

1 **TITLE:**

2

3 **ArchR: An integrative and scalable software package for single-cell chromatin**  
4 **accessibility analysis**

5

6 **AUTHOR LIST AND AFFILIATIONS:**

7

8 Jeffrey M. Granja<sup>1,2,3,†,\*</sup>, M. Ryan Corces<sup>3,4,†</sup>, Sarah E. Pierce<sup>1,5</sup>, S. Tansu Bagdatli<sup>1</sup>, Hani  
9 Choudhry<sup>6</sup>, Howard Y. Chang<sup>1,3,5,7,\*</sup>, William J. Greenleaf<sup>1,3,8,9,\*</sup>

10

11 <sup>1</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA.

12 <sup>2</sup>Program in Biophysics, Stanford University, Stanford, CA, USA.

13 <sup>3</sup>Center for Personal Dynamic Regulomes, Stanford University, Stanford, CA, USA.

14 <sup>4</sup>Department of Pathology, Stanford University School of Medicine, Stanford, CA, USA.

15 <sup>5</sup>Program in Cancer Biology, Stanford University School of Medicine, Stanford, CA, USA.

16 <sup>6</sup>Department of Biochemistry, Faculty of Science, Cancer and Mutagenesis Unit, King Fahd Center  
17 for Medical Research, King Abdulaziz University, Jeddah, Saudi Arabia.

18 <sup>7</sup>Howard Hughes Medical Institute, Stanford University, Stanford, CA, USA.

19 <sup>8</sup>Department of Applied Physics, Stanford University, Stanford, CA, USA.

20 <sup>9</sup>Chan-Zuckerberg Biohub, San Francisco, CA, USA.

21 <sup>†</sup>These authors contributed equally to this work.

22

23 \*Correspondence should be addressed to J.M.G ([jgranja@stanford.edu](mailto:jgranja@stanford.edu)), H.Y.C.  
24 ([howchang@stanford.edu](mailto:howchang@stanford.edu)), or W.J.G ([wjg@stanford.edu](mailto:wjg@stanford.edu))

25

26 **Contact Information**

27 William J. Greenleaf, PhD

28 Stanford University School of Medicine

29 257A Beckman Center, 279 Campus Drive, Stanford, CA 94305-5301

30 Email: [wjg@stanford.edu](mailto:wjg@stanford.edu)

31

32 Howard Y. Chang, MD, PhD

33 Stanford University School of Medicine

34 CCSR 2155c, 269 Campus Drive, Stanford, CA 94305-5168

35 Email: [howchang@stanford.edu](mailto:howchang@stanford.edu)

36

37

38

39

40

## 41 **ABSTRACT**

42 The advent of large-scale single-cell chromatin accessibility profiling has accelerated our ability  
43 to map gene regulatory landscapes, but has outpaced the development of robust, scalable  
44 software to rapidly extract biological meaning from these data. Here we present a software suite  
45 for single-cell analysis of regulatory chromatin in R (ArchR; [www.ArchRProject.com](http://www.ArchRProject.com)) that enables  
46 fast and comprehensive analysis of single-cell chromatin accessibility data. ArchR provides an  
47 intuitive, user-focused interface for complex single-cell analyses including doublet removal,  
48 single-cell clustering and cell type identification, robust peak set generation, cellular trajectory  
49 identification, DNA element to gene linkage, transcription factor footprinting, mRNA expression  
50 level prediction from chromatin accessibility, and multi-omic integration with scRNA-seq. Enabling  
51 the analysis of over 1.2 million single cells within 8 hours on a standard Unix laptop, ArchR is a  
52 comprehensive analytical suite for end-to-end analysis of single-cell chromatin accessibility data  
53 that will accelerate the understanding of gene regulation at the resolution of individual cells.

54

## 55 **INTRODUCTION**

56 Single-cell approaches have revolutionized our understanding of biology, opening the door for a  
57 wide array of applications ranging from interrogation of cellular heterogeneity to identification of  
58 disease-specific processes. The advent of single-cell approaches for the assay for transposase-  
59 accessible chromatin using sequencing (scATAC-seq) has made it possible to study chromatin  
60 accessibility and gene regulation in single cells<sup>1,2</sup>. These chromatin-based assays have  
61 illuminated cell type-specific biology and provided insights into complex biological processes  
62 previously hidden by ensemble averaging<sup>3-7</sup>. Recent methodological advances have increased  
63 the throughput of scATAC-seq, enabling a single lab to generate data from hundreds of thousands  
64 of cells on the timescale of weeks<sup>5,6,8</sup>. These advances have been driven by an increased interest  
65 in chromatin-based gene regulation across a diversity of cellular contexts and biological  
66 systems<sup>1,2,5,6,8</sup>. This capacity for data generation has outpaced the development of intuitive,

67 robust, and comprehensive software for analysis of these scATAC-seq datasets<sup>9</sup> – a crucial  
68 requirement that would facilitate the broad utilization of these methods of investigating gene  
69 regulation at cellular resolution.

70 To this end, we sought to develop a user-oriented software suite for both routine and  
71 advanced analysis of massive-scale single-cell chromatin accessibility data from diverse sources  
72 without the need for high-performance computing environments. This package for single-cell  
73 Analysis of Regulatory Chromatin in R (ArchR; [www.ArchRProject.com](http://www.ArchRProject.com)) provides a facile platform  
74 to interrogate scATAC-seq data from multiple scATAC-seq implementations, including the 10x  
75 Genomics Chromium system<sup>6,7</sup>, the Bio-Rad droplet scATAC-seq system<sup>8</sup>, single-cell  
76 combinatorial indexing<sup>2,5</sup>, and the Fluidigm C1 system<sup>1,4</sup> (**Fig. 1a**). ArchR provides a user-focused  
77 interface for complex scATAC-seq analysis such as marker feature identification, transcription  
78 factor (TF) footprinting, an interactive genome browsing, scRNA-seq integration, and cellular  
79 trajectory analysis (**Fig. 1a**). When compared to other existing tools, such as SnapATAC<sup>10</sup> and  
80 Signac<sup>11</sup>, ArchR provides a more extensive set of features with substantially improved  
81 performance benchmarks (**Supplementary Fig. 1a**). Moreover, ArchR is designed to provide the  
82 speed and flexibility to support interactive analysis, enabling iterative extraction of meaningful  
83 biological interpretations.

84

## 85 **RESULTS**

86

### 87 **The ArchR framework**

88 ArchR takes as input aligned BAM or fragment files, which are first parsed in small chunks per  
89 chromosome, read in parallel to conserve memory, then efficiently stored on disk using the  
90 compressed random-access hierarchical data format version 5 (HDF5) file format. These HDF5  
91 files form the constituent pieces of an ArchR analysis which we call “Arrow” files. Arrow files are  
92 grouped into an “ArchR Project”, a compressed R data file that is stored in memory, which

93 provides an organized, rapid, and low memory-use framework for manipulation of the larger arrow  
94 files stored on disk (**Supplementary Fig. 1b**). Arrow files are always accessed in chunks using  
95 parallel read and write operations that minimize memory while efficiently using the multi-processor  
96 capabilities of most standard computers (**Supplementary Fig. 1c-d**). Moreover, the base file size  
97 of Arrow files remains smaller than the input fragment files across various cellular inputs  
98 (**Supplementary Fig. 2a-b**). These efficiencies provide substantial improvements in speed and  
99 memory usage compared to scATAC-seq software packages such as SnapATAC and Signac.

### 100 **ArchR enables efficient and comprehensive single-cell chromatin accessibility analysis**

101 To benchmark the performance of ArchR, we collected three diverse publicly available datasets  
102 (**Supplementary Table 1**): (i) peripheral blood mononuclear cells (PBMCs) that represent  
103 discrete primary cell types<sup>6,7</sup> (**Supplementary Fig. 2c-e**), (ii) bone marrow stem/progenitor cells  
104 and differentiated cells that represent a continuous cellular hierarchy<sup>7</sup> (**Supplementary Fig. 2f-**  
105 **h**), and (iii) a large atlas of murine cell types from diverse organ systems<sup>5</sup> (**Supplementary Fig.**  
106 **2i-k**). Prior to downstream analysis, we performed rigorous quality control of each dataset to  
107 remove low quality cells. To assess per-cell data quality, ArchR computes TSS enrichment  
108 scores, which have become the standard for bulk ATAC-seq analysis  
109 (<https://www.encodeproject.org/atac-seq/>) and provide clearer separation of low- and high-quality  
110 cells compared to other metrics such as the fraction of reads in promoters<sup>10</sup> (**Supplementary Fig.**  
111 **2c,f**).

112 To quantify the ability of ArchR to analyze large-scale data, we compared the performance  
113 of ArchR to that of SnapATAC and Signac for three of the major scATAC-seq analytical steps  
114 across these three datasets using two different computational infrastructures (**Supplementary**  
115 **Fig 3a and Supplementary Table 2**). We observed that ArchR outperforms SnapATAC and  
116 Signac in speed and memory usage across all comparisons, enabling analysis of 70,000 cell  
117 datasets in under an hour with 32 GB of RAM and 8 cores (**Fig. 1b-c and Supplementary Fig.**  
118 **3b-i**). Additionally, when analyzing a 70,000-cell dataset, SnapATAC exceeded the available



119 memory in the high memory setting (128 GB RAM, 20 cores) (**Fig. 1c**) and both SnapATAC and  
120 Signac exceeded the available memory in the low memory setting (32 GB RAM, 8 cores)  
121 (**Supplementary Fig. 3c**), while ArchR completed these analyses faster and without exceeding  
122 the available memory. Lastly, ArchR can analyze scATAC-seq data directly from BAM files,  
123 enabling the analysis of data from diverse single-cell platforms including the sci-ATAC-seq murine  
124 atlas<sup>5</sup> (**Supplementary Fig. 3j-k**).

125

### 126 **ArchR identifies putative doublets in scATAC-seq data**

127 The presence of so called “doublets” – two cells that are captured within the same nano-reaction  
128 (i.e. a droplet) and thus indexed with the same cellular barcode – often complicate single-cell  
129 analysis. Doublets appear as a superposition of signals from both cells, leading to the false  
130 appearance of distinct clusters or false connections between distinct cell types. To mitigate this  
131 issue, we designed a doublet detection and removal algorithm as part of ArchR. Similar to  
132 methods employed for doublet detection in scRNA-seq<sup>12,13</sup>, ArchR identifies heterotypic doublets  
133 by bioinformatically generating a collection of synthetic doublets, projecting these synthetic  
134 doublets into the low-dimensional data embedding, then identifying the nearest neighbors to these  
135 synthetic doublets as doublets themselves<sup>12,13</sup> (**Fig. 1d-f**). To validate this approach, we carried  
136 out scATAC-seq on a mixture of 10 highly distinct human cell lines (N = 38,072 cells), allowing  
137 for genotype-based identification of doublets via demuxlet<sup>14</sup> as a ground-truth comparison for  
138 computational identification of doublets by ArchR (**Fig. 1g and Supplementary Fig. 4a**). Using  
139 an unbiased optimization for the projection of synthetic doublets, we identified robust parameters  
140 (**Supplementary Fig. 4b**) for doublet prediction (ROC = 0.918) which significantly outperformed  
141 doublet prediction based on the total number of accessible fragments (ROC = 0.641) (**Fig. 1h**  
142 **and Supplementary Fig 4c-h**). With these predicted doublets excluded, the remaining cells  
143 formed 10 large groups according to their cell line of origin (**Fig. 1i**). ArchR’s implementation of

144 heterotypic doublet elimination reduces false cluster identification and thus improves the fidelity  
145 of downstream results.

146

### 147 **ArchR provides high-resolution and efficient dimensionality reduction of scATAC-seq data**

148 ArchR additionally provides methodological improvements over other available software. One of  
149 the most fundamental aspects of ATAC-seq analysis is the identification of a feature set (i.e. a  
150 peak set) for downstream analysis. In the context of single-cell ATAC-seq, identification of peak  
151 regions prior to cluster identification requires peak calling from all cells as a single merged group.  
152 This effectively obscures cell type-specific chromatin accessibility which distorts downstream  
153 analyses. For Signac, a counts matrix is created using a pre-determined peak set, preventing the  
154 contribution of peaks that are specific to lowly represented cell types. Instead of using a pre-  
155 determined peak set, SnapATAC creates a genome-wide tiled matrix of 5-kb bins, allowing for  
156 unbiased genome-wide identification of cell type-specific chromatin accessibility. However, 5-kb  
157 bins are substantially larger than the average regulatory element (~300-500 bp containing TF  
158 binding sites less than 50 bp)<sup>15-17</sup>, thus causing multiple regulatory elements to be grouped  
159 together, again obscuring cell type-specific biology. To avoid both of these pitfalls, ArchR operates  
160 on a genome-wide tiled matrix of 500-bp bins, allowing for the sensitivity to capture cell type-  
161 specific biology at individual regulatory elements across the entire genome. Despite this 10-fold  
162 higher resolution tile matrix, ArchR stores both per-tile accessibility information and all ATAC-seq  
163 fragments in an Arrow file that is smaller than either the original input fragments or the Snap file  
164 from SnapATAC containing the genome-wide tiled matrix at only 5-kb resolution (**Supplementary**  
165 **Fig. 2a-b**).

166 One major application of single-cell analysis is the identification of cellular subsets through  
167 dimensionality reduction and clustering. For dimensionality reduction, ArchR uses an optimized  
168 iterative latent semantic indexing (LSI) method<sup>6,7</sup> (**Supplementary Fig. 5a**), Signac uses an LSI  
169 method, and SnapATAC uses a method based on Jaccard indices. When directly comparing the

170 results from these different dimensionality reduction methods, ArchR identified similar clusters to  
171 other methods while being less biased by low-quality cells and doublets (**Supplementary Fig.**  
172 **5b**). However, when comparing clustering of the bone marrow cell dataset, we found that ArchR  
173 alone maintained the continuous cellular hierarchy expected in this biological system  
174 (**Supplementary Fig. 6a**).

175 To enable the efficient examination of extremely large datasets, ArchR implements a novel  
176 estimated LSI dimensionality reduction by first creating an iterative LSI reduction from a subset  
177 of the total cells, then linearly projecting the remainder of cells into this reduced dimension space  
178 using LSI projection<sup>7</sup> (**Supplementary Fig. 7a**). We compared this approach to the landmark  
179 diffusion map (LDM) estimation method used by SnapATAC which uses a non-linear reduction  
180 based on a subset of cells and then projects the remainder of the cells into this subspace using  
181 LDM projection. When comparing “landmark” subsets of different cell numbers, the estimated LSI  
182 approach implemented by ArchR was more consistent and could recapitulate the clusters called  
183 and the overall structure of the data with as few as 50 cells across both the PBMC (N = 27,845  
184 cells) and bone marrow cell (N = 26,748 cells) datasets (**Supplementary Fig. 7b and 8a-b**). We  
185 speculate that this observed robustness stems from the linearity of the LSI projection as compared  
186 to LDM projection, which occurs in a non-linear subspace. The estimated LSI approach  
187 implemented by ArchR is also faster than the estimated LDM approach implemented by  
188 SnapATAC (**Supplementary Fig. 8c**). Furthermore, the efficiency of the standard iterative LSI  
189 implementation in ArchR limits the requirement for this estimated LSI approach to only extremely  
190 large datasets (>200,000 cells for 32 GB RAM and 8 cores), whereas estimated LDM approaches  
191 are required for comparatively smaller datasets (>25,000 cells for 32 GB and 8 cores) in  
192 SnapATAC. ArchR therefore has the capacity to rapidly and efficiently analyze both large- and  
193 small-scale datasets.

194

195 **Robust inference of gene scores enables accurate cluster identification with ArchR**

196 After clustering, investigators often aim to annotate the biological state related to each cluster.  
197 Methods for inferring gene expression from scATAC-seq data can generate “gene scores” of key  
198 marker genes that can enable accurate cluster annotation<sup>5–8,18</sup>. However, the methods for  
199 integrating chromatin accessibility signal to generate these gene score predictions have not been  
200 extensively optimized. To this end, we used ArchR to benchmark 56 different models for inferring  
201 gene expression from scATAC-seq data using matched scATAC-seq and scRNA-seq data from  
202 PBMCs and bone marrow cells (**Fig. 2a and Supplementary Table 3**). To assess the  
203 performance of each model, we compared the known gene expression from previous methods  
204 integrating scATAC-seq with scRNA-seq<sup>7,11</sup> to the inferred gene scores derived from the model.  
205 By first establishing a rough linkage of ATAC-seq to RNA expression across many relatively  
206 diverse cell types (**Fig. 2a**), we could then determine which method for integrating ATAC-seq  
207 signal to predict gene expression had the best global performance across these data. The 56  
208 gene score models varied by the regions included, the sizes of those regions, and the weights  
209 (based on genomic distance) applied to each region (**Fig. 2b and Supplementary Fig. 9a-h**).  
210 Models that incorporated ATAC-seq signal from the gene-body were more accurate than models  
211 that incorporated signal only from the promoter, likely due to the moderate increase in accessibility  
212 that occurs during active transcription. Moreover, incorporation of distal regulatory elements,  
213 weighted by distance, while accounting for the presence of neighboring genes (see methods)  
214 increased the accuracy of the gene score inference in all cases (**Supplementary Fig. 9a-h**). The  
215 most accurate model across both datasets was Model 42 (a model within the gene body extended  
216 + exponential decay + gene boundary class of models) (**Fig. 2b**) which integrates signal from the  
217 entire gene body, and scales signal with bi-directional exponential decays from the gene TSS  
218 (extended upstream by 5 kb) and the gene transcription termination site (TTS) while accounting  
219 for neighboring genes boundaries (**Fig. 2c**). This model yielded robust genome-wide gene score  
220 predictions in both PBMC and bone marrow cell datasets (**Fig. 2d-f and Supplementary Fig. 9i-**  
221 **j**). We additionally confirmed the efficacy of this class of gene score models using previously

222 published paired bulk ATAC-seq and RNA-seq data from hematopoietic cells (**Supplementary**  
223 **Fig. 9k-m**)<sup>19</sup>. Given this analysis, we implemented this class of gene score models via Model 42  
224 for all downstream analyses involving inferred gene expression in ArchR.

225

### 226 **ArchR enables comprehensive analysis of massive-scale scATAC-seq data**

227 ArchR is designed to handle datasets substantially larger (>1,000,000 cells) than those generated  
228 to date with modest computational resources. To illustrate this, we collected a compendium of  
229 high-quality published scATAC-seq data from immune cells generated with either the 10x  
230 Chromium system or the Fluidigm C1 system (49 samples, ~220k cells; **Supplementary Figure**  
231 **10a-d**). We refer to this compiled dataset as the hematopoiesis dataset. Using both a small-scale  
232 server infrastructure (8 cores, 32 GB RAM, with an HP Lustre file system) and a personal laptop  
233 (MacBook Pro laptop; 8 cores, 32 GB RAM, with an external USB hard drive), ArchR performed  
234 data import, dimensionality reduction, and clustering on ~220k cells in less than three hours (**Fig.**  
235 **3a and Supplementary Fig. 10e**). We next used ArchR to analyze a simulated set of over 1.2  
236 million PBMCs, split into 200 individual samples. Under the same computational constraints,  
237 ArchR performed data import, dimensionality reduction, and clustering of more than 1.2 million  
238 cells in under 8 hours (**Fig. 3a and Supplementary Fig. 10e**).

239 Beyond these straightforward analyses, ArchR also provides an extensive suite of tools  
240 for more comprehensive analysis of scATAC-seq. Here we demonstrate these applications using  
241 the hematopoiesis dataset described above. Estimated LSI of this ~220k-cell dataset  
242 recapitulated the overall structure of the data with a landmark dataset of as few as 500 cells  
243 (**Supplementary Fig. 10f**). Manual inspection of the resultant clusters with our uniform manifold  
244 approximation and projection (UMAP)<sup>20</sup> led us to use the 25,000 cell landmark set (~10% of total  
245 cells), which additionally showed minimal bias due to batch and data quality (**Fig. 3b and**  
246 **Supplementary Fig. 10g-i**). We identified 21 clusters spanning the hematopoietic hierarchy,  
247 calling clusters for even rare cell types such as plasma cells which comprise ~0.1% (265 cells) of

248 the total population. To generate a universal peak set from cluster-specific peaks, ArchR creates  
249 sample-aware pseudo-bulk replicates that recapitulate the biological variability within each cluster  
250 (**Supplementary Fig. 11a**). Peaks are then called from these pseudo-bulk replicates and a set of  
251 reproducible fixed-width non-overlapping peaks are identified using an iterative overlap merging  
252 procedure<sup>21</sup> (**Supplementary Fig. 11b**). Using this approach, we identified 396,642 total  
253 reproducible peaks (**Supplementary Fig. 11c**), of which 215,916 are classified as differentially  
254 accessible peaks across the 21 clusters after bias-matched differential testing (see methods; **Fig.**  
255 **3c**). Motif enrichment within these marker peaks revealed known TF regulators of hematopoiesis  
256 such as GATA1 in erythroid populations, CEBPB in monocytes, and PAX5 in B cell differentiation  
257 (**Fig. 3d**). In addition to motif enrichments, ArchR can calculate peak overlap enrichment with a  
258 compendium of previously published ATAC-seq datasets<sup>19,21–26</sup>, identifying strong enrichment of  
259 peaks consistent with the cell type of each cluster (**Supplementary Fig. 11d**). To further  
260 characterize clusters, ArchR enables the projection of bulk ATAC-seq data into the single-cell-  
261 derived UMAP embedding<sup>7</sup> via a down-sampling approach (**Supplementary Fig. 12a**). This  
262 allows for projection of sorted cell types, facilitating the identification of clusters based on well-  
263 validated bulk ATAC-seq profiles<sup>19</sup> (**Supplementary Fig. 12b**). This projection analysis generates  
264 cell positions from bulk ATAC-seq data consistent with known cell types from a Fluidigm C1  
265 scATAC-seq dataset of sorted hematopoietic cells including highly-similar hematopoietic stem  
266 and progenitor cells<sup>4</sup> (**Supplementary Fig. 12c**) and aligns with inferred gene scores for  
267 canonical hematopoietic marker genes (**Supplementary Fig. 12d**).

268 ArchR also implements a scalable method for determination of transcription factor  
269 deviations from chromVAR<sup>27</sup> in a sample independent manner (**Supplementary Fig. 12e**). TFs  
270 whose expression is highly correlated with their motif accessibility (i.e. putative positive  
271 regulators) can therefore be identified based on the correlation of the inferred gene score to the  
272 chromVAR motif deviation. This analysis identifies known drivers of hematopoietic differentiation  
273 such as GATA1 in erythroid populations, LEF1 in Naive T cell populations, and EOMES in NK/T

274 Cell Memory populations. (**Fig. 3e, Supplementary Fig. 12f, and Supplementary Table 4**).  
275 ArchR also enables rapid footprinting of these TF regulators within clustered subsets while  
276 accounting for Tn5 biases<sup>21</sup> using an improved C++ implementation (**Fig. 3f-h, Supplementary**  
277 **Fig. 12g-i**). Finally, ArchR identifies links between regulatory elements and target genes based  
278 on the co-accessibility of pairs of loci across single cells<sup>1,18</sup> (**Fig. 3i**).

### 279 **The interactive ArchR genome browser**

280 In addition to these robust ATAC-seq analysis paradigms, ArchR provides a fully integrated and  
281 interactive genome browser (**Supplementary Fig. 13a**). The responsive and interactive nature of  
282 the browser is enabled by the optimized storage format within each Arrow file, providing support  
283 for dynamic cell grouping, track resolution, coloration, layout, and more. Launched via a single  
284 command, the ArchR browser enables cell cluster investigations of marker genes such as *CD34*  
285 for early hematopoietic stem and progenitor cells and *CD14* for monocytic populations (**Fig. 3i**  
286 **and Supplementary Fig. 13b-e**) while mitigating the need for external software for visualization  
287 of scATAC-seq data.

288

### 289 **ArchR enables integration of matched scRNA-seq and scATAC-seq datasets**

290 ArchR also provides functionality to integrate scATAC-seq data with scRNA-seq data using  
291 Seurat's infrastructure<sup>11</sup>. In brief, this integration requires matching the chromatin accessibility  
292 profiles and RNA expression for independent heterogeneous cells measured with two different  
293 assays. Single-cell epigenome-to-transcriptome integration is essential for understanding  
294 dynamic gene regulatory processes, and we anticipate this sort of analysis will become even more  
295 prevalent with the advent of platforms for simultaneous scATAC-seq and scRNA-seq. ArchR  
296 efficiently performs this cross-data alignment in parallel using slices of the scATAC-seq data (**Fig.**  
297 **4a**). When performed on the hematopoiesis dataset, this integration enabled accurate scRNA-  
298 seq alignment for >220,000 cells in less than 1 hour (**Fig. 4b**). The alignment showed high  
299 concordance between linked gene expression and inferred gene scores for common



300 hematopoietic marker genes (**Fig. 4c and Supplementary Fig. 14a**). Using this cross-platform  
301 alignment, ArchR also provides methods to identify putative cis-regulatory elements based on  
302 correlated peak accessibility and gene expression<sup>7,21</sup> (**Supplementary Fig. 15a**). In the example  
303 hematopoiesis dataset, this analysis identified 70,239 significant peak-to-gene linkages across  
304 the hematopoietic hierarchy (**Supplementary Fig. 15b and Supplementary Table 5**).

305 Finally, ArchR facilitates cellular trajectory analysis to identify the predicted path of gene  
306 regulatory changes from one set of cells to another, a unique type of insight enabled by single-  
307 cell data. To carry out this analysis, ArchR initially creates a cellular trajectory based on a  
308 sequence of user-supplied clusters or groups. ArchR then identifies individual cell positions along  
309 this trajectory based on Euclidean distance within an N-dimensional subspace<sup>6</sup>. Using B cells as  
310 an example, ArchR traces cells along the B cell differentiation trajectory and identifies 11,999  
311 peak-to-gene links that have correlated regulatory dynamics across the B cell differentiation (**Fig.**  
312 **4e**). Sequencing tracks of the *HMGA1* locus, active in stem and progenitor cells, and the *BLK*  
313 locus, active in differentiated B cells, demonstrate how accessibility at linked peaks correlates  
314 with longitudinal changes in gene expression across pseudo-time (**Fig. 4f-g**). Moreover, using  
315 this same paradigm, ArchR can identify TF motifs with accessibility that are positively correlated  
316 with gene expression of TF genes across the same B cell trajectory (**Fig. 4h**). Transcription factor  
317 footprinting of a subset of these TFs further illustrates the dynamics in the local accessibility at  
318 the binding sites of these lineage-defining TFs across B cell differentiation pseudo-time (**Fig. 4i-**  
319 **k**).

320

## 321 DISCUSSION

322 Chromatin accessibility data provides a lens through which we can observe the gene regulatory  
323 programs that underlie cellular state and identity. The highly cell type-specific nature of cis-  
324 regulatory elements makes profiling of single-cell chromatin accessibility an attractive method to  
325 understand cellular heterogeneity and the molecular processes underlying complex control of



326 gene expression. With the advent of methods to profile chromatin accessibility across thousands  
327 of single cells, scATAC-seq has quickly become a method-of-choice for many single-cell  
328 applications. However, compared to scRNA-seq, analysis of scATAC-seq data remains  
329 comparatively immature with no clear standards, thus dissuading many from adopting this  
330 informative technique.

331 To address this unmet need, we developed ArchR, an end-to-end software solution that  
332 will greatly expedite single-cell chromatin analysis for any biologist. Low memory usage,  
333 parallelized operations, and an intuitive and user-focused, yet extensive and powerful tool suite  
334 make ArchR an ideal platform for scATAC-seq data analysis. In contrast to currently available  
335 software packages, ArchR is designed to handle millions of cells using commonly available  
336 computational resources, such as a laptop running a Unix-based operating system. As such,  
337 ArchR provides the analytical support necessary for the massive scale of ongoing efforts to  
338 catalog the compendium of diverse cell types throughout the body at single-cell resolution<sup>28</sup>. In  
339 addition to the dramatic improvements in run time, memory efficiency, and scale, ArchR supports  
340 state-of-the-art chromatin-based analyses including genome-wide inference of gene activity,  
341 transcription factor footprinting, and data integration with matched scRNA-seq, enabling statistical  
342 linkage of cis- and trans-acting regulatory factors to gene expression profiles. Moreover, the  
343 performance improvements from ArchR enable interactive data analysis whereby end-users can  
344 iteratively adjust analytical parameters and thus optimize identification of biologically meaningful  
345 results. This is especially important in the context of single-cell data where a one-size-fits-all  
346 analytical pipeline is not relevant or desirable. Supervised identification of clusters, resolution of  
347 subtle batch effects, and biology-driven data exploration are intrinsically necessary for a  
348 successful scATAC-seq analysis and ArchR supports these efforts by enabling rapid analytical  
349 processes. ArchR provides an open-source analysis platform with the flexibility, speed, and power  
350 to support the rapidly increasing efforts to understand complex tissues, organisms, and  
351 ecosystems at the resolution of individual cells.

## 352 **Methods**

353

## 354 **Code Availability and Documentation**

355 Extensive documentation and a full user manual are available at [www.ArchRProject.com](http://www.ArchRProject.com). The  
356 software is open-source and all code can be found on GitHub at  
357 <https://github.com/GreenleafLab/ArchR>. Additionally, code for producing the majority of analyses  
358 from this paper is available at the publication page [https://github.com/GreenleafLab/ArchR\\_2020](https://github.com/GreenleafLab/ArchR_2020).

359

## 360 **Data Availability**

361 Bulk and scATAC-seq data from the cell line mixing experiment will be available through GEO  
362 (accession number in progress). All other scATAC-seq data used were from publicly available  
363 sources as outline in **Supplementary Table 1**. We additionally have made available other  
364 analysis files on our publication page [https://github.com/GreenleafLab/ArchR\\_2020](https://github.com/GreenleafLab/ArchR_2020).

365

## 366 **Genome and Transcriptome Annotations**

367 All analyses were performed with the hg19 genome (except the Mouse Atlas with mm9). R-based  
368 analysis used the BSgenome package with “BSgenome.Hsapiens.UCSC.hg19”  
369 (“BSgenome.Mmusculus.UCSC.mm9” for Mouse Atlas) for genomic coordinates and the TxDb  
370 package with “TxDb.Hsapiens.UCSC.hg19.knownGene”  
371 (“TxDb.Mmusculus.UCSC.mm9.knownGene” for Mouse Atlas) gene annotations unless  
372 otherwise stated.

373

## 374 **Cell Type Abbreviations**

375 In many of the figure legends, abbreviations are used for cell types of the hematopoietic system.  
376 HSC, hematopoietic stem cell; LMPP, lymphoid-primed multipotent progenitor cell; CMP, common  
377 myeloid progenitor; CLP, common lymphoid progenitor; GMP, granulocyte macrophage

378 progenitor; CD4 Mem, CD4 memory T cell; CD4 Naive, CD4 naïve T cell; CD8 Naive, CD8 naïve  
379 T cell; CD8 Eff, CD8 effector T cell; CD8 EffMem, CD8 effector memory T cell; CD8 CenMem,  
380 CD8 central memory T cell; Mono, monocyte; pDC, plasmacytoid dendritic cell; NK, natural killer  
381 cell; Ery, erythroid; Baso, basophil.

382

### 383 **scATAC-seq Data Generation – Cell Lines**

384 With the exception of MCF10A, all cell lines were cultured in the designated media containing  
385 10% FBS and penicillin/streptomycin. Jurkat, THP1, and K562 cell lines were ordered from ATCC  
386 and cultured in RPMI-1640. GM12878 cells were ordered from Coriell and cultured in RPMI-1640.  
387 HeLa, HEK-293T, and HT1080 cell lines were ordered from ATCC and cultured in DMEM. T24  
388 cells were ordered from ATCC and cultured in McCoy's 5A. MCF7 cells were ordered from ATCC  
389 and cultured in EMEM containing 0.01 mg/ml of human insulin (Millipore-Sigma 91077C).  
390 MCF10A cells were ordered from ATCC and cultured in DMEM/F12 containing 5% horse serum  
391 (Thermo Fisher 16050130), 0.02 ug/ml human EGF (PeproTech AF-100-15), 0.5 ug/ml  
392 hydrocortisone (Millipore-Sigma H0888), 0.1 ug/ml Cholera toxin (Millipore-Sigma C8052), 10  
393 ug/ml insulin from bovine pancreas (Millipore-Sigma I6634), and penicillin/streptomycin. Cultured  
394 cells were viably cryopreserved in aliquots of 100,000 cells using 100 ul of BAMBANKER freezing  
395 media (Wako Chemicals 302-14681) so that scATAC-seq could be performed on all cells at the  
396 same time. For each cell line, cells were thawed via the addition of 1 mL ice-cold resuspension  
397 buffer (RSB) [10 mM Trish-HCl pH 7.4, 10 mM NaCl, 3 mM MgCl<sub>2</sub>] containing 0.1% Tween-20  
398 (RSB-T). Cells were pelleted in a fixed-angle rotor at 300 RCF for 5 minutes at 4°C. The  
399 supernatant was removed and the pellet was resuspended in 100 uL of ice-cold lysis buffer (RSB  
400 containing 0.1% Tween-20, 0.1% NP-40, and 0.01% digitonin) and incubated on ice for 3 minutes.  
401 To dilute lysis, 1 mL of chilled RSB-T was added to each tube and the cells were pelleted as  
402 before. The supernatant was removed and the pelleted nuclei were resuspended in Diluted Nuclei  
403 Buffer (10x Genomics). The nuclei stock concentration was determined for each cell line using

404 Trypan Blue and a total of 5,000 nuclei from each cell line were pooled together and loaded into  
405 the 10x Genomics scATAC-seq (v1) transposition reaction. The remainder of the scATAC-seq  
406 library preparation was performed as recommended by the manufacturer. Resultant libraries were  
407 sequenced on an Illumina NovaSeq6000 using an S4 flow cell and paired-end 99-bp reads. In  
408 addition to this pooled scATAC-seq, each cell line was used to generate bulk ATAC-seq libraries  
409 as described previously<sup>26</sup>. Bulk ATAC-seq libraries were pooled and purified via PAGE gel prior  
410 to sequencing on an Illumina HiSeq4000 using paired-end 75-bp reads.

411

#### 412 **scATAC-seq Processing – Cell Line Mixing**

413 Raw sequencing data was converted to FASTQ format using cellranger-atac mkfastq (10x  
414 Genomics, version 1.0.0). Single-cell ATAC-seq reads were aligned to the hg19 reference  
415 genome (<https://support.10xgenomics.com/single-cell-atac/software/downloads/latest>) and  
416 quantified using cellranger-count (10x Genomics, version 1.0.0). Genotypes used to perform  
417 demuxlet were determined as follows for each cell line: Bulk ATAC-seq FASTQ files were  
418 processed and aligned using PEPATAC (<http://code.databio.org/PEPATAC/>) as described  
419 previously<sup>21</sup>. Peaks were identified using MACS2 and a union set of variable-width accessible  
420 regions was identified using bedtools merge (version 2.26.0). These accessible regions were  
421 genotyped across all samples using samtools mpileup (version 1.5) and Varscan mpileup2snp  
422 (version 2.4.3) with the following parameters “--min-coverage 5 --min-reads2 2 --min-var-freq 0.1  
423 --strand-filter 1 --output-vcf 1”. All positions containing a single nucleotide variant were compiled  
424 into a master set and then each cell line was genotyped at those specific single-base locations  
425 using samtools mpileup. The allelic depth at each position was converted into a quaternary  
426 genotype (homozygous A, heterozygous AB, homozygous B, or insufficient data to generate a  
427 confident call). Then, for each cell line, inferred genotype probabilities were created based on  
428 those quaternary genotypes and a VCF file was created for input to demuxlet using recommended

429 parameters. Demuxlet was used to identify the cell line of origin for individual cells and to identify  
430 doublets based on mixed genotypes.

431

### 432 **ArchR Methods – Preface**

433 All ArchR features were carefully designed and optimized to enable analysis of 250,000 cells or  
434 greater on a minimal computing environment in R. All ArchR HDF5-formatted processing was  
435 performed with the Bioconductor<sup>29</sup> package “rhdf5”  
436 (<https://www.bioconductor.org/packages/release/bioc/html/rhdf5.html>). All ArchR genomic  
437 coordinate operations were performed with the Bioconductor package “GenomicRanges”  
438 (<https://bioconductor.org/packages/release/bioc/html/GenomicRanges.html>) and “IRanges”  
439 (<https://bioconductor.org/packages/release/bioc/html/IRanges.html>).

440

### 441 **ArchR Methods – scATAC Definitions**

442

443 **Fragments** – In ATAC-seq data analysis, a “fragment” refers to a sequenceable DNA molecule  
444 created by two transposition events. Each end of that fragment is sequenced using paired-end  
445 sequencing. The inferred single-base position of the start and end of the fragment is adjusted  
446 based on the insertion offset of Tn5. As reported previously<sup>30</sup>, Tn5 transposase binds to DNA as  
447 a homodimer with 9-bp of DNA between the two Tn5 molecules. Because of this interaction, each  
448 Tn5 homodimer binding event creates two insertions, separated by 9 bp. Thus, the actual central  
449 point of the “accessible” site is in the very center of the Tn5 dimer, not the location of each Tn5  
450 insertion. To account for this, we apply an offset to the individual Tn5 insertions, adjusting plus-  
451 stranded insertion events by +4 bp and minus-stranded insertion events by -5 bp. This offset is  
452 consistent with the convention put forth during the original description of ATAC-seq<sup>31</sup>. Thus, in  
453 ArchR, “fragments” refers to a table or Genomic Ranges object containing the chromosome,

454 offset-adjusted chromosome start position, offset-adjusted chromosome end position, and unique  
455 cellular barcode ID corresponding to each sequenced fragment.

456  
457 **Tn5 insertions** – In ArchR, “insertions” refers to the offset-adjusted single-base position of Tn5  
458 insertion on either end of the fragment. Insertion positions are accessed in ArchR primarily using  
459 `resize(fragments, 1, “start”)` and `resize(fragments, 1, “end”)`. See the description of “fragments”  
460 above for a detailed description of Tn5 insertion offsets.

461  
462 **Counting Accessibility** – In ArchR, “counting accessibility” refers to counting the number of Tn5  
463 insertions observed within each described feature.

464  
465 **TSS enrichment score** – In ArchR, the “TSS enrichment” refers to the relative enrichment of Tn5  
466 insertions at gene TSS sites genome-wide compared to a local background. This represents a  
467 measure of signal-to-background in ATAC-seq data. See below for how TSS enrichment is  
468 calculated in ArchR. In this work, the TSS enrichment score from ArchR is based on the TSS  
469 regions defined by the `TxDb.Hsapiens.UCSC.hg19.knownGene` (or  
470 `TxDb.Mmusculus.UCSC.mm9.knownGene` for the Mouse Atlas) transcript database object.

471  
472 **ArchR Methods – Arrow Files and ArchRProject**

473 The base unit of an analytical project in ArchR is called an “Arrow file”. Each Arrow file stores all  
474 of the data associated with an individual sample (i.e. metadata, accessible fragments, and data  
475 matrices). Here, an “individual sample” would be the most detailed unit of analysis desired (for  
476 ex. a single replicate of a particular condition). During creation and as additional analyses are  
477 performed, ArchR updates and edits each Arrow file to contain additional layers of information. It  
478 is worth noting that, to ArchR, an Arrow file is actually just a path to an external file stored on disk.  
479 More explicitly, an Arrow file is not an R-language object that is stored in memory but rather an

480 HDF5-format file stored on disk. Because of this, we use an “ArchRProject” object to associate  
481 these Arrow files together into a single analytical framework that can be rapidly accessed in R.  
482 This ArchRProject object is small in size and is stored in memory.

483         Certain actions can be taken directly on Arrow files while other actions are taken on an  
484 ArchRProject which in turn updates each associated Arrow file. Because Arrow files are stored  
485 as large HDF5-format files, “get-er” functions in ArchR retrieve data by interacting with the  
486 ArchRProject while “add-er” functions either (i) add data directly to Arrow files, (ii) add data directly  
487 to an ArchRProject, or (iii) add data to Arrow files by interacting with an ArchRProject.

488

### 489 **ArchR Methods – Reading Input Data into an Arrow File**

490 ArchR can utilize multiple input formats of scATAC-seq data which is most frequently in the format  
491 of fragment files and BAM files. Fragment files are tabix-sorted text files containing each scATAC-  
492 seq fragment and the corresponding cell ID, one fragment per line. BAM files are binarized tabix-  
493 sorted files that contain each scATAC-seq fragment, raw sequence, cellular barcode id and other  
494 information. The input format used will depend on the pre-processing pipeline used. For example,  
495 the 10x Genomics Cell Ranger software returns fragment files while sci-ATAC-seq applications  
496 would use BAM files. Given a specified genome annotation (ArchR has pre-loaded genome  
497 annotations for mm9, mm10, hg19, and hg38 and additional genomes can be added manually),  
498 ArchR reads these input files in sub-chromosomal chunks using Rsamtools. ArchR uses  
499 “scanTabix” to read fragment files and “scanBam” to read BAM files. During this input process,  
500 each input chunk is converted into a compressed table-based representation of fragments  
501 containing each fragment chromosome, offset-adjusted chromosome start position, offset-  
502 adjusted chromosome end position and cellular barcode ID. These chunk-wise fragments are  
503 then stored in a temporary HDF5-formatted file to preserve memory usage while maintaining rapid  
504 access to each chunk. Finally, all chunks associated with each chromosome are read, organized,  
505 and re-written to an “Arrow file” within a single HDF5 group called “fragments”. This pre-chunking

506 procedure enables ArchR to process extremely large input files efficiently and with low memory  
507 usage, enabling full utilization of parallel processing.

508

### 509 **ArchR Methods – QC Based on TSS Enrichment and Unique Nuclear Fragments**

510 Strict quality control (QC) of scATAC-seq data is essential to remove the contribution of low-  
511 quality cells. In ArchR, one characteristic of “low-quality” is a low signal-to-background ratio, which  
512 is often attributed to dead or dying cells which have de-chromatinized DNA which allows for  
513 random transposition genome-wide. Traditional bulk ATAC-seq analysis has used the TSS  
514 enrichment score as part of a standard workflow (<https://www.encodeproject.org/atac-seq/>) for  
515 determination of signal-to-background. We and others have found the TSS enrichment to be  
516 representative across the majority of cell types tested in both bulk ATAC-seq and scATAC-seq.  
517 The idea behind the TSS enrichment score metric is that ATAC-seq data is universally enriched  
518 at gene TSS regions compared to other genomic regions. By looking at per-base-pair accessibility  
519 centered at these TSS regions, we see a local enrichment relative to flanking regions (1900-2000  
520 bp distal in both directions). The ratio between the peak of this enrichment (centered at the TSS)  
521 relative to these flanking regions represents the TSS enrichment score. Traditionally, the per-  
522 base-pair accessibility is computed for each bulk ATAC-seq sample and then this profile is used  
523 to determine the TSS enrichment score. Performing this operation on a per-cell basis in scATAC-  
524 seq is relatively slow and computationally expensive. To accurately approximate the TSS  
525 enrichment score per single cell, we count the average accessibility within a 50-bp region  
526 centered at each single-base TSS position and divide this by the average accessibility of the TSS  
527 flanking positions (+/- 1900 – 2000 bp). This approximation was highly correlated ( $R > 0.99$ ) with  
528 the original method and values were extremely close in magnitude. By default in ArchR, pass-  
529 filter cells are identified as those cells having a TSS enrichment score greater than 4 and more  
530 than 1000 unique nuclear fragments (i.e those fragments that do not map to chrM).

531



## 532 **ArchR Methods – Tile Matrix**

533 Traditional bulk ATAC-seq analysis relies on the creation of a peak matrix from a peak-set  
534 encompassing the precise accessible regions across all samples. This peak set, and thus the  
535 resulting peak matrix, is specific to the samples used in the analysis and must be re-generated  
536 when new samples are added. Moreover, identification of peaks from scATAC-seq data would  
537 optimally be performed after clusters were identified to ensure that cluster-specific peaks are  
538 captured. Thus, the optimal solution for scATAC-seq would be to identify an unbiased and  
539 consistent way to perform analysis prior to cluster identification, without the need for calling peaks.  
540 xbecause this bin size approximates the size of most regulatory elements. To circumvent the  
541 requirement for calling peaks prior to cluster identification, others have tiled the genome into fixed  
542 non-overlapping tiled windows. This method additionally benefits from being stable across  
543 samples and the tiled regions do not change based on inclusion of additional samples. However,  
544 these tiled windows are usually greater than or equal to 5 kb in length, which is more than 10-fold  
545 greater than the size of typical accessible regions containing TF binding sites<sup>15-17</sup>. For this reason,  
546 ArchR uses 500-bp genome-wide tiled windows for all analysis upstream of cluster identification.  
547 To create a tile matrix, ArchR reads in the scATAC-seq fragments for a chromosome and converts  
548 these to insertions. ArchR then floors these insertions to the nearest tile region with  $\text{floor}(\text{insertion}$   
549  $/ \text{tileSize}) + 1$ . The tile regions and cell barcode id (as an integer) are then used as input for  
550 `Matrix::sparseMatrix` which tallies the number of input rows (tiles, denoted as  $i$ ) and columns (cells,  
551 denoted as  $j$ ) and creates a `sparseMatrix`. This analysis is performed for each chromosome and  
552 stored in the corresponding Arrow file. This fast and efficient conversion of scATAC-seq fragments  
553 to a tile matrix, without computing genomic overlaps, facilitates efficient construction of 500-bp  
554 tile matrices for analyses.

555

## 556 **ArchR Methods – Gene Score Matrix**

557 ArchR facilitates the inference of gene expression from chromatin accessibility (called “gene  
558 scores”) by using custom distance-weighted accessibility models. For each chromosome, ArchR  
559 creates a tile matrix (user-defined tile size that is not pre-computed, default is 500 bp), overlaps  
560 these tiles with the gene window (user-defined, default is 100 kb), and then computes the distance  
561 from each tile (start or end) to the gene body (with optional extensions upstream or downstream)  
562 or gene start. We have found that the best predictor of gene expression is the local accessibility  
563 of the gene region which includes the promoter and gene body (**Supplementary Fig. 9**). To  
564 properly account for distal accessibility, for each gene ArchR identifies the subset of tiles that are  
565 within the gene window and do not cross another gene region. This filtering allows for inclusion  
566 of distal regulatory elements that could improve the accuracy of predicting gene expression values  
567 but excludes regulatory elements more likely to be associated with another gene (for ex. the  
568 promoter of a nearby gene). The distance from each tile to the gene is then converted to a  
569 distance weight using a user-defined accessibility model (default is  $e^{-(\text{abs}(\text{distance})/5000)} + e^{-1}$ ). When  
570 the gene body is included in the gene region (where the distance-based weight is the maximum  
571 weight possible), we found that extremely large genes can bias the overall gene scores. In these  
572 cases, the total gene scores can vary substantially due to the inclusion of insertions in both introns  
573 and exons. To help adjust for these large differences in gene size, ArchR applies a separate  
574 weight for the inverse of the gene size ( $1 / \text{gene size}$ ) and scales this inverse weight linearly from  
575 1 to a hard max (which can be user-defined, with a default of 5). Smaller genes thus receive larger  
576 relative weights, partially normalizing this length effect. The corresponding distance and gene size  
577 weights are then multiplied by the number of Tn5 insertions within each tile and summed across  
578 all tiles within the gene window (while still accounting for nearby gene regions as described  
579 above). This summed accessibility is a “gene score” and is depth normalized across all genes to  
580 a constant (user-defined, default of 10,000). Computed gene scores are then stored in the  
581 corresponding Arrow file for downstream analyses.

582

## 583 **ArchR Methods – Iterative LSI Procedure**

584 The default LSI implementation in ArchR is conceptually similar to the method introduced in  
585 Signac (<https://satijalab.org/signac/>), which, for a cell x features matrix (typically tiles or peaks),  
586 uses a term frequency (column sums) that has been depth normalized to a constant (10,000)  
587 followed by normalization with the inverse document frequency (1 / row sums) and then log-  
588 transformed (aka  $\log(\text{TF-IDF})$ ). This normalized matrix is then factorized by singular value  
589 decomposition (SVD) and then standardized across the reduced dimensions for each cell via z-  
590 score. ArchR additionally allows for the use of alternative LSI implementations based on  
591 previously published scATAC-seq papers<sup>5-7</sup>. As mentioned above, the input to LSI-based  
592 dimensionality reduction is the genome-wide 500-bp tile matrix.

593 In scRNA-seq, identifying variable genes is a common way to compute dimensionality  
594 reduction (such as PCA), as these highly variable genes are more likely to be biologically  
595 important, and focusing on these genes likely reduces low-level contributions of variance  
596 potentially due to experimental noise. ScATAC-seq data is binary, precluding the possibility of  
597 identifying variable peaks for dimensionality reduction. Therefore, rather than identifying the most  
598 variable peaks, we initially tried using the most accessible features as input to LSI; however, the  
599 results when running multiple samples exhibited a high degree of noise and low reproducibility.  
600 We therefore moved to our previously described "iterative LSI" approach<sup>6,7</sup>. This approach  
601 computes an initial LSI transformation on the most accessible tiles and identifies lower resolution  
602 clusters that are driven by clear biological differences. For example, when performed on  
603 peripheral blood mononuclear cells, this approach will identify clusters corresponding to the major  
604 cell types (T cells, B cells, and monocytes). Then ArchR computes the average accessibility for  
605 each of these clusters across all features creating "pseudo-bulks". ArchR then identifies the most  
606 variable peaks across these pseudo-bulks to use as features for the second round of LSI. In this  
607 second iteration, the selected variable peaks correspond more similarly to the variable genes  
608 used in scRNA-seq LSI implementations, insofar as they are highly variable across biologically

609 meaningful clusters. We have found this approach can also effectively minimize batch effects and  
610 allows operations on a more reasonably sized feature matrix. Additionally, we observe that this  
611 procedure still allows the identification of rare cell types, such as plasma cells in the bone marrow  
612 cell dataset that exist at ~0.1% prevalence. For larger batch effects, ArchR enables Harmony-  
613 based batch correction on the LSI-reduced coordinates<sup>32</sup>.

614

### 615 **ArchR Methods – Estimated LSI Procedure**

616 For extremely large scATAC-seq datasets, ArchR can estimate the LSI dimensionality reduction  
617 with LSI projection. This procedure is similar to the iterative LSI workflow, however the LSI  
618 procedure differs. First, a subset of randomly selected “landmark” cells is used for LSI  
619 dimensionality reduction. Second, the remaining cells are TF-IDF normalized using the inverse  
620 document frequency determined from the landmark cells. Third, these normalized cells are  
621 projected into the SVD subspace defined by the landmark cells. This leads to an LSI  
622 transformation based on a small set of cells used as landmarks for the projection of the remaining  
623 cells. This estimated LSI procedure is efficient with ArchR because, when projecting the new cells  
624 into the landmark cells LSI, ArchR iteratively reads in the cells from each sample and LSI projects  
625 them without storing them all in memory. This optimization leads to minimal memory usage and  
626 further increases the scalability for extremely large datasets. Even with comparatively small  
627 landmark cell subsets (500-5000 cells), we find that this procedure is able to maintain the global  
628 structure and recapitulates the clusters well; however, the required landmark set size is  
629 dependent on the proportion of different cells within the dataset.

630

### 631 **ArchR Methods – Identification of Doublets**

632 Single-cell data generated on essentially any platform is susceptible to the presence of doublets.  
633 A doublet refers to a single nano-reaction (i.e. a droplet) that received a single barcoded bead  
634 and more than one cell/nucleus. This causes the reads from more than one cell to appear as a

635 single cell. For 10x Genomics applications, the percentage of total "cells" that are actually  
636 doublets is proportional to the number of cells loaded into the reaction. Even at lower cell loadings  
637 as recommended by standard kit use, more than 5% of the data may come from doublets, and  
638 this spurious data exerts substantial effects on clustering. This issue becomes particularly  
639 problematic in the context of developmental/trajectory data because doublets can look like a  
640 mixture between two cell types and this can be confounded with intermediate cell types or cell  
641 states.

642 To predict which "cells" are actually doublets in ArchR, we synthesize *in silico* doublets  
643 from the data by mixing the reads from thousands of combinations of individual cells. Next, we  
644 perform iterative LSI followed by UMAP for each individual sample. We then LSI project the  
645 synthetic doublets into the LSI subspace followed by UMAP projection. ArchR identifies the *k*-  
646 nearest neighbors (user-defined, default 10) to each simulated projected doublet. By iterating this  
647 procedure *N* times (user-defined, default 3 times the total number of cells), we can compute  
648 binomial enrichment statistics (assuming every cell could be a doublet with equal probability) for  
649 each single cell based on the presence of nearby simulated projected doublets (in the LSI or  
650 UMAP subspace defined by the user). This approach is similar to previous approaches<sup>12,13</sup>, but  
651 differs in that LSI is used for dimensionality reduction and UMAP projection is used for  
652 identification. The number of doublets to remove is then determined based on either the number  
653 of cells that pass QC or for the approximate number of cells loaded as defined by the user. While  
654 we have optimized these parameters for general use, users should sensibly check their results  
655 with and without doublet removal.

656

### 657 **ArchR Methods – Identification of Clusters**

658 ArchR uses established scRNA-seq clustering methods that use graph clustering on the LSI  
659 dimensionality reduction coordinates to resolve clusters. By default, ArchR uses Seurat's graph

660 clustering with “Seurat::FindClusters” for identifying high fidelity clusters<sup>11</sup>. ArchR additionally  
661 supports scran<sup>33</sup> for single-cell clustering.

662

663

#### 664 **ArchR Methods – t-SNE and UMAP Embeddings**

665 ArchR supports both t-distributed stochastic neighbor embedding (t-SNE) and uniform manifold  
666 approximation and projection (UMAP) single-cell embedding methodologies. ArchR uses  
667 previously determined reduced dimensions as input for these embeddings. t-SNE analysis is  
668 performed using the “Rtsne” package in R by default. UMAP analysis is performed using the  
669 “uwot” package in R by default. The results are stored within an ArchRProject and then used for  
670 plotting and subsequent analyses.

671

#### 672 **ArchR Methods – Sample-Aware Pseudo-Bulk Replicate Generation**

673 Because of the sparsity of scATAC-seq data, operations are often performed on aggregated  
674 groups of single cells. Most frequently, these groups are defined by clustering, and it is assumed  
675 that each local cluster represents a relatively homogeneous cell type or cell state. This process  
676 of combining data from multiple individual cells creates “pseudo-bulk” data, because it resembles  
677 the data derived from a bulk ATAC-seq experiment.

678 A feature unique to ArchR is the creation of sample-aware pseudo-bulk replicates from  
679 each cell group to use for performing statistical tests (such as reproducible peak identification or  
680 TF footprinting). ArchR does this via a complex decision tree which is dependent upon a user-  
681 specified desired number of replicates and number of cells per replicate as presented in  
682 **Supplementary Fig. 11a**. Briefly, ArchR attempts to create pseudo-bulk replicates in a sample-  
683 aware fashion. This means that each individual pseudo-bulk replicate only contains cells from a  
684 single biological sample. This feature enables the preservation of variability associated with  
685 biological replicates. If the desired number of replicates cannot be created in this fashion, ArchR

686 uses progressively less stringent requirements to create the pseudo-bulk replicates. First, ArchR  
687 attempts to create as many pseudo-bulk replicates in a sample-aware fashion as possible and  
688 then create the remaining pseudo-bulk replicates in a sample-agnostic fashion by sampling  
689 without replacement. If this is not possible, ArchR attempts to create the desired number of  
690 pseudo-bulk replicates in a sample-agnostic fashion by sample without replacement across all  
691 samples. If this is not possible, ArchR attempts the same procedure by sampling without  
692 replacement within a single replicate but with replacement across different replicates without  
693 exceeding a user-specified sampling ratio. If all of these attempts fail, ArchR will create the  
694 specified number of pseudo-bulk replicates by sampling with replacement within a single replicate  
695 and with replacement across different replicates. The fragments from all cells within a pseudo-  
696 bulk replicate are converted to insertions and to a run-length encoding (RLE) coverage object  
697 using the “coverage” function in R. This insertion coverage object (similar to a bigwig) is then  
698 written to a separate HDF5-formatted coverage file. ArchR next identifies single-base resolution  
699 Tn5 insertion sites for each pseudo-bulk replicate, resizes these 1-bp sites to k-bp (user-defined,  
700 default is 6) windows ( $-k/2$  and  $+(k/2 - 1)$  bp from insertion), and then creates a k-mer frequency  
701 table using the “oligonucleotidfrequency(w=k, simplify.as=“collapse”)” function from the  
702 Biostrings package. ArchR then calculates the expected k-mers genome-wide using the same  
703 function with the BSGenome-associated genome file. These Tn5 k-mer values represent the Tn5  
704 bias genome-wide and are then stored in the pseudo-bulk replicate HDF5 coverage file. This  
705 coverage file contains similar information to a bigwig file with Tn5 insertion bias but in a fast-  
706 access HDF5 format. This coverage file can be used for peak-calling and TF footprinting with Tn5  
707 bias correction.

708

### 709 **ArchR Methods – Peak Calling**

710 In ArchR, peak calling is performed on the HDF5-format pseudo-bulk-derived coverage files  
711 described above. By default, ArchR calls peak summits with MACS2 using single-base insertion

712 positions derived from the coverage files (written to a bed file with data.table) with user-specified  
713 values for MACS2 parameters including gsize, shift (default -75), and extsize (default 150) along  
714 with the “nomodel” and “nolambda” flags. These single-base peak summit locations are extended  
715 to a 501-bp width. We use 501-bp fixed-width peaks because they make downstream computation  
716 easier as peak length does not need to be normalized. Moreover, the vast majority of peaks in  
717 ATAC-seq are less than 501-bp wide. Using variable-width peaks also makes it difficult to merge  
718 peak calls from multiple samples without creating extremely large peaks that create confounding  
719 biases.

720 To create a merged non-overlapping fixed-width union peak set, ArchR implements an  
721 iterative overlap removal procedure that we introduced previously<sup>21</sup>. Briefly, peaks are first ranked  
722 by their significance, then the most significant peak is retained and any peak that directly overlaps  
723 with the most significant peak is removed from further analysis. This process is repeated with the  
724 remaining peaks until no more peaks exist. This procedure avoids daisy-chaining and still allows  
725 for use of fixed-width peaks. We use a normalized metric for determining the significance of peaks  
726 because the reported MACS2 significance is proportional to the sequencing depth. This process  
727 is outlined in **Supplementary Fig. 11b**.

728

### 729 **ArchR Methods – Interactive Genome Browser**

730 One challenge inherent to scATAC-seq data analysis is genome-track level visualizations of  
731 chromatin accessibility observed within groups of. Traditionally, track visualization requires  
732 grouping the scATAC-seq fragments, creating a genome coverage bigwig, and normalizing this  
733 track for quantitative visualization. Typically, end-users use a genome browser such as the  
734 WashU Epigenome Browser, the UCSC Genome Browser, or the IGV browser to visualize these  
735 sequencing tracks. This process involves using multiple software and any change to the cellular  
736 groups or addition of more samples requires re-generation of bigwig files etc., which can become  
737 time consuming. For this reason, ArchR has a Shiny-based interactive genome browser that can



738 be launched with a simple line of code “ArchRBrowser(ArchRProj)”. The data storage strategy  
739 implemented in Arrow files allows this interactive browser to dynamically change the cell  
740 groupings, resolution, and normalization, enabling real-time track-level visualizations. The ArchR  
741 Genome Browser also creates high-quality vectorized images in PDF format for publication or  
742 distribution. Additionally, the browser accepts user-supplied input files such as BED files or  
743 GenomicRanges to display features or genomic interaction files that define co-accessibility, peak-  
744 to-gene linkages, or loops from chromatin conformation data.

745 To facilitate this interactive browser, ArchR utilizes the same optimizations described  
746 above for creating a genome-wide TileMatrix to create a TileMatrix for the chosen resolution  
747 specified within the plotting window. Cells corresponding to the same group are summed per tile  
748 and the resulting group matrix represents the accessibility in tiles across the specified window.  
749 This matrix can then be normalized by either the total number of reads in TSS/peak regions, the  
750 total number of cells, or the total number of unique nuclear fragments. By default, ArchR uses the  
751 reads in TSS regions, because this value is computed upon the creation of an Arrow file and is  
752 stable across analyses, unlike the peak regions. Because fragments in Arrow files are split per  
753 chromosome, the low memory cost and high speed of this process enables interactive  
754 visualization of hundreds of thousands of cells in seconds. Additionally, ArchR can plot tracks  
755 without the genome browser using the ArchRBrowserTrack function. ArchR also enables direct  
756 export of group normalized bigwig files using “export.bw” from Rtracklayer that can be directly  
757 used in conventional genome browsers.

758

### 759 **ArchR Methods - Peak Matrix**

760 Once a peak set has been created (see ArchR Methods – Peak Calling), a cell x peak matrix can  
761 readily be made with ArchR. For each Arrow file, ArchR reads in scATAC-seq fragments from  
762 each chromosome and then computes overlaps with the peaks from the same chromosome. A

763 sparse matrix cell x peak matrix is created for these peaks. The matrix is then added to the  
764 corresponding Arrow file. This procedure is iterated across each chromosome.

765

### 766 **ArchR Methods – Creation of Low-Overlapping Aggregates of Cells for Linkage Analysis**

767 ArchR facilitates many integrative analyses that involve correlation of features. Performing these  
768 calculations with sparse single-cell data can lead to substantial noise in these correlative  
769 analyses. To circumvent this challenge, we adopted an approach introduced by Cicero<sup>18</sup> to create  
770 low-overlapping aggregates of single cells prior to these analyses. We filter aggregates with  
771 greater than 80% overlap with any other aggregate in order to reduce bias. To improve the speed  
772 of this approach, we developed an implementation of an optimized iterative overlap checking  
773 routine and a implementation of fast feature correlations in C++ using the “Rcpp” package. These  
774 optimized methods are used in ArchR for calculating peak co-accessibility, peak-to-gene linkage,  
775 and for other linkage analyses.

776

### 777 **ArchR Methods – Peak Co-Accessibility**

778 Co-accessibility analyses have been shown to be useful in downstream applications such as  
779 identifying groups of peaks that are all correlated forming “co-accessible networks”<sup>18</sup>. ArchR can  
780 rapidly compute peak co-accessibility from a peak matrix. These co-accessibility links can  
781 optionally be visualized using the ArchRBrowser. First, ArchR identifies 500+ low-overlapping cell  
782 aggregates (see Creation of Low-Overlapping Aggregates of Cells for Linkage Analysis). Second,  
783 for each chromosome (independently stored within an Arrow file), ArchR reads in the peak matrix  
784 and then creates the cell aggregate x peak matrix. ArchR next identifies all possible peak-to-peak  
785 combinations within a given window (by default 250 kb) and then computes the Pearson  
786 correlation of the log<sub>2</sub>-normalized cell aggregate x peak matrix. In this procedure, column sums  
787 across all chromosomes are used for depth normalization. ArchR iterates through all  
788 chromosomes and then combines the genome-wide results and stores them within the

789 ArchRProject. These can be readily accessed for downstream applications. Additionally, ArchR  
790 enables users to lower the resolution of these interactions to better visualize the main interactors  
791 (keeping the highest correlation value observed in each window).

792

### 793 **ArchR Methods – Motif Annotations**

794 ArchR enables rapid, fine-grained motif analyses. To carry out these analyses, ArchR must first  
795 identify the locations of all motifs in peak regions. ArchR natively supports access to motif sets  
796 curated from chromVAR<sup>27</sup> and JASPAR<sup>34</sup> to be used for these motif analyses. Additionally, ArchR  
797 makes possible the usage of multiple motif databases independently. ArchR first identifies motifs  
798 in peak regions using the matchMotifs function from the “motifmatchr” package  
799 (<https://greenleaflab.github.io/motifmatchr/>) with output being the motif positions within peaks.

800 ArchR then creates a boolean motif overlap sparse matrix for each motif-peak combination that  
801 can be used for downstream applications such as enrichment testing and chromVAR. The motif  
802 positions and motif overlap matrix are stored on disk as an RDS file for later access, which  
803 minimizes the total memory of the ArchRProject, freeing memory for other analyses.

804

### 805 **ArchR Methods – Feature Annotations**

806 ArchR allows for peak overlap analyses with defined feature sets. These feature sets could be  
807 ENCODE CHIP-seq/ATAC-seq peak sets or anything that can be specified as a GenomicRanges  
808 object. To facilitate this operation, we have curated a compendium of previously published ATAC-  
809 seq peak sets<sup>19,21–23,26</sup>, ENCODE CHIP-seq peak sets, and other custom feature sets for end-  
810 users<sup>35</sup>. We believe these custom feature sets will help users better annotate and describe cell  
811 types identified with scATAC-seq. These feature sets are overlapped with the ArchRProject peak  
812 set and then stored as a boolean feature overlap sparse matrix for each feature-peak combination  
813 that can be used for downstream applications such as enrichment testing and chromVAR. This

814 feature overlap matrix is then stored on disk as an RDS file for later access, which minimizes the  
815 total memory of the ArchRProject, freeing memory for other analyses.

816

### 817 **ArchR Methods – Marker Peak Identification with Annotation Enrichment**

818 ArchR allows for robust identification of features that are highly specific to a given group/cluster  
819 to elucidate cluster-specific biology. ArchR can identify these features for any of the matrices that  
820 are created with ArchR (stored in the Arrow files). ArchR identifies marker features while  
821 accounting for user-defined known biases that might confound the analysis (defaults are the TSS  
822 enrichment score and the number of unique nuclear fragments). For each group/cluster, ArchR  
823 identifies a set of background cells that match for the user-defined known biases and weights  
824 each equivalently using quantile normalization. Additionally, when selecting these bias-matched  
825 cells ArchR will match the distribution of the other user-defined groups. For example, if there were  
826 4 equally represented clusters, ArchR will match the biases for a cluster to the remaining 3  
827 clusters while selecting cells from the remaining 3 groups equally. By selecting a group of bias-  
828 matched cells, ArchR can directly minimize these confounding variables during differential testing  
829 rather than using modeling-based approaches. ArchR allows for binomial testing, Wilcoxon testing  
830 (via presto, <https://github.com/immunogenomics/presto/>), and two-sided t-testing for comparing  
831 the group to the bias-matched cells. These p-values are then adjusted for multiple hypothesis  
832 testing and organized across all group/clusters. This table of differential results can then be used  
833 to identify marker features based on user-defined log<sub>2</sub>(Fold Change) and FDR cutoffs.

834

### 835 **ArchR Methods – chromVAR Deviations Matrix**

836 ArchR facilitates chromVAR analysis to identify deviation of accessibility within peak annotations  
837 (i.e. motif overlaps) compared to a controlled background set of bias-matched peaks. A challenge  
838 in using the published version of the chromVAR software is that it requires the full cell x peak  
839 matrix to be loaded into memory in order to compute these deviations. This can lead to dramatic

840 increases in run time and memory usage for moderately sized datasets (~50,000 cells). To  
841 circumvent these limitations, ArchR implements the same chromVAR analysis workflow by  
842 analyzing sample sub-matrices independently (see **Supplementary Fig. 12e**). First, ArchR reads  
843 in the global accessibility per peak across all cells. Second, for each peak, ArchR identifies a set  
844 of background peaks that are matched by GC-content and accessibility. Third, ArchR uses this  
845 background set of peaks and global accessibility to compute bias-corrected deviations with  
846 chromVAR for each sample independently. This implementation requires data from only 5,000-  
847 10,000 cells to be loaded into memory at any given time, minimizing the memory requirements,  
848 enabling scalable analysis with chromVAR, and improving run-time performance.

849

#### 850 **ArchR Methods – Identification of Positive TF-Regulators**

851 ATAC-seq allows for the unbiased identification of TFs that exhibit large changes in chromatin  
852 accessibility at sites containing their DNA binding motifs. However, families of TFs (for ex. GATA  
853 factors) share similar features in their binding motifs when looking in aggregate through position  
854 weight matrices (PWMs). This motif similarity makes it challenging to identify the specific TFs that  
855 might be driving observed changes in chromatin accessibility at their predicted binding sites. To  
856 circumvent this challenge, we have previously used gene expression to identify TFs whose gene  
857 expression is positively correlated to changes in the accessibility of their corresponding motif<sup>21</sup>.  
858 We term these TFs “positive regulators”. However, this analysis relies on matched gene  
859 expression data which may not be readily available in all experiments. To overcome this  
860 dependency, ArchR can identify TFs whose inferred gene scores are correlated to their  
861 chromVAR TF deviation scores. To achieve this, ArchR correlates chromVAR deviation scores of  
862 TF motifs with gene activity scores of TF genes from the low-overlapping cell aggregates (see  
863 above). When using scRNA-seq integration with ArchR, gene expression of the TF can be used  
864 instead of inferred gene activity score.

865

## 866 **ArchR Methods – TF Footprinting**

867 ATAC-seq enables profiling of TF occupancy at base-pair resolution with TF footprinting. TF  
868 binding to DNA protects the protein-DNA binding site from transposition while the displacement  
869 or depletion of one or more adjacent nucleosomes creates increased DNA accessibility in the  
870 immediate flanking sequence. Collectively, these phenomena are referred to as the TF footprint.  
871 To accurately profile TF footprints, a large number of reads are required. Therefore, cells are  
872 grouped to create pseudo-bulk ATAC-seq profiles that can be then used for TF footprinting.

873 One major challenge with TF footprinting using ATAC-seq data is the insertion sequence  
874 bias of the Tn5 transposase<sup>21,36,37</sup> which can lead to misclassification of TF footprints. To account  
875 for Tn5 insertion bias ArchR identifies the k-mer (user-defined length, default length 6) sequences  
876 surrounding each Tn5 insertion site. To do this analysis, ArchR identifies single-base resolution  
877 Tn5 insertion sites for each pseudo-bulk (see above Sample-Aware Pseudo-Bulk Replicate  
878 Generation), resizes these 1-bp sites to k-bp windows ( $-k/2$  and  $+(k/2 - 1)$  bp from insertion), and  
879 then creates a k-mer frequency table using the “oligonucleotidefrequency(w=k,  
880 simplify.as=”collapse”)” function from the Biostrings package. ArchR then calculates the expected  
881 k-mers genome-wide using the same function with the BSgenome-associated genome file. To  
882 calculate the insertion bias for a pseudo-bulk footprint, ArchR creates a k-mer frequency matrix  
883 that is represented as all possible k-mers across a window  $\pm N$  bp (user-defined, default 250 bp)  
884 from the motif center. Then, iterating over each motif site, ArchR fills in the positioned k-mers into  
885 the k-mer frequency matrix. This is then calculated for each motif position genome-wide. Using  
886 the sample’s k-mer frequency table, ArchR can then compute the expected Tn5 insertions by  
887 multiplying the k-mer position frequency table by the observed/expected Tn5 k-mer frequency.  
888 For default TF footprinting with ArchR, motif positions (stored in the ArchRProject) are extended  
889  $\pm 250$  bp centered at the motif binding site. The pseudo-bulk replicates (stored as a HDF5-format  
890 coverage files) are then read into R as a coverage run-length encoding. For each individual motif,  
891 ArchR iterates over the chromosomes, computing a “Views” object using

892 “Views(coverage,positions)”. ArchR uses an optimized C++ function to compute the sum per  
893 position in the Views object. This implementation enables fast and efficient footprinting  
894 (**Supplementary Fig. 12g-i**).

895

### 896 **ArchR Methods – Bulk ATAC-seq LSI projection**

897 ArchR allows for projection of bulk ATAC-seq data into a scATAC-seq subspace as previously  
898 described<sup>7</sup>. ArchR first takes as input a bulk ATAC-seq sample x peak matrix and then identifies  
899 which peaks overlap the features used in the scATAC-seq dimensionality reduction. If there is  
900 sufficient overlap, ArchR estimates a scATAC-seq pseudo-cell x feature matrix within the features  
901 identified to overlap. These pseudo-cells (N = 250) per sample are sampled to be at 0.5x, 1x, 1.5x  
902 and 2x the average accessibility of the cell x feature matrix used. This step prevents unwanted  
903 sampling depth bias for this bulk projection analysis. The pseudo-cell x feature matrix is then  
904 normalized with the term-frequency x inverse document frequency (TF-IDF) method, using the  
905 same inverse document frequency obtained during the scATAC-seq dimensionality reduction.  
906 This normalized pseudo-cell x feature matrix is then projected with singular value decomposition  
907 “t(TF\_IDF) %\*% SVD\$u %\*% diag(1/SVD\$d)” where TF\_IDF is the transformed matrix and SVD  
908 is the previous SVD run using irlba in R. This reduced pseudo-cell x dim matrix can then be input  
909 to “uwot::umap\_transform” which uses the previous scATAC-seq UMAP embedding to project the  
910 pseudo-cells into this embedding.

911

### 912 **ArchR Methods – Data Imputation with MAGIC**

913 ArchR allows for using features such as gene scores and chromVAR deviation scores to assist in  
914 cluster annotation. However, features such as gene scores suffer from dropout noise in single-  
915 cell data. For scRNA-seq there have been many imputation methods developed to remedy this  
916 dropout noise. We have found that an effective method for imputation with scATAC-seq data is  
917 with Markov affinity-based graph imputation of cells (MAGIC)<sup>38</sup>. ArchR implements MAGIC for

918 diffusing single-cell features across similar cells to smooth a single-cell matrix while  
919 simultaneously accounting for drop-out biases. MAGIC creates and stores a cell x cell diffusion  
920 matrix of weights that is then used to smooth the feature matrix with matrix multiplication.  
921 However, this diffusion matrix is dense and scales quadratically with the number of cells. To  
922 circumvent this limitation, ArchR creates equally sized blocks of cells (user-defined, default is  
923 10,000) and then computes the partial diffusion matrix for these cells. These partial diffusion  
924 matrices are then combined to create a blocked diffusion matrix. This blocked diffusion matrix  
925 scales linearly in size leading to more memory efficiency but leads to lower resolution diffusion of  
926 data. To increase the resolution of this blocked diffusion matrix ArchR creates multiple replicates  
927 of the diffusion matrix to independently smooth the data matrix and then takes the average of the  
928 resulting smoothed matrices. ArchR additionally stores these blocked diffusion matrix replicates  
929 on-disk in HDF5-formatted files where each block is stored as its own group for direct access to  
930 specific parts of the matrix. ArchR's MAGIC implementation shifts the memory usage to on-disk  
931 storage and thus enables data diffusion of extremely large datasets ( $N > 200,000$ ) with minimal  
932 computing requirements.

933

### 934 **ArchR Methods – scATAC and scRNA Alignment**

935 ArchR allows for efficient integration with scRNA-seq data utilizing Seurat's integration  
936 infrastructure<sup>11</sup>. When performing this cross-platform alignment across large numbers of cells, we  
937 have found that the required memory and run time increase substantially. Moreover, constraining  
938 this alignment into smaller biologically relevant parts minimizes the alignment space into smaller  
939 alignment "sub-spaces"<sup>7</sup>. Thus, to increase alignment accuracy and improve runtime  
940 performance, ArchR enables the alignment of scATAC-seq and scRNA-seq to be constrained by  
941 user-defined groups of cells from both datasets that define smaller alignment sub-spaces. Within  
942 these sub-spaces, ArchR splits the scATAC-seq cells into equivalent slices of N cells (user-  
943 defined, default is 10,000 cells) and performs alignment with the scRNA-seq cells. This alignment



944 procedure begins with the identification of the top variable genes (user-defined, default is 2,000  
945 genes defined from scRNA-seq) using “Seurat::FindVariableFeatures”. Next, ArchR reads in the  
946 cell x gene scores matrix from the Arrow file for these cells. Then, ArchR imputes these gene  
947 scores using MAGIC and stores this imputed gene score matrix into a Seurat object for integration.  
948 ArchR then uses “Seurat::FindTransferAnchors” with canonical correlation analysis (CCA) to align  
949 this sub-space of cells efficiently. Next, ArchR extracts the aligned scRNA-seq cell, group, and  
950 gene expression profile with “Seurat::TransferData”. These gene expression profiles are stored  
951 in the corresponding Arrow files (stored as “GeneIntegrationMatrix”) for downstream analyses.

952

### 953 **ArchR Methods – scRNA Peak-To-Gene Linkage**

954 We have previously used ATAC-seq peak-to-gene linkages to link putative enhancers and GWAS  
955 risk loci to their predicted target genes<sup>7,21</sup>. ArchR can rapidly compute peak-to-gene links from a  
956 peak matrix and gene expression matrix (see above). These peak-to-gene links can optionally be  
957 visualized using the ArchRBrowser. First, ArchR identifies 500+ low-overlapping cell aggregates  
958 (see Creation of Low-Overlapping Aggregates of Cells for Linkage Analysis). Second, ArchR  
959 reads in the peak matrix and then creates the cell aggregate x peak matrix. Third, ArchR reads in  
960 the gene expression matrix and then creates the cell aggregate x gene matrix. ArchR then  
961 identifies all possible peak-to-gene combinations within a given window of the gene start (user-  
962 defined, default is 250 kb) and then computes the Pearson correlation of the log2-normalized cell  
963 aggregate x peak matrix and cell aggregate x gene matrix across all cell aggregates. ArchR  
964 computes these peak-to-gene links genome-wide and stores them within the ArchRProject, which  
965 can then be accessed for downstream applications. Additionally, ArchR enables users to lower  
966 the resolution of these interactions to better visualize the main interactors (keeping only the  
967 highest correlation value observed in each window).

968

### 969 **ArchR Methods – Cellular Trajectory Analysis**

970 To order cells in pseudo-time, ArchR creates cellular trajectories that order cells across a lower  
971 N-dimensional subspace within an ArchRProject. Previously, we have performed this ordering in  
972 the 2-dimensional UMAP subspace<sup>6</sup> but ArchR has improved upon this methodology to enable  
973 alignment within an N-dimensional subspace (i.e. LSI). First, ArchR requires a user-defined  
974 trajectory backbone that provides a rough ordering of cell groups/clusters. For example, given  
975 user-determined cluster identities, one might provide the cluster IDs for a stem cell cluster, then  
976 a progenitor cell cluster, and then a differentiated cell cluster that correspond to a known or  
977 presumed biologically relevant cellular trajectory (i.e. providing the cluster IDs for HSC, to MPP,  
978 to CMP, to Monocyte). Next, for each cluster, ArchR calculates the mean coordinates for each  
979 cell group/cluster in N-dimensions and retains cells whose Euclidean distance to those mean  
980 coordinates is in the top 5% of all cells. Next, ArchR computes the distance for each cell from  
981 cluster<sub>i</sub> to the mean coordinates of cluster<sub>i+1</sub> along the trajectory and computes a pseudo-time  
982 vector based on these distances for each iteration of i. This allows ArchR to determine an N-  
983 dimensional coordinate and a pseudo-time value for each of the cells retained as part of the  
984 trajectory based on the Euclidean distance to the cell group/cluster mean coordinates. Next,  
985 ArchR fits a continuous trajectory to each N-dimensional coordinate based on the pseudo-time  
986 value using the “smooth.spline” function with df = 250 (degrees of freedom) and spar = 1  
987 (smoothing parameter). Then, ArchR aligns all cells to the trajectory based on their Euclidean  
988 distance to the nearest point along the manifold. ArchR then scales this alignment to 100 and  
989 stores this pseudo-time in the ArchRProject for downstream analyses.

990 ArchR can create matrices that convey pseudo-time trends across features stored within  
991 the Arrow files. For example, ArchR can analyze changes in TF deviations, gene scores, or  
992 integrated gene expression across pseudo-time to identify regulators or regulatory elements that  
993 are dynamic throughout the cellular trajectory. First, ArchR groups cells in small user-defined  
994 quantile increments (default = 1/100) across the cellular trajectory. ArchR then smooths this matrix  
995 per feature using a user-defined smoothing window (default = 9/100) using the

996 “data.table::frollmean” function. ArchR then returns this smoothed pseudo-time x feature matrix  
997 as a SummarizedExperiment for downstream analyses. ArchR additionally can correlate two of  
998 these smoothed pseudo-time x feature matrices using name matching (i.e. positive regulators  
999 with chromVAR TF deviations and gene score/integration profiles) or by genomic position overlap  
1000 methods (i.e. peak-to-gene linkages) using low-overlapping cellular aggregates as described in  
1001 previous sections. Thus, ArchR facilitates integrative analyses across cellular trajectories,  
1002 revealing correlated regulatory dynamics across multi-modal data.

1003

#### 1004 **Benchmarking Analysis – Preface**

1005 For benchmarking analyses, we used one of two computational environments: (1) a MacBook Pro  
1006 laptop containing 32 GB of RAM and a 2.3GHz 8-core Intel Core i9 processor (16 threads) with  
1007 data stored on an external USB hard drive; (2) a large-memory node on a high-performance  
1008 cluster with 128 GB of RAM and two 2.40 GHz 10-core Intel Xeon E5-2640 V4 processors (20  
1009 threads). For benchmarking analyses using more limited compute resources (32 GB and 8 cores)  
1010 we used the same large-memory node configuration but limited the available cores and memories  
1011 using Slurm job submission properties. The main difference between the computational  
1012 environment of the MacBook Pro and the server is the ability of each core on the MacBook Pro  
1013 to use 2 threads whereas hyper-threading is disabled on the server and each core is effectively a  
1014 single thread.

1015 We downloaded scATAC-seq data from previously published and publicly available  
1016 locations. We downloaded the immune cell data fragment files from Satpathy et al. 2019  
1017 (GSE129785), Granja et al. 2019 (GSE139369), and from the 10x Genomics website  
1018 (<https://www.10xgenomics.com/solutions/single-cell-atac/>). For the mouse sci-ATAC-seq data,  
1019 we downloaded the BAM files from <http://atlas.gs.washington.edu/mouse-atac/>. No additional  
1020 steps were used prior to benchmarking analysis. We chose to focus our benchmarking tests  
1021 versus Signac and SnapATAC based on the performance of LSI and LDM shown previously<sup>9</sup>. We

1022 ran all analyses in triplicate using snakemake via a slurm job submission engine on a high-  
1023 performance cluster to accurately limit the available memory and cores. In the case of job failure,  
1024 we allowed for multiple job attempts to ensure that analyses were reproducible. After each failed  
1025 job attempt, the number of parallel threads for each software was lowered to attempt to complete  
1026 the analysis without exceeding the available memory. Unless otherwise stated all analyses were  
1027 run with default parameters for scATAC-seq benchmarking. We provide R markdown html files  
1028 on our publication page [https://github.com/GreenleafLab/ArchR\\_2020](https://github.com/GreenleafLab/ArchR_2020) detailing the exact  
1029 procedures used for all benchmarking analyses.

1030

### 1031 **Benchmarking Analysis – Signac**

1032 Signac (<https://github.com/timoast/signac>) requires a predetermined peak set, thus we  
1033 downloaded the previously published bulk hematopoiesis peak set from Corces et al.  
1034 ([ftp://ftp.ncbi.nlm.nih.gov/geo/series/GSE74nnn/GSE74912/suppl/GSE74912\\_ATACseq\\_All\\_Co](ftp://ftp.ncbi.nlm.nih.gov/geo/series/GSE74nnn/GSE74912/suppl/GSE74912_ATACseq_All_Co)  
1035 [unts.txt.gz](ftp://ftp.ncbi.nlm.nih.gov/geo/series/GSE74nnn/GSE74912/suppl/GSE74912_ATACseq_All_Co)) for all analyses. We first determined which cellular barcodes had more than 1,000  
1036 fragments by using “data.table::fread”. For each individual sample, we created a cell x peak matrix  
1037 with the “FeatureMatrix” function using the fragment files and abundant cell barcodes as input.  
1038 Then, we created a Seurat object from this cell x peak matrix with “CreateSeuratObject”. We then  
1039 determined TSS enrichment scores for each cell across the first 3 chromosomes with the  
1040 “TSSEnrichment” function. Default behavior for the TSSEnrichment function uses the first 2,000  
1041 TSSs; however, we increased this number (to include all TSSs on chr1-3) in order to stabilize the  
1042 TSS enrichment scores for more consistent high-quality cell determination while still minimizing  
1043 the run time. We then kept cells with a TSS enrichment score greater than 2 as high-quality cells  
1044 passing filter. This TSS score cutoff differs from that of ArchR due to differences in the formula  
1045 used for calculating TSS enrichment scores and differences in the gene annotation reference  
1046 used by Signac. We then merged these individual Seurat objects (corresponding to each sample)  
1047 and then performed TF-IDF normalization with “RunTFIDF” and “RunSVD” for LSI dimensionality

1048 reduction. We used the top 25% of features (ranked by accessibility) for LSI to reduce memory  
1049 usage. The first 30 components were used by default for downstream analyses. Clusters were  
1050 identified using “FindClusters” with default parameters. The scATAC-seq embeddings were  
1051 determined using “RunUMAP” for UMAP and “RunTSNE” for tSNE respectively. Lastly, the gene  
1052 score matrix was created using “FeatureMatrix” on the gene start and end coordinates (provided  
1053 from ArchR) extended upstream by 2 kb for each sample and combined afterwards followed by  
1054 log-normalization.

1055

### 1056 **Benchmarking Analysis - SnapATAC**

1057 SnapATAC (<https://github.com/r3fang/SnapATAC>) requires additional preprocessing steps prior  
1058 to creation of a Snap file that can be used for downstream analyses. First, fragment files were  
1059 sorted by their cell barcode with Unix “sort”. Next, these sorted fragment files were converted to  
1060 Snap files by using SnapTools “snap-pre” with parameters “--min-mapq=30 --min-flen=50 --max-  
1061 flen=1000 --keep-chrm=FALSE --keep-single=FALSE --keep-secondary=FALSE --  
1062 overwrite=TRUE --min-cov=1000 --max-num=20000 --verbose=TRUE” as described on the  
1063 GitHub page. A genome-wide tile/bin matrix was then added using “snap-add-bmat” with  
1064 parameters “--bin-size-list 5000” for a 5-kb matrix. To identify high-quality cells, SnapATAC  
1065 computes a promoter ratio score for the fraction of accessible fragments that overlap promoter  
1066 regions. We read in the 5-kb bin matrix into a Snap object using “addBmatToSnap” and then  
1067 created a promoter Genomic Ranges object from the provided transcript annotation file  
1068 (<http://renlab.sdsc.edu/r3fang/share/github/reference/hg19/gencode.v30.annotation.gtf.gz>) and  
1069 then extending the gene start upstream by 2 kb. Next, we overlapped these regions using  
1070 “findOverlaps” and then computed the summed accessibility within these overlapping regions vs  
1071 the total accessibility across all 5-kb bins. We chose a cutoff for promoter ratio as 0.175 by  
1072 manually inspecting the benchmarking dataset total accessibility vs promoter ratio plot as  
1073 described in the GitHub. These high-quality cells were kept for downstream analyses. For

1074 dimensionality reduction, we first filtered bins that were greater than the 95<sup>th</sup> percentile of non-  
1075 zero bins. Next, we ran “runDiffusionMaps” with 30 eigenvectors to be computed (similar in the  
1076 benchmarking analysis of all 3 methods). Clustering was performed with “runKNN” with the first  
1077 20 eigenvectors for a k-nn nearest neighbor search followed by “runCluster” with `louvain.lib = “R-  
1078 igraph”`. The scATAC-seq embeddings were determined using “runViz” with `method = “umap”` for  
1079 UMAP and `method = “Rtsne”` for tSNE for the top 20 eigenvectors. Lastly, the gene score matrix  
1080 was determined by using the gene start and end coordinates (provided from ArchR) as input to  
1081 “createGmatFromMat” with `input.mat = “bmat”` and scaled with “scaleCountMatrix”. For  
1082 comparing estimated dimensionality reduction in SnapATAC (estimated LDM) to estimated LSI in  
1083 ArchR, we first sampled N cells (10,000 or the number of cells specified) based on the inverse of  
1084 their coverage and then computed diffusion maps with “runDiffusionMaps”. The remaining cells  
1085 were projected with “runDiffusionMapsExtension” and the two Snap objects were combined for  
1086 downstream analysis.

1087

### 1088 **Benchmarking Analysis - ArchR**

1089 For analysis with ArchR, we first converted input scATAC-seq data (fragment files or BAM files)  
1090 to Arrow files with “createArrowFiles” with `minFrag = 1000`, `filterTSS = 4`, and `addGeneScoreMat`  
1091 `= FALSE` (`addGeneScoreMat` was set to false to allow for downstream benchmarking of this  
1092 individual step). These Arrow files were then used to create an ArchRProject with the appropriate  
1093 genome annotation. We identified doublet scores for each sample with “addDoubletScores” and  
1094 “filterDoublets” respectively; however, time and memory used for doublet identification were not  
1095 included in the benchmarking results because this step is unique to ArchR and would complicate  
1096 direct comparisons to other software. We then computed the iterative LSI dimensionality reduction  
1097 with “addIterativeLSI” with default parameters (`variableFeatures = 25,000` and `iterations = 2`).  
1098 Clusters were identified using “addClusters” with default parameters. The scATAC-seq  
1099 embeddings were determined using “addUMAP” for UMAP and “addTSNE” for tSNE. Lastly, the

1100 gene score matrix was added by “addGeneScoreMatrix” which stores the depth-normalized cell  
1101 x gene matrix. For comparison of estimated LSI in ArchR to estimated LDM in SnapATAC,  
1102 “addIterativeLSI” was run with an additional parameter for sampling (sampleCellsFinal = 10,000  
1103 or the number of cells specified).

1104

### 1105 **ArchR Analysis – Comparison of Gene Score Methods**

1106 We used ArchR to benchmark 53 models of inferring gene scores to emulate gene expression.  
1107 All models were tested with the same gene annotation reference for direct comparison. We  
1108 additionally used Signac, SnapATAC, and co-accessibility to create gene score models for  
1109 comparison, making a total of 56 models. We used two datasets for evaluation: (1) ~30,000  
1110 PBMCs and (2) ~30,000 bone marrow cells. We first created the gene score models that  
1111 incorporated distance by systematically changing the input parameters for  
1112 “addGeneScoreMatrix”. This parameter sweep included TSS exponential decay functions  
1113 (useTSS = TRUE) and gene body exponential decay functions (useTSS = FALSE). We tried other  
1114 decay functions but saw no appreciable difference so we used exponential decay (this is a user-  
1115 input so any model as a function of relative distance may be inserted). For gene score models  
1116 that were overlap-based (no distance function), we used “addFeatureMatrix” based on a set of  
1117 genomic regions corresponding to either an extended gene promoter [resize(genes, 1, “start”)  
1118 followed by resize(2\*window + 1, “center”)] or an extended gene body [extendGR(genes,  
1119 upstream, downstream)]. For each model, we created a genome-wide gene score matrix and  
1120 extracted these matrices from the Arrow files using “getMatrixFromProject”. We next created 500  
1121 low-overlapping random groupings of 100 cells with ArchR (see above) and took the average  
1122 gene scores for each of these groupings. Next, we collected the gene scores calculated by Signac  
1123 and SnapATAC during our benchmarking tests and averaged the gene scores across the same  
1124 groupings. For co-accessibility, we created gene scores as previously described with Cicero<sup>6,7,18</sup>.  
1125 We first used Cicero to create 5,000 lowly-overlapping cell groupings of 50 cells with “cicero\_cds”.



1126 Next, we calculated the average accessibility for these groupings across all peaks (with  
1127 `getMatrixFromProject`). We correlated all peaks within 250 kb to get peak co-accessibility. We  
1128 annotated the peaks as promoter if within 2.5 kb from the gene start with `“annotate_cds_by_site”`.  
1129 Finally, gene scores for the co-accessibility model were identified with  
1130 `“build_gene_activity_matrix”` with a co-accessibility cutoff of 0.35 followed by  
1131 `“normalize_gene_activities”`. For this co-accessibility model, we tested various parameters such  
1132 as promoter window size, correlation cutoff, and peak-to-peak distance maximums to make sure  
1133 the results were reproducible.

1134 Having a cell aggregate x gene score matrix for all 56 models, we next created a gene  
1135 expression matrix to test these models. We integrated our scATAC-seq (from ArchR’s results)  
1136 with previously annotated scRNA-seq datasets (10k PBMC from 10x website and Bone Marrow  
1137 from Granja et al., 2019) using `“Seurat::FindTransferAnchors”` and `“Seurat::TransferData”` with  
1138 the top 2,000 variable genes from scRNA-seq. This integration was performed for each scATAC-  
1139 seq sample independently and the scRNA-seq data used for each bone marrow alignment was  
1140 constrained to match cell sources together (i.e. BMMC scATAC-seq with BMMC scRNA-seq and  
1141 CD34+ scATAC-seq with CD34+ scRNA-seq)<sup>7</sup>. From this integration, each scATAC-seq cell was  
1142 paired to a matched gene expression profile. We averaged the gene expression profiles for each  
1143 of the 500 lowly-overlapping groups (see above) to create a cell aggregate x gene expression  
1144 matrix.

1145 To benchmark the performance for each gene score model, we identified 2 gene sets: the  
1146 top 2,000 variable genes defined by `“Seurat::FindVariableGenes”` and the top 1,000 differentially  
1147 expressed genes defined by `“Seurat::FindAllMarkers”` (ranking the top N genes for each scRNA-  
1148 seq cluster until 1,000 genes were identified). For these gene sets, we calculated the gene-wise  
1149 correlation (how well do the gene score and gene expression correlate across all genes) and the  
1150 aggregate-wise correlation (how well do the gene score and gene expression correlate across all

1151 cell aggregates). These 4 measures were then ranked across all models, and the average ranking  
1152 was used to score the 56 models.

1153 To orthogonally support this result, we downloaded previously published paired bulk  
1154 ATAC-seq + RNA-seq for hematopoiesis<sup>19</sup>. We then iteratively down-sampled the reads from  
1155 each dataset to create 100 pseudo-cells with 10,000 fragments from each bulk ATAC-seq sample.  
1156 We then created a scATAC-seq fragments file for each pseudo-cell. We performed an identical  
1157 analysis as described above for the 53 ArchR gene score models. For comparing these 53  
1158 models, we used 2 gene sets: the top 2,000 variable genes defined by log<sub>2</sub>-normalized  
1159 expression-ranked variance across each cell type and the top 1,000 marker genes defined by the  
1160 top log<sub>2</sub>(fold change) for each cell type vs the average expression of all cell types. We similarly  
1161 ranked the gene-wise and aggregate-wise correlation across all models, and used the average  
1162 ranking to score each model.

1163

#### 1164 **ArchR Analysis – Large Simulated PBMC ~1.2M Cells**

1165 To further test ArchR's capability to analyze extremely large datasets (N > 200,000), we simulated  
1166 ~1.3M single cells contained within 200 fragment files. We used 4 PBMC samples (2 x 5,000 cells  
1167 and 2 x 10,000 cells from 10x Genomics) for creating this large dataset. We randomly shifted  
1168 each scATAC-seq fragment with a mean difference of +/- 50-100 bp (randomly sampled) and a  
1169 standard deviation of +/- 10-20 bp (randomly sampled). We then sampled the fragments by 80%  
1170 to ensure some differences between simulated cells and then saved these to bg-zipped fragment  
1171 files. We then used ArchR to convert these fragment files to Arrow files with "createArrowFiles"  
1172 with minFrag = 1000, filterTSS = 4 and addGeneScoreMat = TRUE. These Arrow files were then  
1173 assembled into an ArchRProject. We identified doublet scores for each simulated dataset with  
1174 "addDoubletScores" and "filterDoublets" respectively, retaining ~1.2 million cells after doublet  
1175 removal. We then computed the estimated iterative LSI dimensionality reduction with  
1176 "addIterativeLSI" (variableFeatures = 25,000, sampleCellsFinal = 25,000 and 2 iterations).

1177 Estimated clusters were identified using “addClusters” with sampleCells = 50,000. This estimation  
1178 method uses a subset of cells to cluster and then the remaining cells are annotated by their  
1179 nearest neighbors (the maximum annotation observed). An estimated scATAC-seq UMAP was  
1180 created using “addUMAP” with sampleCells = 100,000. This estimation method uses a subset of  
1181 cells to create a UMAP embedding and then the remaining cells are projected into the single-cell  
1182 embedding using “umap::umap\_transform”.

1183

### 1184 **ArchR Analysis – Large Hematopoiesis 220K Cells**

1185 We wanted to test ArchR’s full analysis suite with a large dataset (N > 200,000) comprised of  
1186 previously published immune cell data<sup>6,7</sup>. We additionally grouped all Fluidigm C1-based scATAC-  
1187 seq data from Buenrostro et al. 2018<sup>4</sup> into a fragment file. This amounted to a total of 49 scATAC-  
1188 seq fragment files corresponding to over 200,000 cells. We first used ArchR to convert these  
1189 fragment files to Arrow files using “createArrowFiles” with minFrag = 1000, filterTSS = 8 and  
1190 addGeneScoreMat = TRUE. These Arrow files are then used to create an ArchRProject. We  
1191 identified doublet scores for each simulated dataset with “addDoubletScores” and “filterDoublets”  
1192 respectively. We then computed the estimated iterative LSI dimensionality reduction with  
1193 “addIterativeLSI” (variableFeatures = 25,000, sampleCellsFinal = 25,000 and iterations = 2). A  
1194 scATAC-seq UMAP was then created by using “addUMAP” with minDist = 1 and nNeighbors =  
1195 40. Clusters were initially identified using “addClusters” with default parameters. We re-clustered  
1196 the early progenitor cells (clusters containing CD34+ cells) with a clustering resolution of 0.4 to  
1197 better resolve these cell clusters. We added MAGIC imputation weights with  
1198 “addImputationWeights” for imputing single-cell features that are then overlaid on the UMAP  
1199 embedding. We then manually merged and assigned clusters that correspond to cell types based  
1200 on known marker gene scores and observation of sequencing tracks using the ArchRBrowser.

1201 To identify a union peak set, we created group coverage files, which contain the  
1202 aggregated accessibility of groups of single cells within a cluster, with “addGroupCoverages”. We

1203 then created a reproducible peak set with “addReproduciblePeakSet” and a cell x peak matrix  
1204 with “addPeakMatrix”. Next, we determined background peaks that are matched in GC-content  
1205 and accessibility with “addBgdPeak”. For downstream motif-based analyses we added motif  
1206 overlap annotations with “addMotifAnnotations” for CIS-BP version 1 motifs (version = 1). We  
1207 computed a ChromVAR deviations matrix with “addDeviationsMatrix”. We next identified positive  
1208 TF regulators with “correlateMatrices” where useMatrix1 = “MotifMatrix” and useMatrix2 =  
1209 “GeneScoreMatrix”. To identify which of these correlated TF regulators had strong differential  
1210 motif activity differences we calculated the average motif deviation scores with “exportGroupSE”  
1211 for each cluster and computed the max observed deviation difference between any two clusters.  
1212 This motif difference and the TF-to-gene score correlation were then used to identify positive  
1213 regulators (correlation > 0.5 and a maximum deviation score difference > 50<sup>th</sup> percentile).  
1214 Differential accessibility for each cluster was determined using “markerFeatures” with maxCells =  
1215 1000 and useMatrix = “PeakMatrix”. Marker peaks were defined as peaks with a log<sub>2</sub>(Fold  
1216 Change) > 1.5 and an FDR < 0.01 (Wilcoxon-test with presto,  
1217 <https://github.com/immunogenomics/presto/>). We then determined enriched motifs with  
1218 “peakAnnoEnrichment” in these marker peaks and plotted the motif enrichment p-values for the  
1219 positive TF regulators. ArchR has a curated set of previously published bulk ATAC-seq datasets  
1220 that we used for feature overlap enrichment by computing overlaps with “addArchRAnnotations”  
1221 (collection = “ATAC”) and “peakAnnoEnrichment”. TF footprints, with Tn5-bias correction, were  
1222 calculated by “plotFootprints” with motif positions from “getPosition” and normMethod = subtract.  
1223 Bulk hematopoietic ATAC-seq (GSE74912) was projected into the scATAC-seq subspace using  
1224 “projectBulkATAC” with N = 250 cells. Peak co-accessibility was computed with  
1225 “addCoAccessibility” and accessibility tracks were created with the ArchRBrowser.

1226 We next wanted to integrate our scATAC-seq data with previously published  
1227 hematopoietic scRNA-seq data<sup>7</sup>. To do this analysis, we used “addGeneIntegrationMatrix” with  
1228 sampleCellsATAC = 10,000, sampleCellsRNA = 10,000, and a groupList specifying to group cells

1229 from T/NK clusters and cells from non-T/NK clusters for both scATAC-seq and scRNA-seq prior  
1230 to alignment. This constrained integration improved the alignment accuracy and added a matched  
1231 gene expression profile for each scATAC-seq cell. We overlaid these gene expression profiles on  
1232 the UMAP embedding with “plotEmbedding”. After this integration analysis, we identified peak-to-  
1233 gene links with “addPeak2GeneLinks” and visualized them with “peak2GeneHeatmap”.

1234 To create a cellular trajectory across B cell differentiation, we used “addTrajectory” with  
1235 preFilterQuantile = 0.8, useAll = FALSE, and an initial trajectory of “HSC -> CMP.LMPP -> CLP.1  
1236 -> CLP.2 -> PreB -> B”. We next created trajectory matrices for “MotifMatrix”, “GeneScoreMatrix”,  
1237 “GeneIntegrationMatrix” and “PeakMatrix”. We correlated the deviation score and gene score  
1238 trajectory matrices with “correlateTrajectories”. Additionally, we correlated the deviation score and  
1239 gene expression trajectory matrices with “correlateTrajectories”. We kept TFs whose correlation  
1240 was 0.5 or greater for both of the correlation analyses. We determined these TFs as positive TF  
1241 regulators across the B cell trajectory. We also used ArchR to identify peak-to-gene links across  
1242 the B cell trajectory with “correlateTrajectories” with useRanges = TRUE, varCutOff1 = 0.9, and  
1243 varCutOff2 = 0.9. Lastly, we grouped cells into 5 groups of cells based on pseudo-time across the  
1244 B cell trajectory for track visualization (with the ArchRBrowser) and TF footprinting of the TF  
1245 regulators.

1246

1247

1248

1249

1250

1251

1252

1253

1254

1255 **ACKNOWLEDGEMENTS**

1256 We thank members of the Greenleaf and Chang laboratories for helpful comments. This work was  
1257 supported by NIH RM1-HG007735 and UM1-HG009442 (to H.Y.C. and W.J.G.), R35-CA209919  
1258 (to H.Y.C.), UM1-HG009436 and U19-AI057266 (to W.J.G.), K99-AG059918 and the American  
1259 Society of Hematology Scholar Award (to M.R.C.), and an International Collaborative Award (to  
1260 H.Y.C., H.C.).

1261

1262 **AUTHOR CONTRIBUTIONS**

1263 J.M.G., M.R.C., H.Y.C. and W.J.G conceived the project. J.M.G. and M.R.C. led the design of the  
1264 ArchR software with input from S.E.P and W.J.G. M.R.C. led the scATAC-seq data creation with  
1265 input from S.T.B., H.C. and H.Y.C.. J.M.G. and M.R.C. led the single-cell analysis presented in  
1266 this paper. J.M.G., M.R.C., H.Y.C. and W.J.G wrote the manuscript with input from all authors.

1267

1268 **DECLARATION OF INTERESTS**

1269 W.J.G. and H.Y.C. are consultants for 10x Genomics who has licensed IP associated with ATAC-  
1270 seq. W.J.G. has additional affiliations with Guardant Health (consultant) and Protillion Biosciences  
1271 (co-founder and consultant). H.Y.C. is a co-founder of Accent Therapeutics, Boundless Bio, and  
1272 a consultant for Arsenal Biosciences and Spring Discovery.

1273

1274

1275

1276

1277

1278

1279

1280

1281 **References**

- 1282 1. Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory  
1283 variation. *Nature* **523**, 486–490 (2015).
- 1284 2. Cusanovich, D. A. *et al.* Multiplex single cell profiling of chromatin accessibility by  
1285 combinatorial cellular indexing. *Science* **348**, 910–914 (2015).
- 1286 3. Cusanovich, D. A. *et al.* The cis-regulatory dynamics of embryonic development at single-  
1287 cell resolution. *Nature* **555**, 538–542 (2018).
- 1288 4. Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous Regulatory  
1289 Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535–1548.e16 (2018).
- 1290 5. Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility.  
1291 *Cell* **174**, 1309–1324.e18 (2018).
- 1292 6. Satpathy, A. T. *et al.* Massively parallel single-cell chromatin landscapes of human immune  
1293 cell development and intratumoral T cell exhaustion. *Nat. Biotechnol.* **37**, 925–936 (2019).
- 1294 7. Granja, J. M. *et al.* Single-cell multiomic analysis identifies regulatory programs in mixed-  
1295 phenotype acute leukemia. *Nat. Biotechnol.* **37**, 1458–1465 (2019).
- 1296 8. Lareau, C. A. *et al.* Droplet-based combinatorial indexing for massive-scale single-cell  
1297 chromatin accessibility. *Nat. Biotechnol.* **37**, 916–924 (2019).
- 1298 9. Chen, H. *et al.* Assessment of computational methods for the analysis of single-cell ATAC-  
1299 seq data. *Genome Biol.* **20**, 241 (2019).
- 1300 10. Fang, R. *et al.* Fast and Accurate Clustering of Single Cell Epigenomes Reveals *Cis* -  
1301 Regulatory Elements in Rare Cell Types. *BioRxiv* (2019). doi:10.1101/615179
- 1302 11. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–1902.e21  
1303 (2019).
- 1304 12. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-  
1305 Cell RNA Sequencing Data Using Artificial Nearest Neighbors. *Cell Syst.* **8**, 329–337.e4  
1306 (2019).



- 1307 13. Wolock, S. L., Lopez, R. & Klein, A. M. Scrublet: Computational Identification of Cell  
1308 Doublets in Single-Cell Transcriptomic Data. *Cell Syst.* **8**, 281–291.e9 (2019).
- 1309 14. Kang, H. M. *et al.* Multiplexed droplet single-cell RNA-sequencing using natural genetic  
1310 variation. *Nat. Biotechnol.* **36**, 89–94 (2018).
- 1311 15. Thurman, R. E. *et al.* The accessible chromatin landscape of the human genome. *Nature*  
1312 **489**, 75–82 (2012).
- 1313 16. Andersson, R. & Sandelin, A. Determinants of enhancer and promoter activities of  
1314 regulatory elements. *Nat. Rev. Genet.* **21**, 71–87 (2020).
- 1315 17. Arnosti, D. N. Analysis and function of transcriptional regulatory elements: insights from  
1316 *Drosophila*. *Annu Rev Entomol* **48**, 579–602 (2003).
- 1317 18. Pliner, H. A. *et al.* Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell  
1318 Chromatin Accessibility Data. *Mol. Cell* **71**, 858–871.e8 (2018).
- 1319 19. Corces, M. R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human  
1320 hematopoiesis and leukemia evolution. *Nat. Genet.* **48**, 1193–1203 (2016).
- 1321 20. McInnes, L., Healy, J. & Melville, J. UMAP: Uniform Manifold Approximation and Projection  
1322 for Dimension Reduction. *arXiv* (2018).
- 1323 21. Corces, M. R. *et al.* The chromatin accessibility landscape of primary human cancers.  
1324 *Science* **362**, (2018).
- 1325 22. Calderon, D. *et al.* Landscape of stimulation-responsive chromatin across diverse human  
1326 immune cells. *Nat. Genet.* **51**, 1494–1505 (2019).
- 1327 23. Corces, M. R. *et al.* Single-cell epigenomic identification of inherited risk loci in Alzheimer’s  
1328 and Parkinson’s disease. *BioRxiv* (2020). doi:10.1101/2020.01.06.896159
- 1329 24. Mumbach, M. R. *et al.* Enhancer connectome in primary human cells identifies target genes  
1330 of disease-associated DNA elements. *Nat. Genet.* **49**, 1602–1612 (2017).
- 1331 25. Satpathy, A. T. *et al.* Transcript-indexed ATAC-seq for precision immune profiling. *Nat.*  
1332 *Med.* **24**, 580–590 (2018).

- 1333 26. Corces, M. R. *et al.* An improved ATAC-seq protocol reduces background and enables  
1334 interrogation of frozen tissues. *Nat. Methods* **14**, 959–962 (2017).
- 1335 27. Schep, A. N., Wu, B., Buenrostro, J. D. & Greenleaf, W. J. chromVAR: inferring  
1336 transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods*  
1337 **14**, 975–978 (2017).
- 1338 28. Regev, A. *et al.* The human cell atlas. *Elife* **6**, (2017).
- 1339 29. Huber, W. *et al.* Orchestrating high-throughput genomic analysis with Bioconductor. *Nat.*  
1340 *Methods* **12**, 115–121 (2015).
- 1341 30. Adey, A. *et al.* Rapid, low-input, low-bias construction of shotgun fragment libraries by high-  
1342 density in vitro transposition. *Genome Biol.* **11**, R119 (2010).
- 1343 31. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition  
1344 of native chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-  
1345 binding proteins and nucleosome position. *Nat. Methods* **10**, 1213–1218 (2013).
- 1346 32. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with  
1347 Harmony. *Nat. Methods* **16**, 1289–1296 (2019).
- 1348 33. Lun, A. T. L., McCarthy, D. J. & Marioni, J. C. A step-by-step workflow for low-level analysis  
1349 of single-cell RNA-seq data with Bioconductor. [version 2; peer review: 3 approved, 2  
1350 approved with reservations]. *F1000Res.* **5**, 2122 (2016).
- 1351 34. Fornes, O. *et al.* JASPAR 2020: update of the open-access database of transcription factor  
1352 binding profiles. *Nucleic Acids Res.* **48**, D87–D92 (2020).
- 1353 35. Sheffield, N. C. & Bock, C. LOLA: enrichment analysis for genomic region sets and  
1354 regulatory elements in R and Bioconductor. *Bioinformatics* **32**, 587–589 (2016).
- 1355 36. He, H. H. *et al.* Refined DNase-seq protocol and data analysis reveals intrinsic bias in  
1356 transcription factor footprint identification. *Nat. Methods* **11**, 73–78 (2014).
- 1357 37. Baek, S., Goldstein, I. & Hager, G. L. Bivariate genomic footprinting detects changes in  
1358 transcription factor activity. *Cell Rep.* **19**, 1710–1722 (2017).

1359 38. van Dijk, D. *et al.* Recovering Gene Interactions from Single-Cell Data Using Data Diffusion.  
1360 *Cell* **174**, 716–729.e27 (2018).

1361

1362

1363

1364

1365

1366

1367

1368

1369

1370

1371

1372

1373

1374

1375

1376

1377

1378

1379

1380

1381

1382

1383

1384

1385 **Figure Legends**

1386

1387 **Figure 1. ArchR: A rapid, extensible, and comprehensive scATAC-seq analysis platform.**

1388 **a.** Schematic of the ArchR workflow from input of pre-aligned scATAC-seq data as BAM or  
1389 fragment files to diverse data analysis.

1390 **b-c.** Comparison of run time and memory usage by ArchR, Signac, and SnapATAC for the  
1391 analysis of **(b)** ~20,000 PBMC cells using 32 GB RAM and 8 cores or **(c)** ~70,000 PBMC cells  
1392 using 128 GB RAM and 20 cores. Dots represent individual replicates of benchmarking analysis.

1393 **d.** Initial UMAP embedding of scATAC-seq data from 2 replicates of the cell line mixing experiment  
1394 (N = 38,072 total cells from 10 different cell lines) colored by replicate number.

1395 **e.** Schematic of doublet identification with ArchR.

1396 **f-g.** Initial UMAP embedding of scATAC-seq data from 2 replicates of the cell line mixing  
1397 experiment (N = 38,072 total cells from 10 different cell lines) colored by **(f)** the enrichment of  
1398 projected synthetic doublets or **(g)** the demuxlet identification labels based on genotype  
1399 identification using SNPs within accessible chromatin sites.

1400 **h.** Receiver operating characteristic (ROC) curves of doublet prediction using synthetic doublet  
1401 projection enrichment or the number of nuclear fragments per cell compared to demuxlet as a  
1402 ground truth. The area under the curve (AUC) for these ROC curves are annotated below.

1403 **i.** UMAP after ArchR doublet removal of scATAC-seq data from 2 replicates of the cell line mixing  
1404 experiment (N = 27,220 doublet-filtered cells from 10 different cell lines) colored by demuxlet  
1405 identification labels based on genotype identification using SNPs within accessible chromatin  
1406 sites.

1407

1408 **Figure 2. Optimized gene score inference models provide improved prediction of gene**  
1409 **expression from scATAC-seq data.**

1410 **a.** UMAPs of scATAC-seq data from (top) PBMCs and (bottom) bone marrow cells colored by  
1411 aligned scRNA-seq clusters. This alignment is used for benchmarking all downstream scATAC-  
1412 seq gene score models.

1413 **b.** Heatmaps summarizing the accuracy (measured by Pearson correlation) across 56 gene score  
1414 models for both the top 1,000 differentially expressed and top 2,000 variable genes for both PBMC  
1415 and bone marrow cell datasets. Each heatmap entry is colored by the model rank in the given  
1416 correlation test as described below the heatmap. The model class is indicated to the left of each  
1417 heatmap by color. SA, SnapATAC; SN, Signac; CoA, Co-accessibility.

1418 **c.** Illustration of the gene score Model 42, which uses bi-directional exponential decays from the  
1419 gene TSS (extended upstream by 5 kb) and the gene transcription termination site (TTS) while  
1420 accounting for neighboring gene boundaries (see methods). This model was shown to be more  
1421 accurate than other models such as Model 21 which models an exponential decay from the gene  
1422 TSS.

1423 **d.** Side-by-side UMAPs for PBMCs and bone marrow cells colored by (left) gene scores from  
1424 Model 42 and (right) gene expression from scRNA-seq alignment for key immune cell-related  
1425 marker genes.

1426 **e-f.** Heatmaps of (top) gene expression or (bottom) gene scores for the top 1,000 differentially  
1427 expressed genes (selected from scRNA-seq) across all cell aggregates for (e) PBMCs or (f) bone  
1428 marrow cells. Color bars to the left of each heatmap represent the PBMC or bone marrow cell  
1429 cluster derived from scRNA-seq data.

1430

1431 **Figure 3. ArchR enables comprehensive analysis of massive-scale scATAC-seq data.**

1432 **a.** Run times for ArchR-based analysis of over 220,000 and 1,200,000 single cells respectively  
1433 using a small cluster-based computational environment (32 GB RAM and 8 cores with HP Lustre  
1434 storage) and a personal MacBook Pro laptop (32 GB RAM and 8 cores with an external USB hard  
1435 drive). Color indicates the relevant analytical step.

1436 **b.** UMAP of the hematopoiesis dataset colored by the 21 hematopoietic clusters. UMAP was  
1437 constructed using LSI estimation with 25,000 landmark cells.

1438 **c.** Heatmap of 215,916 ATAC-seq marker peaks across all hematopoietic clusters identified with  
1439 bias-matched differential testing. Color indicates the column Z-score of normalized accessibility.

1440 **d.** Heatmap of motif hypergeometric enrichment adjusted p-values within the marker peaks of  
1441 each hematopoietic cluster. Color indicates the motif enrichment ( $-\log_{10}(\text{p-value})$ ) based on the  
1442 hypergeometric test.

1443 **e.** Side-by-side UMAPs of (left) gene scores and (right) motif deviation scores for ArchR-identified  
1444 TFs where the inferred gene expression is positively correlated with the chromVAR TF deviation  
1445 across hematopoiesis.

1446 **f-h.** Tn5 bias-adjusted transcription factor footprints for GATA, SPI1, and EOMES motifs,  
1447 representing positive TF regulators of hematopoiesis. Lines are colored by the 21 clusters shown  
1448 in **Figure 3c**.

1449 **i.** Genome accessibility track visualization of marker genes with peak co-accessibility. (Left) *CD34*  
1450 genome track (chr1:208,034,682-208,134,683) showing greater accessibility in earlier  
1451 hematopoietic clusters (1-5, 7-8 and 12-13). (Right) *CD14* genome track (chr5:139,963,285-  
1452 140,023,286) showing greater accessibility in earlier monocytic clusters (13-15).

1453

1454 **Figure 4. Integration of scATAC-seq and scRNA-seq data by ArchR identifies gene**  
1455 **regulatory trajectories of hematopoietic differentiation.**

1456 **a.** Schematic of scATAC-seq alignment with scRNA-seq data in M slices of N single cells. These  
1457 slices are independently aligned to a reference scRNA-seq dataset and then the results are  
1458 combined for downstream analysis. This integrative design facilitates rapid large-scale integration  
1459 with low-memory requirements.

1460 **b-d.** UMAP of scATAC-seq data from the hematopoiesis dataset colored by **(b)** alignment to  
1461 previously published hematopoietic scRNA-seq-derived clusters, **(c)** integrated scRNA-seq gene

1462 expression for key marker TFs and genes, or (d) cell alignment to the ArchR-defined B cell  
1463 trajectory. In (d), the smoothed arrow represents a visualization of the interpreted trajectory  
1464 (determined in the LSI subspace) in the UMAP embedding.

1465 e. Heatmap of 11,999 peak-to-gene links identified across the B cell trajectory with ArchR.

1466 f-g. Genome track visualization of the (f) *HMGA1* locus (chr6:34,179,577-34,249,577) and (g)  
1467 *BLK* locus (chr8:11,301,521-11,451,521). Single-cell gene expression across pseudo-time in the  
1468 B cell trajectory is shown to the right. Inferred peak-to-gene links for distal regulatory elements  
1469 across the hematopoiesis dataset is shown below.

1470 h. Heatmap of positive TF regulators whose gene expression is positively correlated with  
1471 chromVAR TF deviation across the B cell trajectory.

1472 i-k. Tn5 bias-adjusted transcription factor footprints for (i) NFE2, (j) EBF1, and (k) IRF8 motifs,  
1473 representing positive TF regulators across the B cell trajectory. Lines are colored by the position  
1474 in pseudo-time of B cell differentiation.

1475

## 1476 **Supplementary Figure Legends**

1477

### 1478 **Supplementary Fig. 1. ArchR infrastructure and supported analyses.**

1479 a. Comparison of supported scATAC-seq analysis features across ArchR, Signac and  
1480 SnapATAC.

1481 b. (Left) Schematic of the ArchR Arrow file format where accessible reads and arrays are  
1482 organized within. Arrow files can then be used as input for an ArchRProject (Right). The  
1483 ArchRProject stores the locations of these Arrow files and extracts their cell-centric metadata. All  
1484 analysis with ArchR operates through this ArchRProject which can readily access data from Arrow  
1485 files stored on disk.



1486 **c.** Schematic demonstrating how ArchR operations that involve using Arrow fragments (i.e.  
1487 addTileMatrix) operate on each chromosome independently in parallel for many Arrow files and  
1488 then add the resulting matrix back to the corresponding Arrow files again in parallel.

1489 **d.** Schematic demonstrating how ArchR operations that use Arrow matrices (i.e. addIterativeLSI)  
1490 access a subset of each chromosome's matrix from each Arrow file in parallel that are then  
1491 merged to create a filtered matrix for subsequent analysis.

1492

1493 **Supplementary Fig. 2.**

1494 **a-b.** File sizes of storage formats (for both accessible fragments and counts matrix) for ArchR and  
1495 SnapATAC compared to **(a)** the total number of cells they represent or **(b)** the total number of  
1496 fragments corresponding to the cells represented in each file. Line colors represent the different  
1497 software used or the original fragment files.

1498 **c.** QC filtering plots for the PBMCs dataset from (left) ArchR, showing the TSS enrichment score  
1499 vs unique nuclear fragments per cell, or (right) SnapATAC, showing the promoter ratio / fraction  
1500 of reads in promoters (FIP) vs unique nuclear fragments per cell. Dot color represents the density  
1501 in arbitrary units of points in the plot.

1502 **d-e.** Aggregate **(d)** TSS insertion profiles centered at all TSS regions or **(e)** fragment size  
1503 distributions for the cells passing ArchR QC thresholds for each sample in the PBMCs dataset.  
1504 Line color represents the sample from the dataset as indicated below the plot.

1505 **f.** QC filtering plots for the bone marrow cell dataset from (left) ArchR, showing the TSS  
1506 enrichment score vs unique nuclear fragments per cell, or (right) SnapATAC, showing the  
1507 promoter ratio / fraction of reads in promoters (FIP) vs unique nuclear fragments per cell. Dot  
1508 color represents the density in arbitrary units of points in the plot.

1509 **g-h.** Aggregate **(g)** TSS insertion profiles centered at all TSS regions or **(h)** fragment size  
1510 distributions for the cells passing ArchR QC thresholds for each sample in the bone marrow cell  
1511 dataset. Line color represents the sample from the dataset as indicated below the plot.

1512 i. QC filtering plots from ArchR for each individual organ type from the mouse atlas dataset  
1513 showing the TSS enrichment score vs unique nuclear fragments per cell. Dot color represents the  
1514 density in arbitrary units of points in the plot.

1515 j-k. Aggregate (j) TSS insertion profiles centered at all TSS regions or (k) fragment size  
1516 distributions for the cells passing ArchR QC thresholds for each sample in the mouse atlas  
1517 dataset. Line colors represent different samples as indicated to the left of the plot.

1518

### 1519 **Supplementary Fig. 3.**

1520 a. Schematic describing the individual benchmarking steps compared across ArchR, Signac, and  
1521 SnapATAC for (1) Data Import, (2) Dimensionality Reduction and Clustering, and (3) Gene Score  
1522 Matrix Creation.

1523 b-i. Comparison of ArchR, Signac, and SnapATAC for run time and peak memory usage for the  
1524 analysis of (b) ~20,000 cells from the PBMCs dataset using 128 GB of RAM and 20 cores (plot  
1525 corresponds to **Figure 1b**), (c) ~70,000 cells from the PBMCs dataset using 32 GB of RAM and  
1526 8 cores (plot corresponds to **Figure 1c**), (d-e) ~10,000 cells from the PBMCs dataset using (d)  
1527 32 GB of RAM and 8 cores or (e) 128 GB of RAM and 20 cores, (f-g) ~30,000 cells from the  
1528 PBMCs dataset using (f) 32 GB of RAM and 8 cores or (g) 128 GB of RAM and 20 cores, and (h-  
1529 i) ~30,000 cells from the bone marrow dataset using (h) 32 GB of RAM and 8 cores or (i) 128 GB  
1530 of RAM and 20 cores. Dots represent individual replicates of benchmarking analysis.

1531 j. Benchmarks from ArchR for run time and peak memory usage for the analysis of ~70,000 cells  
1532 from the sci-ATAC-seq mouse atlas dataset for (left) 32 GB of RAM with 8 cores and (right) 128  
1533 GB of RAM with 20 cores. Dots represent individual replicates of benchmarking analysis.

1534 k. t-SNE of mouse atlas scATAC-seq data (N = 64,286 cells) colored by individual samples.

1535

### 1536 **Supplementary Fig. 4.**

1537 **a.** QC filtering plots from ArchR for (top) replicate 1 and (bottom) replicate 2 from the cell line  
1538 mixing dataset showing the TSS enrichment score vs unique nuclear fragments per cell. Dot color  
1539 represents the density in arbitrary units of points in the plot.

1540 **b.** Accuracy of various doublet prediction methods for (top) replicate 1 and (bottom) replicate 2  
1541 from the cell line mixing dataset, measured by the area under the curve (AUC) of the receiver  
1542 operating characteristic (ROC), across different in silico cell loadings. Accuracy is determined with  
1543 respect to genotype-based identification of doublets using demuxlet. Above each plot, “KNN”  
1544 represents the number of cells nearby each projected synthetic doublet to record when calculating  
1545 doublet enrichment scores. The distance for KNN recording is determined in the LSI subspace  
1546 for LSI projection and in the UMAP embedding for UMAP projection parameters.

1547 **c-h.** UMAP of scATAC-seq data showing the (**c-d**) simulated doublet density, (**e-f**) simulated  
1548 doublet enrichment, or (**g-h**) cell line identity based on genotyping information and demuxlet for  
1549 (**c,e,g**) replicate 1 (N = 15,345 cells) and (**d,f,h**) replicate 2 (N = 22,727 cells) of the cell line mixing  
1550 dataset.

1551

#### 1552 **Supplementary Fig. 5.**

1553 **a.** Schematic of the iterative LSI procedure implemented in ArchR for dimensionality reduction.

1554 **b.** UMAPs of scATAC-seq data from ~30,000 cells from the PBMCs dataset to compare clustering  
1555 results across ArchR with doublet removal, ArchR without doublet removal, Signac, SnapATAC,  
1556 and SnapATAC with estimated LDM. Each UMAP is colored by (left) sample, (middle) clusters as  
1557 defined by ArchR with doublet removal, and (right) the number of unique nuclear fragments.

1558

#### 1559 **Supplementary Fig. 6.**

1560 **a.** UMAPs of scATAC-seq data from ~30,000 cells from the bone marrow dataset to compare  
1561 clustering results across ArchR with doublet removal, ArchR without doublet removal, Signac,  
1562 SnapATAC, and SnapATAC with estimated LDM. Each UMAP is colored by (left) sample, (middle)

1563 clusters as defined by ArchR with doublet removal, and (right) the number of unique nuclear  
1564 fragments.

1565

1566 **Supplementary Fig. 7.**

1567 **a.** Schematic of the estimated LSI framework implemented by ArchR. Briefly, a subset of cells,  
1568 referred to as “landmark” cells, are used for LSI dimensionality reduction. The remaining cells are  
1569 then linearly projected with LSI projection into this landmark-defined LSI subspace. This method  
1570 enables massive-scale analysis of scATAC-seq data with ArchR.

1571 **b.** UMAPs of scATAC-seq data from ~30,000 cells from the PBMCs dataset showing the results  
1572 of dimensionality reduction from (left) estimated LSI with ArchR after doublet removal or (right)  
1573 estimated LDM with SnapATAC. For each analytical case, a range of cell numbers is used for the  
1574 landmark cell subset (top to bottom). Within each analytical case, two UMAPs are presented,  
1575 colored by the clusters identified without estimation from (left) ArchR or (right) SnapATAC.

1576

1577 **Supplementary Fig. 8.**

1578 **a.** UMAPs of scATAC-seq data from ~30,000 cells from the bone marrow cell dataset showing  
1579 the results of dimensionality reduction from (left) estimated LSI with ArchR after doublet removal  
1580 or (right) estimated LDM with SnapATAC. For each analytical case, a range of cell numbers is  
1581 used for the landmark cell subset (top to bottom). Within each analytical case, two UMAPs are  
1582 presented, colored by the clusters identified without estimation from (left) ArchR or (right)  
1583 SnapATAC.

1584 **b.** Comparison of clustering fidelity based on adjusted Rand index in ArchR by estimated LSI or  
1585 in SnapATAC by estimated LDM across multiple landmark subset sizes.

1586 **c.** Benchmarking of run time for ArchR estimated LSI and SnapATAC estimated LDM for ~30,000  
1587 cells from (left) the PBMCs dataset and (right) the bone marrow cell dataset for (top) 128 GB of  
1588 RAM with 20 cores and (bottom) 32 GB of RAM with 8 cores.

1589

1590 **Supplementary Fig. 9.**

1591 **a-h.** Distribution of Pearson correlations of inferred gene score and aligned gene expression for  
1592 **(a,c,e,g)** each gene or **(b,d,f,h)** each cell group across groups of 100 cells (N = 500 groups).

1593 Distributions are either presented for **(a,b,e,f)** the top 1,000 differentially expressed genes or  
1594 **(c,d,g,h)** the top 2,000 most variable genes for each of the 56 gene score models tested. In each  
1595 plot, the red dotted line represents the median value of the best-performing model. Violin plots  
1596 represent the smoothed density of the distribution of the data. In box plots, the lower whisker is  
1597 the lowest value greater than the 25% quantile minus 1.5 times the interquartile range, the lower  
1598 hinge is the 25% quantile, the middle is the median, the upper hinge is the 75% quantile and the  
1599 upper whisker is the largest value less than the 75% quantile plus 1.5 times the interquartile range.

1600 SA, SnapATAC; SN, Signac; CoA, Co-accessibility.

1601 **i-j.** UMAPs of scATAC-seq data from **(i)** cells from the PBMCs dataset (N = 27,845 cells) or **(j)**  
1602 cells from the bone marrow cell dataset (N = 26,748 cells) colored by **(top)** inferred gene scores  
1603 or **(bottom)** gene expression for several marker genes.

1604 **k.** Schematic illustrating the methodology used to assess the accuracy of inferred gene scores  
1605 based on orthogonal matched bulk ATAC-seq and bulk RNA-seq data of various sorted  
1606 hematopoietic cell types.

1607 **l.** Heatmaps summarizing the accuracy (measured by Pearson correlation) across 56 gene score  
1608 models for both the top 1,000 differentially expressed and top 2,000 variable genes for bulk ATAC-  
1609 seq and RNA-seq data from sorted hematopoietic cell types. Each heatmap entry is colored by  
1610 the model rank in the given correlation test as described below the heatmap. The model class is  
1611 indicated to the left of each heatmap by color. SA, SnapATAC; SN, Signac; CoA, Co-accessibility.

1612 **m.** Heatmaps of **(left)** gene expression or **(right)** gene scores for the top 1,000 differentially  
1613 expressed genes (selected from bulk RNA-seq) across all cell types from the matched bulk ATAC-  
1614 seq and RNA-seq data.

1615

1616 **Supplementary Fig. 10.**

1617 **a.** Bar plot showing the number of cells passing ArchR QC thresholds from each of the immune  
1618 cell scATAC-seq datasets used for the ~220k cell hematopoiesis dataset.

1619 **b-c.** Aggregate **(b)** TSS insertion profiles centered at all TSS regions or **(c)** fragment size  
1620 distributions for the cells passing ArchR QC thresholds for each sample in the hematopoiesis  
1621 dataset. Line color represents the sample from the dataset as indicated in **Supplementary Figure**  
1622 **10a.**

1623 **d.** Summary of quality control information for each cell from the hematopoiesis dataset. The  
1624 distribution of (left) TSS enrichment scores, (middle) the number of unique nuclear fragments,  
1625 and (right) the fraction of reads in peak regions (FRiP) are shown for each single cell passing  
1626 filter.

1627 **e.** Benchmarking of peak memory usage for analysis of (top) the ~220,000 cells from the  
1628 hematopoiesis dataset and (bottom) ~1,200,000 simulated PBMCs using a computational  
1629 infrastructure with 32 GB of RAM and 8 cores with an HP Lustre file storage system.

1630 **f.** UMAPs of scATAC-seq data derived from estimated LSI of the hematopoiesis dataset using  
1631 different numbers of landmark cells. These UMAPs are colored by the clusters identified from the  
1632 25,000-cell estimated LSI shown in **Figure 3b.**

1633 **g-i.** UMAPs of scATAC-seq data as shown in **Figure 3b**, colored by **(g)** the different experimental  
1634 samples (as shown in **Supplementary Figure 10a**), **(h)** the number of unique nuclear fragments,  
1635 or **(i)** the per-cell TSS enrichment score.

1636

1637 **Supplementary Fig. 11.**

1638 **a.** Schematic for the generation of sample-aware pseudo-bulk replicates in ArchR for downstream  
1639 analyses. Briefly, for each cell grouping (in most cases identified by clusters), cells are split per  
1640 sample of origin. Next, for each cell grouping these sample-aware cell groups are tested for being

1641 larger than a specified minimum number of cells to create a specified minimum number of sample-  
1642 aware replicates. If these requirements are not met with a simple splitting, ArchR accounts for  
1643 each different case by using sub-sampling approaches (see methods).

1644 **b.** Schematic for iterative overlap peak merging in ArchR to identify non-overlapping fixed-width  
1645 peaks. Briefly, peaks (peak summits that are extended to yield fixed-width peaks) are called per  
1646 sample and then ranked by significance. Next, for all peaks across multiple samples, the peak  
1647 with the highest significance is kept. Peaks directly overlapping this most-significant peak are  
1648 discarded and then this procedure is repeated until all peaks have either been kept or discarded,  
1649 thus converging upon a non-overlapping fixed-width peak set.

1650 **c.** Bar plot showing the number of final peaks identified across all clusters (“Union Peaks”) and  
1651 within each cluster from the hematopoiesis dataset. Bars are colored by peak annotation relative  
1652 to a supplied gene set.

1653 **d.** Heatmap of hypergeometric enrichment testing the overlap of curated peak sets from  
1654 previously published bulk ATAC-seq data (provided by ArchR) with the marker peak sets identified  
1655 for each cluster in the hematopoiesis dataset in **Figure 3c**.

1656

### 1657 **Supplementary Fig. 12.**

1658 **a.** Schematic for the projection of bulk ATAC-seq data into an existing single-cell embedding using  
1659 LSI projection. Briefly, bulk ATAC-seq data is deeply sequenced (10-20 million fragments), down  
1660 sampled to a fragment number corresponding to the average single-cell experiment, and LSI-  
1661 projected into the single-cell subspace.

1662 **b.** LSI projection of bulk ATAC-seq data from diverse hematopoietic cell types into the scATAC-  
1663 seq embedding of the hematopoiesis dataset.

1664 **c-d.** UMAP of scATAC-seq data from the hematopoiesis dataset (N = 215,031 cells) colored by  
1665 (c) sorted cells processed with the Fluidigm C1 system or (d) inferred gene scores for marker  
1666 genes of hematopoietic cells.



1667 **e.** Schematic of the scalable chromVAR method implemented in ArchR. Briefly, ArchR computes  
1668 global accessibility within each peak and then computes chromVAR deviations for each sample  
1669 independently. This design facilitates large-scale chromVAR analysis with minimal memory usage  
1670 for massive-scale scATAC-seq datasets.

1671 **f.** Dot plot showing the identification of positive TF regulators through correlation of chromVAR  
1672 TF deviation scores and inferred gene scores in cell groups (Correlation > 0.5 and Deviation  
1673 Difference in the top 50<sup>th</sup> percentile). These TFs were additionally filtered by the maximum  
1674 observed deviation score difference observed across each cluster average. This additional filter  
1675 removes TFs that are correlated but do not have large accessibility changes in hematopoiesis.

1676 **g.** Schematic of TF footprinting with Tn5 bias correction in ArchR. Briefly, base-pair resolution  
1677 insertion coverage files from sample-aware pseudo-bulk replicates are used to compute the  
1678 insertion frequency around each motif for each replicate. For each motif, the total observed k-  
1679 mers relative to the motif center per bp are identified. This k-mer position frequency table can  
1680 then be multiplied by the individual sample Tn5 k-mer frequencies to compute the Tn5 insertion  
1681 bias per replicate.

1682 **h.** TF footprint for the NFIA motif. Lines are colored by cluster identity from the hematopoiesis  
1683 dataset shown in **Figure 3b**.

1684 **i.** Benchmarking of run time for TF footprinting with ArchR for the 102 sample-aware pseudo-bulk  
1685 replicates from the hematopoiesis dataset.

1686

1687 **Supplementary Fig. 13.**

1688 **a.** Schematic of the ArchR integrative genome browser. Briefly, the ArchR integrative browser is  
1689 launched with a single command into an interactive Shiny session. From there, users can select  
1690 any gene to visualize the accessibility genome track. Additionally, users can change cell  
1691 groupings, resolution, layout and more with an intuitive user interface. Lastly, users can supply

1692 custom feature regions (such as peak sets) or looping/linkage sets (such as peak co-  
1693 accessibility).

1694 **b-e.** Genome accessibility track visualization of marker genes with peak co-accessibility for **(b)**  
1695 the *CD1C* locus (chr1:158,209,562-158,299,563), **(c)** the *AVP* locus (chr20:3,040,369-  
1696 3,090,370), **(d)** the *RORC* locus (chr1:151,764,347-151,819,348), and **(e)** the *SDC1* locus  
1697 (chr2:20,400,193-20,450,194).

1698

1699 **Supplementary Fig. 14.**

1700 **a.** Side-by-side UMAPs for the hematopoiesis dataset cells colored by (top) gene expression  
1701 ( $\log_2(\text{Normalized Counts} + 1)$ ) from scRNA-seq alignment or (bottom) inferred gene scores  
1702 ( $\log_2(\text{Gene Score} + 1)$ ) from gene score Model 42 (see **Figure 2c**) for key immune marker genes.

1703

1704 **Supplementary Fig. 15.**

1705 **a.** Schematic of identification of peak-to-gene links with ArchR. First, all combinations of peak-to-  
1706 gene linkages are identified. Second, the peak accessibility and gene expression for cell groups  
1707 are calculated. Finally, all potential peak-to-gene linkages are tested and significant links ( $R >$   
1708  $0.45$  and  $\text{FDR} < 0.1$ ) are kept.

1709 **b.** Heatmap of 70,239 peak-to-gene links identified across the hematopoiesis dataset with ArchR.

1710

1711

1712

1713

1714

1715

1716

1717

1718 **Supplementary Tables**

1719

1720 **Supplementary Table 1. scATAC-seq Data Sets**

1721 This table contains information about each scATAC-seq data set used in this study including QC  
1722 statistics, scATAC platform and source.

1723

1724 **Supplementary Table 2. scATAC-seq Benchmarking Results**

1725 This table contains information corresponding to benchmarking results of Signac, SnapATAC and  
1726 ArchR for the benchmarking data sets used in this study. Information such as run time and  
1727 maximum memory usage are present in this table.

1728

1729 **Supplementary Table 3. Gene Score Models**

1730 This table contains information for each of the Gene Score models used in **Figure 2**. Descriptions  
1731 of each model are provided in this table.

1732

1733 **Supplementary Table 4. Positive Hematopoietic Regulators**

1734 This table contains information for the identification of positively correlated Hematopoietic TFs.  
1735 Information such as Pearson correlation, linkage statistics and motif are located in this table.

1736

1737 **Supplementary Table 5. Hematopoiesis Peak To Gene Linkages**

1738 This table contains information corresponding to the peak to gene linkages in Hematopoiesis.  
1739 Information such as peak coordinate, gene coordinate and Pearson correlation can be found in  
1740 this table.

1741

1742

1743

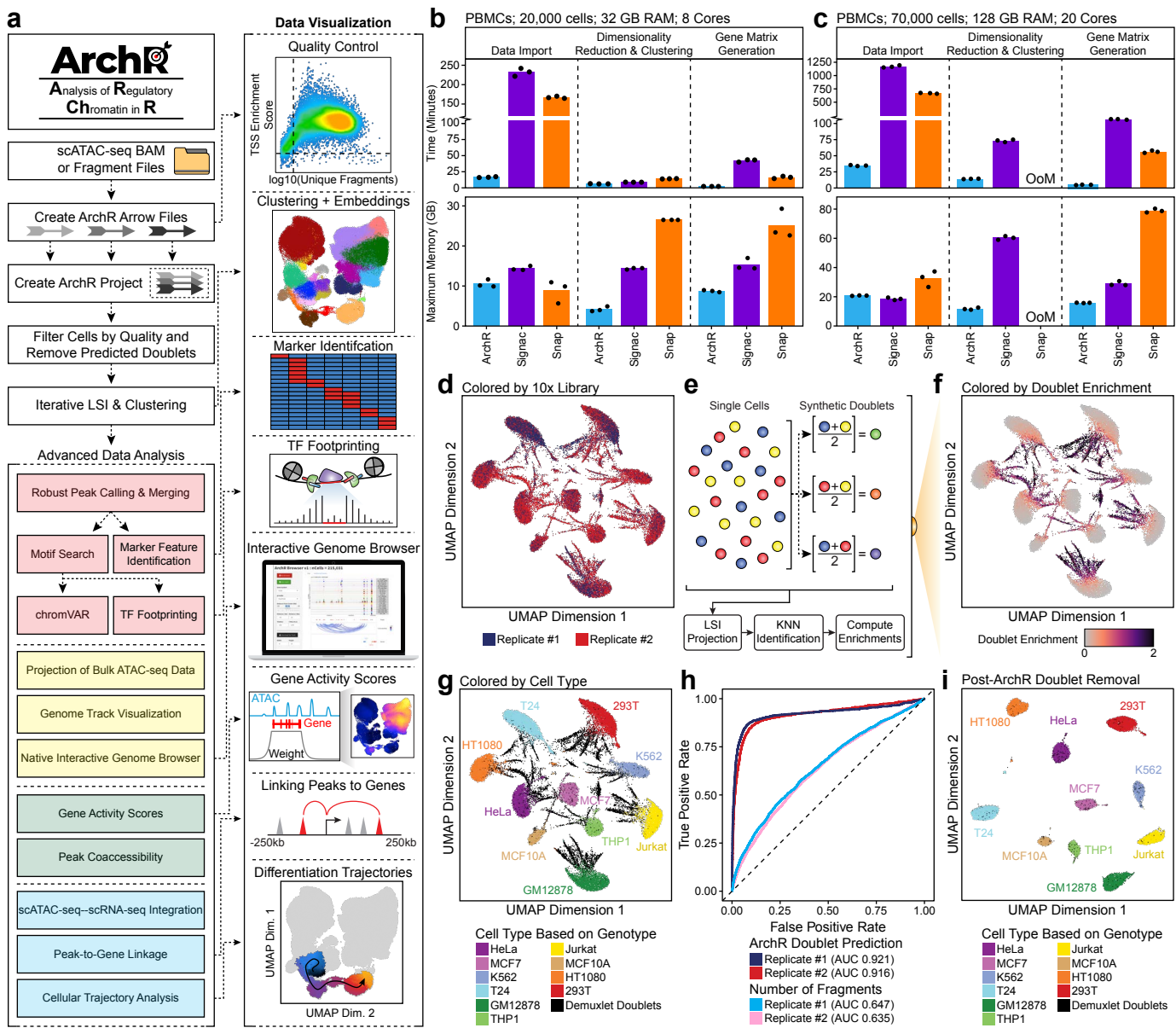


Figure 1

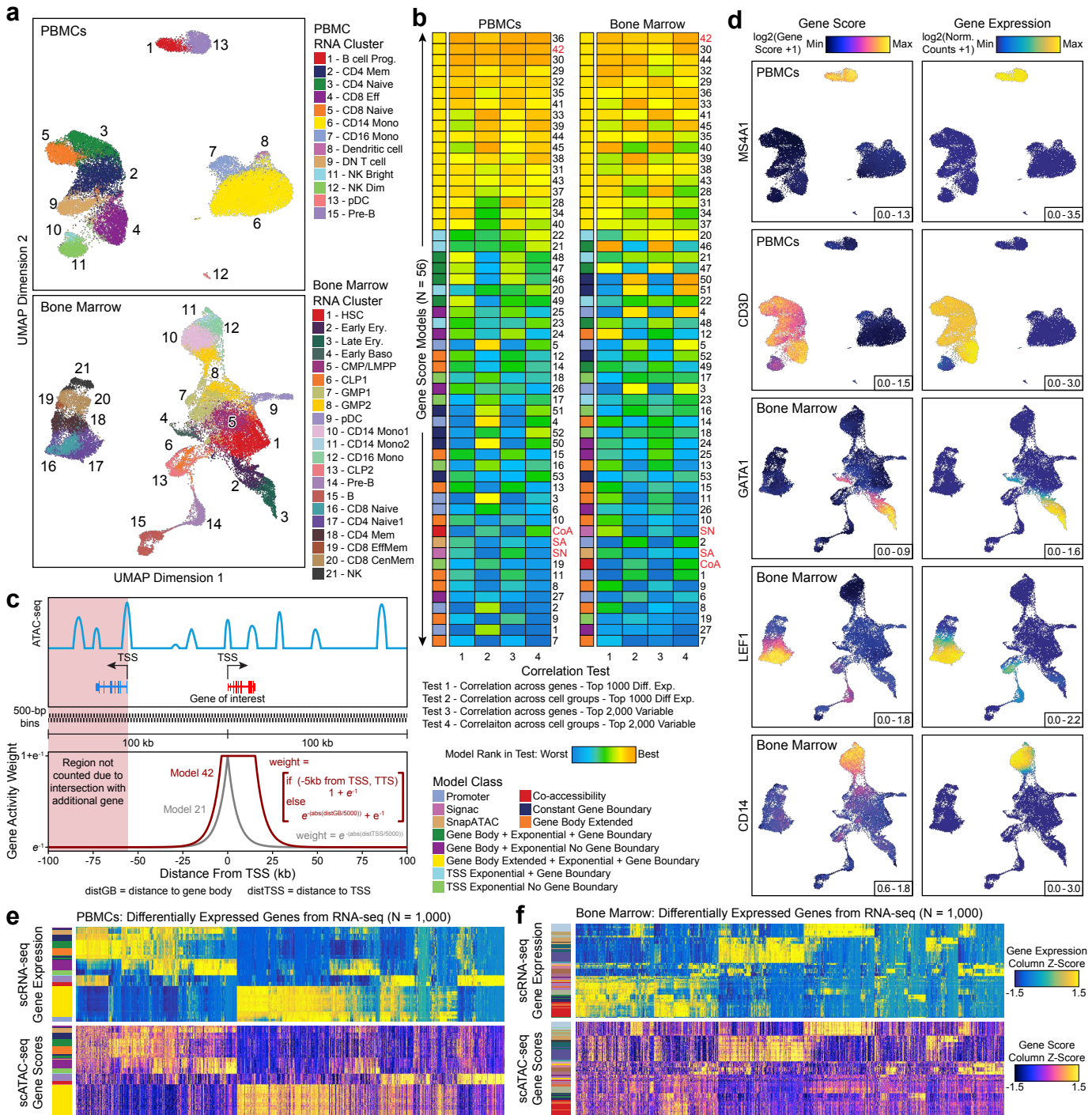


Figure 2





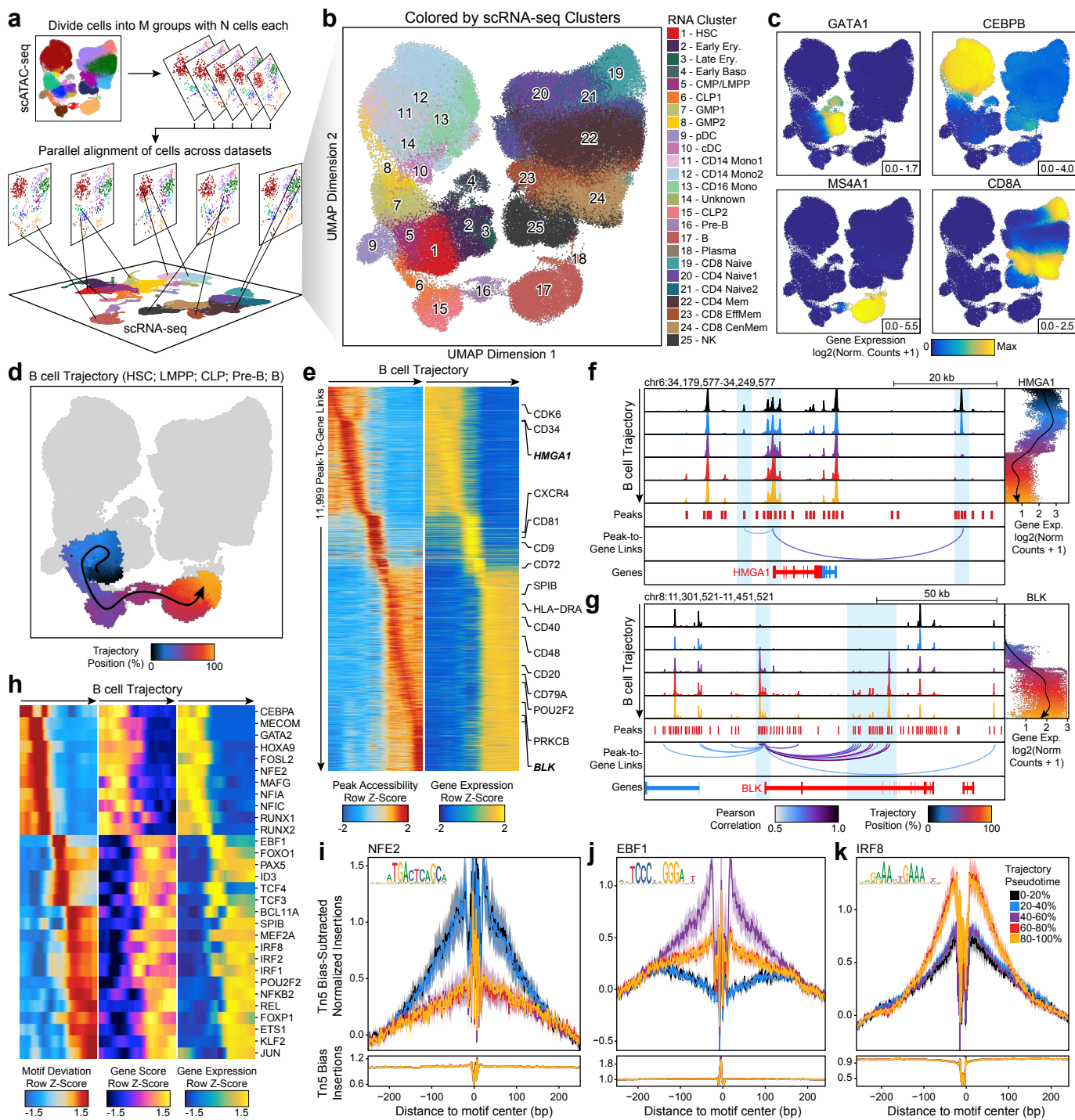


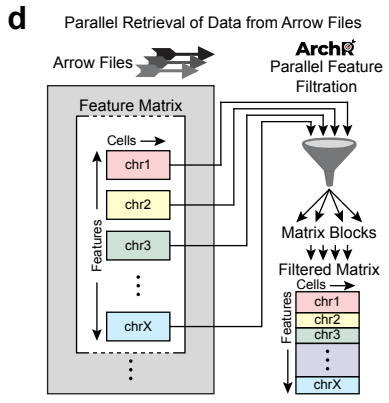
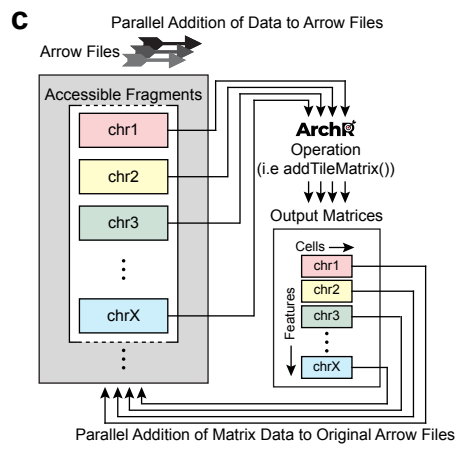
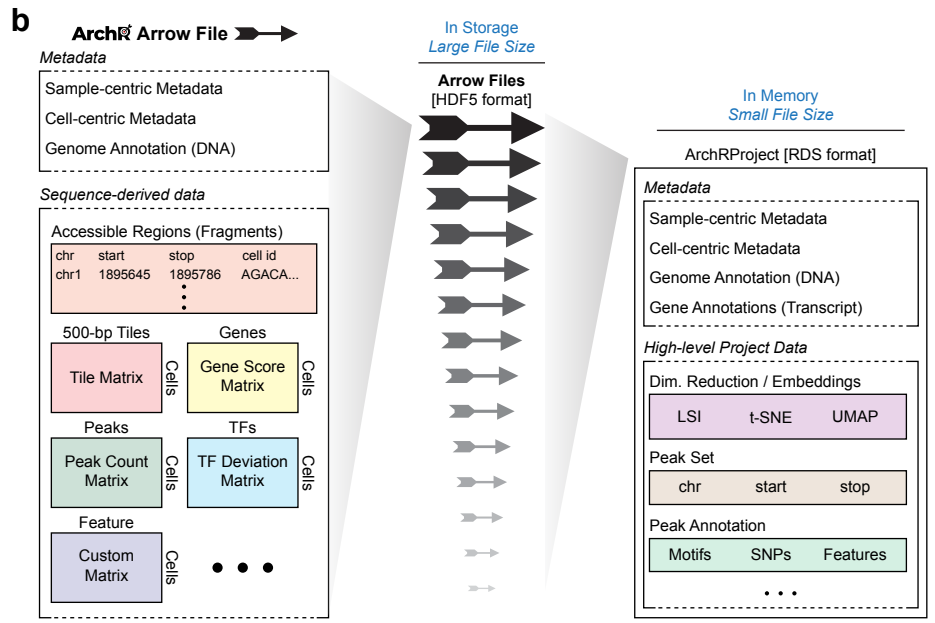
Figure 4



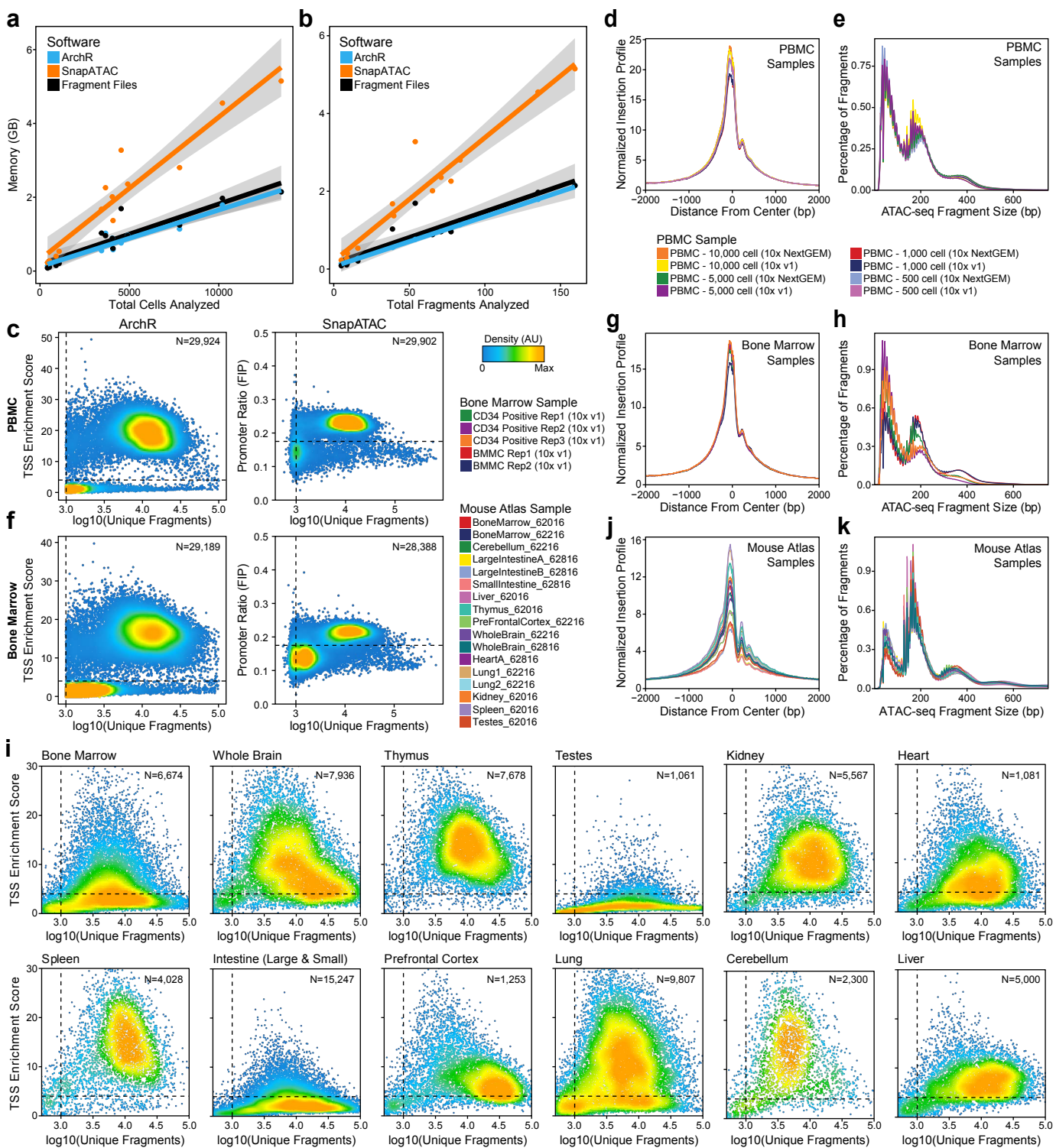
**a**

	ArchR	Signac	SnapATAC	
Pre-processing	NR	NA	✓	
Data import / base file type creation	✓	NA	✓	Data Import
QC filter cells	✓	✓	✓	
Matrix creation	✓ (Tile)	✓ (Peak)	✓ (Tile)	
Doublet removal	✓	NP	NP	Doublet Removal
Data imputation with MAGIC	✓	NP	✓	Gene Scores
Genome-wide gene score matrix	✓	✓	✓	
Dimensionality reduction and clustering	✓	✓	✓	Clustering
UMAP and tSNE plotting	✓	✓	✓	
Cluster peak calling	✓	NP	✓	Standard ATAC-seq Analyses
Cluster-based peak matrix creation	✓	NP	✓	
Motif enrichment	✓	✓	✓	
chromVAR motif deviations	✓	✓	✓	
Footprinting	✓	NP	NP	
Feature set annotation	✓	NP	NP	
Track plotting	✓	✓	NP	
Co-accessibility	✓	NP	NP	Data Visualization
Interactive genome browser	✓	NP	NP	Advanced ATAC-seq Analyses
Cellular trajectory analysis	✓	NP	NP	
Project bulk data into scATAC embedding	✓	NP	NP	
Integration of RNA-seq and ATAC-seq	✓	✓	✓	Integration of RNA-seq and ATAC-seq
Genome-wide peak-to-gene links	✓	NP	NP	

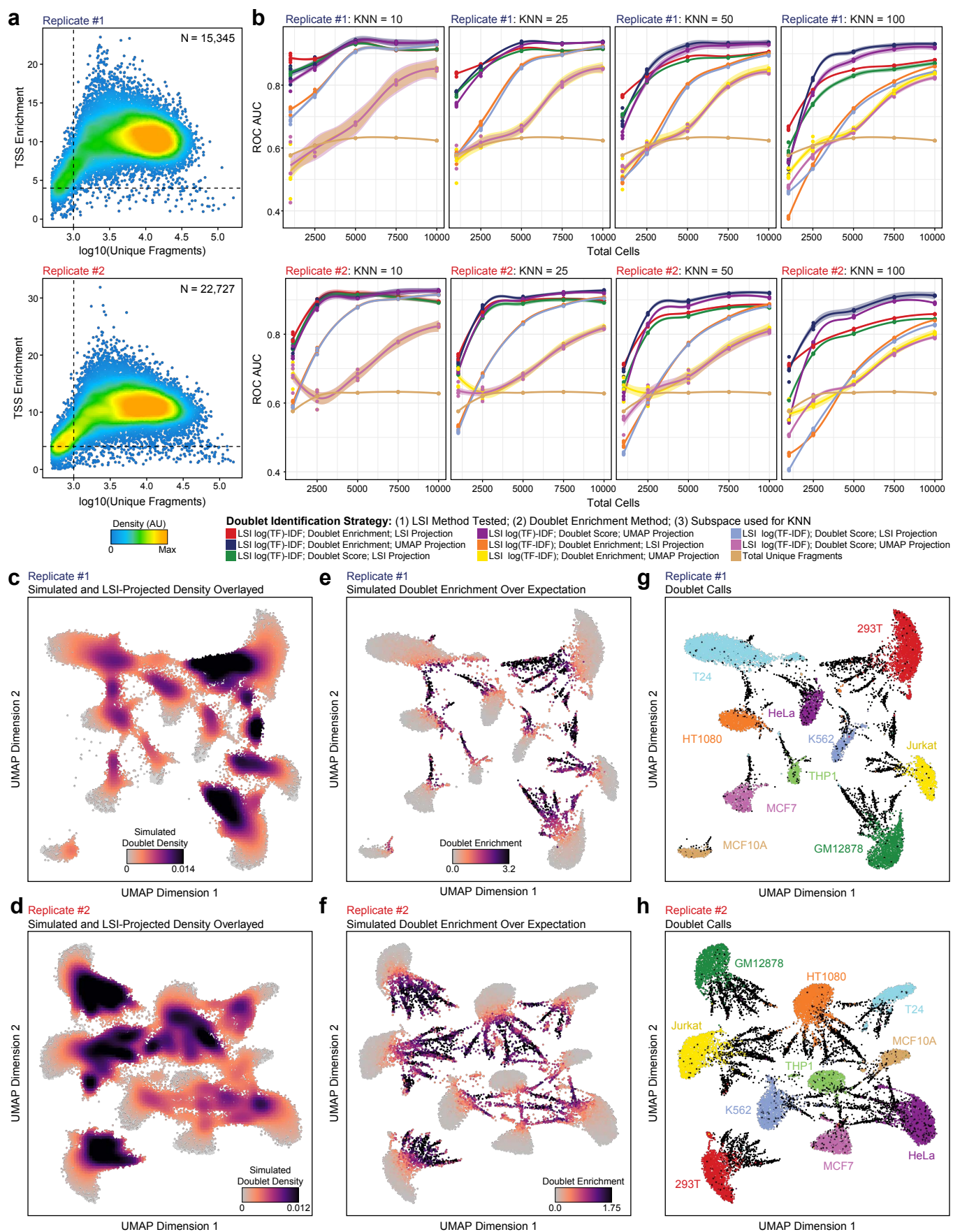
NR = Not Required NA = Not Applicable NP = Not Possible



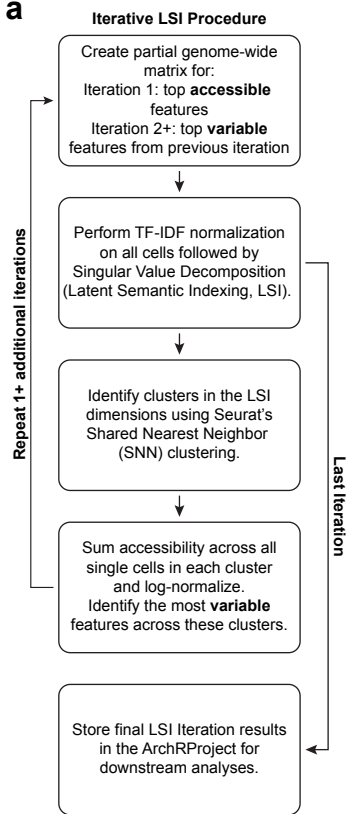
Supplementary Figure 1



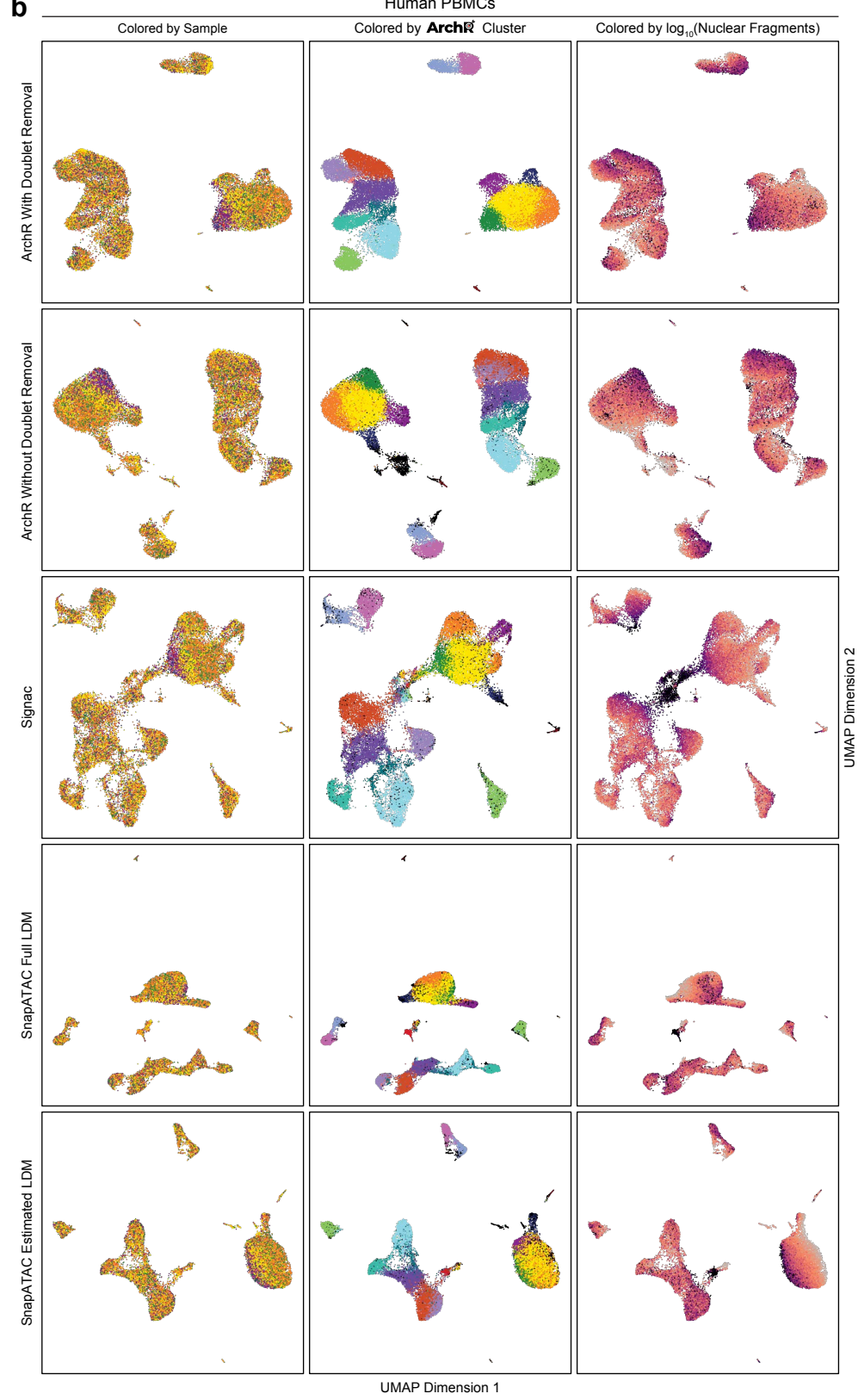








- Sample**
- PBMC - 1,000 cell (10x NextGEM)
  - PBMC - 1,000 cell (10x v1)
  - PBMC - 5,000 cell (10x NextGEM)
  - PBMC - 5,000 cell (10x v1)
  - PBMC - 10,000 cell (10x NextGEM)
  - PBMC - 10,000 cell (10x v1)
  - PBMC - 500 cell (10x NextGEM)
  - PBMC - 500 cell (10x v1)
- ArchR Cluster**
- Low Quality Cells (ArchR Filtered)
  - Doubles (ArchR Identified)
  - 1 - Dendritic cell
  - 2 - CD14 Mono1
  - 3 - CD14 Mono2
  - 4 - CD14 Mono3
  - 5 - CD16 Mono
  - 6 - Pre-B
  - 7 - Pro-B
  - 8 - CD4 Naive1
  - 9 - CD4 Naive2
  - 10 - CD4 Mem
  - 11 - CD8 Naive
  - 12 - CD8 EffMem
  - 13 - CD8 CenMem
  - 14 - NK
  - 15 - Double Negative
  - 16 - Other
- log<sub>10</sub>(Nuclear Fragments)**
- 3.25 4.50
- 



Supplementary Figure 5

**a**

## Human Bone Marrow

Colored by Sample

Colored by **ArchR** ClusterColored by  $\log_{10}$ (Nuclear Fragments)

ArchR With Doublet Removal

ArchR Without Doublet Removal

Signac

SnapATAC Full LDM

SnapATAC Estimated LDM

Sample

- BMCC Rep1 (10x v1)
- BMCC Rep2 (10x v1)
- CD34 Positive Rep1 (10x v1)
- CD34 Positive Rep2 (10x v1)
- CD34 Positive Rep3 (10x v1)

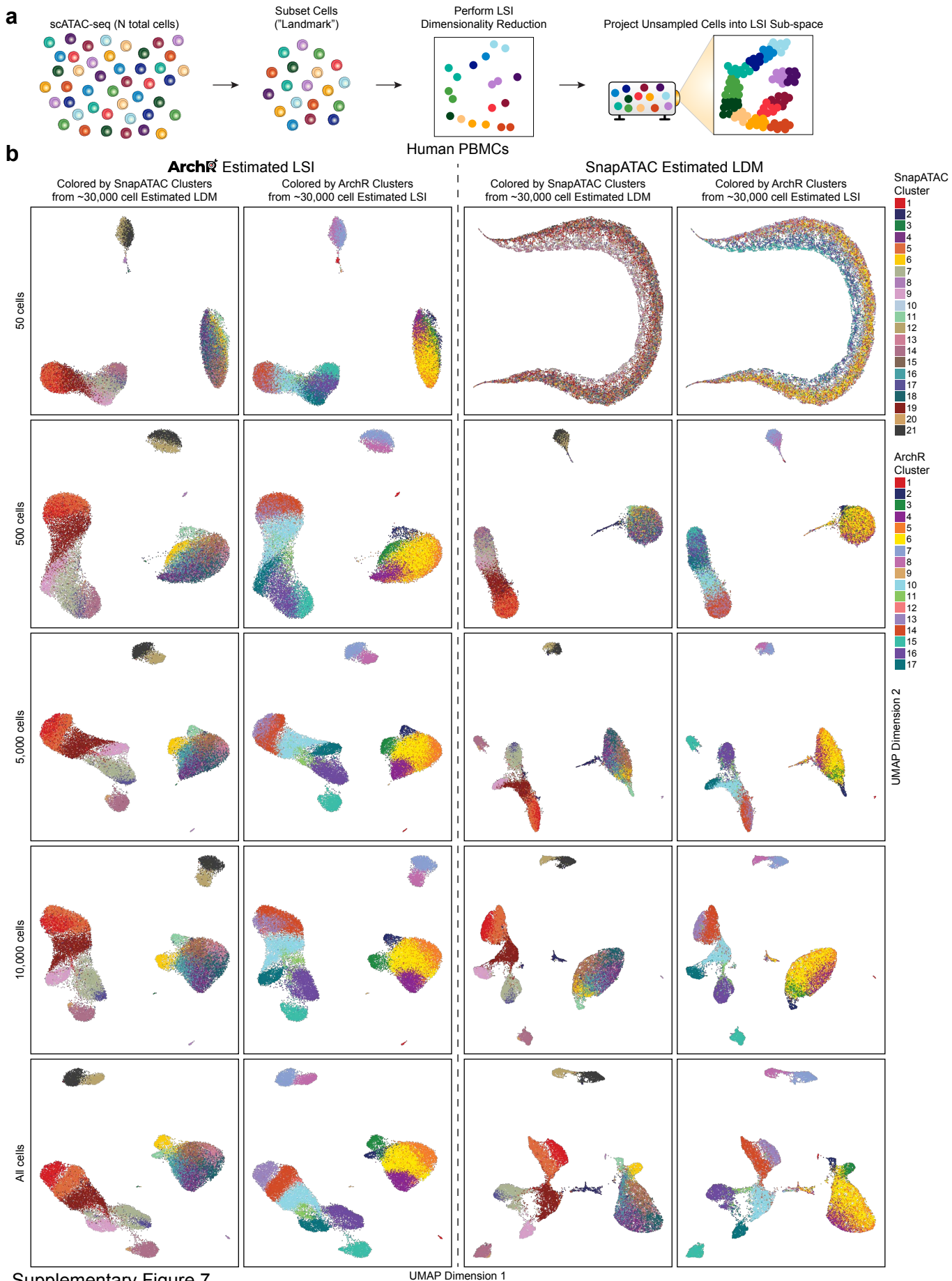
ArchR Cluster

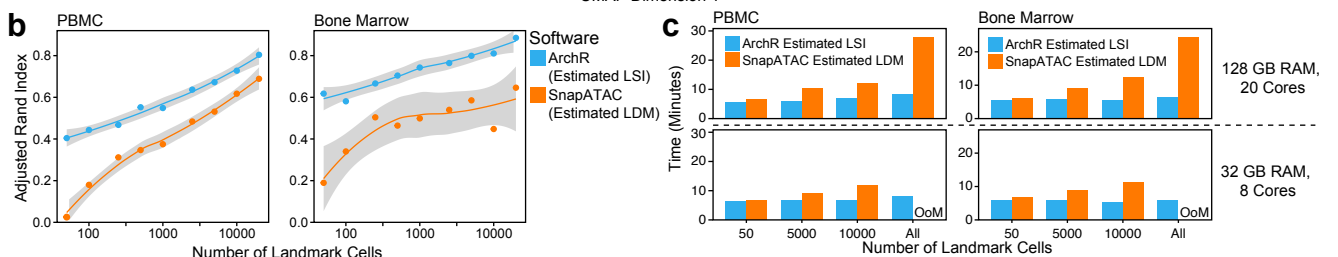
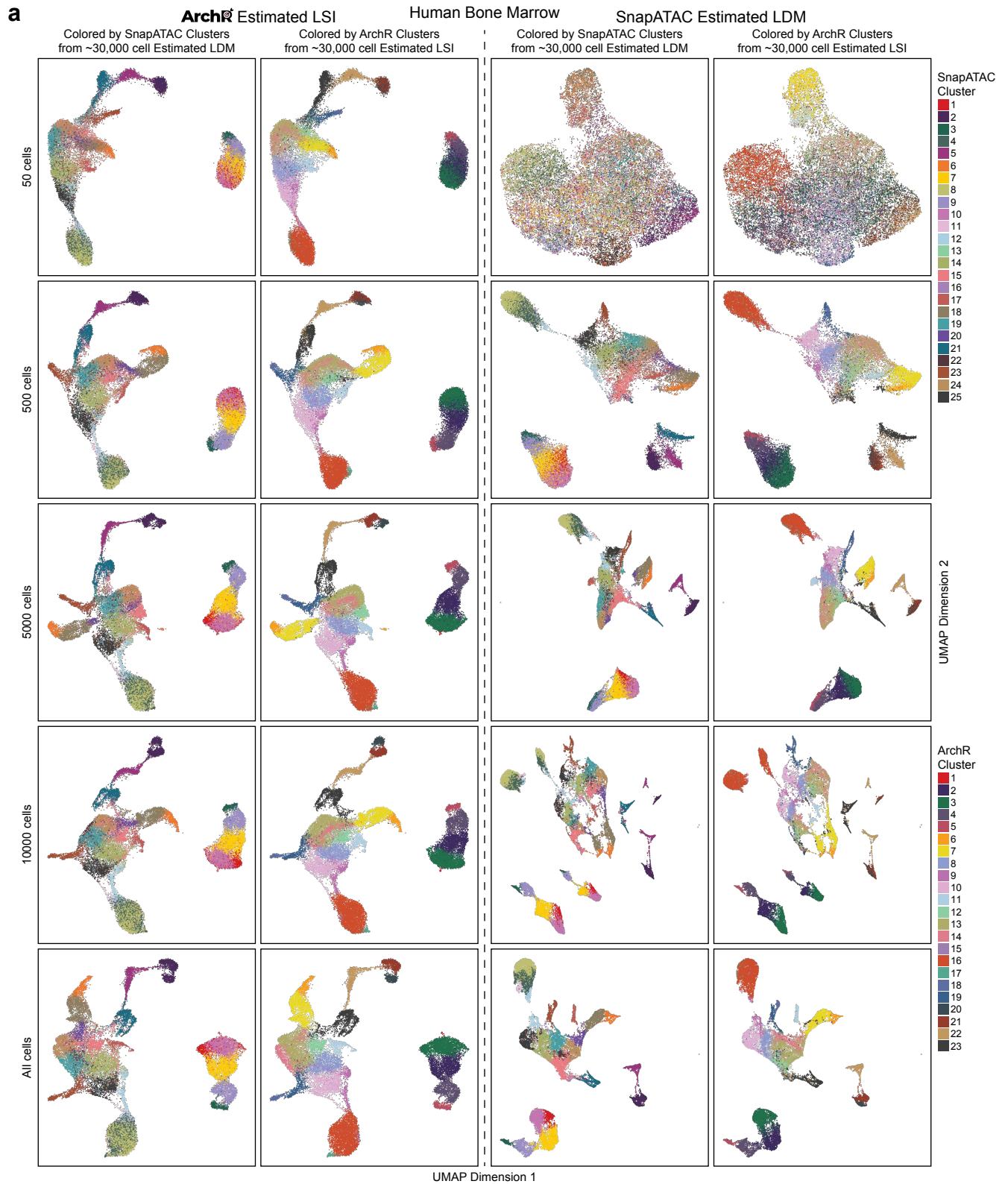
- Low Quality Cells (ArchR Filtered)
- Doublets (ArchR Identified)
- 1 - HSC
- 2 - CMP
- 3 - Early Ery.
- 4 - Late Ery.
- 5 - LMPP
- 6 - CLP
- 7 - Pre-B
- 8 - B1
- 9 - B2
- 10 - GMP1
- 11 - GMP2
- 12 - GMP3
- 13 - Neutrophil
- 14 - Early Baso.
- 15 - Monocyte
- 16 - pDC1
- 17 - pDC2
- 18 - CD8 CenEff Mem
- 19 - NK
- 20 - CD4 Mem
- 21 - CD4/CD8 Naive
- 22 - Other T cell

 $\log_{10}$ (Nuclear Fragments)

UMAP Dimension 2

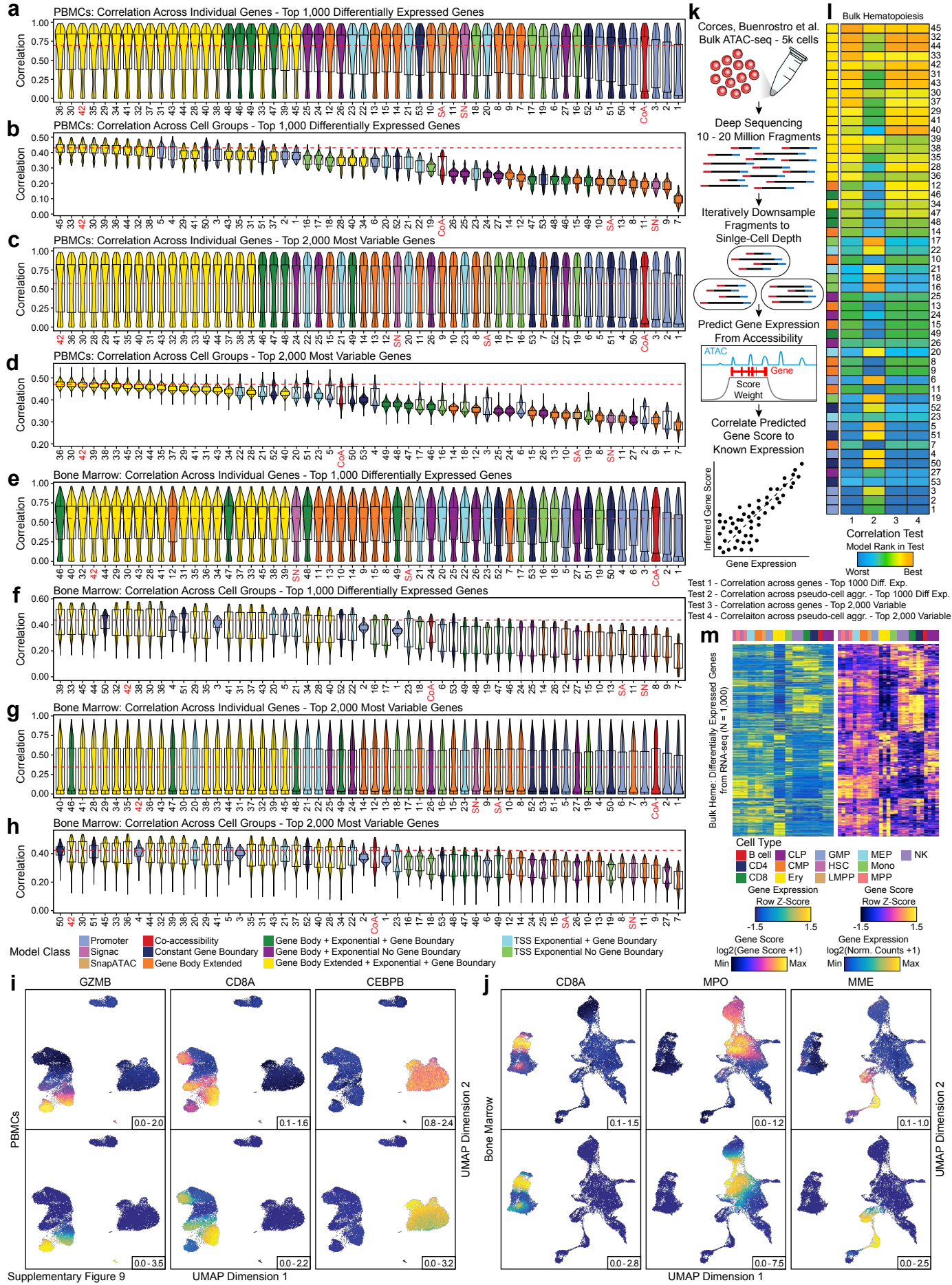
UMAP Dimension 1

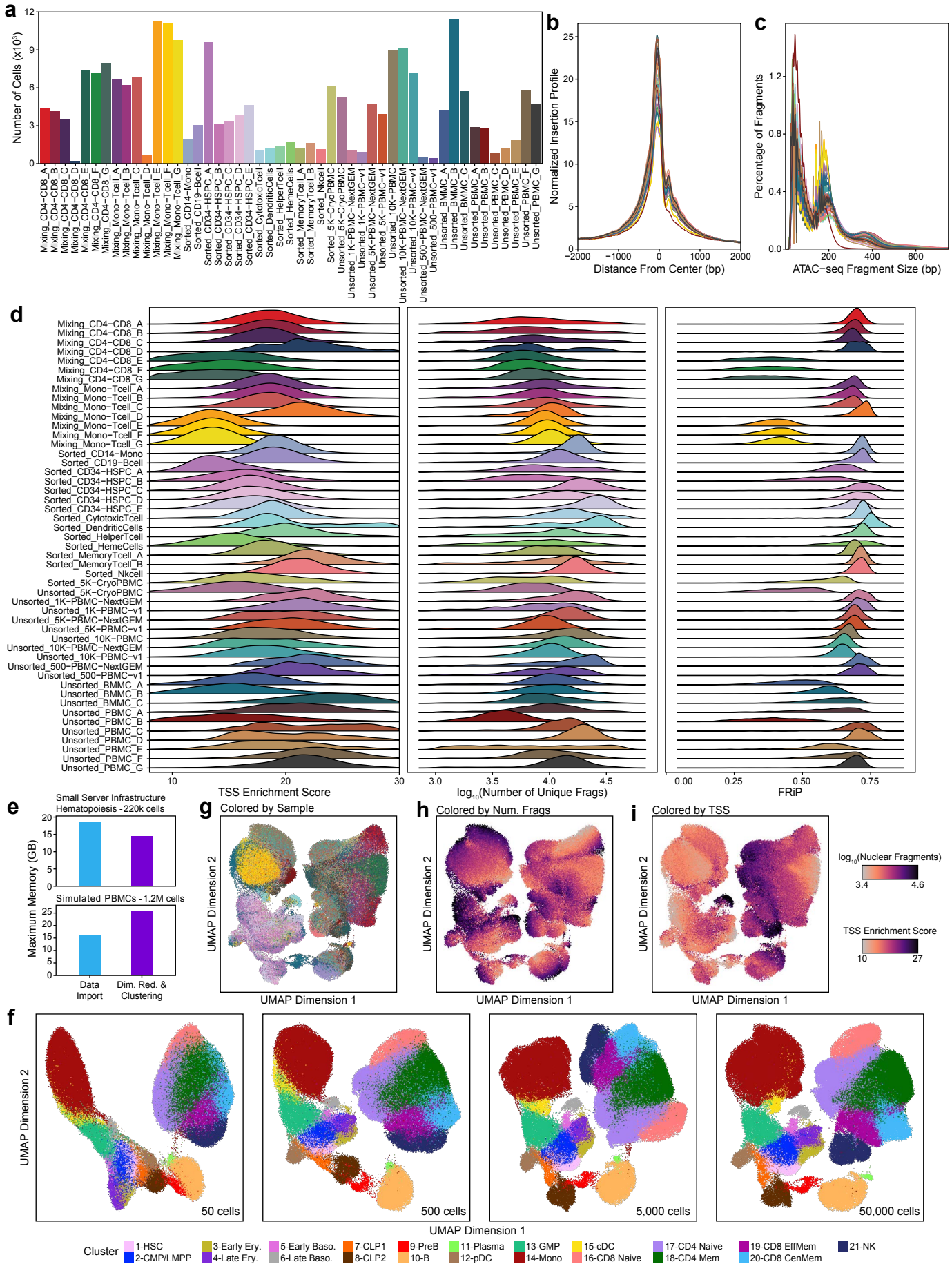




Supplementary Figure 8





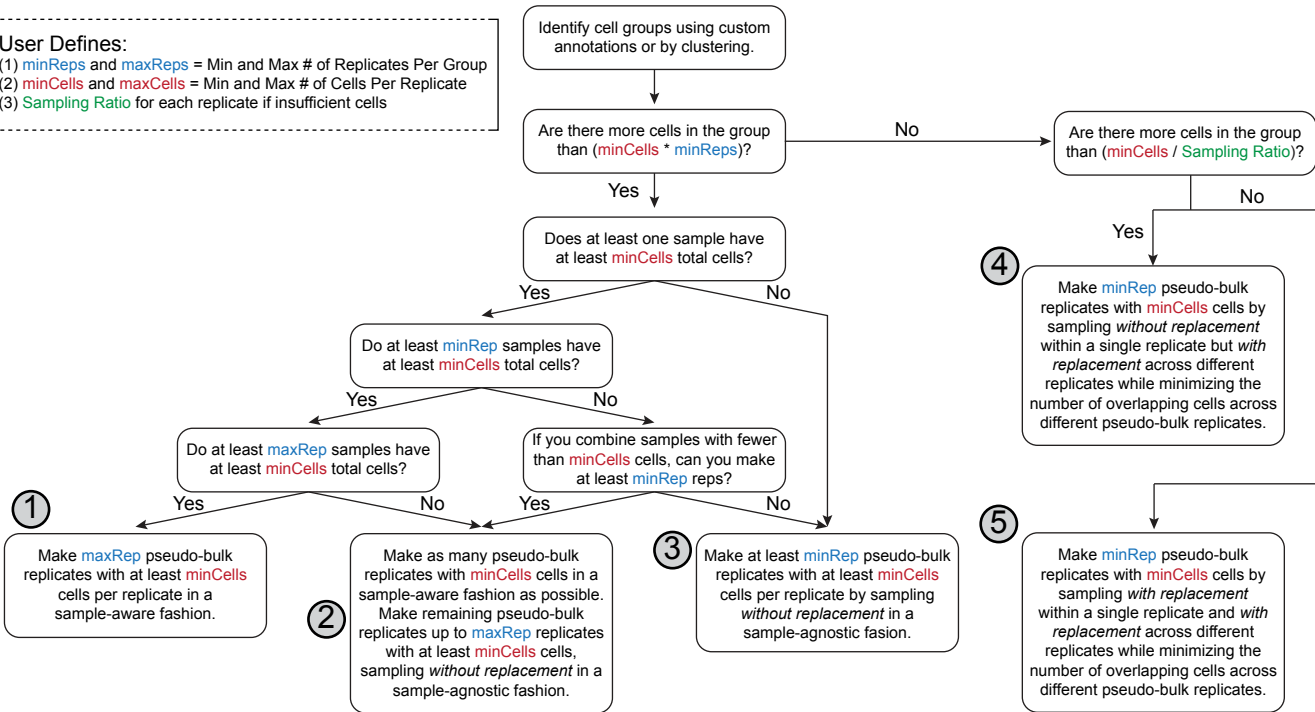


a

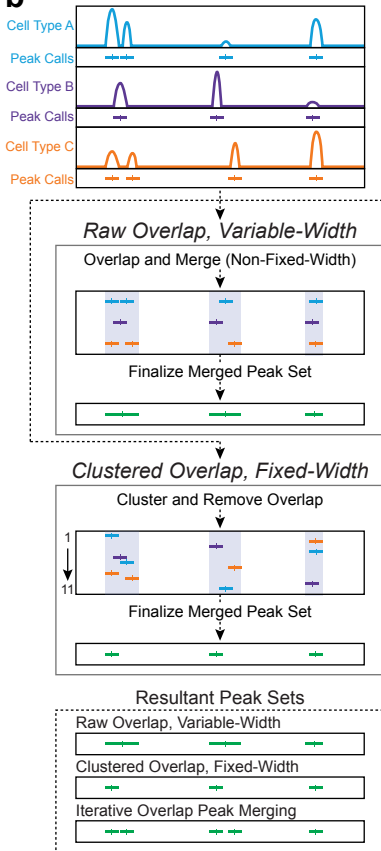
## ArchR Sample-Aware Pseudo-Bulk Replicates

## User Defines:

- (1) **minReps** and **maxReps** = Min and Max # of Replicates Per Group
- (2) **minCells** and **maxCells** = Min and Max # of Cells Per Replicate
- (3) **Sampling Ratio** for each replicate if insufficient cells



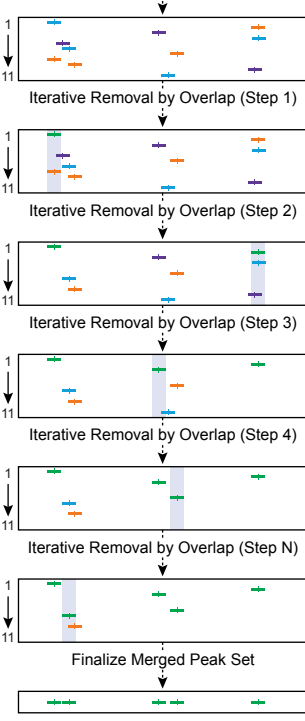
b



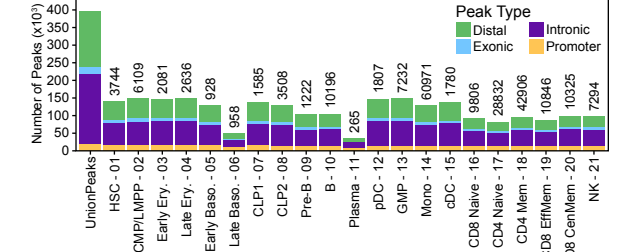
## ArchR

## Iterative Overlap Peak Merging

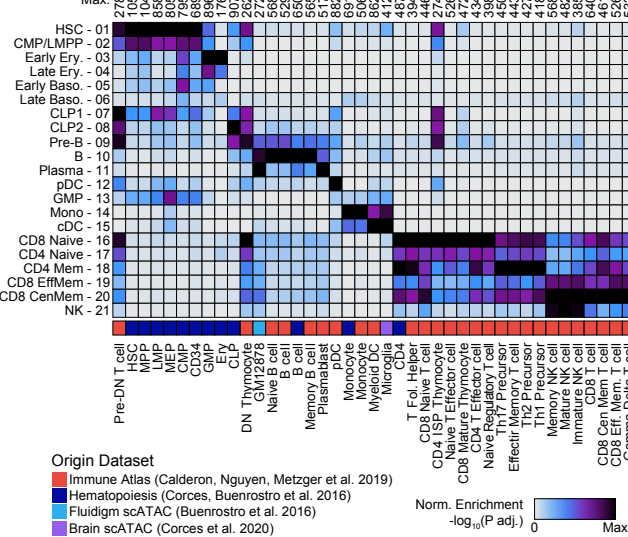
## Rank Peaks by Normalized Significance



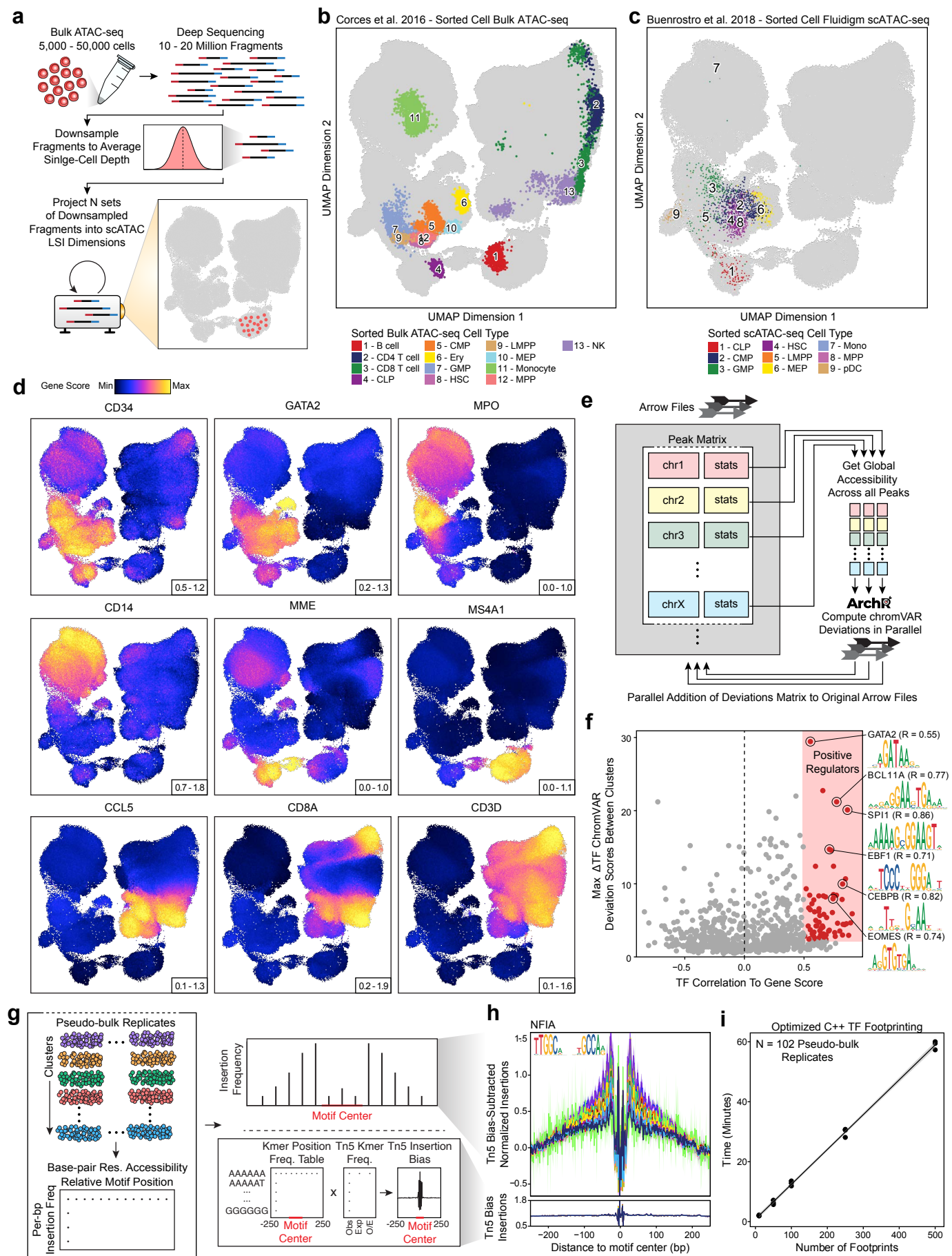
c



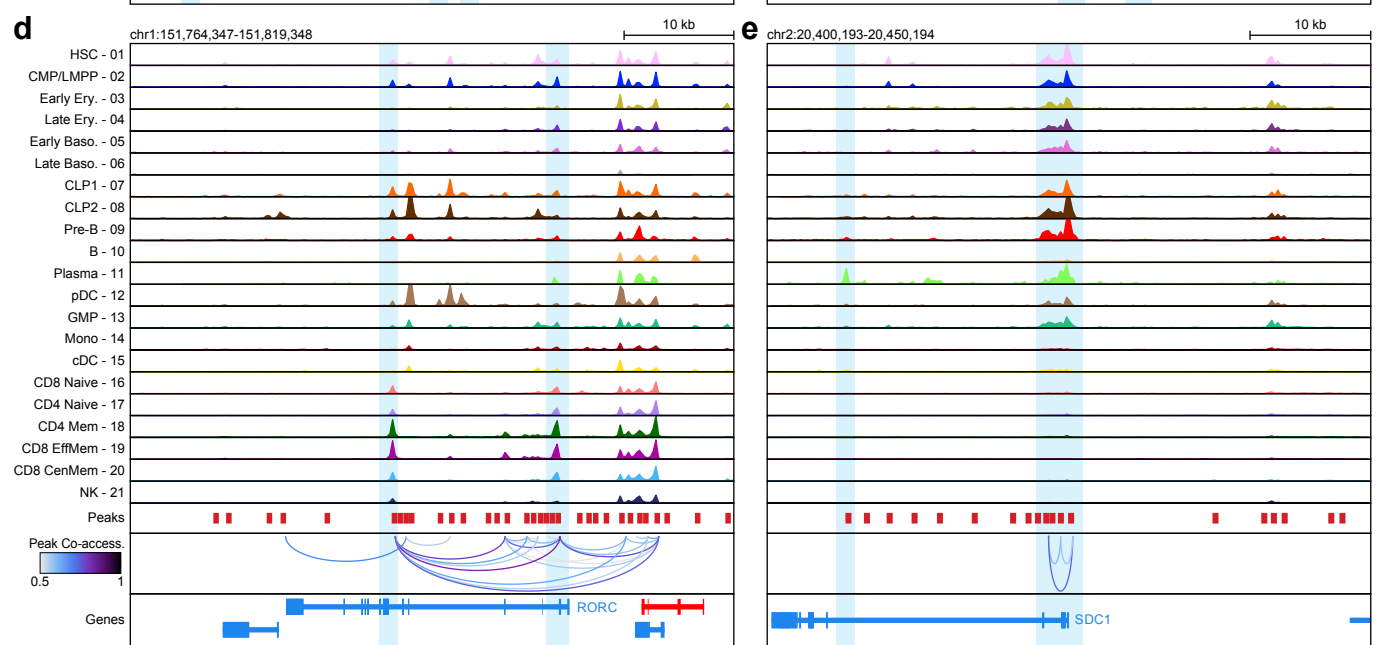
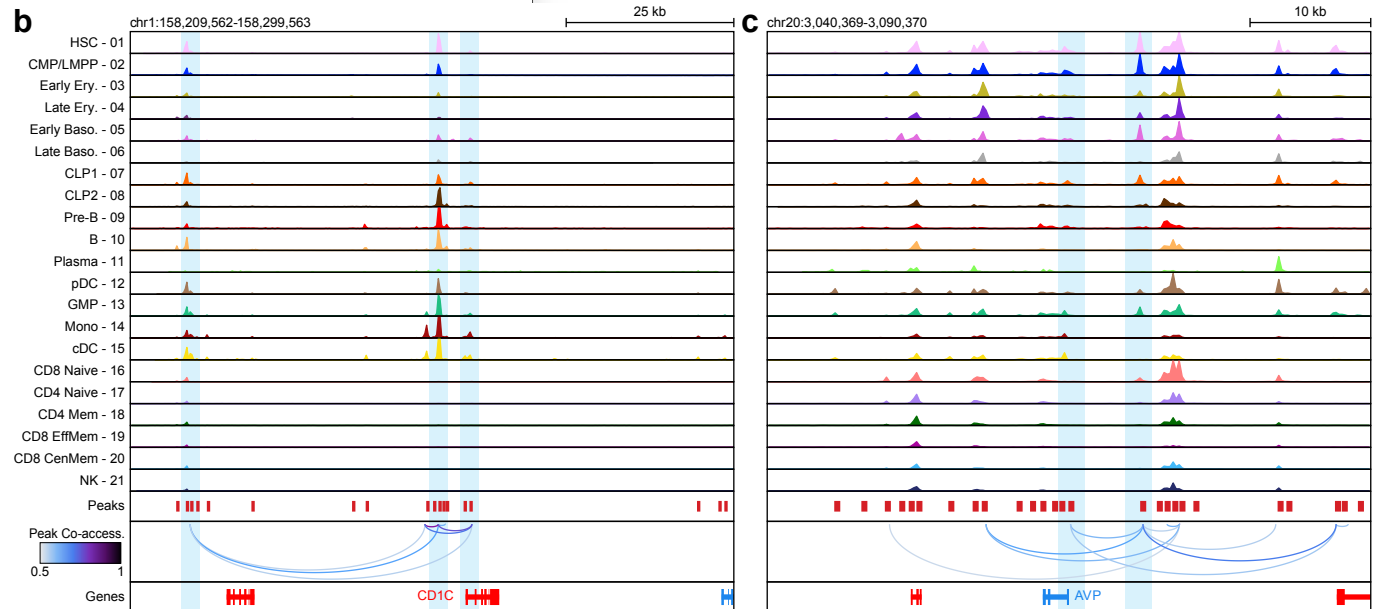
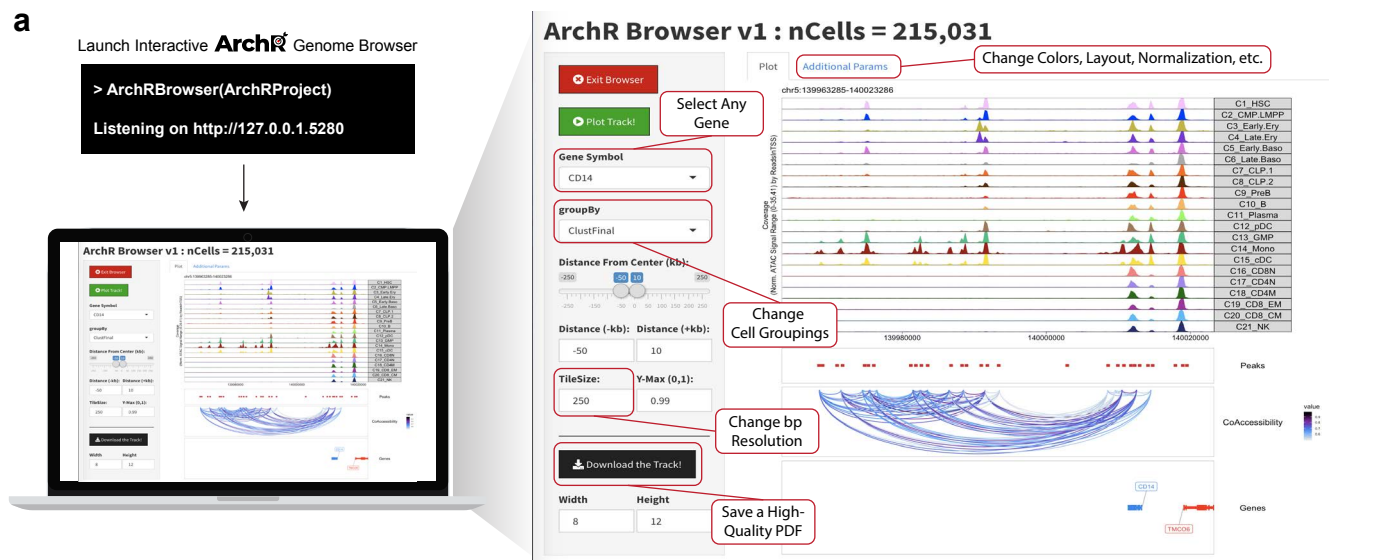
d







Supplementary Figure 12



Supplementary Figure 13

**a**



