

Benchmarking Automated Cell Type Annotation Tools for Single-cell ATAC-seq Data

1 Yuge Wang¹, Xingzhi Sun², Hongyu Zhao^{1,3,*}

2 ¹Department of Biostatistics, Yale School of Public Health, New Haven, CT, United States

3 ²Department of Statistics and Data Science, Yale University, New Haven, CT, United States

4 ³Program of Computational Biology and Bioinformatics, Yale University, New Haven, CT, United
5 States

6 * **Correspondence:**

7 Hongyu Zhao

8 hongyu.zhao@yale.edu

9 **Keywords:** label transfer, scATAC-seq, scRNA-seq, machine learning, benchmark.

10 Abstract

11 As single-cell chromatin accessibility profiling methods advance, scATAC-seq has become ever
12 more important in the study of candidate regulatory genomic regions and their roles underlying
13 developmental, evolutionary and disease processes. At the same time, cell type annotation is critical
14 in understanding the cellular composition of complex tissues and identifying potential novel cell
15 types. However, most existing methods that can perform automated cell type annotation are designed
16 to transfer labels from an annotated scRNA-seq data set to another scRNA-seq data set, and it is not
17 clear whether these methods are adaptable to annotate scATAC-seq data. Several methods have been
18 recently proposed for label transfer from scRNA-seq data to scATAC-seq data, but there is a lack of
19 benchmarking study on the performance of these methods. Here, we evaluated the performance of
20 five scATAC-seq annotation methods on both their classification accuracy and scalability using
21 publicly available single-cell datasets from mouse and human tissues including brain, lung, kidney,
22 PBMC and BMMC. Using the BMMC data as basis, we further investigated the performance of these
23 methods across different data sizes, mislabeling rates, sequencing depths and the number of cell types
24 unique to scATAC-seq. Bridge integration, which is the only method that requires additional
25 multimodal data and does not need gene activity calculation, was overall the best method and robust
26 to changes in data size, mislabeling rate and sequencing depth. Conos was the most time and memory
27 efficient method but performed the worst in terms of prediction accuracy. scJoint tended to assign
28 cells to similar cell types and performed relatively poorly for complex datasets with deep annotations
29 but performed better for datasets only with major label annotations. The performance of scGCN and
30 Seurat v3 was moderate, but scGCN was the most time-consuming method and had the most similar
31 performance to random classifiers for cell types unique to scATAC-seq.

32 1 Introduction

33 With the advancement of single-cell sequencing technologies, researchers not only can profile single-
34 cell transcriptomes by scRNA-seq, but can also measure multiple modalities at the single-cell level
35 (Packer and Trapnell, 2018; Carter and Zhao, 2021), among which scATAC-seq is probably the most
36 widely used sequencing technology (Buenrostro et al., 2015; Cusanovich et al., 2015). scATAC-seq
37 can quantify chromatin accessibility across tens of thousands of single cells and is an important tool

38 to study gene regulation accompanied with scRNA-seq (Buenrostro et al., 2018; Fiers et al., 2018; Jia
39 et al., 2018; Wang et al., 2022). Single-cell studies usually start with cell type annotations and
40 accurate and robust annotations are crucial for downstream functional analyses that are often
41 conducted in a cell-type-specific manner. Cell type annotation is often laborious and involves
42 automated annotations from computational tools followed by verification and manual annotations
43 from experts (Clarke et al., 2021). Although there are many tools designed for automated cell type
44 annotations for scRNA-seq data (Abdelaal et al., 2019; Pasquini et al., 2021), only a limited number
45 of tools are available and suitable for scATAC-seq data. As scATAC-seq becomes more mature and
46 widely adopted in single-cell studies, there is a need to comprehensively evaluate their performance
47 on annotating scATAC-seq data.

48 Currently, there are two types of annotation tools that can be applied to scATAC-seq data. The first
49 category includes those originally designed for scRNA-seq data (intra-modality annotation), such as
50 Seurat v3 (Stuart et al., 2019), Conos (Barkas et al., 2019) and scGCN (Song et al., 2021). The
51 second category includes tools designed specifically for scATAC-seq data or for cross-modality
52 annotation. The two representative methods in the second category are scJoint (Lin et al., 2022) and
53 Bridge integration (Hao et al., 2022). Unlike the other methods that directly transfer labels from
54 scRNA-seq to scATAC-seq after unifying the feature set through gene activity calculation, Bridge
55 integration leverages a multimodal data as a bridge, avoiding potential loss of information and
56 incorrectness of assumptions on feature relationships when calculating gene activities.

57 In this study, we benchmark these scATAC-seq annotation tools using real single-cell datasets from
58 various tissues with available cell type annotations as the ground truth. The real data we considered
59 included both paired data where scATAC-seq and scRNA-seq were simultaneously measured in each
60 single cell and unpaired data where scATAC-seq and scRNA-seq were separately measured from the
61 same tissue. We evaluated the performance of different methods on both annotation accuracy and
62 scalability. For accuracy, we considered both the overall accuracy as well as accuracy on ATAC-
63 specific cell types. For scalability, we compared running time and memory usage across different
64 datasets. Apart from evaluating real data across different tissues, we also investigated the model
65 performance across different cell numbers, mislabeling proportions, sequencing depths and number
66 of unique cell types using a well-annotated human bone marrow mononuclear cell (BMMC)
67 multimodal data (Luecken et al., 2021). The results of our study offer a basis for future methodology
68 development and provide a reference for users to choose appropriate tools for cell type annotation
69 from scATAC-seq data.

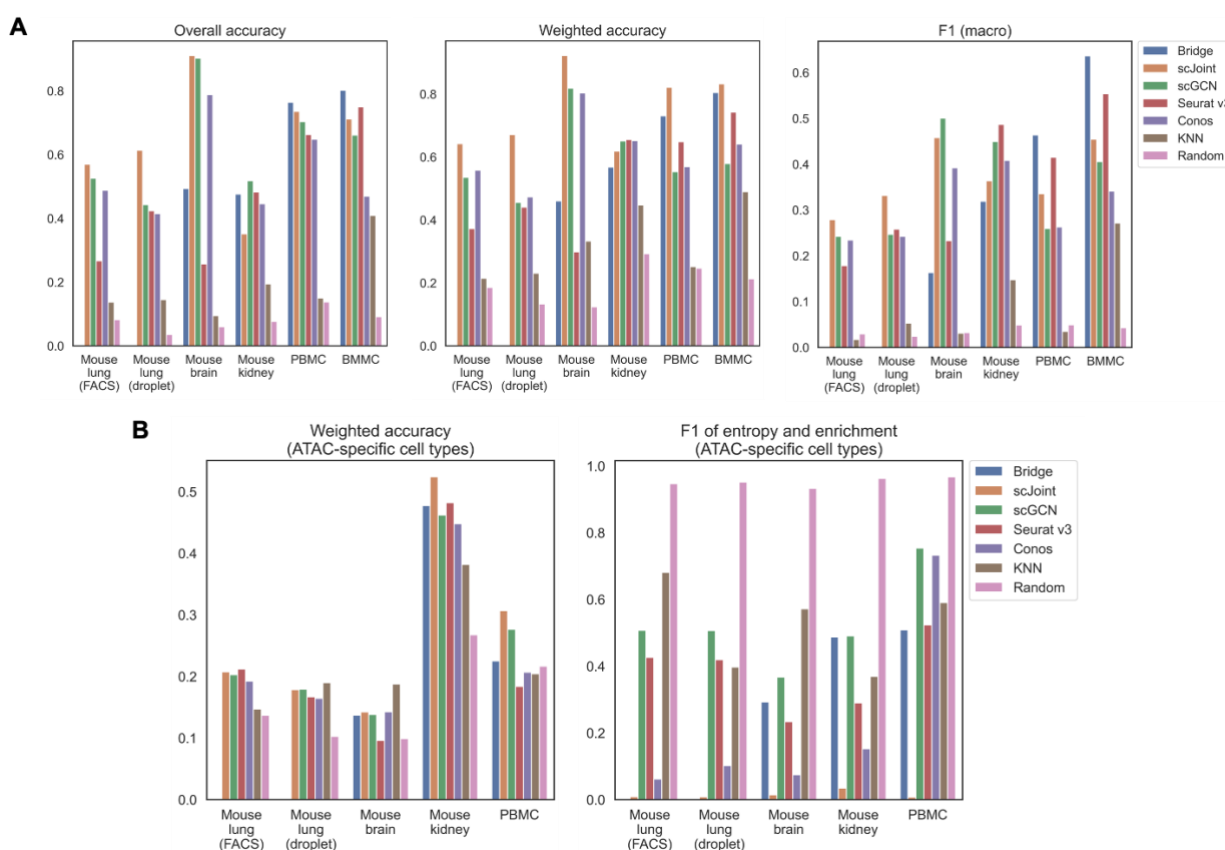
70 **2 Results**

71 **2.1 Performance across Different Tissues**

72 In this study, we used data from five different tissues, including mouse lung (Consortium, 2018;
73 Cusanovich et al., 2018), mouse brain (Chen et al., 2019; Ma et al., 2020), mouse kidney (Cao et al.,
74 2018; Miao et al., 2021), human peripheral blood mononuclear cell (PBMC) (Granja et al., 2019) and
75 human bone marrow mononuclear cells (BMMC) (Luecken et al., 2021) to benchmark five methods
76 for automated scATAC-seq label annotation, including Conos, Seurat v3, scGCN, scJoint and Bridge
77 integration. For mouse lung, scRNA-seq data from both 10x Chromium (droplet-based) and Smart-
78 seq2 (FACS-based) were collected. Among all the methods, only Bridge integration required
79 multimodal data where scATAC-seq and scRNA-seq were simultaneously measured. Therefore, we
80 collected multimodal data for each tissue except for mouse lung (the SHARE-seq data for mouse

81 lung were sequenced too shallowly to be used). For the mouse brain, both SHARE-seq and SNARE-
 82 seq data were used as the multimodal data to benchmark Bridge integration separately.

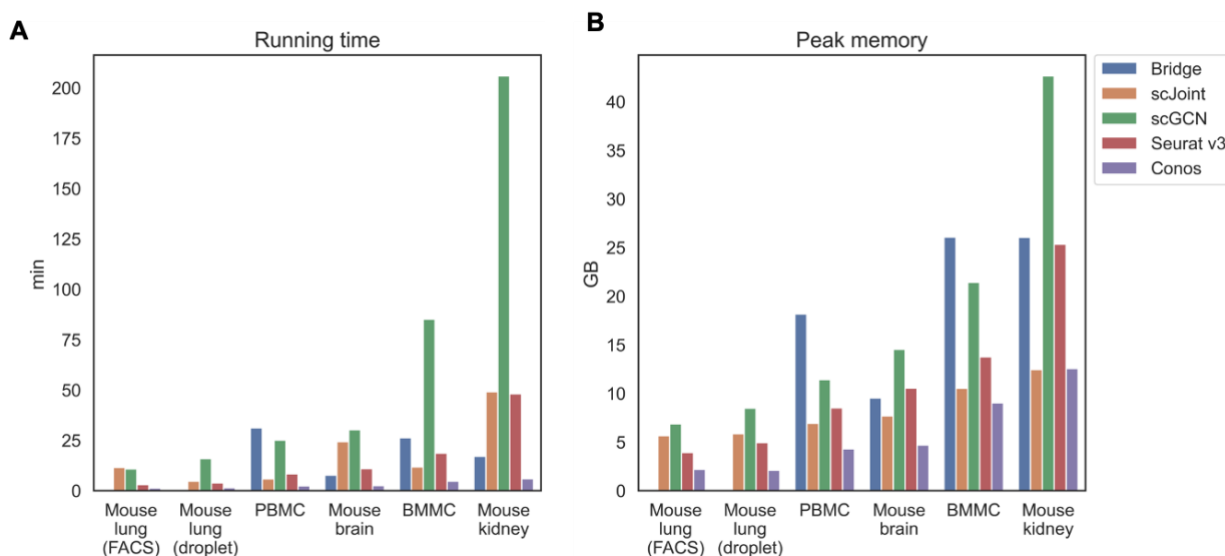
83 We calculated three accuracy-related metrics on all ATAC cells, namely overall accuracy, weighted
 84 accuracy and F1 (macro) of precision and recall. For the first and third metrics, they were calculated
 85 based on predicted label, which was the cell type whose predicted probability was the largest. For
 86 weighted accuracy, we considered the similarity among cell types by calculating the weighted
 87 average of the entire predicted probability vector of each cell. Therefore, even though a predicted
 88 label was false, the score could be high if similar cell types had higher predicted probabilities. As can
 89 be seen from Figure 1A, all the five methods had better performance than plain K nearest neighbor
 90 (KNN) and the random classifiers. For mouse lung (both FACS and droplet) and mouse brain, scJoint
 91 had consistent and leading performance across all the three metrics, with only slightly lower F1
 92 (macro) than scGCN on mouse brain. For the two human tissues (PBMC and BMMC), Bridge
 93 integration achieved the highest overall accuracy and F1 (macro); while for weighted accuracy,
 94 Bridge integration was the second best performer, following scJoint. For mouse kidney, there was no
 95 leading method across all three metrics, but scGCN and Seurat v3 had overall better performance.



96

97 **Figure 1.** Performance of label transfer methods on single-cell data from selected mouse and human
 98 tissues. (A) Overall metrics considering performance on all scATAC-seq cells. (B) Metrics calculated
 99 on scATAC-seq cells labeled with ATAC-specific cell types. The Bridge results shown here for the
 100 mouse brain used SNARE-seq as the multimodal ‘bridge’. Comparison of results using SNARE-seq
 101 and SHARE-seq can be found in Supplementary Figure S1. For mouse lung (both FACS and
 102 droplet), Bridge integration was not considered because of no available multimodal data.

103 Apart from the three metrics assessing all ATAC cells, we designed two additional metrics for cell
104 types that uniquely existed in ATAC data, namely weighted accuracy and F1 of entropy and
105 enrichment (details in Materials and Methods). For ATAC-specific cell types, they could never be
106 correctly classified because their labels did not exist in the reference RNA data. Then, there are two
107 expected patterns for the predicted probability vectors of these cells. One is having predicted
108 probability vectors close to the background distribution of cell types, and the other is having higher
109 predicted probabilities for similar cell types in the RNA data. Both can have their own benefits in real
110 practice. For example, for the first case, one can perform another round of manual annotations on
111 cells with predicted probabilities close to the background distribution; while for the second case, one
112 can tell from the predicted probabilities which existing cell types are the closest to the unknown cell
113 type, however, this might suffer from misclassifying novel cell types due to biological similarity to
114 known cell types. F1 (entropy and enrichment) and weighted accuracy were calculated over cells
115 with unique cell types in the ATAC data to cover the first and the second cases, respectively (Figure
116 1B). Although scJoint consistently had relatively high weighted accuracy across tissues, there were
117 not significant differences in weighted accuracy among all the five methods. For F1 (entropy and
118 enrichment), the scores of scJoint and Conos were extremely low, while scGCN achieved the best
119 scores among the five methods, followed by Bridge integration and Seurat v3.



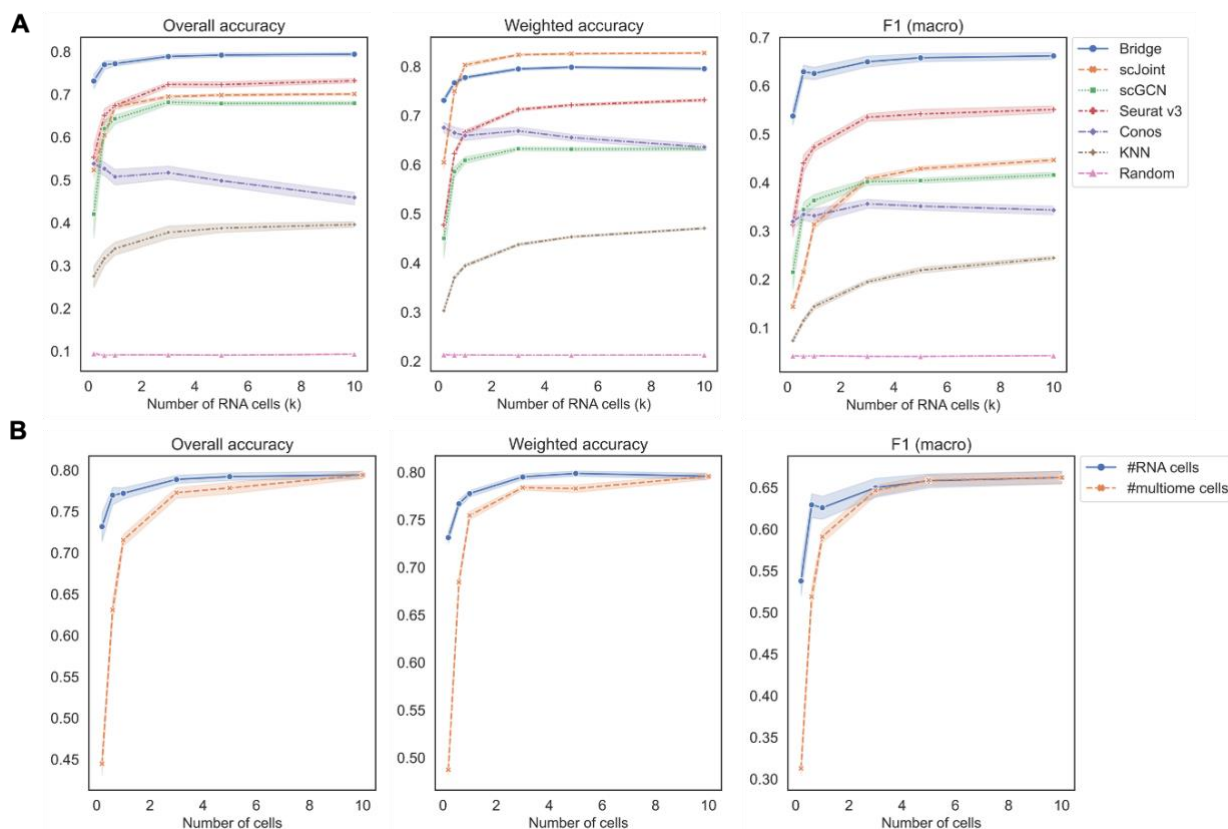
120

121 **Figure 2.** Running time (A) and peak memory usage (B) of different methods on selected tissues.
122 Tissues are placed in the increasing order of their scales from left to right. For mouse lung (both
123 FACS and droplet), Bridge was not considered because of no available multimodal data.

124 Apart from the prediction accuracy, we evaluated the efficiency and scalability of the five methods
125 by recording their running time and peak memory usage on each tissue (Figure 2). scGCN was the
126 most time-consuming method and took the largest memory on mouse kidney, where there were about
127 71,000 cells in total. Conos was the most time and memory efficient method and remained nearly
128 constant as the data scale increased. For the remaining three methods (Bridge integration, scJoint and
129 Seurat v3), their running time didn't differ significantly, but Bridge integration consumed more
130 memory than others.

131 2.2 Performance across Different Data Scales

132 The BMMC data is a first-of-its-kind single-cell multimodal dataset which consists of about 70,000
 133 cells with paired scRNA-seq and scATAC-seq measurements from 10 diverse donors at four
 134 sequencing sites. This dataset contains the largest number of cell types (22) among all selected
 135 tissues and captures both developmental and differentiated cell types. This dataset is the most
 136 comprehensive multi-modal benchmark dataset to date as far as we know, so we designed several
 137 experiments using the BMMC data to investigate the performance of different methods across
 138 diverse data characteristics. For all of the following experiments on the BMMC data, we manually
 139 separated all donors into three groups and used them as unimodal RNA data, unimodal ATAC data,
 140 and multimodal data, respectively (see Materials and Methods).



141

142 **Figure 3.** Performance of methods on different data scales of BMMC. (A) Change the number of
 143 cells in scRNA-seq only, while keeping scATAC-seq and multimodal (for Bridge only) cell numbers
 144 as 10k. Results shown here for Conos were paired with Pagoda2 for data processing. (B)
 145 Performance of Bridge while changing number of cells in scRNA-seq and in the multimodal data,
 146 respectively. The error band shows the 95% confidence interval.

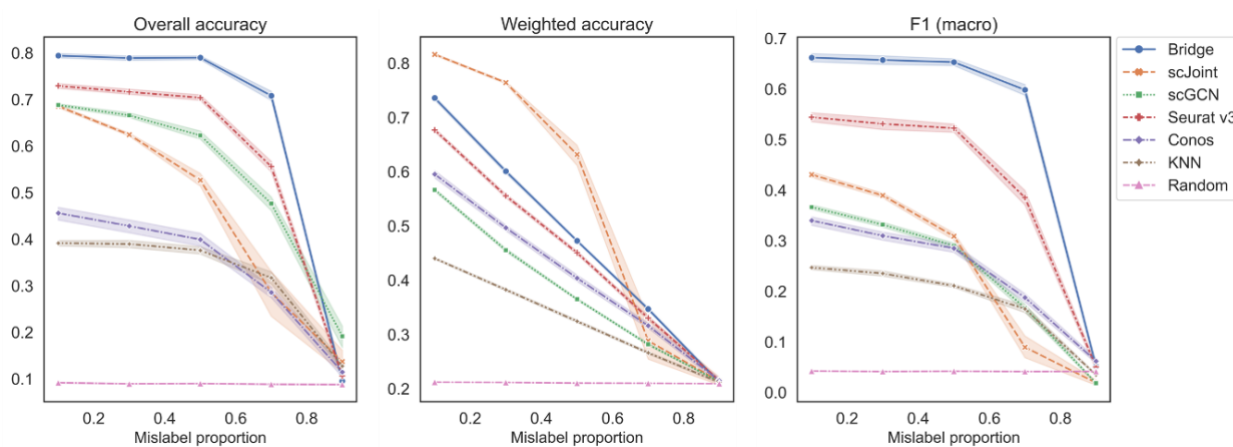
147 Figure 3A shows the performance of different methods across an increasing number of RNA cells,
 148 where KNN classifier and random classifiers were used as baseline references. We can observe that
 149 the value of three metrics didn't further increase when the cell number reached 3k, which is a
 150 relatively small number given the current high-throughput sequencing technologies. In terms of
 151 overall accuracy and F1 (macro) of precision and recall, the order of the five methods from the best
 152 to the worst were the same, which was Bridge > Seurat v3 > scJoint > scGCN > Conos. For weighted
 153 accuracy, which took into consideration the similarity among cell types (see Materials and Methods
 154 for details) when assessing the predicted probability matrix, scJoint achieved the highest score and
 155 Conos was slightly better than scGCN, while the order of the rest of the methods remained the same.

156 Conos is a graph-based method and either Seurat or Pagoda2 is recommended for data processing
157 before constructing the cell graph. We found the performance of Conos was worse when paired with
158 Seurat (Supplementary Figure S2A), resulting in both lower values of the three metrics and higher
159 instability. For Bridge integration, since it requires additional multimodal data as the ‘Bridge’, we
160 performed another set of experiments specifically for Bridge by varying the number of cells in the
161 multimodal data. We found the performance also stabilized when the cell number reached 3k and
162 Bridge was more sensitive to the smaller number of cells in the multimodal data than in the unimodal
163 RNA data (Figure 3B).

164 We recorded the running time and peak memory usage of the five methods when increasing the
165 number of RNA cells (Supplementary Figure S3A). scGCN was the most time-consuming method
166 and the second most memory-consuming method. Most of the time of running scGCN was spent on
167 processing the data where intra-data and inter-data graphs were constructed. Bridge integration
168 required the largest memory usage among all the methods because it involved additional multimodal
169 data as the bridge, while its running time was close to that of scJoint. Conos and Seurat v3 were the
170 two fastest methods and Conos was the least memory-consuming method.

171 2.3 Performance across Different Mislabeling proportions

172 The second set of experiments was designed to study the performance across different mislabeling
173 proportions of the RNA data (Figure 4). For overall accuracy and F1 (macro), their scores remained
174 constant for Bridge integration and Seurat v3 until mislabeling proportion reached 50% and
175 decreased sharply when the proportion exceeded 70%. For scGCN, scJoint and Conos, their scores
176 decreased slowly when the proportion was less than 50% and decreased faster after that. For
177 weighted accuracy, almost all methods except scJoint decreased linearly as the mislabeling
178 proportion increased. The order of the five methods was similar to the previous experiment, with
179 Bridge and Seurat v3 being the top two methods in terms of overall accuracy and F1 (macro), and
180 Conos and scGCN being the two worst-performers. scJoint was still the best method when
181 considering the weighted accuracy. We also compared the performance of Conos when paired with
182 Seurat and Pagoda2 separately and found that Conos (Seurat) was significantly worse than Conos
183 (Pagoda2) across all metrics, especially when the mislabeling proportion was low (Supplementary
184 Figure S2B).



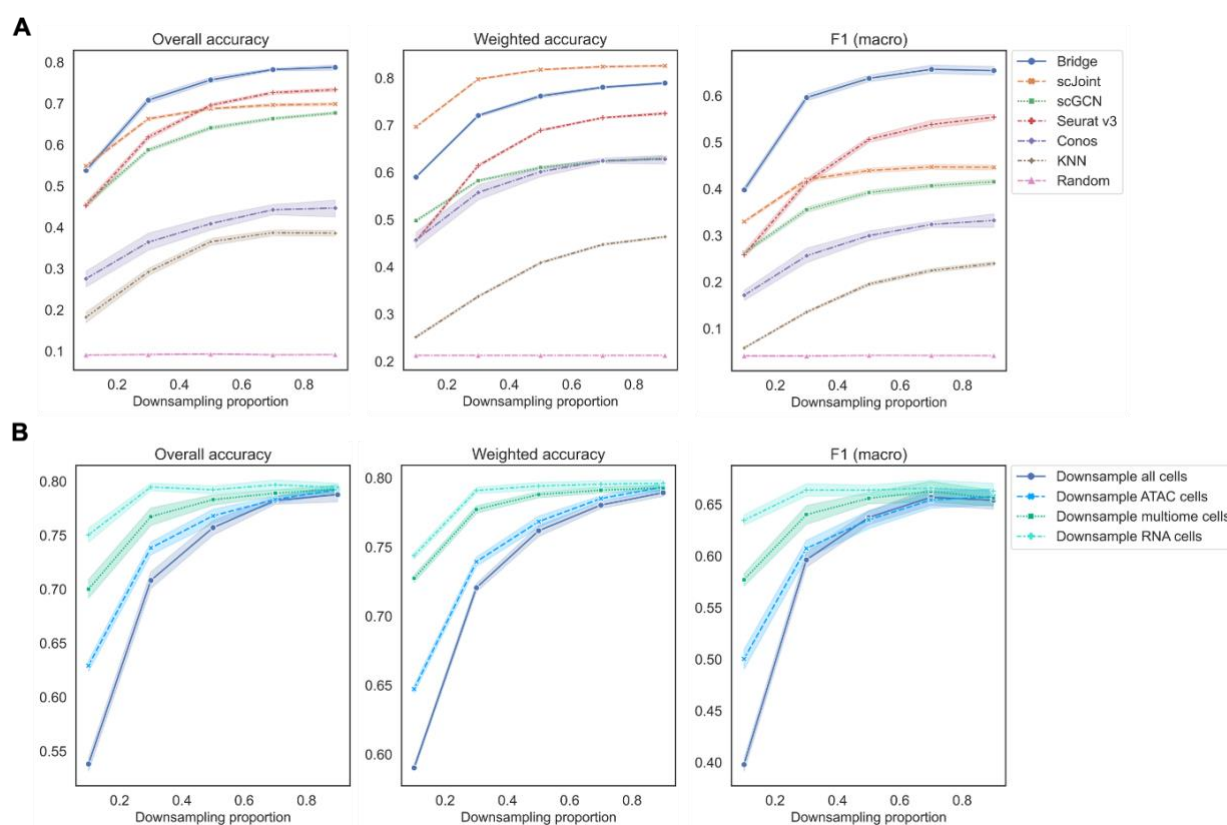
185

186 **Figure 4.** Performance of methods on different mislabeling proportions of BMMC. Results shown
187 here for Conos were paired with Pagoda2 for data processing. The error band shows the 95%
188 confidence interval.

189 2.4 Performance across Different Downsampling Proportions

190 In the third set of experiments, we downsampled each count matrix to some proportions to mimic
 191 different levels of sequencing depth. As we can observe from Figure 5A, all methods had a
 192 decreasing trend in their performance as the downsampling proportion decreased (lower sequencing
 193 depth). When the downsampling proportion was no less than 50%, the order of methods in terms of
 194 overall accuracy and F1 (macro) was the same, which was Bridge > Seurat v3 > scJoint > scGCN >
 195 Conos. When the sequencing depth was extremely low (downsampling proportion < 50%), Bridge
 196 integration was still the best performer, but the performance of scJoint became better than Seurat v3.
 197 As for weighted accuracy, scJoint was the best performer across all the methods. For Conos, its
 198 performance was worse when using Seurat for data processing compared to using Pagoda2
 199 (Supplementary Figure S2C).

200 If we downsampled cells in the RNA or ATAC data separately, similar patterns can be observed
 201 (Supplementary Figure S4). Still, Bridge integration and Seurat v3 had the highest overall accuracy
 202 and F1 (macro), and scJoint had the highest weighted accuracy. For Bridge integration, we found that
 203 at the same downsampling rate, downsampling all cells resulted in the worst performance, followed
 204 by downsampling ATAC cells, multimodal cells and RNA cells only (Figure 5B). Therefore, Bridge
 205 integration was the most sensitive to the sequencing depth of ATAC cells and least sensitive to the
 206 sequencing depth of RNA cells of the BMMC data. For other methods, they were more sensitive to
 207 the sequencing depth of ATAC cells than that of the RNA cells as well (Supplementary Figure S4).



208

209 **Figure 5.** Performance of methods on different downsampling proportions of BMMC. (A)
 210 Downsample scRNA-seq, scATAC-seq and multimodal data (for Bridge only) at the same time.
 211 Results shown here for Conos were paired with Pagoda2 for data processing. Figures for scenarios
 212 where only scRNA-seq or scATAC-seq were downsampled can be found in Supplementary Figure S4.

213 (B) Performance of Bridge integration under different downsampling scenarios. The error band
214 shows the 95% confidence interval.

215 **2.5 Performance When There Exist ATAC-Specific Cell Types**

216 The last set of experiments was designed to investigate the performance of methods when there were
217 ATAC-specific cell types by manually removing some cell types in the reference scRNA-seq data.
218 For overall accuracy and F1 (macro), Bridge integration achieved the highest scores followed by
219 scJoint and Seurat v3 with close performances, and Conos was still the worst performer
220 (Supplementary Figure S5A). The overall accuracy didn't change much as the number of removed
221 cell types in RNA increased, while F1 (macro) decreased linearly when the number of removed cell
222 types increased. For weighted accuracy, scJoint was the best method followed by Bridge integration
223 and Seurat v3. For Conos, its performance became worse when Seurat was used for data processing
224 (Supplementary Figure S5B). For Bridge integration, we found that the values of overall accuracy,
225 weighted accuracy and F1 (macro) were smaller when removing cell types in the RNA data
226 compared to removing cell types in the multimodal data given the same number of removed cell
227 types (Supplementary Figure S5C).

228 Since after removing cell types in scRNA-seq, there existed ATAC-specific cell types, we also
229 calculated the metrics designed for assessing performance of methods on these cell types. As shown
230 in the last two plots in Supplementary Figure S5A, scJoint had the highest weighted accuracy
231 followed by Bridge integration and Conos with close performances; while scGCN was the best
232 performer in terms of F1 (entropy and enrichment) and scJoint performed worst. Therefore, scJoint
233 tended to classify ATAC-specific cell types to their similar cell types in the reference data.

234 **3 Discussion**

235 We performed a comprehensive benchmarking study on five automated scATAC-seq label
236 annotations methods across five different tissues using both unimodal and multimodal single-cell
237 data. By conducting experiments on the well-annotated BMMC data, we also studied the
238 performance across different cell numbers, mislabeling proportions, sequencing depths and number
239 of unique cell types. We designed three overall metrics and two metrics for ATAC-specific cell types
240 to evaluate the prediction accuracy. In addition, we assessed the running time and memory usage of
241 each method.

242 Through the designed experiments on BMMC, we found that lower number of RNA cells, higher
243 mislabeling proportions, and lower sequencing depth could lead to worse performance of all
244 methods. When changing the number of RNA cells, we found that all methods were not sensitive to
245 the data scale when the cell number was larger than 3k. When changing the mislabeling proportion,
246 most methods had a significant decrease in overall accuracy and F1 (macro) only after the
247 mislabeling proportion reached 50%. Bridge integration was able to maintain accuracy at a high level
248 even when the mislabeling proportion was 70%. In contrast, all methods were sensitive to lower
249 sequencing depth. Across all the experimental scenarios, we found Bridge integration was
250 consistently the best performer in terms of overall accuracy and F1 (macro), and the second-best
251 performer in terms of weighted accuracy. scJoint was found to always achieve the highest weighted
252 accuracy across all experiments, suggesting it did a good job in relating similar cell types. In contrast,
253 Conos performed the worst regardless of the processing pipeline used (Seurat or Pagoda2).
254 Additionally, for Bridge integration, we found that the sequencing depth of scATAC-seq and
255 multimodal data played a more important role than the sequencing depth of scRNA-seq. This might

256 be because scATAC-seq is known to be sparser than scRNA-seq due to the limitation of current
257 sequencing technologies (Minnoye et al., 2021).

258 By benchmarking across different tissues, we found that all methods had better performance than
259 KNN and random classifiers when considering all cells. On human PBMC and BMMC where all data
260 were measured by 10x and were published no earlier than 2019, Bridge was the leading method.
261 However, for mouse lung and mouse brain, scJoint was the best performer. Note that the sequencing
262 depth of SHARE-seq mouse lung data was too low so that we were not able to assess the
263 performance of Bridge integration on mouse lung (Ma et al., 2020). For mouse brain, when applying
264 Bridge integration, we had to remap the original fastq data of unimodal ATAC data to mm10 because
265 there was inconsistency between the reference genome used for the provided unimodal ATAC (mm9)
266 and multimodal ATAC data (both SHARE-seq and SNARE-seq used mm10). After remapping, we
267 found the sequencing depth of the unimodal ATAC data was extremely low, with median count sum
268 per cell being 78 (mapped to the peak set of SNARE-seq) and 94 (mapped to the peak set of SHARE-
269 seq). While for other tissues, there were usually thousands of counts per cell (Supplementary Table
270 S3). Such high sparsity might cause the poor performance of Bridge integration on mouse brain,
271 which was consistent with the finding in BMMC experiments of changing sequencing depth. For
272 mouse kidney, Bridge integration performed relatively badly but the difference between it and other
273 methods was not significant, and the bad performance might also result from the low sequencing
274 depth of multimodal RNA data (Supplementary Table S3).

275 For performance on ATAC-specific cell types, we found scJoint consistently had the highest
276 weighted accuracy but the lowest F1 (entropy and enrichment), suggesting that it tended to classify
277 unique cell types to existing cell types that were the most similar to them. On the contrary, scGCN
278 was the best method in terms of F1 (entropy and enrichment), followed by Bridge and Seurat v3.

279 In terms of efficiency and scalability, scGCN was both time and memory consuming, and Conos was
280 the most efficient algorithm. Bridge integration required additional multimodal data, so it consumed
281 more memory than others, but its memory usage didn't increase sharply when the data scale
282 increased because it utilized dictionary learning and only performed heavy computation on a subset
283 of data (Hao et al., 2022).

284 Our study had some limitations. First, the conclusions are tissue and technology specific. Second, the
285 granularity of cell types was coarse for most tissues, like the three mouse tissues after unifying
286 annotations across datasets. The performance of methods might change if finer cell annotations were
287 provided.

288 Based on the findings in our benchmarking study, we have the following recommendations. If all
289 data are from 10x and multimodal data from the same tissue are available, Bridge integration is likely
290 the best method for label transfer; otherwise, scJoint is the to-go method. For scJoint, the caveat is
291 that it tends to misclassify ATAC-specific cell types to the biologically similar cell types in RNA. If
292 one cares about ATAC-specific cell types, a better strategy might be using scGCN or Seurat v3 and
293 another method in two separate rounds. For scGCN or Seurat v3, manual annotations can be
294 performed on cells that have high entropy and low enrichment.

295 **4 Materials and Methods**

296 **4.1 Single-cell Data Preprocessing**

297 A full list of data used in this study can be found in the Supplementary Table S1. Descriptions of
298 preprocessing pipelines specific to each dataset are provided below. Moreover, to facilitate the
299 evaluation of label prediction performance, we manually unified the naming conventions of cell
300 labels provided in the scRNA-seq and scATAC-seq (Supplementary Table S2). Details for data
301 preprocessing can be found in our GitHub repository.

302 Human BMMC. This is so far the largest single-cell multimodal RNA and ATAC dataset with well-
303 annotated labels and hierarchical batch structures. To mimic the case where scRNA-seq, scATAC-
304 seq and multimodal data were measured separately, we manually separated all batches to three
305 groups without any overlaps. Specifically, batches s1d2, s1d3, s3d3, s4d9, and s3d10 were used as
306 scRNA-seq (26,450 cells), s2d4, s2d5, s3d7, and s4d8 were used as scATAC-seq (22,653 cells), and
307 s1d1, s2d1, s4d1 were used as multimodal data (18,467 cells). Since the raw gene activity matrix was
308 not provided, the gene activity matrix for cells assigned to the scATAC-seq group was obtained using
309 Signac (Stuart et al., 2021).

310 Human PBMC. The reference genomes used for scATAC-seq (hg19) and 10x multiome ATAC-seq
311 (hg38) were different and only the latter had public raw sequence data in fastq formats. We remapped
312 the 10x multiome data using cellranger-arc to get the peak count matrix and fragment files. Since
313 Bridge integration requires that the peak sets of count matrices in scATAC-seq and multimodal
314 ATAC data are the same, we requantified the abundance of scATAC-seq peaks on the multimodal
315 peak set using the FeatureMatrix function in Signac. For the gene activity matrix, we used Signac to
316 do the calculation.

317 Mouse kidney. To unify the feature set as required by Bridge integration, we requantified the
318 scATAC-seq peaks on the multimodal peak set as what we did for human PBMC data. In addition,
319 since the gene activity matrix for mouse kidney scATAC-seq was not provided, we calculated it
320 using the GeneActivity function in Signac.

321 Mouse brain. The reference genome used for scATAC-seq (mm9) was different from that used for
322 ATAC in the two brain multimodal data (mm10). To correct the inconsistency, we used the provided
323 bam files of scATAC-seq data to map it to mm10 in three steps. First, samtools was used to convert
324 bam to fastq files. Second, fastq files were mapped to mm10 to get new bam files using bowtie2 and
325 samtools sequentially. Last, sinto was used to get fragment files from bam files. After getting
326 fragment files, Signac was used to obtain the count matrix using the peak set in the multimodal
327 ATAC data (SNARE-seq and SHARE-seq separately) and the fragment files. For the scATAC-seq
328 gene activity matrix, we used the provided one.

329 Mouse lung. We did not find an appropriate multimodal data for mouse lung, so the data for this
330 tissue were only used to benchmark methods that do not require multimodal data (only Bridge
331 integration requires). For the gene activity matrix, we used the one provided by the original paper.

332 **4.2 Description and Implementation of Methods**

333 Conos. Conos is designed as a graph-based batch effect removal method. The joint graph embedding
334 using nearest neighbors and Pearson correlation is constructed as the first step to connect all cells.
335 Then, the label transfer from reference data to query data can be implemented by information
336 propagation between graph vertices through an iterative diffusion process.

337 Seurat v3. Seurat first identifies a set of anchors between the reference and the query data through
338 canonical correlation analysis (CCA) and mutual nearest neighbors (MNNs). Then, a weight matrix

339 is constructed to quantify the distance between each query cell and anchor cell in the query data by a
340 Gaussian kernel. Last, the prediction score of any cell in the query data is calculated as a weighted
341 average of labels of anchor cells in the reference data.

342 scGCN. The first step of scGCN is to build a hybrid graph of all cells using MNNs approach and
343 CCA. Based on the constructed graph, a semi-supervised graph convolutional neural network is
344 trained to embed cells from both reference and query data on the same latent space and predict cell
345 type labels for cells in the query data.

346 scJoint. Like scGCN, a semi-supervised neural network with cross entropy loss is trained to jointly
347 embed cells from both scRNA-seq and scATAC-seq. Different from scGCN that directly utilizes the
348 trained network to predict probability vectors through Softmax layers, scJoint performs label transfer
349 by training an additional kNN classifier in the embedding space.

350 Bridge integration. This method utilizes multimodal data as a bridge to transfer labels from scRNA-
351 seq to scATAC-seq. The multimodal dataset is treated as a dictionary and each cell is an atom, on
352 which dictionary representations of both unimodal scRNA-seq and scATAC-seq are constructed.
353 After dimensionality reduction of multimodal cells via Laplacian Eigendecompositions, unimodal
354 cells can be embedded on the same space by the dictionary representations. Then, the final label
355 transfer can be achieved by any single-cell integration techniques and Bridge integration chooses
356 mnnCorrect.

357 For Conos, Seurat v3, scGCN and scJoint, the raw count matrix of scRNA-seq and gene activity
358 score matrix of scATAC-seq were provided as inputs. In addition, the raw count matrix of scATAC-
359 seq was provided for Seurat v3 to perform dimension reduction. For Bridge integration, since the
360 information transfer was realized by using the multimodal data as a bridge, the gene activity matrix
361 was not needed. Instead, we provided raw count matrices of scRNA-seq, scATAC-seq (mapped to
362 the same peak set of multimodal ATAC data) and multimodal data for Bridge integration. The
363 implementation of each method followed the instructions on their websites. Details can be found in
364 the scripts on our GitHub repository and package versions can be found in Supplementary Table S4.

365 **4.3 Benchmarking Design**

366 To investigate the model performance across different cell numbers, mislabeling proportions,
367 sequencing depths and number of unique cell types, we designed the following set of experiments
368 based on the human BMMC multimodal data. For each specific setting, 20 replicates were generated
369 using unique random seeds.

370 Change data scale. This was separated into three sub-experimental designs. (1) Change the cell
371 numbers in scRNA-seq (reference) while keeping the scATAC-seq and the multimodal cell numbers
372 (for Bridge integration) as 10k. The chosen numbers were 0.2k, 0.6k, 1k, 3k, 5k, and 10k. (2) Change
373 the cell numbers in the multimodal data while keeping the scRNA-seq and scATAC-seq cell numbers
374 as 10k. The chosen numbers were 0.2k, 0.6k, 1k, 3k, 5k, and 10k. This setting was used for Bridge
375 integration specifically.

376 Change mislabeling proportion. The mislabeling proportions for scRNA-seq cells were chosen as
377 10%, 30%, 50%, 70%, and 90%. Mislabeled cells were randomly selected and assigned wrong labels
378 based on the background compositions of other cell labels.

379 Change sequencing depth. The sequencing depths were manually changed by downsampling reads to
380 10%, 30%, 50%, 70%, and 90% of the original number of reads using R package DropletUtils
381 (Griffiths et al., 2018; Lun et al., 2019). We set four different scenarios under this experiment, which
382 are changing sequencing depth in (1) all cells, (2) RNA cells, (3) ATAC cells, and (4) multiome cells
383 (for Bridge integration).

384 Change the number of unique cell types. We randomly removed 2, 4 or 6 selected cell types in the
385 scRNA-seq data. Candidate cell types were those whose cell numbers were between 200 and 1,000.

386 4.4 Evaluation Metrics

387 Accuracy. After getting the predicted probability matrix across all cells in scATAC-seq, the cell type
388 that had the highest predicted probability was assigned to each cell as the predicted label. Then the
389 overall accuracy was calculated using the predicted labels and true labels.

390 Weighted accuracy. To account for the prediction uncertainty and similarity across cell types. We
391 proposed a weighted accuracy (WACC) by taking the average of the predicted probability vector
392 weighted by cell type similarities.

$$393 \quad WACC = 1/N \sum_i \sum_{j \in C_R} S_{c(i),j} P_{i,j}$$

394 In the equation above, P is the predicted probability matrix with each row as a cell in scATAC-seq
395 and each column as a cell type observed in scRNA-seq reference data. C_R is the set of all cell types in
396 scRNA-seq and N is the total number of scATAC-seq cells. S is a cross-modality cell type similarity
397 matrix with each row as a cell type in scATAC-seq and each column as a cell type in scRNA-seq and
398 $c(i)$ is a function mapping cell i to its true cell type label.

399 The similarity matrix was calculated in three steps. First, partition-based graph abstraction (PAGA)
400 (Wolf et al., 2019) was performed on the normalized count matrix of scRNA-seq and gene activity
401 matrix of scATAC-seq separately. Then, the within-modality similarity matrix was calculated based
402 on the Euclidean distance of each pair of cell types using the PAGA positions. Specifically, we
403 applied \exp to the negative of the calculated distance matrix. Last, we calculated the cross-modality
404 similarity matrix using the two within-modality matrices by considering three scenarios. If two cell
405 types existed in both modalities, their similarity was calculated as the average of two within-modality
406 similarities. If one cell type is modality-specific, its similarity with any common cell type would be
407 the similarity calculated using the modality that contained the two cell types. If a cell type l only
408 existed in scATAC-seq and the other cell type k was only observed in scRNA-seq, their similarity
409 was calculated as

$$410 \quad S_{l,k} = [S^{ATAC}_{l,common} \circ \mathbf{1}\{S^{ATAC}_{l,common} \geq .5\}] S^{RNA}_{common,k} / \sum_{i \in common} \mathbf{1}\{S^{ATAC}_{l,i} \geq .5\}$$

411 where S^{ATAC} and S^{RNA} are within-modality similarity matrix for ATAC and RNA, respectively and
412 $common$ is the set of all common cell types. The first product is Hadamard product which is element
413 wise and the second product is matrix multiplication.

414 Precision, recall and F1 score. Precision is defined as true positive (TP) over the summation of TP
415 and false positive (FP) and recall is defined as TP over the summation of TP and false negative (FN).
416 F1 score is the harmonic mean of precision and recall,

$$417 \quad F_1 = 2 \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}.$$

418 Since this is a multi-class classification problem, we need to specify whether we want macro or
419 micro level metrics. It's easy to show that overall accuracy is equivalent to micro precision, recall
420 and F1 score under the multi-class scenario. Therefore, we calculated macro level precision and
421 recall in this study, which is the average of precisions and recalls obtained for each class. Then,
422 macro F1 score is calculated based on macro precision and recall.

423 Entropy and enrichment. To evaluate the performance of methods on cell types unique to scATAC-
424 seq data, we borrowed the two metrics proposed in scGCN which are scaled entropy and enrichment
425 (Song et al., 2021). Scaled entropy is defined as

$$426 \quad NE = -\frac{1}{M \log_2 |C_R|} \sum_i \sum_{j \in C_R} \frac{S_{i,j}}{\sum_{j \in C_R} S_{i,j}} \log_2 \frac{S_{i,j}}{\sum_{j \in C_R} S_{i,j}}, \text{ where } S_{i,j} = \frac{P_{i,j}}{Q_j}.$$

427 $P_{i,j}$ is the predicted probability for cell i with unique cell type label in scATAC-seq and cell type j ,
428 and Q_j is the proportion of cell type j in scRNA-seq as the background probability. C_R is the set of all
429 cell types in scRNA-seq and M is the total number of scATAC-seq cells with unique cell labels. The
430 final score is normalized by $\log_2 |C_R|$ to make it in the range of [0, 1]. Another metric is enrichment
431 score,

$$432 \quad ES = \frac{1}{M} \sum_i \max_{j \in C_R} \frac{S_{i,j}}{\sum_{j \in C_R} S_{i,j}}.$$

433 The enrichment score is also bounded within 0 and 1. For cell types only observed in scATAC-seq,
434 an ideal method should deliver high normalized entropy and low enrichment score. Therefore, we
435 also calculated an F1 score to combine these two

$$436 \quad F_1 = 2 \frac{NE \cdot (1-ES)}{NE + (1-ES)}.$$

437 Running time and memory. All methods were run on Yale's high performance computing clusters
438 with one computing core. For neural network methods scGCN and scJoint, they were run using
439 GPUs; and for the rest methods, they were run using CPUs. The CPU of our device is Intel® Xeon
440 ® Gold 6240, 2.6 GHz, and the GPU is NVIDIA RTX 3090 with 25 GB RAM. When evaluating
441 running time, we did not count the time used for data preprocessing (e.g. remap to alternative
442 reference genome, requantify scATAC-seq peaks, and calculate gene activity matrix) because the
443 needed steps for different tissues were different. For memory assessment, we used the recorded peak
444 memory usage of each method.

445 **5 Conflict of Interest**

446 The authors declare that the research was conducted in the absence of any commercial or financial
447 relationships that could be construed as a potential conflict of interest.

448 **6 Data Availability Statement**

449 All the single-cell data used in this manuscript are publicly available. Detailed information of each
450 data and their downloadable links can be found in Supplementary Table S1. The related scripts for
451 reproducing results in this manuscript are available on GitHub at
452 <https://github.com/AprilYuge/ATAC-annotation-benchmark>.

453 **7 Author Contributions**

454 YW collected data, performed label unification and similarity matrix calculation, designed the
455 benchmarking pipeline and evaluation metrics, assisted in preparing scripts for running each method,
456 evaluated the model performance, and wrote the manuscript. XS wrote scripts for running each
457 method, performed data processing and gene activity calculation, assisted in model evaluation, and
458 provided feedback to the manuscript. HZ supervised the entire project, contributed to the design of
459 the benchmarking pipeline and evaluation metrics, revised the manuscript critically for important
460 intellectual content and provided approval for the publication of this manuscript.

461 **8 Funding**

462 This study was supported in part by NIH grants R56 AG074015 and P50 CA196530.

463 **9 Supplementary Material**

464 The Supplementary Material for this article can be found online at the bioRxiv website.

465 **10 References**

- 466 Abdelaal, T., Michielsen, L., Cats, D., Hoogduin, D., Mei, H., Reinders, M.J., et al. (2019). A
467 comparison of automatic cell identification methods for single-cell RNA sequencing data.
468 *Genome biology* 20(1), 1-19.
- 469 Barkas, N., Petukhov, V., Nikolaeva, D., Lozinsky, Y., Demharter, S., Khodosevich, K., et al. (2019).
470 Joint analysis of heterogeneous single-cell RNA-seq dataset collections. *Nature methods*
471 16(8), 695-698.
- 472 Buenrostro, J.D., Corces, M.R., Lareau, C.A., Wu, B., Schep, A.N., Aryee, M.J., et al. (2018).
473 Integrated single-cell analysis maps the continuous regulatory landscape of human
474 hematopoietic differentiation. *Cell* 173(6), 1535-1548. e1516.
- 475 Buenrostro, J.D., Wu, B., Litzenburger, U.M., Ruff, D., Gonzales, M.L., Snyder, M.P., et al. (2015).
476 Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*
477 523(7561), 486-490.
- 478 Cao, J., Cusanovich, D.A., Ramani, V., Aghamirzaie, D., Pliner, H.A., Hill, A.J., et al. (2018). Joint
479 profiling of chromatin accessibility and gene expression in thousands of single cells. *Science*
480 361(6409), 1380-1385.
- 481 Carter, B., and Zhao, K. (2021). The epigenetic basis of cellular heterogeneity. *Nature Reviews*
482 *Genetics* 22(4), 235-250.
- 483 Chen, S., Lake, B.B., and Zhang, K. (2019). High-throughput sequencing of the transcriptome and
484 chromatin accessibility in the same cell. *Nature biotechnology* 37(12), 1452-1457.
- 485 Clarke, Z.A., Andrews, T.S., Atif, J., Pouyababar, D., Innes, B.T., MacParland, S.A., et al. (2021).
486 Tutorial: guidelines for annotating single-cell transcriptomic maps using automated and
487 manual methods. *Nature protocols* 16(6), 2749-2764.

- 488 Consortium, T.M. (2018). Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris.
489 *Nature* 562(7727), 367-372.
- 490 Cusanovich, D.A., Daza, R., Adey, A., Pliner, H.A., Christiansen, L., Gunderson, K.L., et al. (2015).
491 Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing.
492 *Science* 348(6237), 910-914.
- 493 Cusanovich, D.A., Hill, A.J., Aghamirzaie, D., Daza, R.M., Pliner, H.A., Berletch, J.B., et al. (2018).
494 A single-cell atlas of in vivo mammalian chromatin accessibility. *Cell* 174(5), 1309-1324.
495 e1318.
- 496 Fiers, M.W., Minnoye, L., Aibar, S., Bravo González-Blas, C., Kalender Atak, Z., and Aerts, S.
497 (2018). Mapping gene regulatory networks from single-cell omics data. *Briefings in*
498 *functional genomics* 17(4), 246-254.
- 499 Granja, J.M., Klemm, S., McGinnis, L.M., Kathiria, A.S., Mezger, A., Corces, M.R., et al. (2019).
500 Single-cell multiomic analysis identifies regulatory programs in mixed-phenotype acute
501 leukemia. *Nature biotechnology* 37(12), 1458-1465.
- 502 Griffiths, J.A., Richard, A.C., Bach, K., Lun, A.T., and Marioni, J.C. (2018). Detection and removal
503 of barcode swapping in single-cell RNA-seq data. *Nature communications* 9(1), 1-6.
- 504 Hao, Y., Stuart, T., Kowalski, M., Choudhary, S., Hoffman, P., Hartman, A., et al. (2022). Dictionary
505 learning for integrative, multimodal, and scalable single-cell analysis. *bioRxiv*.
- 506 Jia, G., Preussner, J., Chen, X., Guenther, S., Yuan, X., Yekelchik, M., et al. (2018). Single cell
507 RNA-seq and ATAC-seq analysis of cardiac progenitor cell transition states and lineage
508 settlement. *Nature communications* 9(1), 1-17.
- 509 Lin, Y., Wu, T.-Y., Wan, S., Yang, J.Y., Wong, W.H., and Wang, Y. (2022). scJoint integrates atlas-
510 scale single-cell RNA-seq and ATAC-seq data with transfer learning. *Nature Biotechnology*
511 40(5), 703-710.
- 512 Luecken, M.D., Burkhardt, D.B., Cannoodt, R., Lance, C., Agrawal, A., Aliee, H., et al. (Year). "A
513 sandbox for prediction and integration of dna, rna, and proteins in single cells", in: *Thirty-fifth*
514 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*
515 *(Round 2)*.
- 516 Lun, A.T., Riesenfeld, S., Andrews, T., Gomes, T., and Marioni, J.C. (2019). EmptyDrops:
517 distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data.
518 *Genome biology* 20(1), 1-9.
- 519 Ma, S., Zhang, B., LaFave, L.M., Earl, A.S., Chiang, Z., Hu, Y., et al. (2020). Chromatin potential
520 identified by shared single-cell profiling of RNA and chromatin. *Cell* 183(4), 1103-1116.
521 e1120.
- 522 Miao, Z., Balzer, M.S., Ma, Z., Liu, H., Wu, J., Shrestha, R., et al. (2021). Single cell regulatory
523 landscape of the mouse kidney highlights cellular differentiation programs and disease
524 targets. *Nature communications* 12(1), 1-17.
- 525 Minnoye, L., Marinov, G.K., Krausgruber, T., Pan, L., Marand, A.P., Secchia, S., et al. (2021).
526 Chromatin accessibility profiling methods. *Nature Reviews Methods Primers* 1(1), 1-24.
- 527 Packer, J., and Trapnell, C. (2018). Single-cell multi-omics: an engine for new quantitative models of
528 gene regulation. *Trends in Genetics* 34(9), 653-665.

- 529 Pasquini, G., Arias, J.E.R., Schäfer, P., and Busskamp, V. (2021). Automated methods for cell type
530 annotation on scRNA-seq data. *Computational and Structural Biotechnology Journal* 19,
531 961-969.
- 532 Song, Q., Su, J., and Zhang, W. (2021). scGCN is a graph convolutional networks algorithm for
533 knowledge transfer in single cell omics. *Nature communications* 12(1), 1-11.
- 534 Stuart, T., Butler, A., Hoffman, P., Hafemeister, C., Papalexi, E., Mauck III, W.M., et al. (2019).
535 Comprehensive integration of single-cell data. *Cell* 177(7), 1888-1902. e1821.
- 536 Stuart, T., Srivastava, A., Madad, S., Lareau, C.A., and Satija, R. (2021). Single-cell chromatin state
537 analysis with Signac. *Nature methods* 18(11), 1333-1341.
- 538 Wang, Y., Chen K., Cai Z., and Zhao H. (2022). Gene regulatory network inference using single-cell
539 multiome ATAC-seq and RNA-seq data (Abstract). Presented at the Annual Meeting of The
540 American Society of Human Genetics, October 26, 2022 in Los Angeles, CA.
- 541 Wolf, F.A., Hamey, F.K., Plass, M., Solana, J., Dahlin, J.S., Göttgens, B., et al. (2019). PAGA: graph
542 abstraction reconciles clustering with trajectory inference through a topology preserving map
543 of single cells. *Genome biology* 20(1), 1-9.

544