define cell-cell
1041 relatedness. Specifically, we anticipate autoencoder-like models that integrate genomic
1042 sequence context, gene body positions, and precise accessible chromatin information
1043 will yield information-rich features and that more advanced manifold learning methods
1044 will help to remove redundancy and better preserve heterogeneity within single cell
1045 populations. Such achievements may enable us to overcome the inherent sparsity and
1046 high dimensionality that characterizes scATAC-seq data.

1047 **Conclusions**

1048 Our benchmarking results highlight SnapATAC, cisTopic, and Cusanovich2018 as the
1049 top performing scATAC-seq data analysis methods to perform clustering across all
1050 datasets and different metrics. Methods that preserve information at the peak-level
1051 (cisTopic, Cusanovich2018, Scasat) or bin-level (SnapATAC) generally outperform those
1052 that summarize accessible chromatin regions at the motif/k-mer level (chromVAR,
1053 BROCKMAN, SCRAT) or over the gene-body (Cicero, Gene Scoring). In addition,
1054 methods that implement a dimensionality reduction step (BROCKMAN, cisTopic,
1055 Cusanovich2018, Scasat, SnapATAC) generally show advantages over the other
1056 methods without this important step. SnapATAC is the most scalable method; it was
1057 the only method capable of processing more than 80,000 cells. Cusanovich2018 is the
1058 method that best balances analysis performance and running time.

1059 Taken together, our manuscript provides a framework for evaluating and
1060 benchmarking new and existing methodologies as well as provides important
1061 guidelines for the analysis of scATAC-seq data. Importantly, we provide more than 100
1062 well organized and documented Jupyter notebooks to illustrate and reproduce all the
1063 analyses performed in this benchmarking work. We believe our systematic analysis
1064 could guide the development of computational approaches aimed at solving the
1065 remaining challenges associated with analyzing scATAC-seq datasets.

1066 **Methods**

1067 Our assessment of methods was based on public scATAC-seq datasets made available
1068 in public repositories by the respective authors (see **Data and code availability**). As
1069 such, we refer to the original publications for further details on experimental design
1070 and data pre-processing/alignment. For peak calling, we used the ENCODE pipeline
1071 (<https://www.encodeproject.org/atac-seq/>) except for the 10X PBMCs data for which
1072 peaks were already available through the Cell Ranger pipeline optimized for this
1073 technology. Whenever changes were required for running a given method, those are
1074 noted in the respective sections.

1075 *Datasets*

1076 **Human hematopoiesis I** (Buenrostro et al. 2018)

1077 This dataset comprised of 10 FACS-sorted cell populations from CD34⁺ human bone
1078 marrow, namely, hematopoietic stem cells (HSCs), multipotent progenitors (MPPs),
1079 lymphoid-primed multipotent progenitors (LMPPs), common-myeloid progenitors
1080 (CMPs), granulocyte-macrophage progenitors (GMPs), megakaryocyte-erythrocyte
1081 progenitors (MEPs), common-lymphoid progenitor (CLPs), plasmacytoid dendritic cells
1082 (pDCs), monocytes, and an uncharacterized CD34⁺ CD38⁻ CD45RA⁺ CD123⁻ cell
1083 population. A total of 2,034 cells from 6 human donors were used for analysis. A peak
1084 file (including 491,437 peaks) obtained from bulk ATAC-seq dataset was provided.

1085 **sci-ATAC-seq mouse tissues** (Cusanovich et al. 2018)

1086 This dataset comprises cells from 13 tissues of adult mouse, namely, bone marrow,
1087 cerebellum, heart, kidney, large intestine, liver, lung, prefrontal cortex, small intestine,

1088 spleen, testes, thymus, and whole brain, with over 2,000 cells per tissue. A total of
1089 81,173 cells from 5 mice were used for analysis. A subset was obtained by randomly
1090 down-sampling 15% cells from each tissue and was comprised of 12,178 cells.

1091 **Human hematopoiesis II (10X PBMCs)**

1092 This dataset is composed of peripheral blood mononuclear cells (PBMCs) from one
1093 healthy donor. A total of 5,335 cells were used for analysis.

1094 **Simulated scATAC-seq datasets**

1095 In order to evaluate and benchmark various approaches, we generated synthetic
1096 (labeled) data from down-sampling 18 FACS-sorted bulk populations that were
1097 previously described [28]. For ease of interpretation, we considered only 6 isolated
1098 populations (HSC, CMP, NK, CD4, CD8, Erythroblast). For the erythropoiesis
1099 simulation, eight additional populations (P1-P8) originally described in [17] were also
1100 considered.

1101

1102 Our simulation framework starts with a peak x cell type counts matrix (from bulk
1103 ATAC-seq) and generates a single-cell counts matrix (C) for an arbitrary number of
1104 synthetic single cells. Explicitly, for a simulated single cell j and corresponding peak i
1105 from bulk cell type t , we seek to generate $c_{i,j}$ where $c_{i,j} \in \text{Error! Bookmark not defined.}$,
1106 noting that these values correspond to possible observations in a diploid genome. Next,
1107 we define the rate (r_i^t) at which the peak i is prevalent in the bulk ATAC-seq data for
1108 cell type t . This rate is determined by the ratio of reads observed in peak i over the sum
1109 of all reads. Assuming a total of k peaks for the matrix C and for user-defined
1110 parameters q (noise parameter; $q \in [0,1]$) and n (number of simulated fragments), we
1111 define $c_{i,j}$ as follows:

1112

$$1113 \quad c_{i,j} \sim \text{rbinom}(2, p_i^t)$$

1114 where

1115

$$1116 \quad p_i^t = (r_i^t) \binom{1}{2} n (1 - q) + (1/k) \binom{1}{2} n (q)$$

1117

1118 Intuitively, the parameter p_i^t defines the probability that a count will be observed in
1119 peak i for a single cell. Additionally, p_i^t can be decomposed into the sum two terms. As
1120 $q \rightarrow 0$, the first term dominates, and the probability of observing a count in peak i is
1121 simply the scaled probability of the ratio of reads for that peak from the bulk ATAC-seq
1122 data (r_i^t). Thus, when $q = 0$, the simulated data has no noise. Conversely as $q \rightarrow 1$, the
1123 second term dominates, and p_i^t reduces to a flat probability that is no longer
1124 parameterized by the peak i or cell type t and thus represents a random distribution of
1125 n fragments into k peaks.

1126
1127 For bone marrow-based simulations we simulated 200 cells per labeled cell type while
1128 for erythropoiesis-based simulation we simulated 100 cells per labeled cell type.
1129 Eventually we have 1,200 cells for each simulated dataset. In the base simulations, we
1130 parametrized $n = 2,500$ fragments in peaks in expectation for all cells. For additional
1131 simulations that compared different data coverages, we set n to various values (5000,
1132 2500, 1000, 500, 250 respectively) to benchmark this effect. To evaluate the effect of noise
1133 in our simulation, we set q to three values (0, 0.2, 0.4) to benchmark the robustness to
1134 noise. At values of $q > 0.4$, no method could reliably separate all the subpopulations.
1135 Finally, since our simulation started at the reads in peaks level, for some methods, the
1136 core algorithm associated with the method was extracted in order to benchmark it in
1137 this setting. Additionally, full code to reproduce these simulated dataset matrices has
1138 been made available with our online code resources.

1139 *Peak calling*

1140 For real datasets, peaks were called using the ENCODE ATAC-seq processing pipeline
1141 (<https://www.encodeproject.org/atac-seq>). Briefly, single-cells were aggregated into cell
1142 populations according to cell type, obtained either by FACS sorting or by tissue of
1143 origin. Peaks were called for each cell population and merged into a single file with
1144 bedtools [30].

1145 *Building the features matrix*

1146 **BROCKMAN** This method starts by defining regions of interest, which will be scanned
1147 for k -mer content, as 50 bp windows around each transposon integration site and
1148 merging overlapping regions. Then, a frequency matrix of k -mers-by-cells is built by

1149 counting all possible gapped k -mers (for k from 1 to 8) within the previously defined
1150 windows. This frequency matrix is scaled so that each k -mer has mean 0 and standard
1151 deviation 1. Principal component analysis (PCA) is applied to the scaled k -mers-by-cells
1152 frequency matrix, and significant principal components (PCs) as estimated with the
1153 jackstraw method are selected to build a final features matrix for downstream analyses.

1154 **ChromVAR** This method starts by counting reads under chromatin-accessible peaks in
1155 order to build a count matrix of peaks-by-cells (X). Then, a set of chromatin features
1156 such as transcription factor (TF) motifs or k -mers are considered. Reads mapping to
1157 each peak that contains a given TF motif (or k -mer) are counted in order to build a
1158 count matrix of motifs-by-cells or k -mers-by-cells (M). Moreover, a raw accessibility
1159 deviation matrix of motifs (or k -mers)-by-cells (Y) is generated by calculating the
1160 difference between M and the expected number of fragments based on X . Then,
1161 background peak sets are created for each motif to remove technical confounders.
1162 Background motifs-by-cells raw accessibility deviations are then used to calculate a bias
1163 corrected deviation matrix and to compute a deviation z -score used for downstream
1164 analyses.

1165 **cisTopic** This method starts by building a peaks-by-cells binary matrix by checking if a
1166 peak region is accessible, i.e., at least one read falls within the peak region. Then, latent
1167 Dirichlet allocation (LDA) is performed on this binary matrix and two probability
1168 distributions are generated, a topics-by-cells probability matrix and a regions-by-topics
1169 probability matrix. The former is the final features matrix for downstream analyses.

1170 **Cicero** This method defines promoter peaks as the union of annotated TSS minus 500
1171 base pairs and macs2 defined peaks around the TSS. It takes as input the peaks-by-cells
1172 binary matrix. It also requires either pseudo temporal ordering or coordinates in a low
1173 dimensional space (t-SNE) so that cells can be readily grouped. It then computes the co-
1174 accessibility scores between sites using Graphical Lasso. To get the gene activity scores,
1175 it selects sites that are proximal to gene TSS or distal sites linked to them and weight
1176 them by their co-accessibility. Then all the sites are summed and weighted according to
1177 their co-accessibility to produce a genes-by-cells feature matrix that is used in this
1178 benchmarking analysis.

1179 **Gene Scoring** This method first constructs a peaks-by-cells count matrix and defines
1180 regions of interest as the 50kb upstream and downstream of gene TSSs. Then it finds the
1181 overlap between ATAC-seq peaks and TSS regions and the peaks are weighted by a
1182 function of the distance to the linked genes. Finally, the peaks-by-cells count matrix is
1183 converted into genes-by-cells weighted count matrix by multiplying the weighted peaks
1184 by genes matrix. The genes-by-cells weighted count matrix is the final features matrix
1185 for downstream analyses.

1186 *Cusanovich2018* This method starts by binning the genome into fixed-size windows (by
1187 default, 5kbp), and building a binary matrix from evaluating whether any reads map to
1188 each bin. Bins that overlap ENCODE-defined blacklist regions are filtered out, and the
1189 top 20,000 most commonly used bins are retained. Then, the bins-by-cells binary matrix
1190 is normalized and rescaled using the term frequency-inverse document frequency (TF-
1191 IDF) transformation. Next, singular value decomposition (SVD) is performed to
1192 generate a PCs-by-cells LSI score matrix, which is used to group cells by hierarchical
1193 clustering into different clades. Within each clade, peak calling is performed on the
1194 aggregated scATAC-seq profiles, and identified peaks are combined into a new peaks-
1195 by-cells binary matrix. Finally, the new peaks-by-cells matrix is transformed with TF-
1196 IDF and SVD as before to get a matrix of PCs-by-cells, which is the final features matrix
1197 for downstream analyses.

1198 **scABC** This method starts by building a peaks-by-cells count matrix of read coverage
1199 within peak regions. Then, the weights of cells are calculated by a nonlinear
1200 transformation of the read coverage within the peaks background, defined as a 500 kb
1201 region around peaks. Since the weights will be used as part of weighted K-medoids
1202 clustering to define cell landmarks and further perform finer re-clustering instead of
1203 normalizing the peaks-by-cells matrix, the feature matrix in scABC is defined as the
1204 peaks-by-cells count matrix.

1205 **Scasat** This method first constructs a peaks-by-cells binary accessibility matrix by
1206 checking if at least one read overlaps with the peak region. Then Jaccard distance is
1207 computed based on the binary matrix to get a cells-by-cells dissimilarity matrix.
1208 Multidimensional scaling (MDS) is further performed to reduce the dimension and to
1209 generate the final feature matrix for downstream analysis.

1210 **SCRAT** This method starts by aggregating reads from each cell according to different
1211 features (such as TF motifs or region of interest of each gene), and then building a count
1212 matrix of features-by-cells. The features-by-cells count matrix is normalized by library
1213 and region size to get the final feature matrix for downstream analyses.

1214 **SnapATAC** This method starts by binning the genome into fixed-size windows (by
1215 default 5kb) and estimating read coverage for each bin to build a bins-by-cells binary
1216 count matrix. Bins that overlap ENCODE-defined blacklist regions are filtered out, as
1217 well as those with exceedingly high or low z-scored coverage. Then, the bins-by-cell
1218 matrix is transformed into a cells-by-cells Jaccard index similarity matrix, which is
1219 further transformed by normalization and regressing out coverage bias between cells.
1220 Finally, PCA is applied to the normalized similarity matrix, and the top PCs are used to
1221 build a PCs-by-cells matrix that is the final features matrix for downstream analyses.

1222 *Clustering*

1223 For this study we used three commonly used clustering methods: k-means, hierarchical
1224 clustering (with default ward linkage) as implemented in the scikit-learn library [31]
1225 and Louvain clustering (a community-detection-based method) [32, 33] as implemented
1226 in Scanpy [34], For both hierarchical clustering and k-means, we set the number of
1227 clusters to the number of unique FACS-sorted labels or known tissues. In the 10X
1228 PBMCs scATAC-seq dataset, which lacks the FACS-sorted labels, we instead set the
1229 number of clusters to 8 since this is the expected number of populations based on
1230 previous studies [19]. For the Louvain algorithm, we set the size of local neighborhood
1231 to 15 for all the datasets. Since Louvain method requires ‘resolution’ instead of the
1232 number of clusters and different number of clusters will affect the clustering evaluation,
1233 to make the comparison fair, we use the binary search algorithm on the ‘resolution’
1234 (ranging from 0.0 to 3.0) to find the same number of clusters as the other two clustering
1235 methods. If the precise number of clusters did not match the desired value, the
1236 ‘resolution’ value inducing the closest number of clusters to the desired value was used.

1237 *Metrics for evaluating clustering results*

1238 To evaluate clustering solutions for datasets with a known ground truth (i.e. for each
1239 cell we have a label that indicated the cell type) we used three well-established metrics:
1240 Adjusted Rand Index (ARI), Mutual information and Homogeneity. Briefly, for the

1241 Adjusted Rand Index (ARI) first the Random Index (RI) is defined as a similarity
1242 measure between two clusters considering all pairs of samples assigned in the same or
1243 different clusters in the predicted and true clustering. Then, the raw RI score is adjusted
1244 for chance in the ARI score as described in the following formula:

1245

$$1246 \quad ARI = \frac{RI - E(RI)}{\max(RI) - E(RI)}$$

1247

1248

1249 Where RI is the pre-computed random index and E is the expected random index.

1250 Mutual Information is a measure of the mutual dependence between two variables. The
1251 Mutual Information value is computed according to the following formula, where $|U_i|$
1252 is the number of the samples in cluster U_i and $|V_j|$ is the number of the samples in
1253 cluster V_j :

$$1254 \quad MI(U, V) = \sum_{i=1}^{|U|} \sum_{j=1}^{|V|} \frac{|U_i \cap V_j|}{N} \log \frac{N|U_i \cap V_j|}{|U_i||V_j|}$$

1255

1256

1257 The homogeneity score is used to check if the algorithm used for the clustering can
1258 assign to each cluster only samples belonging to a single class. Its value h is bounded
1259 between 0 and 1, and a low value indicates low homogeneity and vice versa. The score
1260 is computed as follow:

$$1261 \quad h = 1 - \frac{H(Y_{true}|Y_{pred})}{H(Y_{true})}$$

1262

1263

1264 where $H(Y_{true} | Y_{pred})$ is the probability to assign true samples to a set of predicted
1265 samples, while $H(Y_{true})$ are the labels of the samples.

1266 To evaluate clustering solutions for the 10X PBMCs dataset we proposed a simple score
1267 called the Residual Average Gini Index (RAGI) and compared the accessibility of
1268 housekeeping genes with previously characterized marker genes [19]. We reasoned that
1269 a good clustering solution should contain clusters that are enriched for accessibility of
1270 different marker genes, and each marker gene should be highly accessible in only one or
1271 a few clusters. First, to quantify the accessibility of each gene in each cell we used the
1272 Gene Scoring approach described above. Briefly, the accessibility at each TSS is the
1273 distance-weighted sum of reads within or near the region. Second, to quantify the
1274 enrichment of each gene in each cluster of cells, we computed the mean of the
1275 accessibility values in all cells for each cluster. Third, based on the vector of mean
1276 accessibility values (one per cluster), we computed the Gini Index [35] for each marker
1277 gene. The Gini Index measures how imbalanced the accessibility of a gene is across
1278 clusters. This score is bound by [0,1] where 1 means total imbalance (i.e. a gene is
1279 accessible in one cluster only) and 0 means no enrichment. This score has been
1280 previously used on scRNA-seq to perform clustering [36, 37]. As a control, we also
1281 calculated the Gini Index for a set of annotated housekeeping genes reported in
1282 (https://m.tau.ac.il/~elieis/HKG/HK_genes.txt). Housekeeping genes should show
1283 minimal specificity for any given cluster since, by definition, they are highly expressed
1284 in all cells. Based on the set of Gini Index values for marker and housekeeping genes we
1285 calculated several metrics: (i) the mean Gini Index for the two groups; (ii) the difference
1286 in means to assess the average residual specificity that a clustering solution has with
1287 respect to marker genes (this is our proposed RAGI metric); and (iii) the Kolmogorov
1288 Smirnov statistic and its p-value comparing the two groups of Gini Indices for marker
1289 and house-keeping genes. We sorted the methods based on the descending order of the
1290 differences in means (**Sup Table 13**); a positive value indicates that the marker genes on
1291 average separate the clusters better than uninformative housekeeping genes.

1292 *Rare cell type-specific peak analysis*

1293 FACS-sorted bulk ATAC-seq data was downloaded and processed from a previously
1294 described resource [5]. For each simulation, we created a randomly-sampled set of 200
1295 million unique (PCR-deduplicated) reads, which roughly represents a complexity
1296 similar to recommendations from the 10X Chromium scATAC-seq solution. Cell type-
1297 specific peaks were defined using the full dataset for each of the three cell types. Peaks
1298 were called using macs2 callpeak with custom parameters as in the ENCODE pipeline,
1299 i.e. “--nomodel --shift - 100 --extsize 200” to account for Tn5 insertions rather than read
1300 abundance when inferring peaks. Overlaps between the isolated minor population and
1301 the synthetic mixtures were computed using GenomicRanges[38] findOverlaps
1302 function, which is equivalent to bedtools[30] overlap. For each minor population (B-cell,
1303 CD4+ T-cell, Monocyte) and each prevalence (1, 5, 10, 20, 30%), each simulation was
1304 repeated 5 times for a total of 75 simulations. Reads from the other two (major)
1305 populations were sampled equivalently to make up the synthetic mixture for
1306 comparison.

1307 *Data and code availability*

1308 All the results presented in this manuscript can be reproduced using the Jupyter
1309 notebooks available both at <https://github.com/pinelloab/scATAC-benchmarking/> and
1310 in the supplementary material (**Sup Data**). For the analyzed real datasets, the
1311 *Buentrostro2018* dataset was downloaded from GEO accession GSE96772, the 10X
1312 PBMCs dataset was downloaded from https://support.10xgenomics.com/single-cell-atac/datasets/1.0.1/atac_v1_pbmc_5k, and the sci-ATAC-seq mouse dataset was
1313 downloaded from
1314 [http://krishna.gs.washington.edu/content/members/ajh24/mouse_atlas_data_release/ba](http://krishna.gs.washington.edu/content/members/ajh24/mouse_atlas_data_release/bams)
1315 [ms](http://krishna.gs.washington.edu/content/members/ajh24/mouse_atlas_data_release/bams). For the simulated bone marrow dataset, data for the FACS-sorted bulk ATAC-seq
1316 populations were downloaded from GEO accession GSE119453. For the simulated
1317 erythropoiesis dataset, the additional populations were downloaded from GEO
1318 accession GSE115672.

1320 *Author Contributions*

1321 H.C. and L.P. conceived this project and designed the framework with input from all
1322 the authors. H.C. and T.A. preprocessed data. H.C., C.L., T.A., M.E.V., S.P.G. and K.C.
1323 implemented scATAC-seq methods. H.C. and K.C. performed clustering analysis. H.C.,

1324 T.A., L.P. performed clustering validation. C.L. simulated data. H.C. and C.L. analyzed
1325 simulated data. L.P. and J.D.B. provided guidance. All the authors wrote the
1326 manuscript.

1327 *Acknowledgments*

1328 This project has been made possible in part by grant number 2018- 182734 to L.P. from
1329 the Chan Zuckerberg Initiative DAF, an advised fund of Silicon Valley Community
1330 Foundation. L.P. is also partially supported by a National Human Genome Research
1331 Institute (NHGRI) Career Development Award (R00HG008399).

1332 *Conflicts of interest*

1333 J.D.B. holds a patent for the invention of ATAC-seq.

1334

1335

1336 **References**

1337

- 1338 1. Cusanovich, D.A., et al., *A Single-Cell Atlas of In Vivo Mammalian Chromatin*
1339 *Accessibility*. Cell, 2018. **174**(5): p. 1309-1324 e18.
- 1340 2. Schep, A.N., et al., *chromVAR: inferring transcription-factor-associated accessibility from*
1341 *single-cell epigenomic data*. Nat Methods, 2017. **14**(10): p. 975-978.
- 1342 3. de Boer, C.G. and A. Regev, *BROCKMAN: deciphering variance in epigenomic regulators*
1343 *by k-mer factorization*. BMC Bioinformatics, 2018. **19**(1): p. 253.
- 1344 4. Ji, Z., W. Zhou, and H. Ji, *Single-cell regulome data analysis by SCRAT*. Bioinformatics,
1345 2017. **33**(18): p. 2930-2932.
- 1346 5. Corces, M.R., et al., *Lineage-specific and single-cell chromatin accessibility charts human*
1347 *hematopoiesis and leukemia evolution*. Nat Genet, 2016. **48**(10): p. 1193-203.
- 1348 6. Kiselev, V.Y., T.S. Andrews, and M. Hemberg, *Challenges in unsupervised clustering of*
1349 *single-cell RNA-seq data*. Nat Rev Genet, 2019.
- 1350 7. McInnes, L., J. Healy, and J. Melville, *Umap: Uniform manifold approximation and*
1351 *projection for dimension reduction*. arXiv preprint arXiv:1802.03426, 2018.
- 1352 8. Pliner, H.A., et al., *Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell*
1353 *Chromatin Accessibility Data*. Mol Cell, 2018. **71**(5): p. 858-871 e8.
- 1354 9. Bravo González-Blas, C., et al., *cisTopic: cis-regulatory topic modeling on single-cell*
1355 *ATAC-seq data*. Nature Methods, 2019.
- 1356 10. Cusanovich, D.A., et al., *The cis-regulatory dynamics of embryonic development at*
1357 *single-cell resolution*. Nature, 2018. **555**(7697): p. 538-542.

- 1358 11. Cusanovich, D.A., et al., *Multiplex single cell profiling of chromatin accessibility by*
1359 *combinatorial cellular indexing*. Science, 2015. **348**(6237): p. 910-4.
- 1360 12. Lareau, C.A., et al., *Droplet-based combinatorial indexing for massive scale single-cell*
1361 *epigenomics*. bioRxiv, 2019: p. 612713.
- 1362 13. Zamanighomi, M., et al., *Unsupervised clustering and epigenetic classification of single*
1363 *cells*. Nat Commun, 2018. **9**(1): p. 2410.
- 1364 14. Baker, S.M., et al., *Classifying cells with Scasat, a single-cell ATAC-seq analysis tool*.
1365 Nucleic Acids Res, 2019. **47**(2): p. e10.
- 1366 15. Fang, R., et al., *Fast and Accurate Clustering of Single Cell Epigenomes Reveals Cis-*
1367 *Regulatory Elements in Rare Cell Types*. BioRxiv, 2019.
- 1368 16. Mathelier, A., et al., *JASPAR 2016: a major expansion and update of the open-access*
1369 *database of transcription factor binding profiles*. Nucleic acids research, 2015. **44**(D1): p.
1370 D110-D115.
- 1371 17. Leif S. Ludwig, et al., *Transcriptional States and Chromatin Accessibility Underlying*
1372 *Human Erythropoiesis*. Cell Reports, 2019.
- 1373 18. Buenrostro, J.D., et al., *Integrated Single-Cell Analysis Maps the Continuous Regulatory*
1374 *Landscape of Human Hematopoietic Differentiation*. Cell, 2018. **173**(6): p. 1535-1548
1375 e16.
- 1376 19. Pliner, H.A., J. Shendure, and C. Trapnell, *Supervised classification enables rapid*
1377 *annotation of cell atlases*. BioRxiv, 2019.
- 1378 20. Consortium, E.P., *An integrated encyclopedia of DNA elements in the human genome*.
1379 Nature, 2012. **489**(7414): p. 57.
- 1380 21. Bernstein, B.E., et al., *The NIH roadmap epigenomics mapping consortium*. Nature
1381 biotechnology, 2010. **28**(10): p. 1045.
- 1382 22. Yoshida, H., et al., *The cis-Regulatory Atlas of the Mouse Immune System*. Cell, 2019.
1383 **176**(4): p. 897-912. e20.
- 1384 23. Stark, R. and G. Brown, *DiffBind: differential binding analysis of ChIP-Seq peak data*. R
1385 package version, 2011. **100**: p. 4-3.
- 1386 24. Amemiya, H.M., A. Kundaje, and A.P. Boyle, *The ENCODE Blacklist: Identification of*
1387 *Problematic Regions of the Genome*. Sci Rep, 2019. **9**(1): p. 9354.
- 1388 25. Satpathy, A.T., et al., *Massively parallel single-cell chromatin landscapes of human*
1389 *immune cell development and intratumoral T cell exhaustion*. BioRxiv, 2019: p. 610550.
- 1390 26. Chen, H., et al., *Single-cell trajectories reconstruction, exploration and mapping of omics*
1391 *data with STREAM*. Nature Communications, 2019. **10**(1).
- 1392 27. Qiu, X., et al., *Reversed graph embedding resolves complex single-cell trajectories*. Nat
1393 Methods, 2017. **14**(10): p. 979-982.
- 1394 28. Ulirsch, J.C., et al., *Interrogation of human hematopoiesis at single-cell and single-*
1395 *variant resolution*. Nature genetics, 2019: p. 1.
- 1396 29. Quinlan, A.R. and I.M. Hall, *BEDTools: a flexible suite of utilities for comparing genomic*
1397 *features*. Bioinformatics, 2010. **26**(6): p. 841-842.
- 1398 30. Pedregosa, F., et al., *Scikit-learn: Machine learning in Python*. Journal of machine
1399 learning research, 2011. **12**(Oct): p. 2825-2830.
- 1400 31. Blondel, V.D., et al., *Fast unfolding of communities in large networks*. Journal of
1401 statistical mechanics: theory and experiment, 2008. **2008**(10): p. P10008.

- 1402 32. Levine, J.H., et al., *Data-driven phenotypic dissection of AML reveals progenitor-like cells*
1403 *that correlate with prognosis*. Cell, 2015. **162**(1): p. 184-197.
- 1404 33. Wolf, F.A., P. Angerer, and F.J. Theis, *SCANPY: large-scale single-cell gene expression*
1405 *data analysis*. Genome biology, 2018. **19**(1): p. 15.
- 1406 34. Gini, C., *Concentration and dependency ratios*. Rivista di politica economica, 1997. **87**: p.
1407 769-792.
- 1408 35. Jiang, L., et al., *GiniClust: detecting rare cell types from single-cell gene expression data*
1409 *with Gini index*. Genome Biol, 2016. **17**(1): p. 144.
- 1410 36. Tsoucas, D. and G.C. Yuan, *GiniClust2: a cluster-aware, weighted ensemble clustering*
1411 *method for cell-type detection*. Genome Biol, 2018. **19**(1): p. 58.
- 1412 37. Lawrence, M., et al., *Software for computing and annotating genomic ranges*. PLoS
1413 computational biology, 2013. **9**(8): p. e1003118.

1414

1415

1416

1417

1418

1419

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439 **Supplementary Notes**

1440

1441 **Supplementary Note 1: Analysis of the simulated datasets**

1442 For all the synthetic datasets, the input is a peaks-by-cells raw count matrix generated
1443 as described in the Methods section. For all methods, we first order peaks based on the
1444 number of cells in which the peak is observed and select the top 8,000 peaks (making
1445 sure each of these peaks appear at least in one cell).

1446

1447 For BROCKMAN, we scanned for gapped k-mers (the default setting is used, i.e. length
1448 1–8, all possible gaps) within peaks to calculate the scaled k-mer frequencies for each
1449 cell. For chromVAR, we used both TF binding motifs from the JASPAR database
1450 (human) or short k-mers (k=6) within peaks to score the accessibility deviation across
1451 cells. For Cicero, we run it with the default parameters to calculate gene activity scores.
1452 For cisTopic, we run it with the same parameters shown in their online tutorial
1453 ([https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db594](https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db5940f7/vignettes/CompleteAnalysis.html)
1454 [0f7/vignettes/CompleteAnalysis.html](https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db5940f7/vignettes/CompleteAnalysis.html)). For *Cusanovich2018* we first binarize the count
1455 matrix and then perform the proposed TF-IDF transformation and SVD. For Gene
1456 Scoring, we select peaks overlapping with the regions of 50,000 bp upstream and
1457 downstream of TSSs as described in [1]. For scABC, since its feature matrix is the same
1458 as input matrix of peaks-by-cells, we instead run the steps of calculating the weights of
1459 cells that are used later for their proposed clustering approach. For Scasat, we first
1460 binarize the count matrix and then calculate Jaccard distance, followed by Multi
1461 Dimensional Scaling (MDS) with 10 dimensions (the same number of components as
1462 used for the Control-Naive). For SCRAT, the accessibility of TF binding motifs is
1463 summarized within peaks. We attempted to adjust for the library size and peak region
1464 length as suggested in the original study, however we noticed that this step
1465 dramatically penalizes this method performance in all the tested conditions (**Sup Fig. 1**).
1466 This step was therefore disabled for all the analyses performed with SCRAT. For
1467 SnapATAC, we use the fixed-size peaks as its bins. The Jaccard Index is normalized
1468 with the authors' proposed method, *normOVE*. For methods that implement PCA step,
1469 we use the elbow plot to decide the optimal number of PCs. For methods that do not
1470 implement a step of dimensionality reduction, we use the R package *irlba* [2] to perform
1471 PCA.

1472

1473 All the notebooks detailing the exact procedures are available at
1474 https://github.com/pinellolab/scATAC-benchmarking/tree/master/Synthetic_Data.

1475

1476 **Supplementary Note 2: Analysis of the Buenrostro2018 dataset**

1477

1478 For this dataset we started with aligned files in bam format (one per cell). We removed
1479 duplicated reads using the function *MarkDuplicates* version 2.20.2 with the option
1480 *REMOVE DUPLICATES = TRUE* from Picard (<https://broadinstitute.github.io/picard/>).

1481

1482 For the methods that do not provide an explicit function to read in bam files and count
1483 reads under peaks, including Cicero, Cusanovich2018, GeneScoring, Scasat, and
1484 Control-Naïve, we used a simple script to obtain a common peaks-by-cells raw count
1485 matrix (e.g. [https://github.com/pinellolab/scATAC-
1486 benchmarking/tree/master/Real_Data/Buenrostro_2018/run_methods/Cusanovich2018/c
1487 ount_reads_peaks.sh](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real_Data/Buenrostro_2018/run_methods/Cusanovich2018/count_reads_peaks.sh)). For the methods that implement the same strategy to filter peaks
1488 based on their frequency, including Cicero, Control-Naive, Cusanovich2018,
1489 GeneSoring, Scasat, and scABC, we filter out peaks that are observed in less than 1% of
1490 cells. For chromVAR, we run its function *filterPeaks* with the default setting to filter out
1491 peaks based on the minimum number of fragments and merge overlapping peaks. For
1492 the methods that implement a PCA step, including BROCKMAN, Control-Naïve,
1493 Cusanovich2018, and SnapATAC, we decided the number of PCs based on the elbow
1494 plot. For Scasat, which implements MDS, we set the number of dimension as 15
1495 according to its tutorial
1496 [https://github.com/ManchesterBioinference/Scasat/blob/master/ScAsAT_functions_Bue
1497 nrostro_All_Bam_Together.ipynb](https://github.com/ManchesterBioinference/Scasat/blob/master/ScAsAT_functions_Buenrostro_All_Bam_Together.ipynb). For cisTopic, the number of topics (dimensions) is
1498 decided by its function *selectModel* with default settings.

1499

1500 For the clustering analysis, we set the expected number of clusters as the number of
1501 FACS-sorting labels (10 in this case). For k-means, we use the *k-means++* to select the initial
1502 cluster centers. For hierarchical clustering, we use the *Ward* linkage based on Euclidean
1503 distance. Both k-means and hierarchical clustering are implemented in scikit-learn
1504 package[3]. For Louvain, we set the number of neighbors to 15 and the resolution is
1505 decided using a binary search with 20 steps that explores values of the resolution
1506 parameter in the interval 0~3 . The Louvain algorithm used is implemented in Scanpy[4].

1507

1508 For the UMAP visualization, we run the function 'umap' from the R package *umap* with
1509 default settings.

1510

1511 All the notebooks for this analysis are available at

1512 <https://github.com/pinellolab/scATAC->

1513 [benchmarking/tree/master/Real Data/Buenrostro 2018](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real%20Data/Buenrostro%202018) and

1514 <https://github.com/pinellolab/scATAC->

1515 [benchmarking/tree/master/Real Data/Buenrostro 2018 bulkpeaks](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real%20Data/Buenrostro%202018%20bulkpeaks).

1516

1517 **Supplementary Note 3: Analysis of 10x PBMCs dataset**

1518

1519 For this dataset, we started with a single merged bam file downloaded from the 10x
1520 website and preprocessed with Cell Ranger: [https://support.10xgenomics.com/single-](https://support.10xgenomics.com/single-cell-atac/datasets/1.0.1/atac_v1_pbmc_5k)

1521 [cell-atac/datasets/1.0.1/atac v1 pbmc 5k](https://support.10xgenomics.com/single-cell-atac/datasets/1.0.1/atac_v1_pbmc_5k). We noticed that all the methods except

1522 SnapATAC don't support this format i.e. a single bam file for multiple cells. Therefore,

1523 using the cell barcodes passing quality filtering from Cell Ranger, we split this file in

1524 multiple bam files, one per cell recovering 5,335 single-cell bam files. We also removed

1525 duplicate reads using Picard and performed UMAP visualization as discussed in

1526 **Supplementary Note 2** . For the clustering analysis, we set the expected number of

1527 clusters as the number of putative cell types (8 in this case) as previous studies suggested

1528 [5, 6].

1529

1530 All the notebooks are available at <https://github.com/pinellolab/scATAC->

1531 [benchmarking/tree/master/Real Data/10x PBMC 5k](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real%20Data/10x%20PBMC%205k).

1532

1533 **Supplementary Note 4: Analysis of the sci-ATAC-seq mouse dataset**

1534

1535 For this dataset, we started with multiple merged bam file from 17 samples across 13

1536 tissues downloaded from

1537 [http://krishna.gs.washington.edu/content/members/ajh24/mouse atlas data release/ba](http://krishna.gs.washington.edu/content/members/ajh24/mouse_atlas_data_release/ba)

1538 [ms](http://krishna.gs.washington.edu/content/members/ajh24/mouse_atlas_data_release/ba). For each tissue we performed the same steps as in 10x PBMCs dataset to decompose

1539 the single merged bam file to multiple single cell bam files (81,173 bam files). The

1540 downloaded bam files were already deduplicated. The downsampled dataset of 12,178
1541 cells is generated by randomly selecting 15% from each sample.

1542

1543 The scATAC-seq methods and UMAP visualization are implemented as in
1544 **Supplementary Note 2**. For the clustering analysis, we set the expected number of
1545 clusters as the number of tissues (13 in this case).

1546

1547 All the notebooks are available at [https://github.com/pinellolab/scATAC-](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real_Data/Cusanovich_2018)
1548 [benchmarking/tree/master/Real_Data/Cusanovich_2018](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real_Data/Cusanovich_2018) and
1549 [https://github.com/pinellolab/scATAC-](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real_Data/Cusanovich_2018_subset)
1550 [benchmarking/tree/master/Real_Data/Cusanovich_2018_subset](https://github.com/pinellolab/scATAC-benchmarking/tree/master/Real_Data/Cusanovich_2018_subset).

1551

1552 **Supplementary Note 5: Memory requirements and implementation choices**

1553

1554 As mentioned in the main text, SnapATAC is the only methods that allows to process
1555 successfully large datasets, as the sciATAC-seq mouse dataset with ~80000 cells. Here we
1556 investigate why the other methods failed to analyze this large dataset. We hypothesize
1557 that main reason is related to the way the methods load/process the data in memory. In
1558 fact, we discovered that several methods require to load the entire dataset in the central
1559 memory (RAM).

1560

1561 BROCKMAN, Cicero and Gene Scoring try to load the entire dataset in memory using
1562 the *read.table* function or the *fread* function within the *data.table* package in R. Other
1563 methods such as: Cusanovic, Scrat, chromVAR, scABC and Scasat, store the entire
1564 dataset in memory within a *Matrix* object in R. CisTopic, has an optimized step to map
1565 the reads into the genome using the *Rsubread* function. This function creates a hash table
1566 of the entire genome and allows the user to select the amount of memory to use. At the
1567 end, the entire dataset is stored in the computer memory in a *CisTopicObject* data
1568 structure.

1569

1570 SnapATAC, preprocess the entire dataset and store it a *.snap* file. This file is based on the
1571 HDF5 technology that allows out of core computation. In SnapATAC the Python library
1572 h5py (a wrapper for HDF5 core library) is used to create the custom snap file

1573 format. More information about this custom file are available here :
1574 https://github.com/r3fang/SnapTools/blob/master/docs/snap_format.docx .

1575

1576 **Supplementary Note 6: End-to-end user-perspective clustering analysis**

1577

1578 For the methods that explicitly implement the step of clustering in their tutorials,
1579 including Cusanovich2018, cisTopic, SnapATAC, scABC, Cicero, and Scasat, in addition
1580 to the three clustering methods used in this benchmark framework, we also performed
1581 the clustering analysis as shown in each tutorial. For Cusanovich2018, we followed the
1582 tutorial at <http://atlas.gs.washington.edu/fly-atac/docs/> and used density peak
1583 algorithm [7] to identify clusters. For cisTopic, we followed the tutorial at
1584 [https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db5940](https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db5940f7/vignettes/CompleteAnalysis.html)
1585 [f7/vignettes/CompleteAnalysis.html](https://rawcdn.githack.com/aertslab/cisTopic/f628c6f60918511ba0fa4a85366ebf52db5940f7/vignettes/CompleteAnalysis.html) and used ward hierarchical clustering to cluster
1586 cells. For SnapATAC, we followed tutorial at
1587 https://github.com/r3fang/SnapATAC/blob/master/examples/10X_P50/README.md
1588 and used Leiden algorithm to cluster cells. For scABC we followed the tutorial at
1589 <https://github.com/SUwonglab/scABC/blob/master/vignettes/ExampleWorkflow.html>
1590 and used weighted k-medoids clustering. For Cicero we followed the tutorial at
1591 [https://www.bioconductor.org/packages/devel/bioc/vignettes/cicero/inst/doc/website.ht](https://www.bioconductor.org/packages/devel/bioc/vignettes/cicero/inst/doc/website.html)
1592 [ml](https://www.bioconductor.org/packages/devel/bioc/vignettes/cicero/inst/doc/website.html) . To be consistent with the feature matrix used in the benchmarking framework,
1593 instead of using its default peaks-by-cells count matrix, we used gene activity scores as
1594 the input of clustering analysis. After reducing the dimensionality with tSNE, density
1595 peak clustering algorithm is used to cluster cells. For Scasat, we follow the tutorial at
1596 [https://github.com/ManchesterBioinference/Scasat/blob/master/ScAsAT_functions_Bue](https://github.com/ManchesterBioinference/Scasat/blob/master/ScAsAT_functions_Buenroostro_All_Bam_Together.ipynb)
1597 [nroostro_All_Bam_Together.ipynb](https://github.com/ManchesterBioinference/Scasat/blob/master/ScAsAT_functions_Buenroostro_All_Bam_Together.ipynb) and use *ward.D2* hierarchical clustering for clustering.

1598

1599 We run all the six methods on three real datasets, Buenroostro2018, 10x PBMCs 10x
1600 dataset, sci-ATAC-seq mouse dataset. For Buenroostro2018 and sci-ATAC-seq mouse
1601 dataset, we specified the number of clusters as the number of FACS-sorted labels and
1602 the number of tissues respectively. For 10x PBMCs, we specified the number of clusters
1603 as 8 as suggested by the previous studies [5, 6].

1604

1605

1606 **Supplementary Note 7: Running time**

1607

1608 For the real datasets, we recorded the execution time of each method to generate a
1609 feature matrix starting from an aligned and deduplicated bam file. We noticed that not
1610 all the methods provide specific functions to read in bam files. Some methods only start
1611 with features by cells raw matrix (e.g. Cicero). In addition, the functions to count reads
1612 of some methods were not generalizable across the different scATAC-seq techniques
1613 (e.g. *Cusanovich2018*). Therefore, to make a fair comparison we used a common script
1614 ([https://github.com/pinelloab/scATAC-](https://github.com/pinelloab/scATAC-benchmarking/blob/master/Real%20Data/Buenrostro%202018/run_methods/Control/count_reads_peaks.sh)
1615 [benchmarking/blob/master/Real Data/Buenrostro 2018/run_methods/Control/count reads](https://github.com/pinelloab/scATAC-benchmarking/blob/master/Real%20Data/Buenrostro%202018/run_methods/Control/count_reads_peaks.sh)
1616 [peaks.sh](https://github.com/pinelloab/scATAC-benchmarking/blob/master/Real%20Data/Buenrostro%202018/run_methods/Control/count_reads_peaks.sh)) to obtain the peaks by cells matrix starting from bam files for the following 4
1617 methods: Control-Naïve, *Cusanovich2018*, Gene Scoring, Scasat. BROCKMAN, perform
1618 two steps to obtain the final feature matrix (q bash script to count k-mer frequency and
1619 a R function to assemble the matrix), so we are considering the sum of their running
1620 times. Similarly, the running time for SnapATAC is based on two steps: the *snaptools*
1621 utility that converts a bam to the required *.snap* format and the R function that generates
1622 the feature matrix.

1623

1624 For the simulated datasets, we recorded the execution time of generating feature
1625 matrices starting from a simulated peaks-by-cell count matrix. For scABC, since its
1626 feature matrix is the same as the input, to have a useful running time, we instead record
1627 the time to calculate the cells weights, which are necessary for downstream analysis.

1628

1629 We also assessed the scalability of the methods with respect to the read coverage (250,
1630 500, 1000, 2500 and 5000 fragments per peaks). We observed that the running time of most
1631 methods is not affected by the read coverage. This is not surprising given that our
1632 simulation them number of peaks is fixed, so the dimensionality of the matrix is
1633 unchanged. However, for cisTopic, we noticed an exponential increase in running times
1634 as we increase the number of fragments (**Sup Fig. 22**). We assume this might be due to
1635 the topic modelling approach used by cisTopic since it tries to learn the probability
1636 distribution over the regions for each topic while high coverage will result in the increase
1637 in the number of accessible regions.

1638

1639

1640 1. Lareau, C.A., et al., *Droplet-based combinatorial indexing for massive scale single-cell*
1641 *epigenomics*. bioRxiv, 2019: p. 612713.

- 1642 2. Baglama, J. and L. Reichel, *Augmented implicitly restarted Lanczos bidiagonalization*
1643 *methods*. SIAM Journal on Scientific Computing, 2005. **27**(1): p. 19-42.
- 1644 3. Pedregosa, F., et al., *Scikit-learn: Machine learning in Python*. Journal of machine
1645 learning research, 2011. **12**(Oct): p. 2825-2830.
- 1646 4. Wolf, F.A., P. Angerer, and F.J. Theis, *SCANPY: large-scale single-cell gene expression*
1647 *data analysis*. Genome Biol, 2018. **19**(1): p. 15.
- 1648 5. Pliner, H.A., J. Shendure, and C. Trapnell, *Supervised classification enables rapid*
1649 *annotation of cell atlases*. BioRxiv, 2019.
- 1650 6. Bravo González-Blas, C., et al., *cisTopic: cis-regulatory topic modeling on single-cell*
1651 *ATAC-seq data*. Nature Methods, 2019.
- 1652 7. Rodriguez, A. and A. Laio, *Clustering by fast search and find of density peaks*. Science,
1653 2014. **344**(6191): p. 1492-1496.

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674

1675

1676

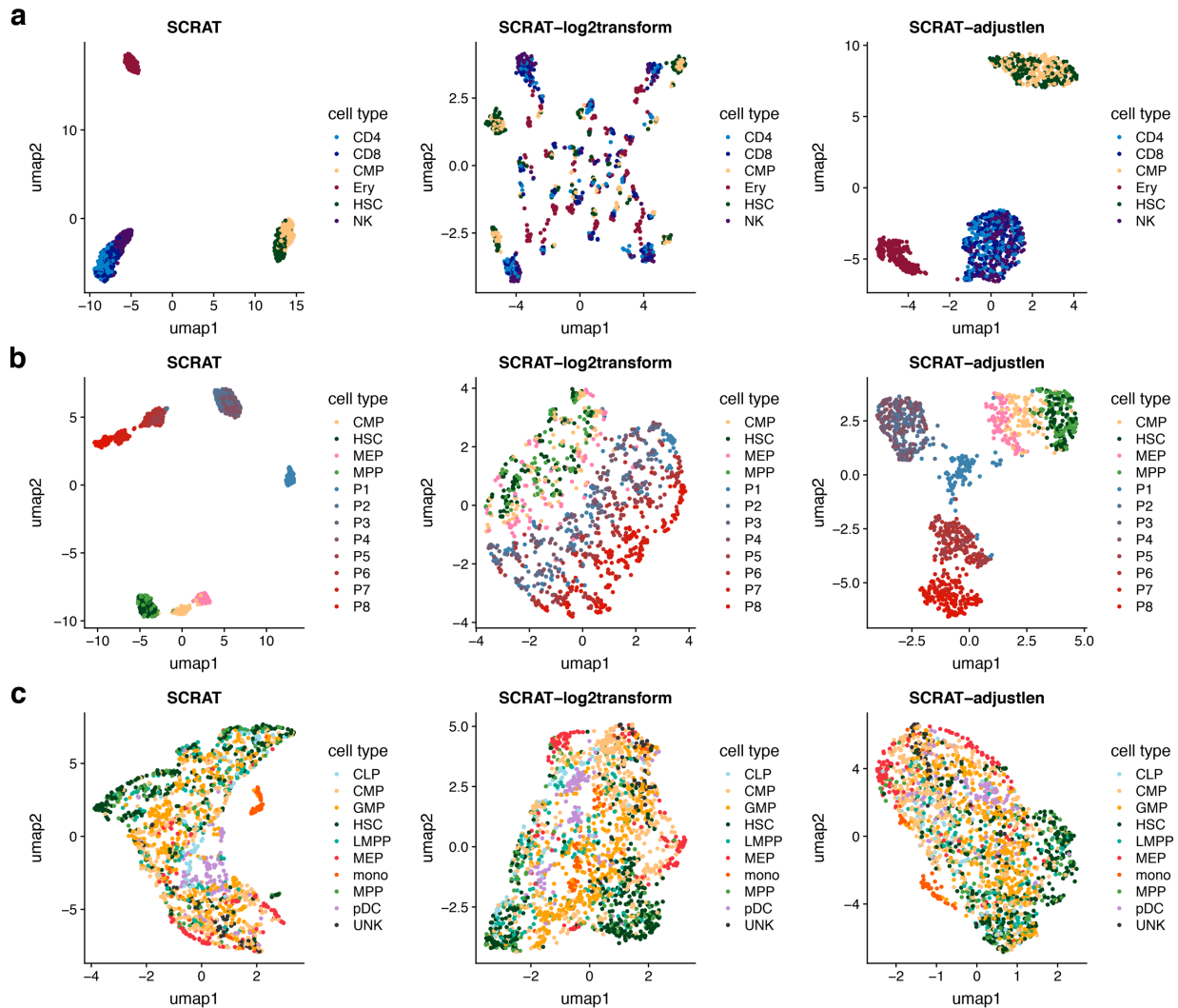
1677

1678

1679 **Supplementary Figures**

1680

1681



1682

1683

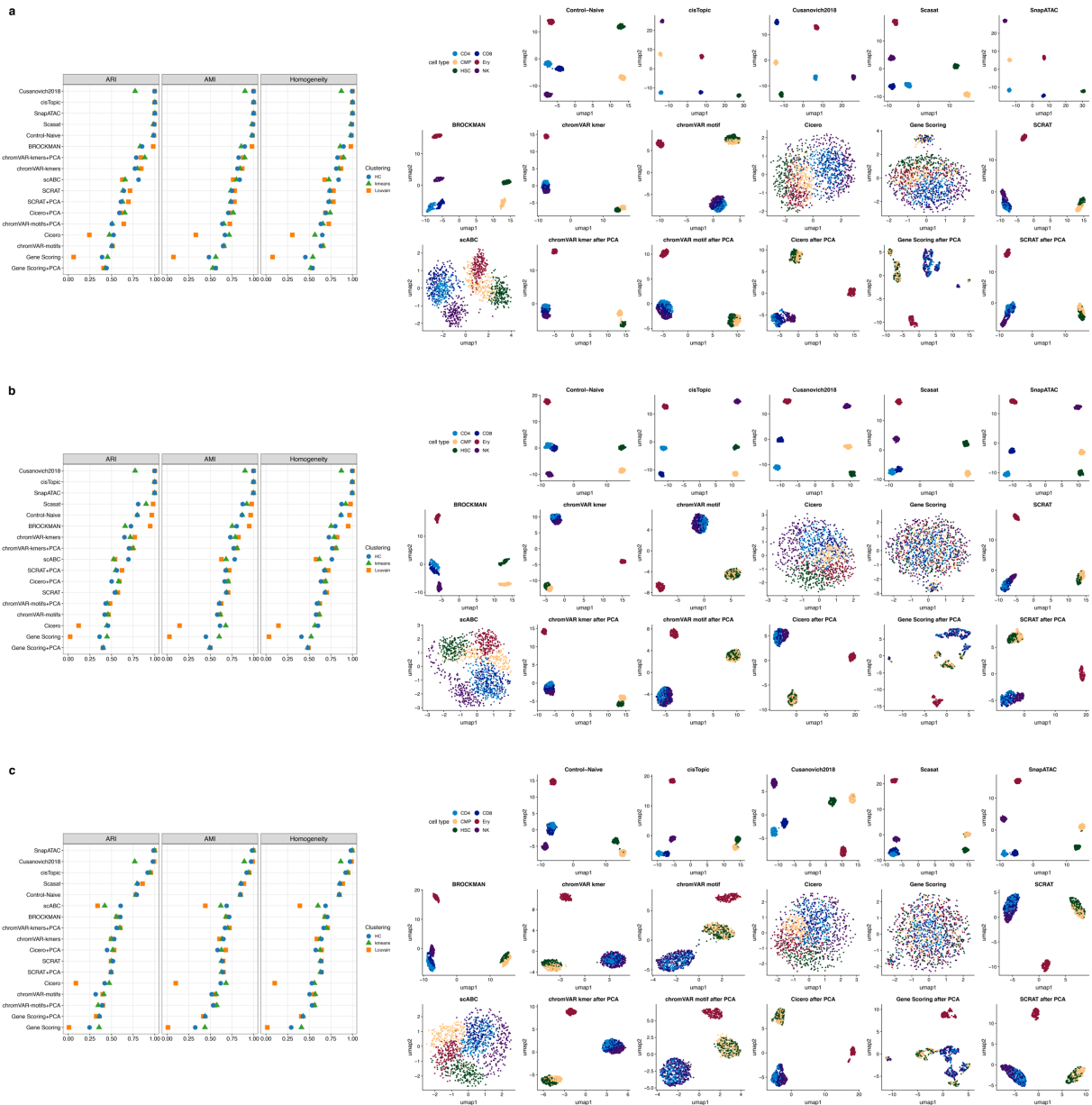
1684 **Figure S1.** UMAP visualization of cells based on SCRAT feature matrix with different parameter
 1685 settings (Left: log2transform=FALSE, adjustlen=FALSE. Middle: log2transform=TRUE,
 1686 adjustlen=FALSE. Right: log2transform=FALSE, adjustlen=TRUE) in three datasets. **(a)** simulated
 1687 bone marrow dataset at a noise level of 0.2 with a coverage of 2,500 fragments **(b)** simulated
 1688 erythropoiesis dataset at a noise level of 0.2 with a coverage of 2,500 fragments **(c)** *Buenrostro*
 1689 *2018* dataset.

1690

1691

1692

1693



1694

1695

1696 **Figure S2.** Clustering evaluation according to AMI, ARI and Homogeneity metrics (*left*) and UMAP

1697 visualization of cells colored by known cell labels (*right*) in simulated bone marrow datasets with

1698 a coverage of 2,500 fragments at **(a)** no noise (0), **(b)** moderate noise (0.2) and **(c)** high noise

1699 (0.4).

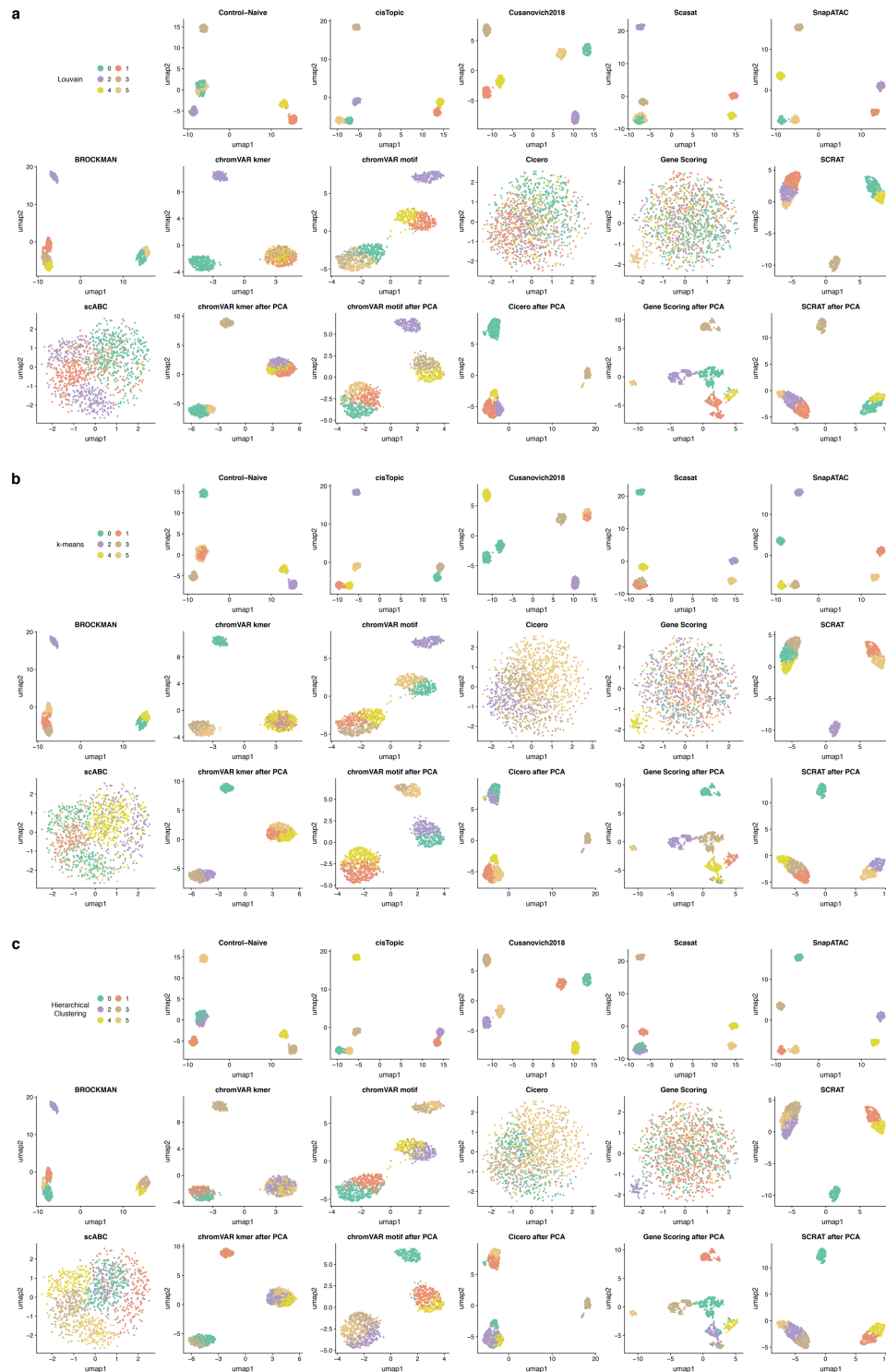
1700

1701

1702

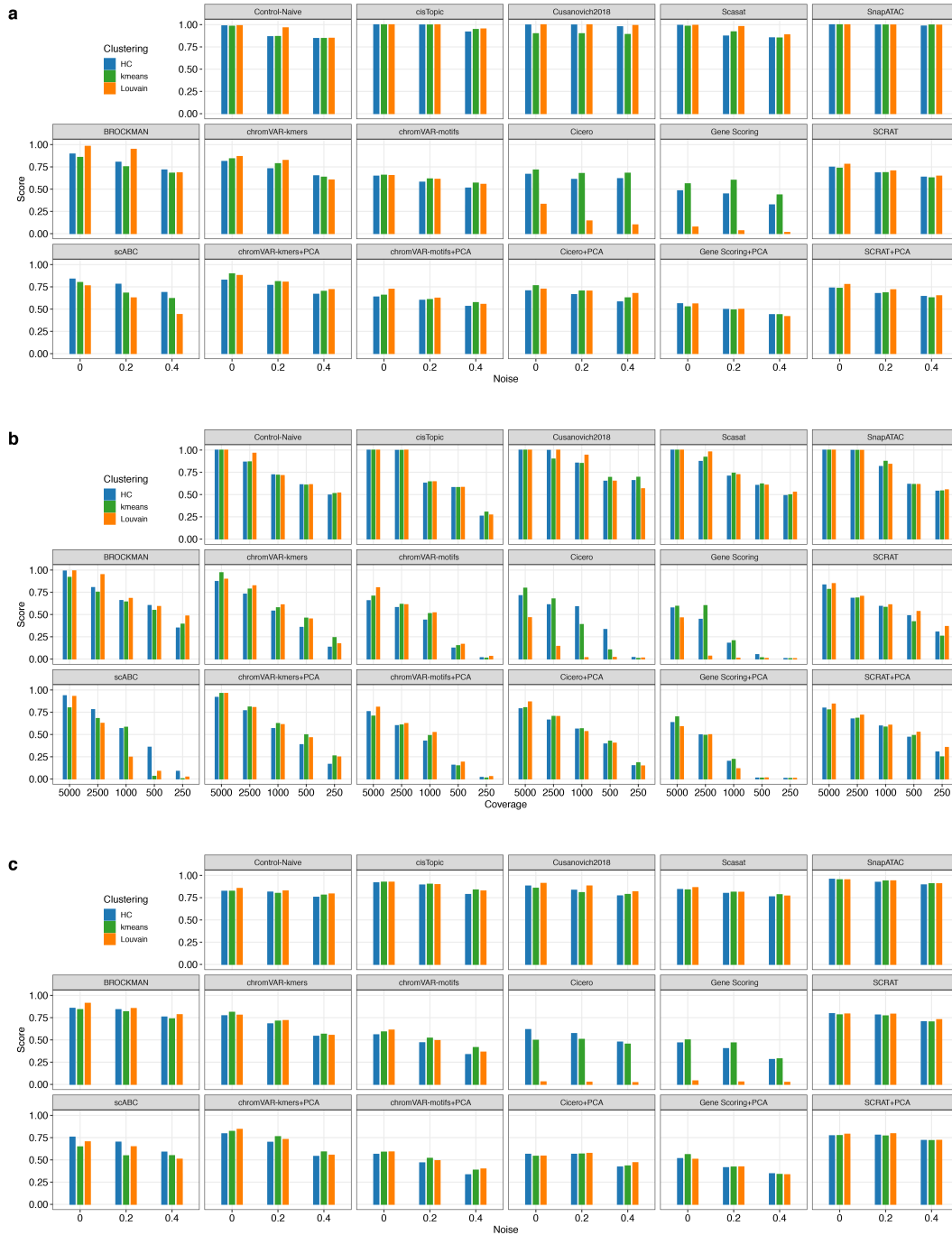
1703

1704



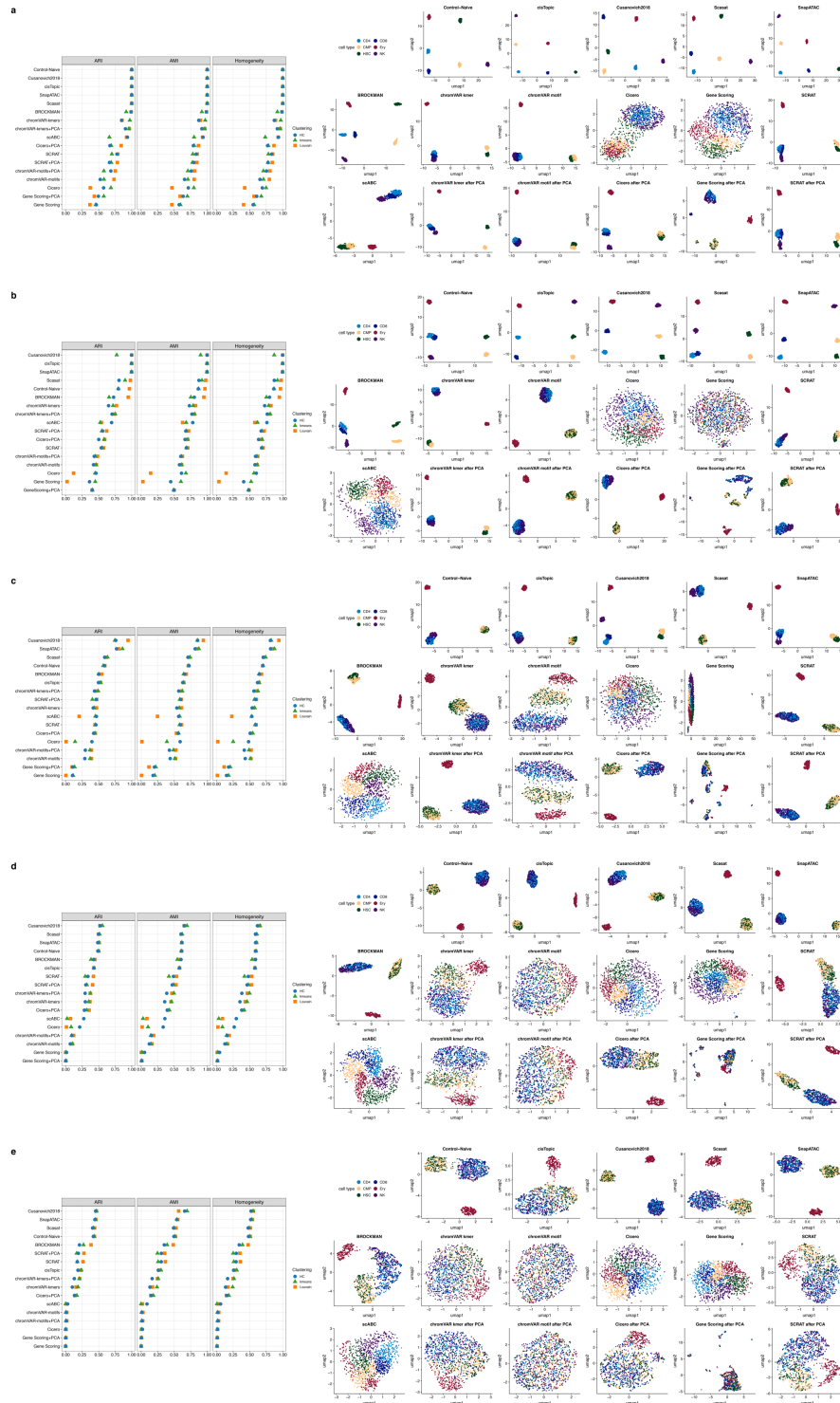
1705
1706
1707
1708
1709
1710

Figure S3. UMAP visualization of cells colored by clustering solution on the simulated bone marrow dataset with a noise level of 0.4 and a coverage of 2,500 fragments using **(a)** Louvain algorithm, **(b)** k-means clustering, and **(c)** hierarchical clustering (HC).

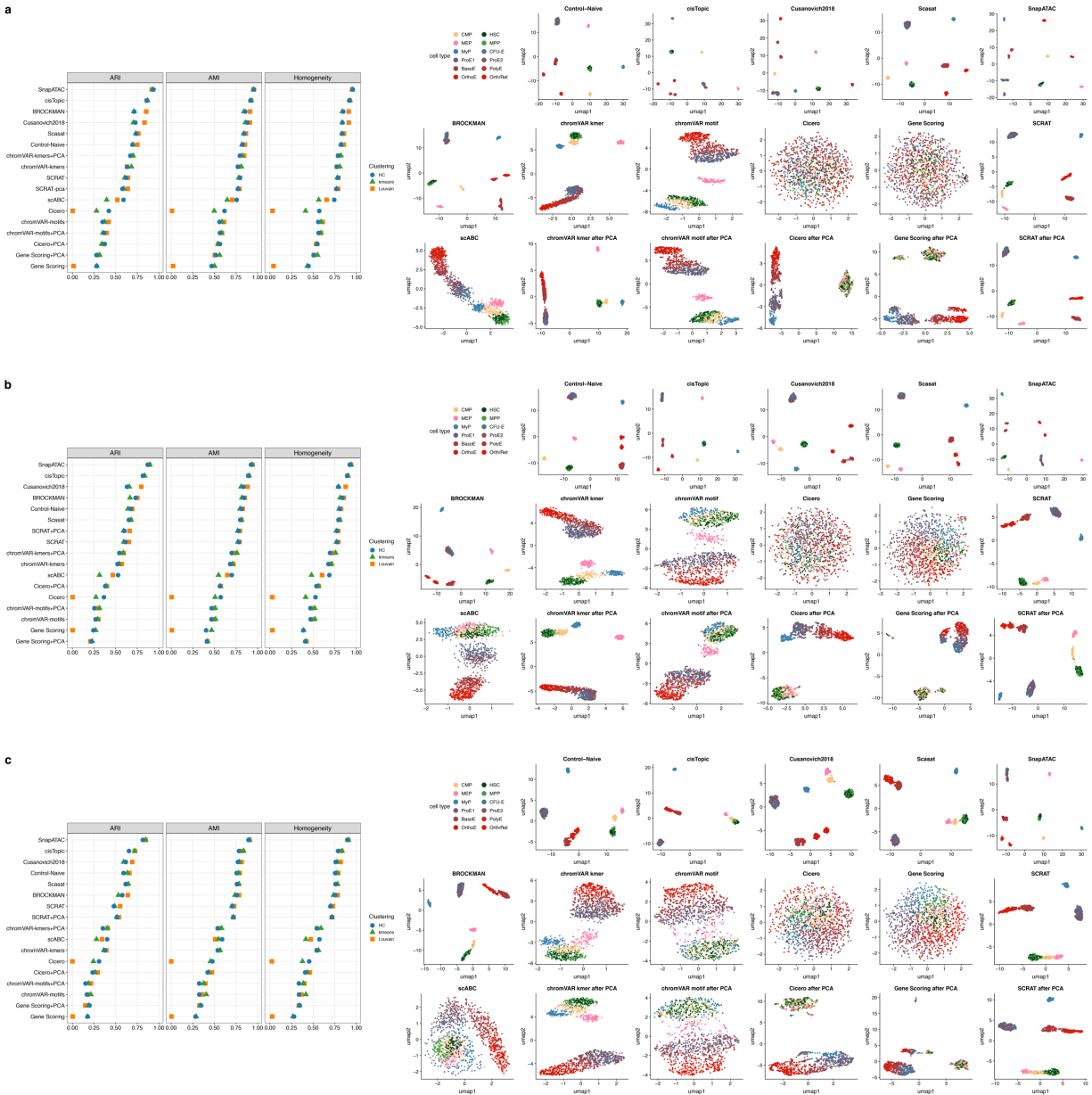


1711
1712
1713
1714
1715
1716
1717
1718

Figure S4. Summary of clustering scores at different noise levels and coverages based on three different clustering methods including hierarchical clustering (HC), k-means clustering and the Louvain algorithm. **(a)** clustering scores at noise levels of 0, 0.2, and 0.4 for the simulated bone marrow dataset with a coverage of 2,500. **(b)** clustering scores at the coverages of 5000, 2500, 1000, 500, 250 in the simulated bone marrow dataset at the noise level of 0.2. **(c)** clustering scores at the noise levels of 0, 0.2, and 0.4 for the simulated erythropoiesis dataset with a coverage of 2,500.



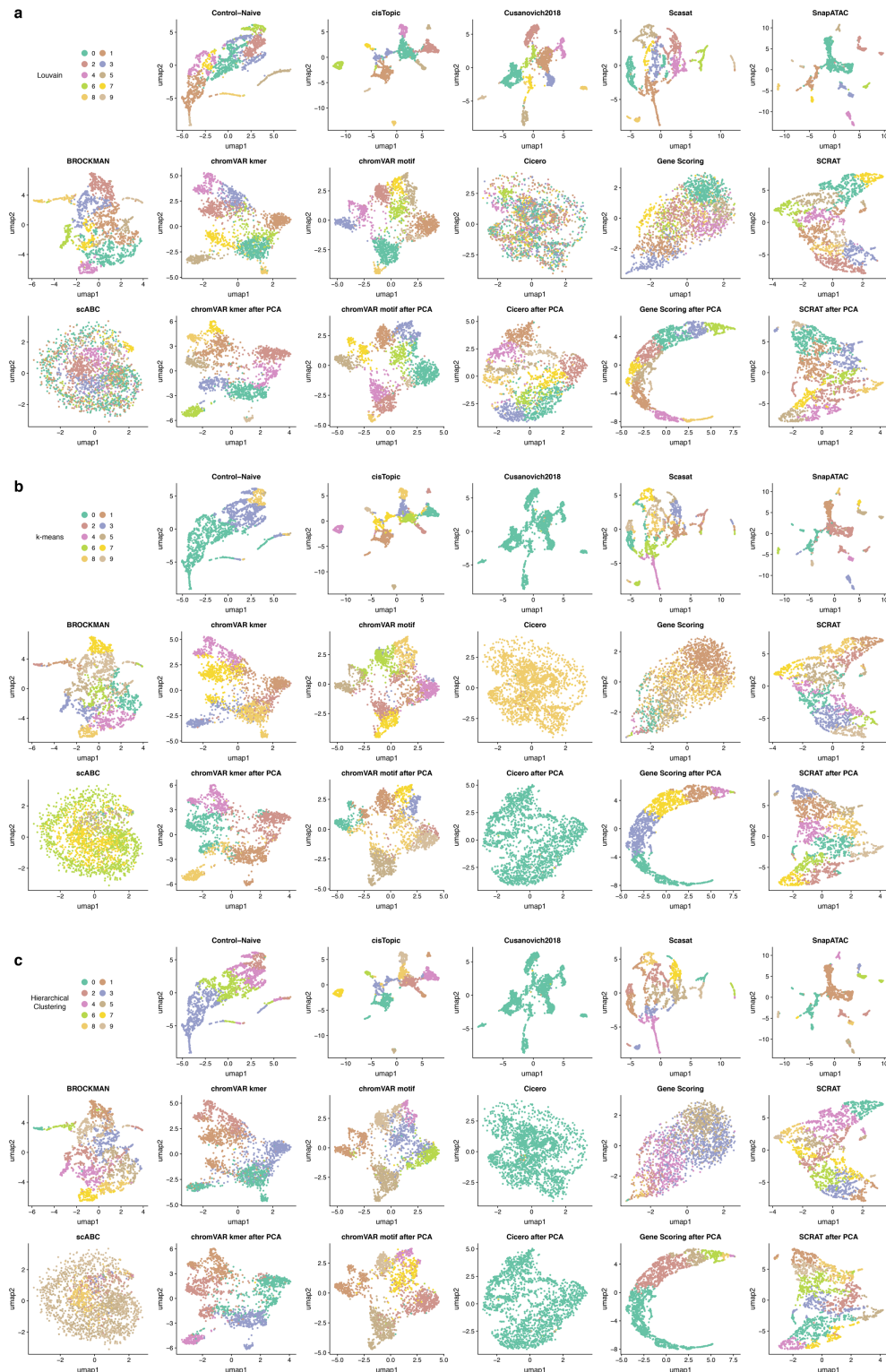
1719
 1720 **Figure S5.** Clustering evaluation according to AMI, ARI and Homogeneity metrics (*left*) and
 1721 UMAP visualization of cells colored by known cell labels (*right*) for the simulated bone marrow
 1722 dataset with a noise level of 0.2 and varying coverages: **(a)** 5000 reads, **(b)** 2500 reads, **(c)** 1000
 1723 reads, **(d)** 500 reads, and **(e)** 250 reads.
 1724



1725
1726

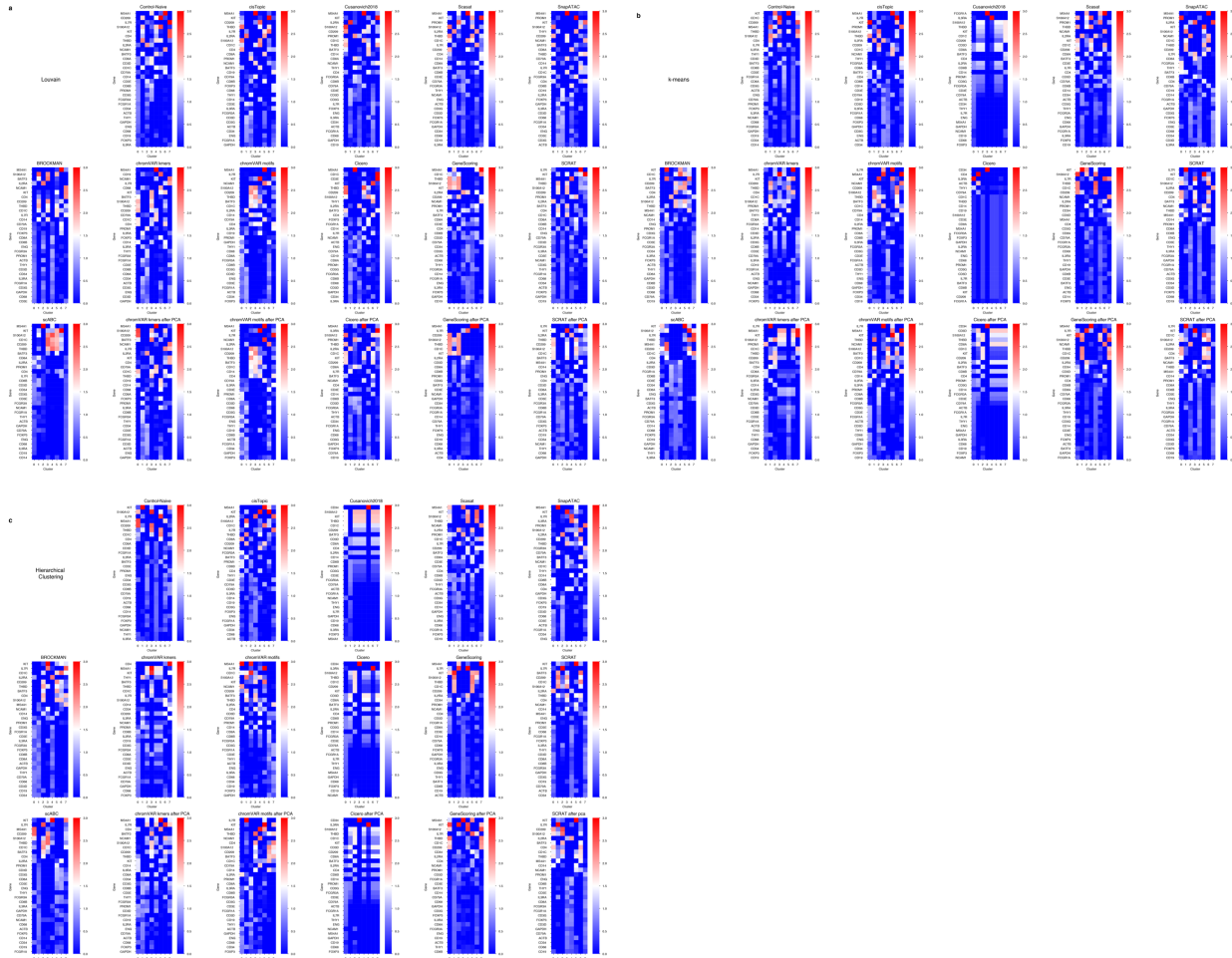
1727 **Figure S6.** Clustering evaluation according to AMI, ARI and Homogeneity metrics (*left*) and
 1728 UMAP visualization of cells colored by known cell labels (*right*) for the simulated erythropoiesis
 1729 datasets with a coverage of 2,500 fragments and **(a)** no noise (0), **(b)** moderate noise (0.2) or **(c)**
 1730 high noise (0.4).

1731
1732
1733
1734
1735
1736



1737
1738

1739 **Figure S7.** UMAP visualization of cells colored by the clustering solution on the *Buenrostro2018*
1740 dataset using **(a)** the Louvain algorithm, **(b)** k-means clustering, and **(c)** hierarchical clustering
1741 (HC).



1742

1743

1744 **Figure S8.** Heatmap for the average accessibility across clusters (columns) and the marker genes
1745 (rows) that are used to calculate the RAGI metric on the 10X PBMCs dataset. **(a)** Louvain
1746 clustering solution **(b)** k-means clustering solution **(c)** hierarchical clustering (HC) clustering
1747 solution.

1748

1749

1750

1751

1752

1753

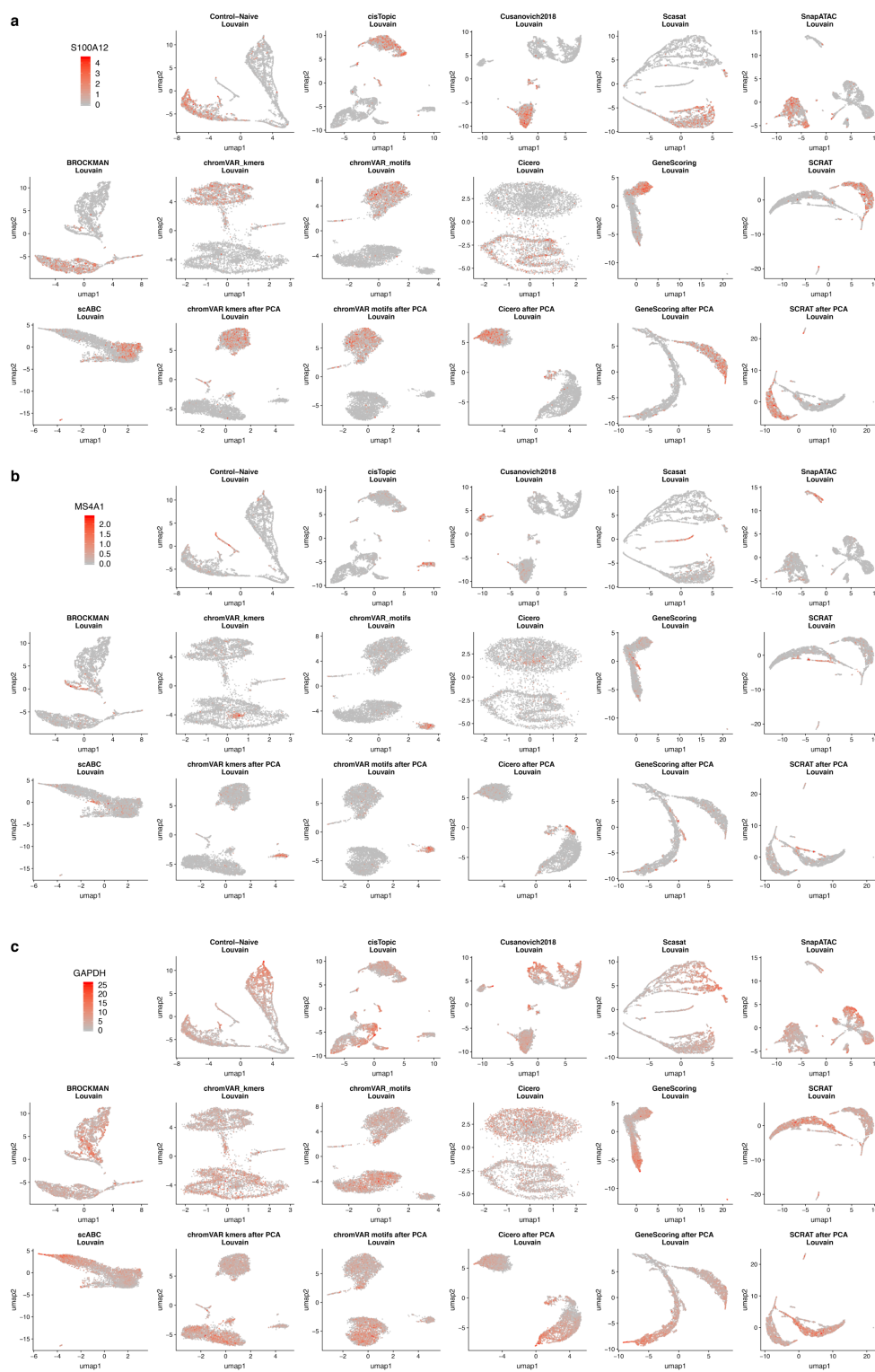
1754

1755

1756

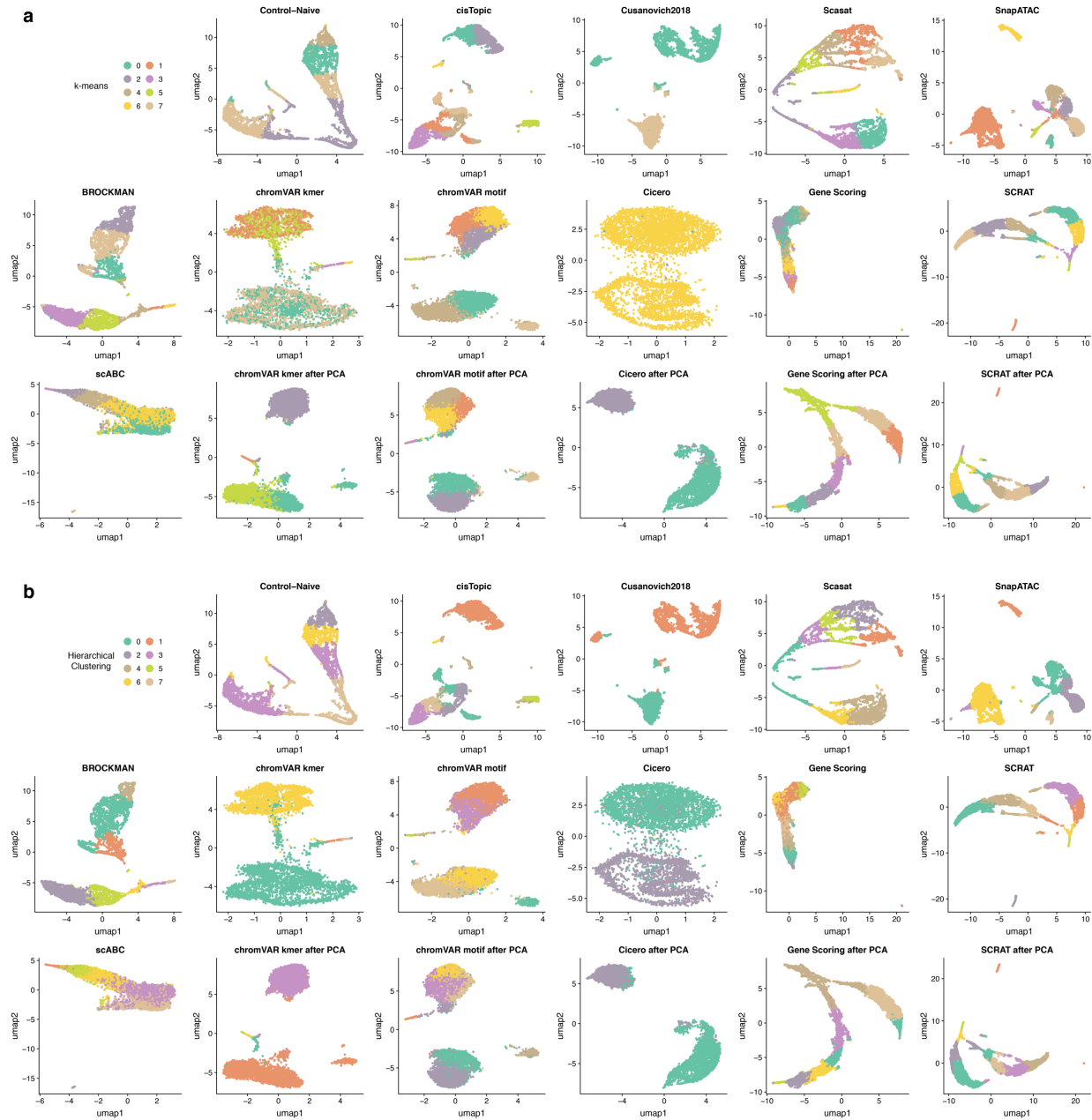
1757

1758



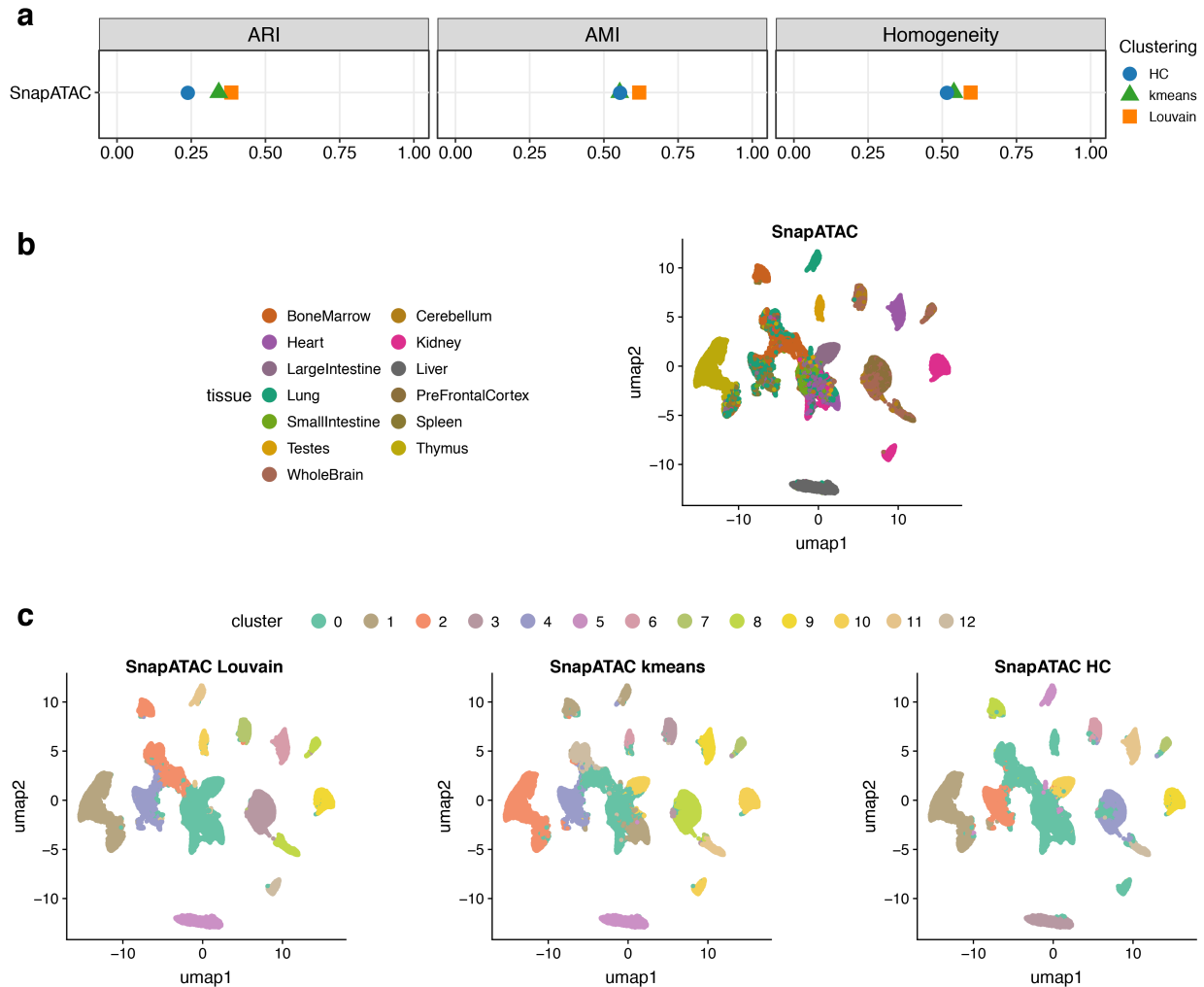
1759
1760
1761
1762
1763

Figure S9. UMAP visualization of cells colored by the accessibility of marker genes: **(a)** S100A12 and **(b)** MS4A1 and **(c)** GAPDH (housekeeping gene) and on the 10X PBMCs dataset.



1764
1765
1766
1767
1768
1769
1770
1771
1772
1773

Figure S10. UMAP visualization of cells colored by the clustering solution on 10X PBMCs dataset using **(a)** k-means clustering and **(b)** hierarchical clustering (HC).



1774

1775

1776 **Figure S11.** Assessment of SnapATAC on the full sci-ATAC-seq mouse dataset. **(a)** Clustering

1777 scores according to AMI, ARI and Homogeneity metrics **(b)** UMAP visualization of cells colored

1778 by the known tissues. **(c)** UMAP visualization of cells colored by three clustering solutions: the

1779 Louvain algorithm, k-means clustering, and hierarchical clustering (HC).

1780

1781

1782

1783

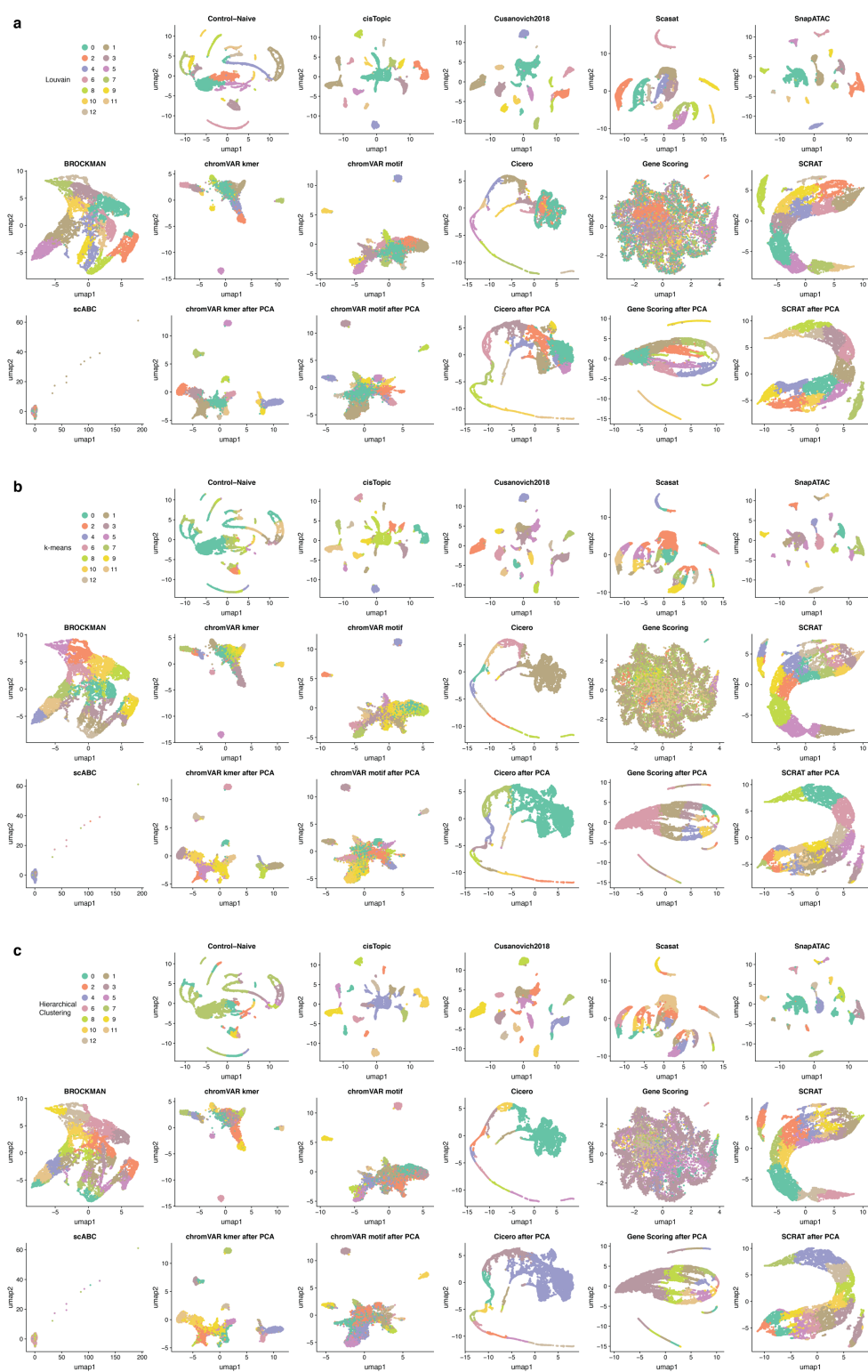
1784

1785

1786

1787

1788



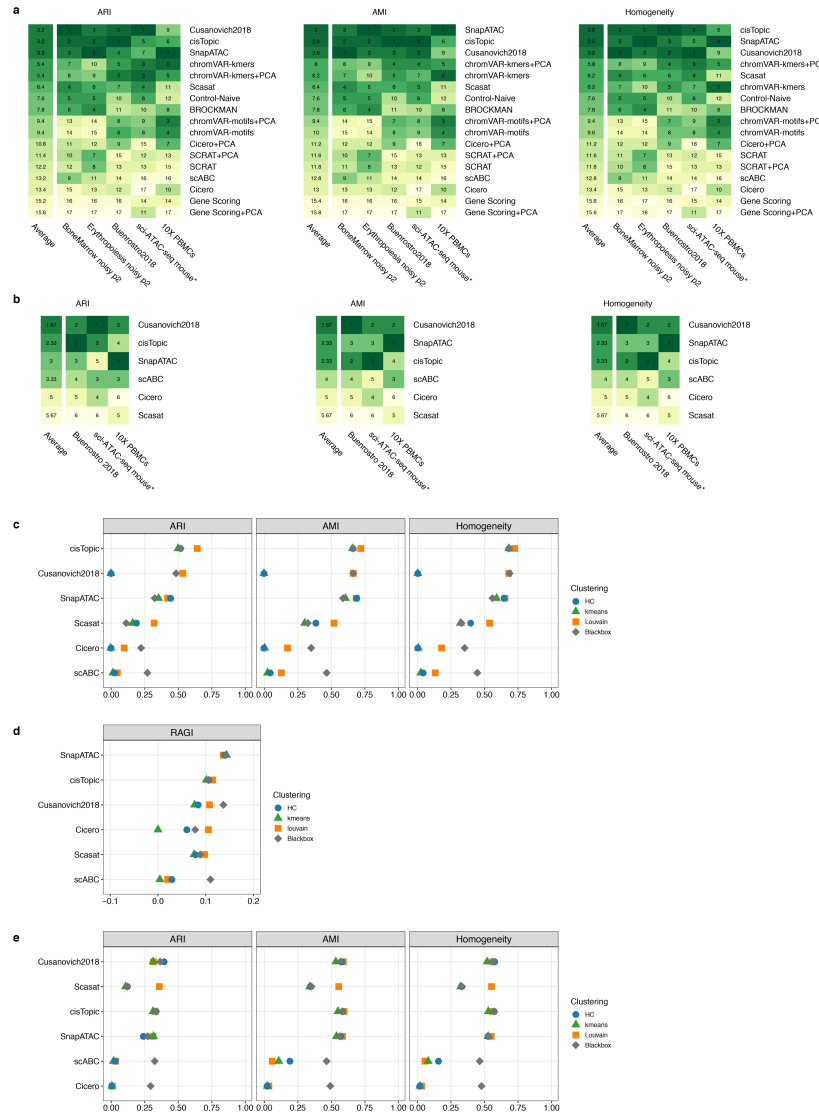
1789

1790

1791 **Figure S12.** UMAP visualization of cells colored by the clustering solution on the downsampled
1792 sci-ATAC-seq mouse dataset using (a) the Louvain algorithm, (b) k-means clustering, and (c)

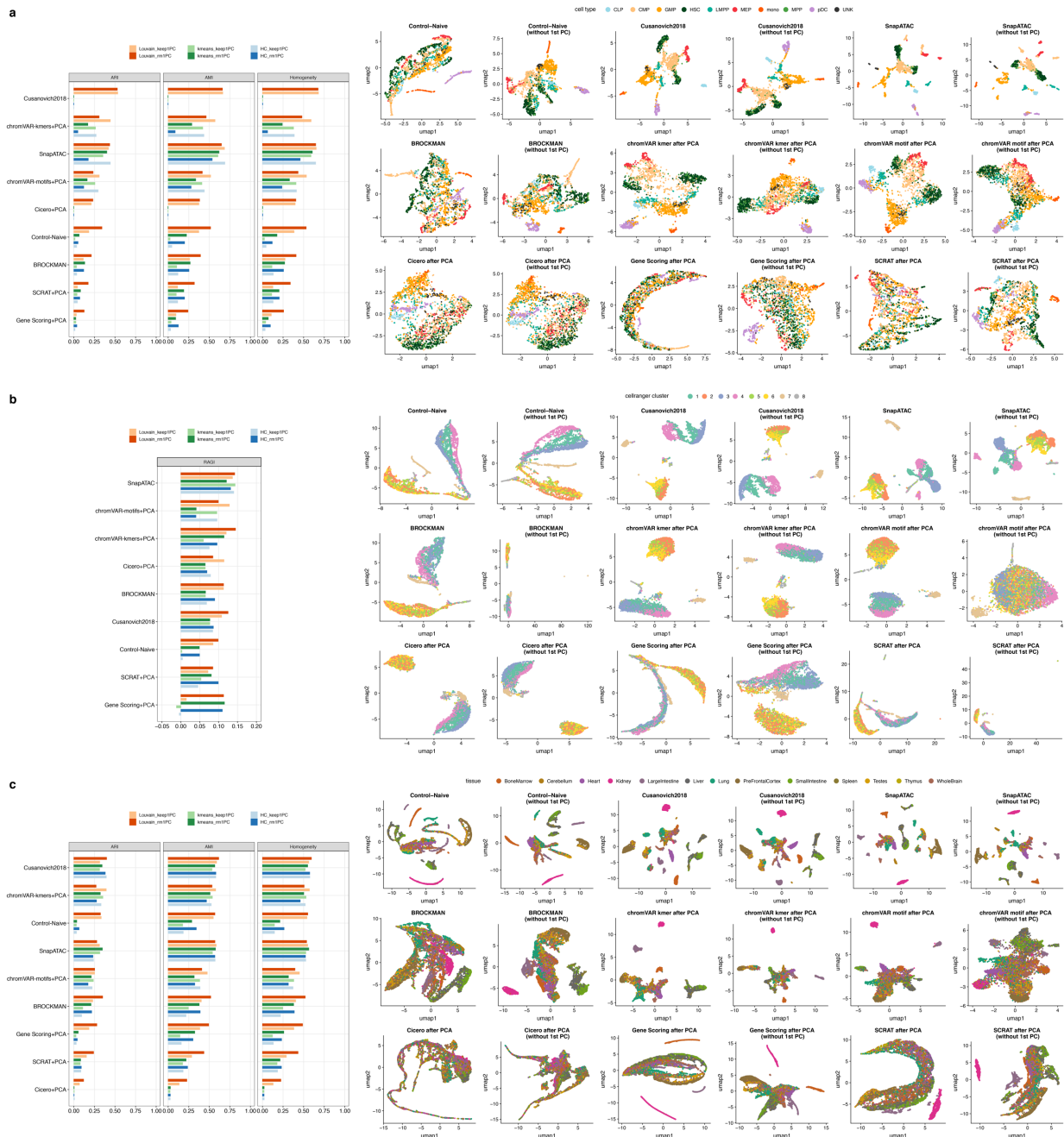
1793

hierarchical clustering (HC).



1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808

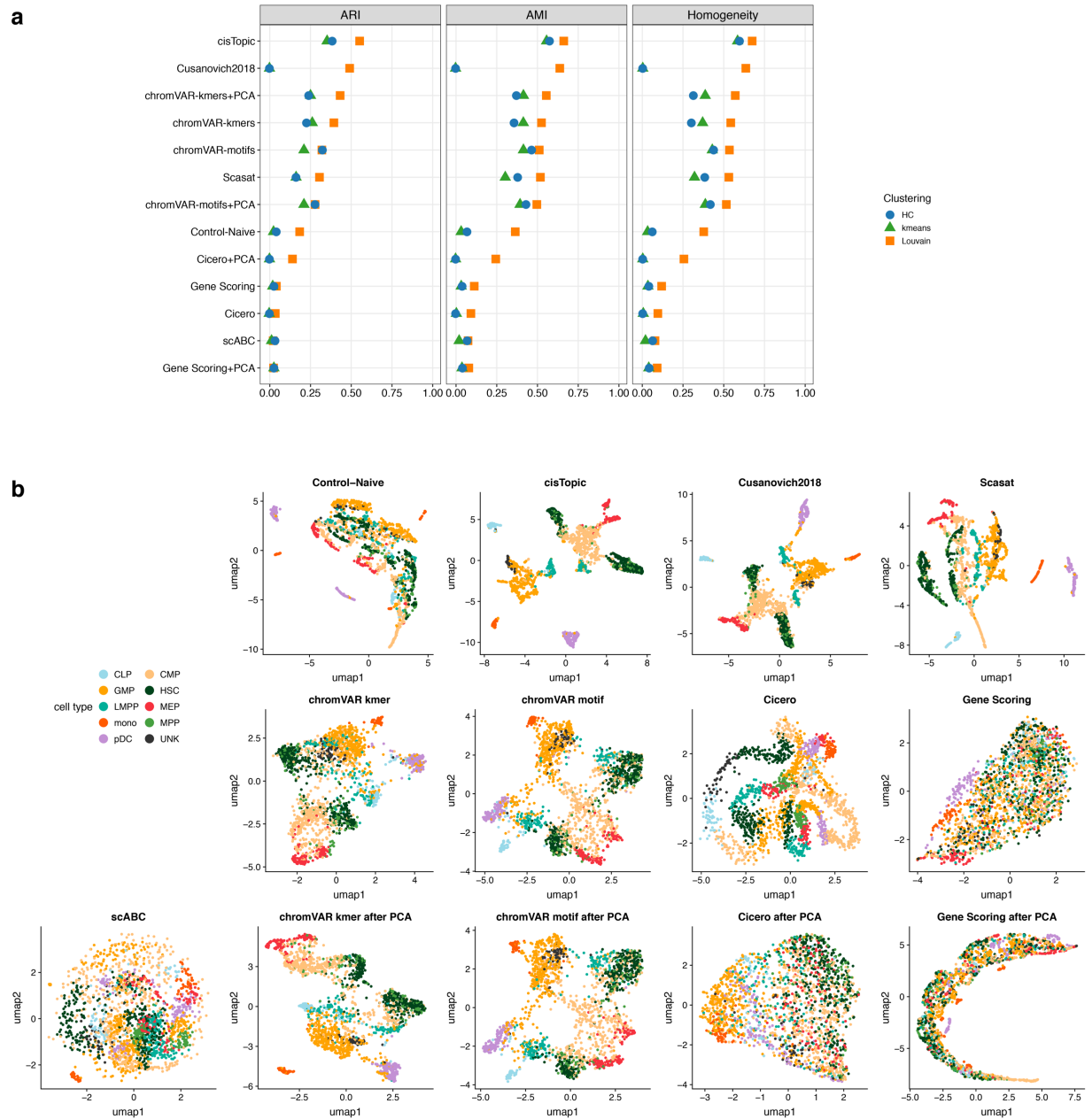
Figure S13. Ranking of method performance. **(a)** Rank was based on the best-performing clustering method for each metric on all methods and datasets. The column on the left shows the averaged rank per method across all datasets. * indicates a downsampled dataset of the indicated original dataset. **(b)** Rank of each method based on the best-performing clustering approach for each metric on methods assessed with an end-to-end clustering pipeline (termed as a ‘blackbox’) applied to the *Buenrostro2018*, downsampled sci-ATAC-seq mouse and 10X PBMCs datasets. The column on the left shows the averaged rank per method over these three datasets. * indicates a downsampled dataset of the indicated original dataset. **(c)** Dot plot of clustering scores for each metric applied to the *Buenrostro2018* dataset, including the ‘blackbox’ approach. **(d)** Dot plot of clustering scores for each metric applied to the 10X PBMCs dataset, including the ‘blackbox’ approach. **(e)** Dot plot of scores for each metric applied to the downsampled sci-ATAC-seq mouse dataset, including the ‘blackbox’ approach.



1809
1810

1811 **Figure S14.** Comparison between keeping the first PC and removing the first PC. **Left:** Clustering
1812 scores when the first PC is kept and for removal of the first PC, for each metric. **Right:** UMAP
1813 visualization of cells colored by known cell labels. The analyses are performed on **(a)** the
1814 *Buenrostro2018* dataset. **(b)** the 10X PBMCs dataset. **(c)** the downsampled sci-ATAC-seq mouse
1815 dataset.

1816
1817



1818
1819

1820 **Figure S15.** Assessment of methods using the peaks called from bulk ATAC-seq on the
1821 *Buenrostro2018* dataset. Only the methods that rely on peaks are included. **(a)** Clustering
1822 evaluation according to AMI, ARI and Homogeneity metrics **(b)** UMAP visualization of cells
1823 colored by known cell labels.

1824

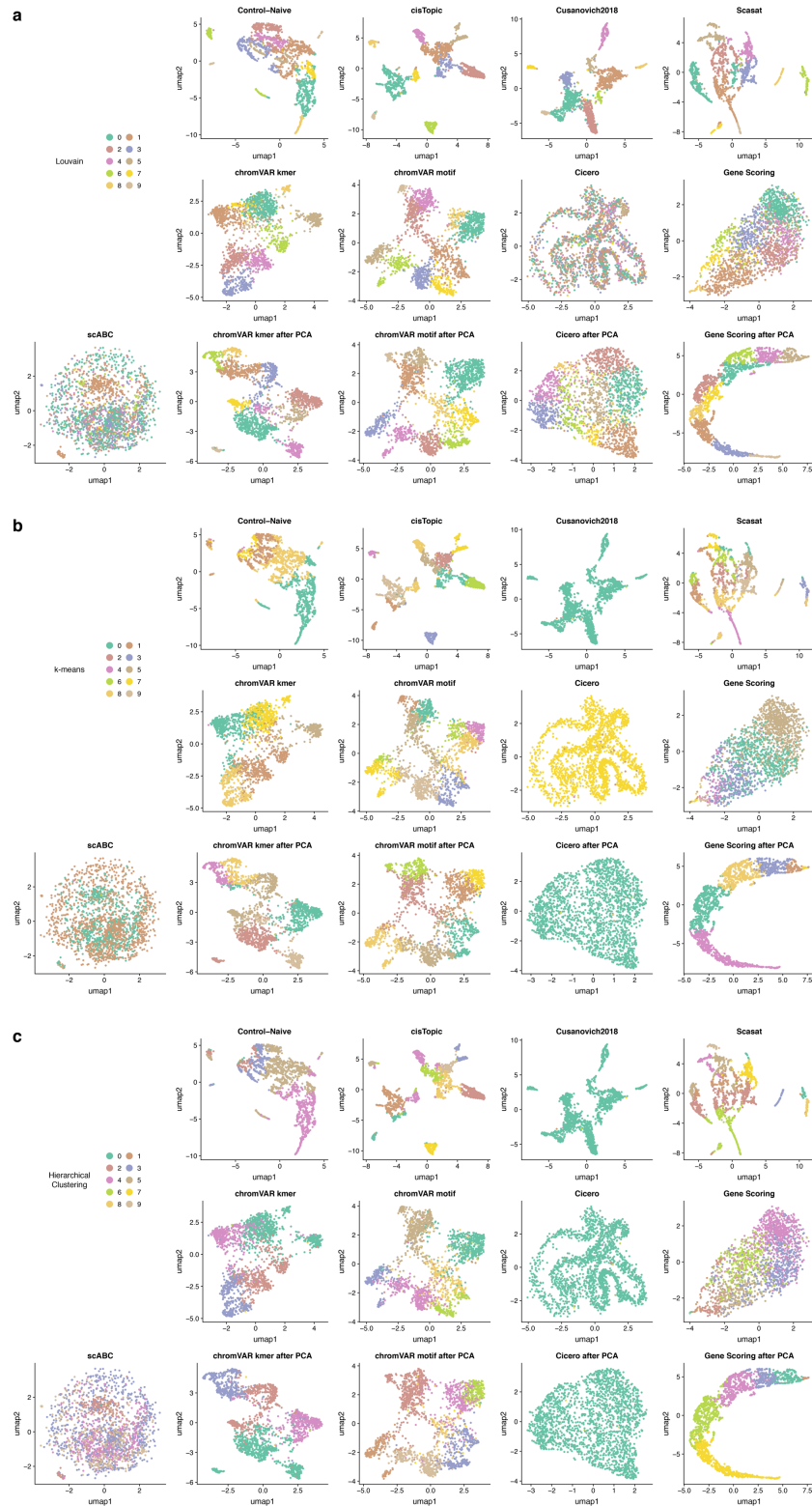
1825

1826

1827

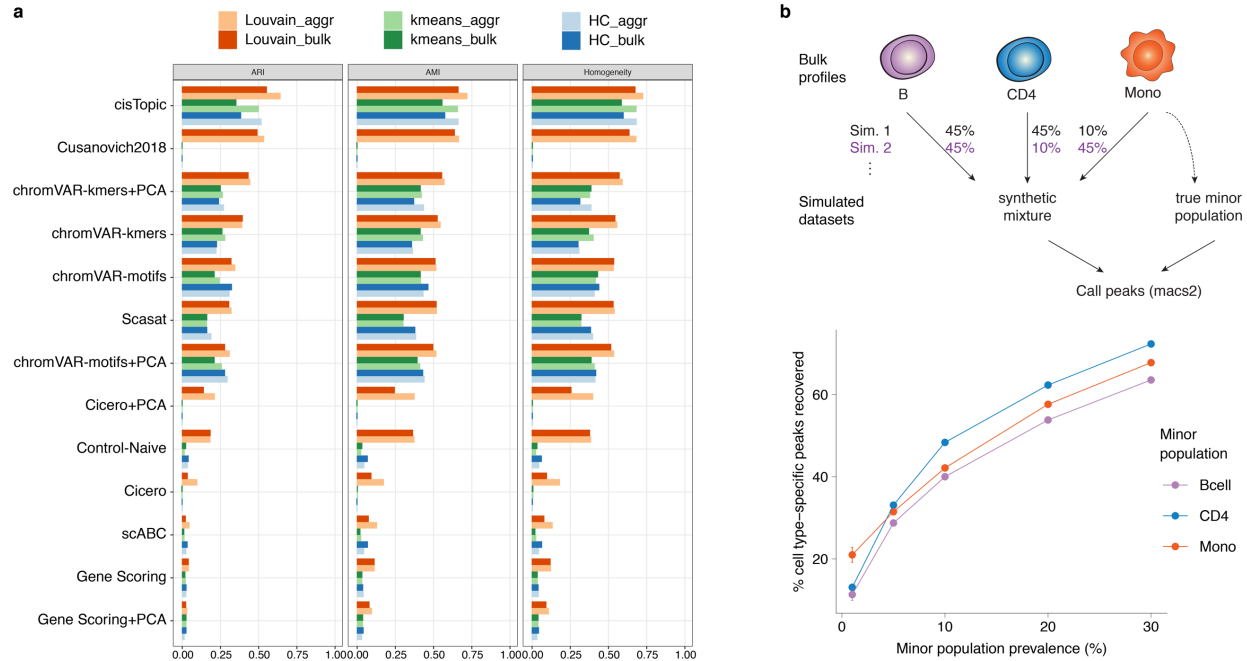
1828

1829



1830
1831
1832
1833

Figure S16. UMAP visualization of cells colored by the clustering solution on the *Buenrostro2018* dataset paired with bulk peaks using **(a)** k-means clustering and **(b)** hierarchical clustering (HC).



1834

1835

1836

1837

Figure S17. (a) Comparison of clustering scores between bulk ATAC-seq peaks and aggregated

1838

scATAC-seq peaks for each metric on the *Buenrostro2018* dataset. **(b) Top:** Simulation

1839

procedure from bulk ATAC-seq data. The three cell types (B-cells, CD4+ T-cells, and monocytes)

1840

are mixed in various proportions for each synthetic mixture. **Bottom:** The results of simulation

1841

in (b) **Top:** The x-axis reflects the proportion of the minor population. The y-axis reflects the

1842

percentage of recovered cell-type-specific peaks after performing peak calling on each mixture

1843

of single cells.

1844

1845

1846

1847

1848

1849

1850

1851

1852

1853

1854

1855

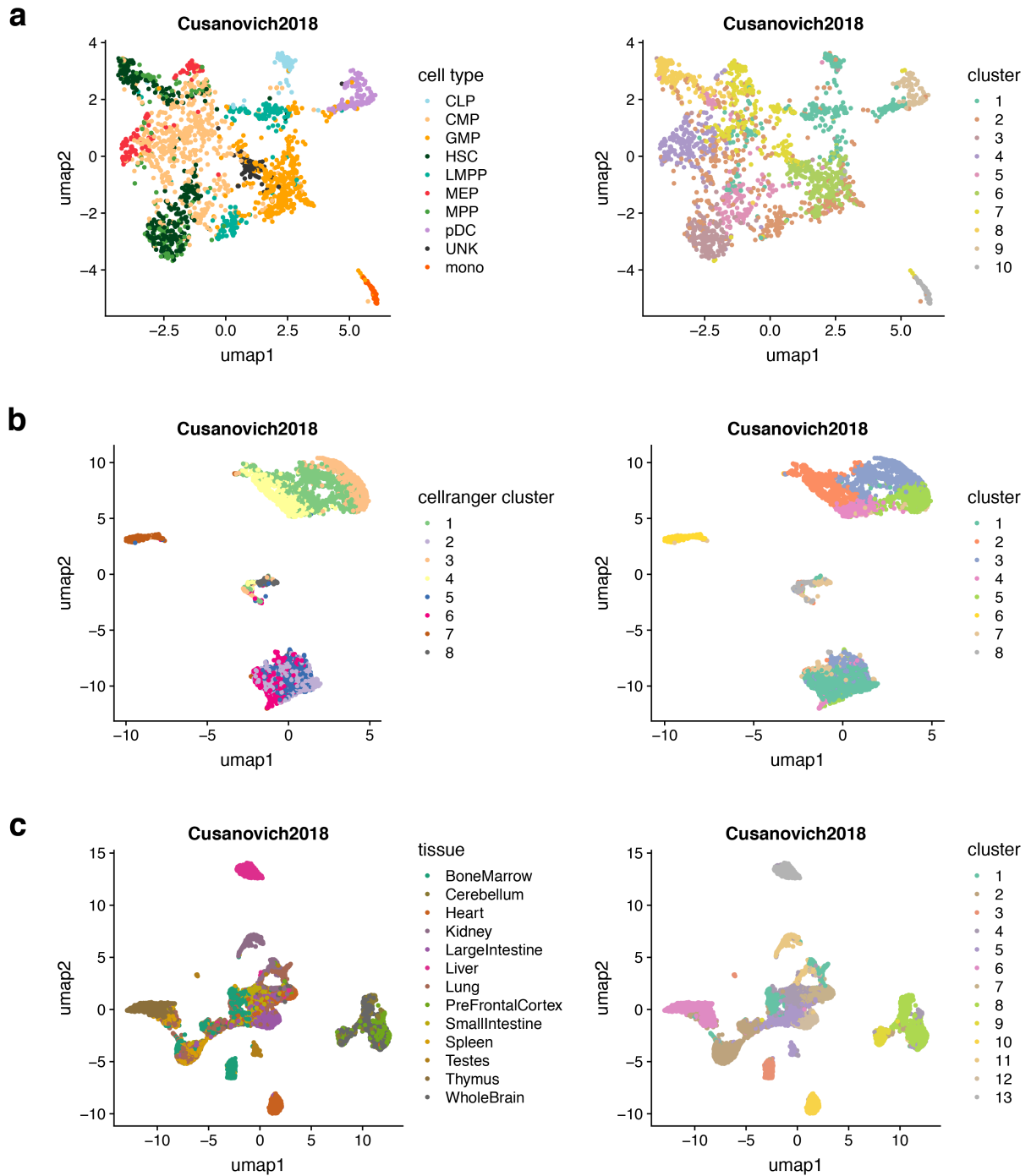
1856

1857

1858

1859

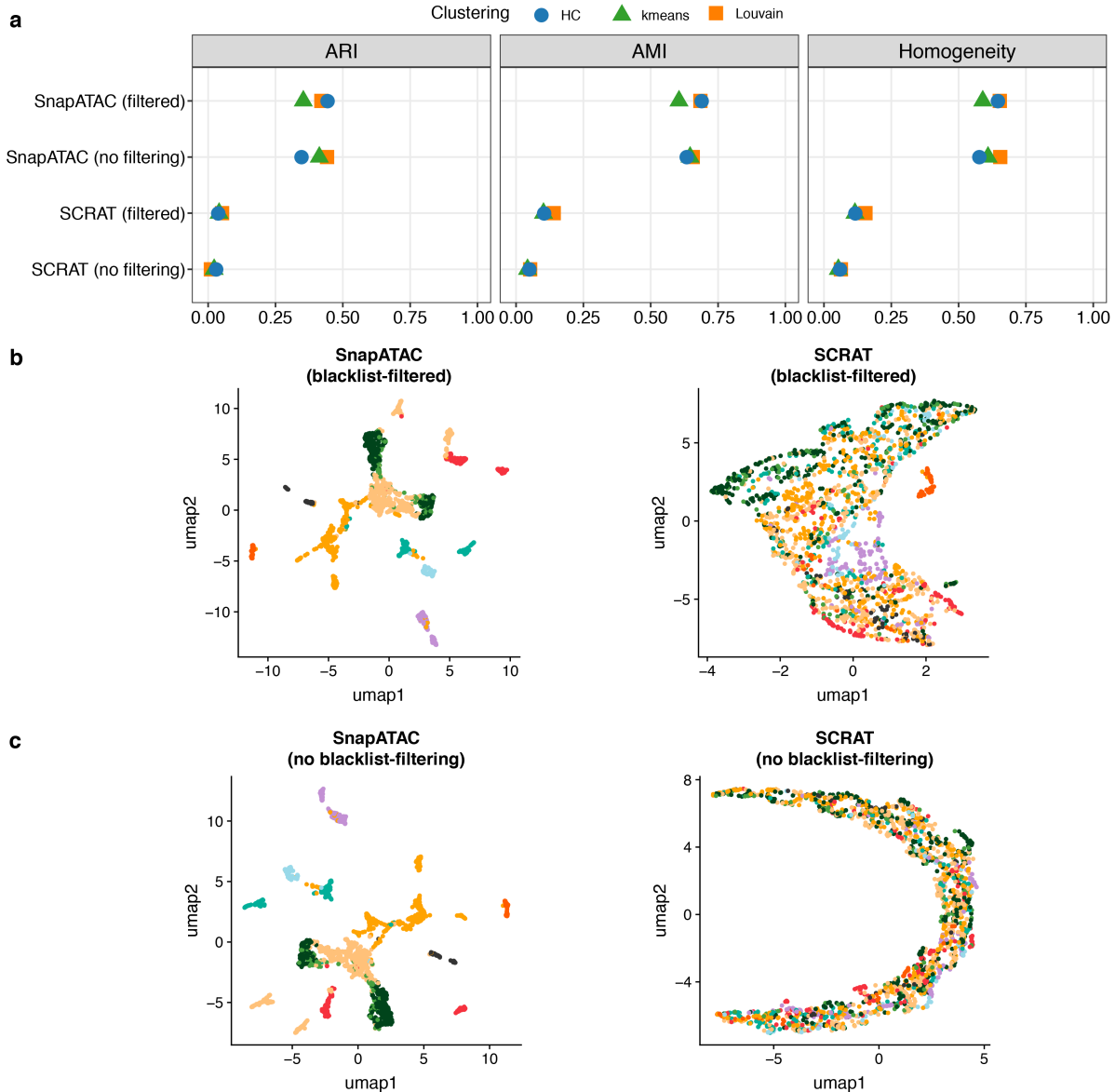
1860



1861

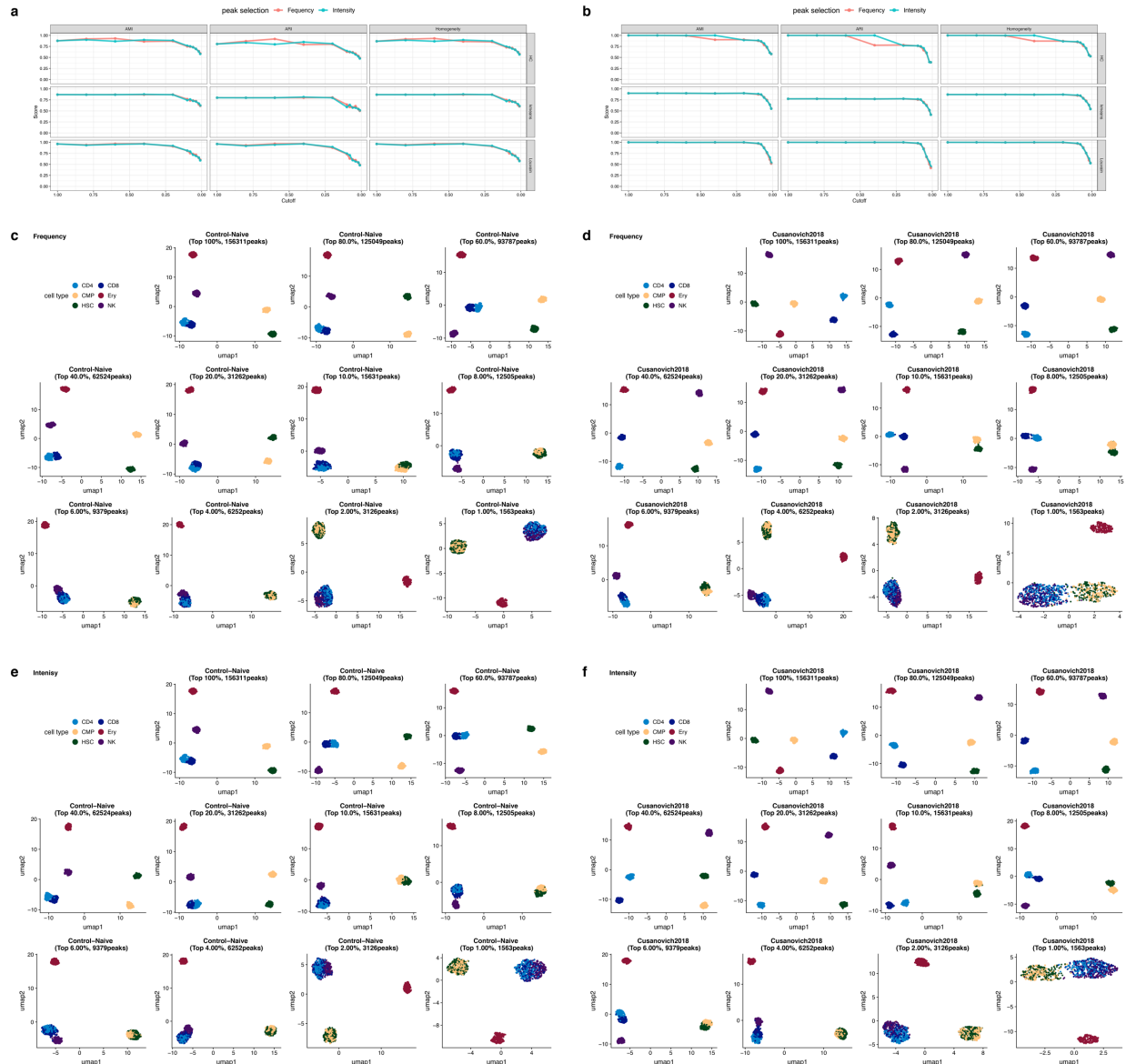
1862 **Figure S18.** Comparison between the known populations and the identified clades (pseudo-
1863 bulk) using *Cusanovich2018*. **Left:** UMAP visualization of cells colored by the known labels.
1864 **Right:** UMAP visualization of cells colored by the identified clades using *Cusanovich2018*. The
1865 analyses are performed on **(a)** the *Buenrostro2018* dataset. **(b)** the 10X PBMCs dataset. **(c)** the
1866 downsampled sci-ATAC-seq mouse dataset.

1867



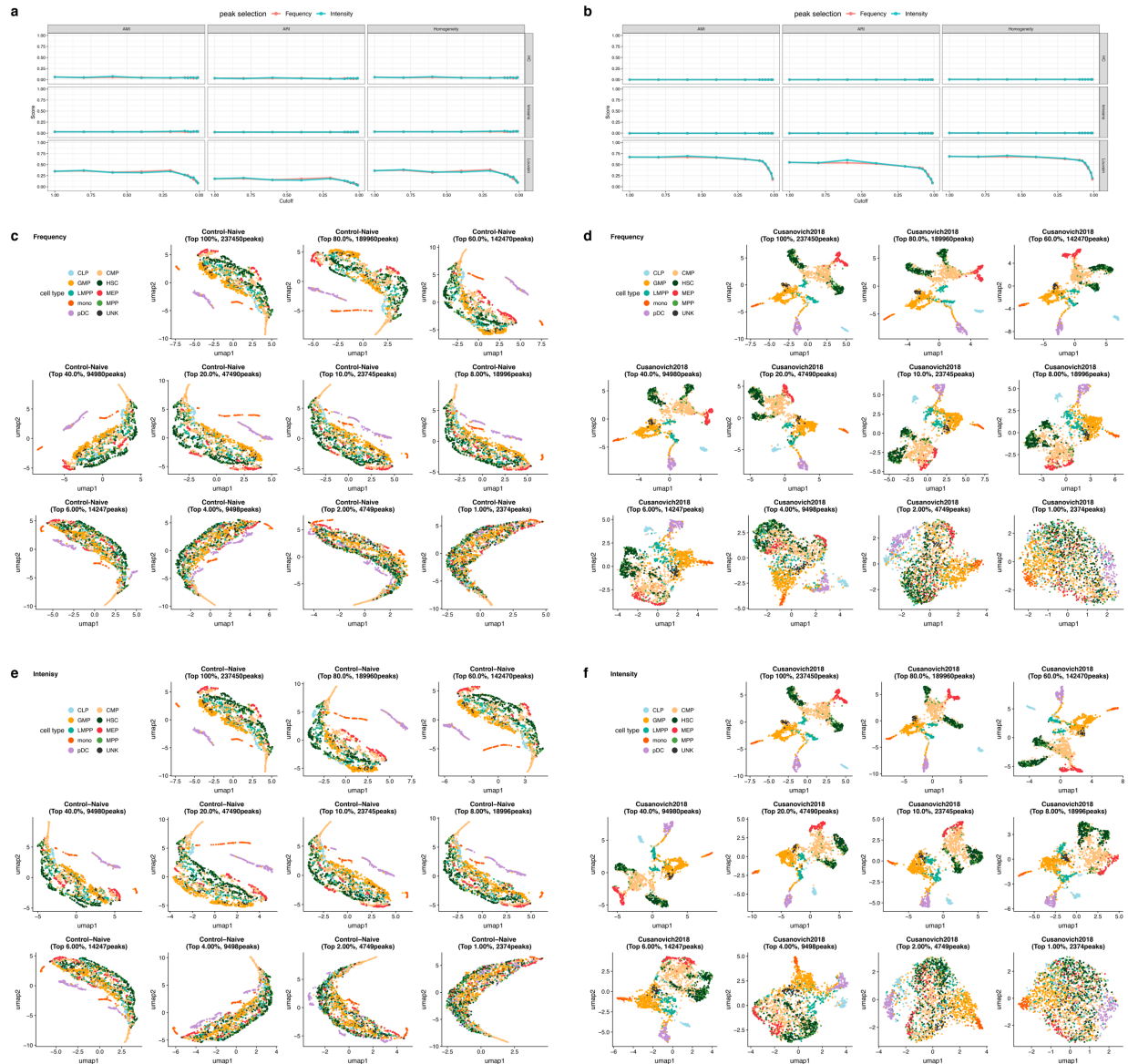
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881

Figure S19. Assessment of the effect of ENCODE blacklisted regions on the benchmarking results in the *Buenrostro2018* dataset. **(a)** Comparison of clustering scores between filtering or not filtering the blacklisted regions **(b)** UMAP visualization based on SnapATAC (*left*) and SCRAT (*right*) feature matrices after filtering the ENCODE blacklisted regions. Cell are colored by the FACS-sorting labels. **(c)** UMAP visualization based on SnapATAC (*left*) and SCRAT (*right*) feature matrices without filtering the ENCODE blacklisted regions. Cell are colored by the FACS-sorting labels.



1882
1883

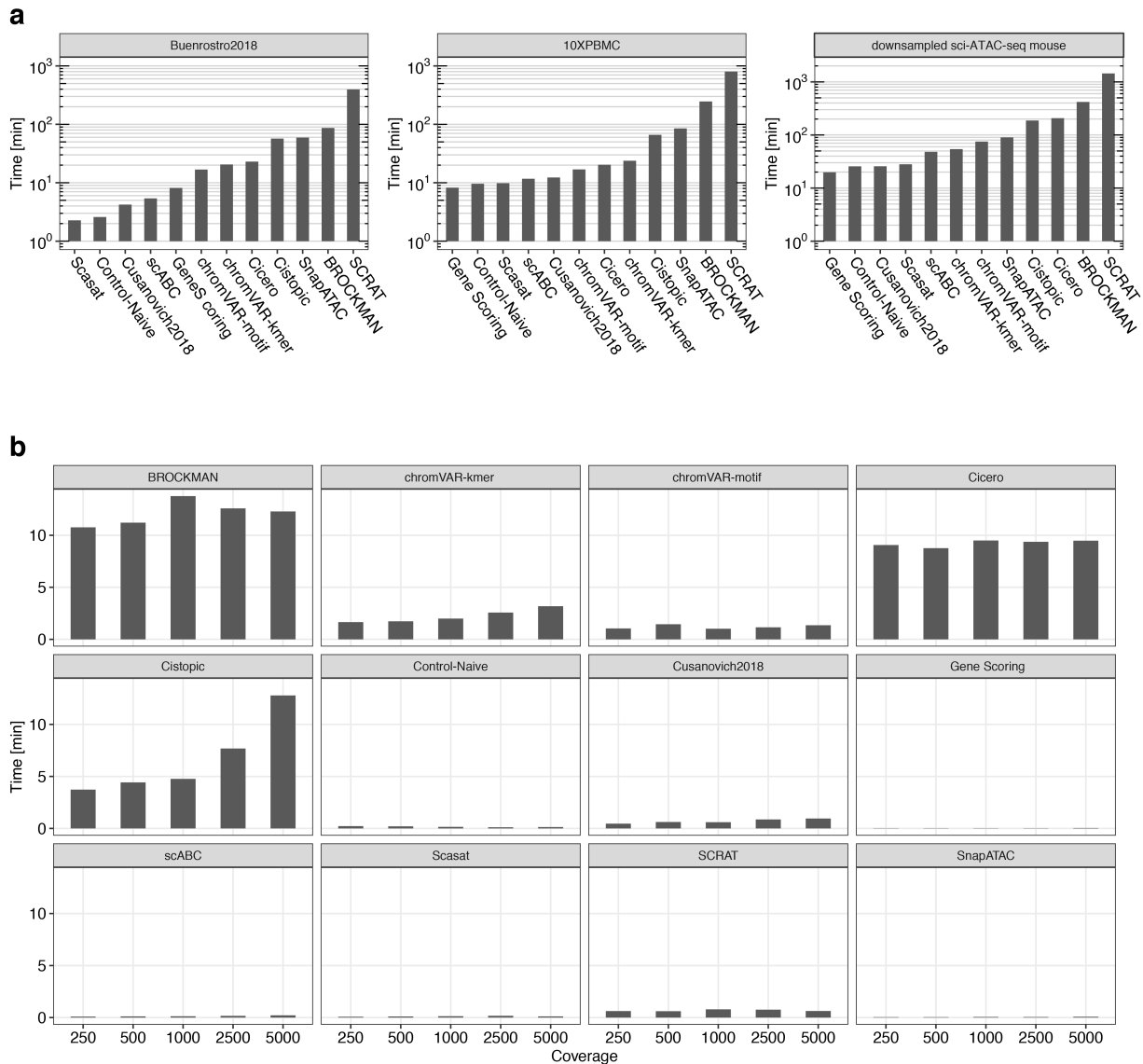
1884 **Figure S20.** Comparison between frequency-based and intensity-based peak selection for each
1885 metric on the simulated bone marrow dataset with a noise level of 0.2 with a coverage of 2,500
1886 fragments. **(a)** Clustering scores for each metric and clustering method across different cutoffs
1887 for the Control-Naïve method. **(b)** Clustering scores for each metric and clustering method
1888 across different cutoffs for the *Cusanovich2018* method. **(c)** UMAP visualization of cells colored
1889 by the known labels using a frequency-based peak selection for Control-naïve method. **(d)**
1890 UMAP visualization of cells colored by the known labels using a frequency-based peak selection
1891 for the *Cusanovich2018* method. **(e)** UMAP visualization of cells colored by the known labels
1892 using an intensity-based peak selection for Control-naïve method. **(f)** UMAP visualization of
1893 cells colored by the known labels using an intensity-based peak selection for the
1894 *Cusanovich2018* method.



1895
1896

1897 **Figure S21.** Comparison between frequency-based and intensity-based peak selection for each
1898 metric on the *Buenrostro2018* dataset. **(a)** Clustering scores for each metric and clustering
1899 method across different cutoffs for the Control-Naïve method. **(b)** Clustering scores for each
1900 metric and clustering method across different cutoffs for the *Cusanovich2018* method. **(c)**
1901 UMAP visualization of cells colored by FACS-sorting labels using a frequency-based peak
1902 selection for the Control-naïve method. **(d)** UMAP visualization of cells colored by FACS-sorting
1903 labels using a frequency-based peak selection for the *Cusanovich2018* method. **(e)** UMAP
1904 visualization of cells colored by FACS-sorting labels using an intensity-based peak selection for
1905 the Control-naïve method. **(f)** UMAP visualization of cells colored by the FACS-sorting labels
1906 using an intensity-based peak selection for the *Cusanovich2018* method.

1907



1908
1909

1910 **Figure S22.** Running time results. **(a)** Running time, in minutes for each method applied to the
1911 *Buenrostro2018*, 10X PBMCs, and downsampled sci-ATAC-seq mouse datasets. **(b)** Running
1912 time, in minutes for each method on the simulated bone marrow dataset at a noise level of 0.2
1913 with read coverages of 250, 500, 1000, 2500, and 5000 fragments.
1914