# Benchmarking Algorithms for Gene Set Scoring of Single-cell ATAC-seq Data

Xi Wang[1,2,#], Qiwei Lian[1,2,#], Haoyu Dong[1], Shuo Xu[2], Yaru Su[3], Xiaohui Wu[1,*]

*[1] Pasteurien College, Suzhou Medical College of Soochow University, Soochow University, Suzhou 215000, China*

*[2] Department of Automation, Xiamen University, Xiamen 361005, China*

*[3] College of Mathematics and Computer Science, Fuzhou University, Fuzhou 350116, China*

[#] Equal contribution.

* Corresponding author.

E-mail: xhwu@suda.edu.cn (Wu X).

**Running title:** *Wang X et al / Benchmarking for Gene Set Scoring of scATAC-seq*

Word count: 7937

Keyword count: 5

Figure count: 9

Table count: 0

Supplementary Figures: 1

Supplementary Tables: 3

Supplementary Files: 0

Reference Count: 58

Article Title: 64 characters

Running Title: 41 characters

Abstract: 238 words

## Abstract

Gene set scoring (GSS) has been routinely conducted for gene expression analysis of bulk or single-cell RNA-seq data, which helps to decipher single-cell heterogeneity and cell-type-specific variability by incorporating prior knowledge from functional gene sets. Single-cell assay for transposase accessible chromatin using sequencing (scATAC-seq) is a powerful technique for interrogating single-cell chromatin-based gene regulation, and genes or gene sets with dynamic regulatory potentials can be regarded as cell-type specific markers as if in scRNA-seq. However, there are few GSS tools specifically designed for scATAC-seq, and the applicability and performance of RNA-seq GSS tools on scATAC-seq data remain to be investigated. We systematically benchmarked ten GSS tools, including four bulk RNA-seq tools, five single-cell RNA-seq (scRNA-seq) tools, and one scATAC-seq method. First, using matched scATAC-seq and scRNA-seq datasets, we find that the performance of GSS tools on scATAC-seq data is comparable to that on scRNA-seq, suggesting their applicability to scATAC-seq. Then the performance of different GSS tools were extensively evaluated using up to ten scATAC-seq datasets. Moreover, we evaluated the impact of gene activity conversion, dropout imputation, and gene set collections on the results of GSS. Results show that dropout imputation can significantly promote the performance of almost all GSS tools, while the impact of gene activity conversion methods or gene set collections on GSS performance is more GSS tool or dataset dependent. Finally, we provided practical guidelines for choosing appropriate pre-processing methods and GSS tools in different scenarios.

**Keywords:** Single-cell ATAC-seq; Gene set scoring; Pathway analysis; Single-cell RNA-seq; Benchmark

## Introduction

Assay for Transposase-Accessible Chromatin using sequencing (ATAC-seq) is a powerful and the most widely used epigenomic technique for interrogating chromatin accessibility on a genome-wide scale [1]. In particular, the advent of single-cell ATAC-seq (scATAC-seq) has made it possible to profile chromatin-accessibility variations in single cells, which allows to illuminate chromatin-based gene regulation with an unprecedented cellular resolution and discover new cell subpopulations [2, 3]. One of the ultimate goals for analyzing single-cell chromatin accessibility data is to quantitatively understand the relationship between the variation of chromatin accessibility and that of the expression of nearby genes [4]. A first step toward this goal is to link regulatory DNA elements with their target genes on a genome-wide scale and predict gene activity (GA) score by modelling the chromatin accessibility at the gene

59   level. Several tools are currently in progress to convert chromatin accessibility signals to GA

60   scores, including Cicero [4], MAESTRO [5], ArchR [6], snapATAC [7], and Signac [8]. The

61   inferred GA scores facilitate the integrative analysis of single-cell RNA-seq (scRNA-seq) and

62   scATAC-seq data, and the scores of key marker genes can be used for accurate annotation of

63   cell types as if in scRNA-seq [4, 6, 9].

64       In addition to single gene analysis, gene set analysis, analogue to pathway analysis, has

65   become a routine step for analyzing gene expression data, which has proven to be effective in

66   estimating the activity of pathways or transcription factors (TFs) for uncovering transcriptional

67   heterogeneity and disease subtypes [10-12]. In single-cell RNA-seq studies, gene set scoring

68   (GSS), or commonly referred to as pathway activity transformation, has been broadly

69   conducted to quantify the enrichment and relevance of gene sets in individual cells. GSS

70   converts the gene-level data into gene set-level information; gene sets contain genes

71   representing distinct biological processes (e.g., the same Gene Ontology annotation) or

72   pathways (e.g., the Molecular Signature Database (MSigDB) [13]). Therefore, GSS helps to

73   decipher single-cell heterogeneity and cell-type-specific variability by incorporating prior

74   knowledge from functional gene sets or pathways [14, 15]. A wide spectrum of GSS tools have

75   been designed for scRNA-seq data, such as Pagoda2 [16], Vision [17], and AUCell [18], which

76   infer pathway-level information from the gene expression profile for the characterization of

77   transcriptional heterogeneity of cell populations. Similarly, gene sets with dynamic regulatory

78   potentials inferred from scATAC-seq can also be regarded as cell-type specific markers as if

79   in scRNA-seq [5].

80       Single-cell ATAC-seq data and RNA-seq data have analogous characteristic structures,

81   both of which suffer from similar sparsity and noise. In recent years, great breakthroughs have

82   been made in the computational modelling of scRNA-seq data, such as dropout imputation,

83   dimensionality reduction, cell type identification, GSS, and regulatory networks inference [19-

84   22]. In contrast, the progress on computational modelling in the field of scATAC-seq lags far

85   behind that of scRNA-seq [23, 24]. As a compromise, many scRNA-seq analysis methods are

86   directly applied to scATAC-seq data. For example, Liu et al. [25] benchmarked tools dedicated

87   to imputing scRNA-seq data (e.g., MAGIC [26] and SAVER [27]) for recovering dropout

88   peaks in scATAC-seq data and found that most scRNA-seq imputation tools can be readily

89   applied to scATAC-seq data. Tools for alignment, quality control, peak calling, and differential

90   peak analysis for RNA-seq and/or ChIP-seq data are widely used for ATAC-seq data [23]. This

91   series of evidence indicates that GSS tools for scRNA-seq could in principle be applicable to

92   scATAC-seq as well. However, due to the close-to-binary nature and extreme sparsity of the

93    scATAC-seq data, it remains elusive whether these limitations would distort or confound the

94    results produced by the direct application of RNA-seq methods to scATAC-seq. To the best of

95    our knowledge, currently only one tool, UniPath [28], provides a function dedicated to scoring

96    gene sets for scATAC-seq, therefore, it is timely and imperative to further investigate the

97    applicability and performance of more GSS tools designed for bulk or single-cell RNA-seq on

98    scATAC-seq data.

99        Currently the performance of GSS tools designed for bulk or single-cell RNA-seq on

100   scRNA-seq data sequenced with diverse scRNA-seq protocols has been comprehensively

101   evaluated. Zhang et al. [15] evaluated the performance of eleven pathway activity

102   transformation tools on 32 scRNA-seq datasets and found Pagoda2 [16] exhibited the best

103   overall performance. Holland et al. [14] compared the performance of six TFs or pathway

104   activity estimators on simulated and real scRNA-seq data, which found that bulk tools can be

105   applied to scRNA-seq, partially outperforming scRNA-seq tools. These studies focused only

106   on scRNA-seq, to the best of our knowledge, there has been no systematic benchmark study to

107   evaluate the performance of GSS tools on scATAC-seq data. Here we systematically evaluated

108   the performance of ten GSS tools using ten scATAC-seq datasets, including four tools designed

109   for bulk RNA-seq, five tools designed for scRNA-seq, and one method proposed for scATAC-

110   seq. The performance was quantitively evaluated under four scenarios of dimensionality

111   reduction, clustering, classification, and cell type determination, which are critical steps of

112   single-cell analysis in most scRNA-seq and scATAC-seq studies. Our benchmark results

113   provide abundant evidence that GSS tools designed for RNA-seq are also applicable to

114   scATAC-seq. Using three matched scATAC-seq and scRNA-seq datasets, results showed that

115   the performance of GSS tools for scATAC-seq data on clustering cells or distinguishing cell

116   types was comparable to that for scRNA-seq. In particular, the performance of several GSS

117   tools designed for RNA-seq exceeds the current only method dedicated to scATAC-seq, under

118   diverse evaluation scenarios. Moreover, we evaluated the impact of data preprocessing of

119   scATAC-seq on GSS, including dropout imputation and GA transformation. Benchmark results

120   show that dropout imputation can significantly promote the performance of almost all GSS

121   tools. In contrast, the performance of different GA transformation methods varies greatly across

122   different GSS tools and different datasets. In addition, we also evaluated the performance of

123   GSS tools using different gene set collections in the context of clustering and found that

124   different GSS tools and different datasets have different degrees of robustness to different gene

125   collections. Our benchmark results provide practical guidelines for choosing appropriate GSS

126 tools for raw scATAC-seq data or data after dropout imputation, and also provide important
127 clues on how to preprocess the scATAC-seq data for more effective GSS.

128 **Results**

129 **Overview of the benchmark workflow**

130 We benchmarked ten GSS tools, including four tools for bulk RNA-seq (PLAGE [29], z-score
131 [30], ssGSEA [31], and GSVA [32]), five tools for scRNA-seq (AUCell [18], Pagoda2 [16],
132 Vision [17], VAM [33], and UniPath [28]) and one function provided in the UniPath for scoring
133 gene sets from scATAC-seq (hereinafter called UniPathATAC), using ten real scATAC-seq
134 datasets with different number of cells and cell types (Figure 1). UniPathATAC can score gene
135 sets directly from scATAC-seq data, using the peak-cell matrix as the input to obtain the gene
136 set score matrix. In contrast, the input of RNA-seq GSS tools is the gene-cell matrix, thus the
137 peak-level profile obtained from scATAC-seq data needs to be converted into the GA matrix,
138 using a GA transformation tool. Four GA tools, including MAESTRO [5], Signac [8], ArchR
139 [6], and snapATAC [7], were examined. MAESTRO obtains the GA matrix from the peak-cell
140 matrix, while other three GA tools from the fragment file (Materials and methods). Unless
141 otherwise specified, Signac was used as the default GA conversion tool as it runs fast and has
142 good performance in our preliminary test. But we also conducted in-depth evaluation on the
143 impact of different GA tools on GSS. Moreover, the pipeline for evaluating GSS tools involves
144 an additional preprocessing step -- imputation of dropout peaks. We adopted three popular
145 imputation tools developed for scRNA-seq (MAGIC [26], DrImpute [34], and SAVER [27])
146 and one tool designed for scATAC-seq (SCALE [35]). It should be noted that the imputation is
147 performed on the peak-cell matrix rather than the fragment file, therefore, only MAESTRO [5]
148 can be used for GA conversion from the imputed data. In addition, we examined six gene set
149 collections from MSigDB (version 7.1), including KEGG (Kyoto Encyclopedia of Genes and
150 Genomes), GO:BP (GO Biological Process), GO:MF (GO Molecular Function), GO:CC (GO
151 Cellular Component), REACTOME, and TFT (Transcription Factor Target) (Table S1). Unless
152 otherwise specified, KEGG that contains 186 gene sets in MSigDB was used as the default
153 prior information. We benchmarked GSS tools under diverse scenarios of dimensionality
154 reduction, clustering, classification and cell type determination. Each GSS tool was used to
155 obtain the gene set score matrix from each scATAC-seq dataset (hereafter called GSS-ATAC),
156 which was then evaluated in the context of each evaluation scenario.

157 **GSS tools designed for RNA-seq are applicable to scATAC-seq**

158 We used three matched datasets of scATAC-seq and scRNA-seq that are derived from the same

159 cells, including Brain, PBMC3K, and PBMC10K (Table S2), to examine whether GSS tools

160 designed for RNA-seq are applicable to scATAC-seq data. First, we used Signac to convert the

161 peak-cell matrix to the GA matrix and then performed each GSS tool to obtain the GSS-ATAC

162 matrix. We also used the nine RNA-seq GSS tools to score gene sets for the matched scRNA-

163 seq data to obtain the corresponding gene set score matrix for scRNA-seq (hereafter called

164 GSS-RNAseq). Then the performances of different GSS tools were evaluated by

165 dimensionality reduction measured by Silhouette, clustering measured by ARI and

166 classification measured by accuracy based on the GSS-ATAC or the GSS-RNAseq matrix

167 obtained by different tools. We conducted the pipeline for each dataset and then calculated the

168 average value of each performance indicator of the three datasets. For both scRNA-seq and

169 scATAC-seq data, two methods, Pagoda2 and PLAGE, generally provide better performance

170 than other methods in terms of all the three performance indicators (Figure 2A). Other GSS

171 tools exhibit comparable and moderate performance. Although the performance of GSS tools

172 on scRNA-seq and scATAC-seq is comparable, most GSS tools provide slightly better

173 performance on scRNA-seq than on scATAC-seq. This is not unexpected because that these

174 tools, except for UniPathATAC, were designed for RNA-seq and the reference cell types of

175 scATAC-seq datasets were determined by the scRNA-seq data rather than scATAC-seq. Still,

176 the consistency between clustering results obtained by GSS-ATAC and the reference cell types,

177 measured by ARI, is even slightly higher than that of GSS-RNAseq obtained by several tools,

178 including GSVA, VAM, and Vision (GSVA: 0.51 *vs.* 0.47; VAM: 0.38 *vs.* 0.36; Vision: 0.50

179 *vs.* 0.49). In particular, the performance of the two tools with the best performance, Pagoda2

180 and PLAGE, is higher than UniPathATAC, a method designed specifically for scATAC-seq,

181 under all evaluation schemes (e.g., ARI of Pagoda2 = 0.60, PLAGE = 0.57, UniPathATAC =

182 0.55). Moreover, 2D embeddings of both GSS-ATAC and GSS-RNAseq matrices obtained by

183 different GSS tools show comparable discrimination of the cell types (Figure 2B).

184 In addition to Signac, we also used three other GA tools for transforming the peak profile

185 to the gene-level activity scores and then calculated the GSS-ATAC matrix using different GSS

186 tools. Results on the PBMC10K data show that the GSS-ATAC matrix based on the GA matrix

187 obtained by different GA tools yields comparable ARI score to that using scRNA-seq data

188 (Figure 2C), demonstrating again the applicability of RNA-seq GSS tools to scATAC-seq.

189 Among the four GA tools, ArchR is less robust than other three GA tools for the PBMC10K

190 data (Figure 2C). Taken together, these results preliminarily show that GSS tools designed for

191 RNA-seq have comparable performance on both scRNA-seq and scATAC-seq data and thus

192   are applicable to scATAC-seq data. In the following benchmark evaluation, we used more

193   scATAC-seq datasets and considered different factors, including pre-processing steps, gene set

194   collections, and GA methods, to evaluate different GSS tools more comprehensively.

195   **Evaluation of GSS tools using different scATAC-seq datasets**

196   Having preliminarily demonstrated that GSS tools designed for RNA-seq are applicable to

197   scATAC-seq, next we used eight scATAC-seq datasets (Table S2), which are from human and

198   mouse with number of cells ranging from 500 to 10K, to further evaluate the performance of

199   different GSS tools.  Generally, the performance of GSS tools is highly dependent on datasets

200   (Figure 3). Regardless of the evaluation scenario, the performance of all tools on

201   Hematopoiesis, Leukemia, and SNAREmix is extremely poor, significantly lower than that on

202   other five datasets. We then examined the raw scATAC-seq data to check whether the datasets

203   with generally poor GSS results have low data quality. Indeed, we found significantly lower

204   consistency between the clusters and the reference cell types of the three datasets with poor

205   GSS results than the other five datasets (Figure S1). Although different GSS tools have varied

206   performance on different datasets, Pagoda2 and PLAGE perform overall better than other tools.

207   For example, the average ARI scores of all the eight datasets of Pagoda2 and PLAGE are much

208   higher than that of the third tool UniPathATAC (Pagoda2 = 0.32, PLAGE = 0.30,

209   UniPathATAC = 0.24). Of note, UniPathATAC is specially designed for scATAC-seq.

210   Similarly, according to the scenario of classification, the average accuracy of PLAGE and

211   Pagoda2 is also much higher than other tools (PLAGE=0.72, Pagoda2=0.67, other tools=0.62).

212   These results revealed that the performance of the scATAC-seq specific tool, UniPathATAC,

213   is only moderate, which is generally lower than that of two GSS tools for RNA-seq, Pagoda2

214   and PLAGE, suggesting again the feasibility of applying RNA-seq GSS tools to scATAC-seq

215   data.

216   **Evaluation of the impact of dropout imputation on GSS**

217   Similar to scRNA-seq, scATAC-seq is plagued by extremely high sparsity and noise, therefore

218   single-cell dropout peaks are usually recovered before downstream analysis. In contrast to the

219   considerable progress that has been made in dropout imputation of scRNA-seq data, much

220   fewer imputation tools for scATAC-seq are available. Till now, SCALE [35] is the only

221   imputation method specially designed for scATAC-seq. A previous benchmark study [25]

222   suggested that imputation tools designed for scRNA-seq are also applicable to scATAC-seq.

223   Therefore, in addition to SCALE, we also considered three widely used scRNA-seq imputation

224   tools, including MAGIC, DrImpute, and SAVER. Of note, the recovered peak-cell matrix can

225   only be transformed into gene-cell activity matrix by MAESTRO, whereas the other three GA

226 tools cannot because they use the fragment file for GA conversion. The performance of
227 different GSS tools was compared under three evaluation scenarios -- dimensional reduction,
228 clustering, and classification, using nine scATAC-seq datasets.

229     In general, regardless of imputation methods or GSS tools used, the performance of GSS
230 using recovered peak profile is significantly improved compared with that using the raw peak
231 profile (Figure 4A). Among the four imputation methods, SCALE that is designed for scATAC-
232 seq provides the overall best performance, ranking first or second in almost all comparisons.
233 Among the three scRNA-seq imputation methods, the overall performance of DrImpute is the
234 best, followed by MAGIC. Except that the performance of SAVER is apparently the worst in
235 most cases, the performance of the other three tools is relatively close. Moreover, the impact
236 of the same imputation tool on the performance of different GSS tools is quite consistent, and
237 no GSS tool relies on a specific imputation method. Next, we examined in detail the change of
238 ARI scores of different GSS tools before and after imputation by SCALE under the clustering
239 scenario (Figure 4B). In almost all cases, regardless of datasets or GSS tools, ARI scores based
240 on recovered data are increased significantly. However, the performance improvement of
241 different datasets after imputation varies greatly; the increase of ARI value under Leukemia,
242 Hematopoiesis, and Brain is much slighter than that under other six datasets. Moreover, after
243 imputation, the performance of different GSS tools on the same dataset also varies greatly. For
244 example, after imputation, the ARI score of different tools on InSilico varies from 0.41 by
245 UniPathATAC to 0.88 by Vision; the ARI score on GM12878HL varies from 0.02 by VAM to
246 0.79 by Pagoda2. In addition, the performance ranking of these tools changes after imputation.
247 Pagoda2 and PLAGE are top performers using the raw data (Figures 2 & 3), while their ranking
248 falls to a medium level after imputation. The performance of almost all tools has been greatly
249 improved using data after imputation, but none is obviously the best -- several tools, including
250 GSVA, Vision, Pagoda2, ssGSEA, and AUCell, achieve comparably good performance.
251 Interestingly, the ARI score of Pagoda2 on raw data of InSilico and PBMC3K is much higher
252 than that of other tools, however, the performance after imputation is even lower than that
253 before imputation or most other tools. This result indicates that the impact of data imputation
254 for a tool that already performs well on the raw data may be limited. In contrast, some GSS
255 tools have very poor performance before imputation, while a substantial improvement was
256 obtained after imputation. For example, the ARI score of Vision on the InSilico raw data is only
257 0.13, while it is increased greatly to 0.88 using data after imputation. The UMAP visualization
258 of the GSS-ATAC matrix obtained from the InSilico data shows significantly more
259 distinguishable cell types using data after imputation (Figure 4C). These results demonstrate

260 that the performance of GSS tools can be significantly improved by the incorporation of the
261 imputation step in data preprocessing, particularly for those GSS tools having poor
262 performance on the raw data.

263 **Evaluation of GSS tools by the enrichment analysis of marker gene sets**

264 Next, we used marker genes of known cell types as the reference to further evaluate the
265 accuracy of cell type recognition using gene sets quantified by different GSS tools (Materials
266 and Methods, Table S3). Considering the abundance of cell types and the availability of cell
267 marker information in the CellMarker database [36], here we used the two PBMC datasets with
268 25 sub-types for evaluation. ssGSEA has the highest accuracy of cell type recognition when
269 only the top one to three gene sets were used (Figure 5). For example, when identifying cell
270 types only based on the top one gene set, the accuracy of ssGSEA is ~71%, which is much
271 higher than other tools (Vision = 51% in the second place). Several other tools also achieve
272 comparable accuracy to ssGSEA when using ≥ 3 top gene sets, including VAM, Pagoda2, and
273 Vision, which reach an accuracy of > 82% using top five gene sets.  Surprisingly, for PLAGE
274 which has comparable performance with Pagoda2 in other evaluation scenarios (Figures 2 &
275 3), none of the top gene sets identified by PLAGE is enriched on correct cell types. In particular,
276 although UniPathATAC is designed purposely for scATAC-seq, its performance is
277 consistently lower than several other GSS tools for RNA-seq. Taken together, among the ten
278 GSS tools, six tools, including ssGSEA, VAM, Pagoda2, Vision, AUCell, and z-score, provide
279 overall better performance than other tools. UniPathATAC and GSVA rank at the second level,
280 while UniPath and PLAGE perform the worst.

281 **Evaluation of the impact of GA transformation on GSS**

282 GA conversion is a necessary step before using RNA-seq GSS tools on scATAC-seq data. Here
283 we evaluated the performance of different GA tools by calculating the correlation between the
284 GA profile from scATAC-seq and the gene expression profile from scRNA-seq, using three
285 matched scRNA-seq and scATAC-seq datasets (Brain, PBMC3K, and PBMC10K). Generally,
286 Signac and snapATAC provide better consistency between GA inferred from scATAC-seq and
287 gene expression level from scRNA-seq than MAESTRO and ArchR (Figure 6A). Using the
288 SCALE-imputed data for GA conversion by MAESTRO, the consistency measured by
289 correlation is increased (*P* value < 5.8e-108 between MAESTRO/SCALE and MAESTRO/raw
290 for each dataset), suggesting that imputation could increase the performance of GA conversion.
291 Next, we compared the effect of GA tools on GSS using more scATAC-seq datasets. Since GA
292 tools except for MAESTRO are only applicable to the raw scATAC-seq peak profile, we used
293 the raw data without imputation for evaluation. GA matrix obtained by GA tools were used as

294    the input for the ten GSS tools to score gene sets. There is no clear consensus on which

295    approach is the best; no GA method has significantly higher impact on the performance of all

296    GSS tools than other methods (Figure 6B). Among the ten GSS tools, the performance variation

297    of different GA methods on AUCell and UniPath is greater than that on other GSS tools.

298    Among the four GA tools, the performance of different GSS tools on GA matrix obtained by

299    Signac and snapATAC is more robust and relatively higher than that by MAESTRO or ArchR.

300    Moreover, different from other GSS tools for RNA-seq that can only score gene sets from the

301    GA matrix, UniPathATAC can score gene sets directly from the peak profile without GA

302    transformation, while its performance is inferior than several GSS tools designed for RNA-seq,

303    such as Pagoda2 and PLAGE. Collectively, Signac and snapATAC provide relatively better

304    results than MAESTRO and ArchR in both evaluation scenarios, whereas MAESTRO has the

305    unique ability to obtain GA from imputed data.

306    **Evaluation of the impact of different gene set collections on GSS**

307    Next, we investigated the impact of six gene set collections from MSigDB (Table S1) on the

308    performance of GSS tools, using nine scATAC-seq datasets. In the evaluation pipeline, we

309    used SCALE for dropout imputation, followed by MAESTRO for GA transformation. Then

310    we applied different GSS tools to each GA matrix to calculate the GSS-ATAC matrix based

311    on each gene set collections, and evaluated the performance in the context of clustering. The

312    impact of different gene set collections on GSS performance is not as evident as that of

313    imputation tools (Figure 7A *vs.* Figure 4A). The average ARI score using TFT or GO:BP is

314    slightly lower than that using other four gene set collections (TFT = 0.367; GO:BP = 0.389;

315    others: 0.4 to 0.419). Moreover, different GSS tools have different degree of robustness to

316    different gene set collections on different datasets (Figure 7B). For four datasets (Brain,

317    Hematopoiesis, Leukemia, and PBMC3K), the performance of all GSS tools is relatively stable,

318    regardless of which gene set collection is used (Figure 7C). In contrast, for the other five

319    datasets, the performance of different GSS tools is more affected by gene set collections. For

320    example, for InSilico which shows overall high performance, AUCell, GSVA, and Vision are

321    much less sensitive to gene sets than other tools (Figure 7B). Among the ten GSS tools, the

322    performance of Vision and UniPath is the least affected by gene sets, while UniPathATAC is

323    the most sensitive to gene sets (Figure 7C). In particularly, Pagoda2 is the top performer on

324    raw scATAC-seq data according to our evaluation (Figures 2A & 3), however, its robustness

325    to different gene sets is only moderate (Figures 7B & C). Overall, Vision has relatively more

326    robust and generally high performance across different gene set collections.

327    **Running time evaluation**

328     The computing speed of Vision and z-score is significantly faster than that of other tools. Even
329     when the number of cells and gene sets increases, the running time only increases slightly
330     (Figure 8). In contrast, GSVA and VAM run fast when the data size is small, while the running
331     time increases significantly with the increase of data size. Among these tools, ssGSEA and
332     UniPath take significantly more computing time than other tools. Nevertheless, among these
333     experiments, it only takes up to six hours (Unipath: 328.84 min) even for the longest case by
334     these two tools. PLAGE and Pagoda2, which show the best performance on the raw data, are
335     quite efficient, which are second in line to the fastest tools, Vision and z-score. However,
336     Pagoda2 failed to complete calculation in some cases, which needs to be used with caution.
337     According to the calculation speed, the ranking for the top three tools with overall high
338     performance on data after imputation is Vision > Pagoda2 > ssGSEA. In addition,
339     UniPathATAC, a tool specially designed for scATAC-seq, has a medium computing speed,
340     which is close to Pagoda2.

341     **Practical guidelines for choosing GSS tools**

342     Here we summarized the performance of different GSS tools on ten scATAC-seq datasets in
343     various evaluation pipelines in the context of clustering, considering different GA tools,
344     imputation tools, and gene set collections (Figure 9A). For the preprocessing of scATAC-seq
345     data in the GSS pipeline, our results showed that dropout imputation can significantly improve
346     the GSS results, and SCALE or DrImpute provide overall better performance than the other
347     two imputation tools. In contrast, using different GA tools or gene set collections has much
348     less impact on GSS results. Regardless of gene set collections, for peak-cell data after dropout
349     imputation by SCALE (only MAESTRO can be used for GA transformation in this case),
350     Vision and GSVA show an overall better performance on the SCALE-recovered data than other
351     GSS tools (average ARI: GSVA = 0.47, Vision = 0.46, others = 0.29 to 0.44). For raw peak-
352     cell data, Pagoda2 in conjunction with snapATAC (ARI = 0.31) or Signac (ARI = 0.29)
353     performs the best, followed by PLAGE. In particular, it is worth noting that RNA-seq GSS
354     tools are only applicable to scATAC-seq when the peak-level open-chromatin profile of
355     scATAC-seq has been converted into gene-level activity scores by GA tools. Although our
356     benchmark demonstrates that dropout imputation greatly improves the performance of GSS
357     tools, only MAESTRO can be applied to the recovered peak-cell matrix for GA transformation,
358     while other GA tools cannot due to that the fragment file needed for GA conversion cannot be
359     imputed.

360     Based on our comprehensive evaluation and unique features of different tools, we propose
361     some practical guidelines for choosing appropriate tools for GSS (Figure 9B). For GSS from

362    raw scATAC-seq data without dropout imputation, we recommend two tools with overall best

363    performance and high speed, PLAGE and Pagoda2, combined with snapATAC or Signac for

364    GA transformation (Figures 2 and 8). Meanwhile, users can also use SCALE to recover the

365    peak-cell profile, followed by GA conversion with MAESTRO, and then adopt Vision, which

366    has relatively good performance (Figures 4, 5, and 7) and speed (Figure 8) for data after

367    imputation. Since the performance of different GSS tools on data after imputation is greatly

368    improved and becomes closer (Figure 4), users can also try multiple GSS tools with comparable

369    performance to Vision, such as GSVA, Pagoda2, ssGSEA, and AUCell, for comparative

370    analysis, especially when the data size is small. If users want to perform GSS without GA

371    conversion, then UniPathATAC is the only tool available at present. In addition, considering

372    that different gene set collections have relatively limited and uncertain impact on the

373    performance of GSS tools (Figure 7) but are important for biological interpretation, it is

374    recommended to try different gene set collections in the GSS pipeline.

## Discussion

376    GSS has been widely conducted in bulk or single-cell RNA-seq studies, which helps to

377    decipher single-cell heterogeneity and cell-type-specific variability by incorporating prior

378    knowledge from functional gene sets or pathways. ScATAC-seq is a powerful epigenetic

379    technique for interrogating single-cell chromatin-based gene regulation, and genes or gene sets

380    with dynamic regulatory potentials can be regarded as cell-type specific markers as if in

381    scRNA-seq. The GA score transformed from the chromatin accessibility profile of scATAC-

382    seq is potentially a reliable predictor of gene expression and can be used for cell type annotation

383    [4-8]. GA scores facilitate the use of RNA-seq GSS tools to score gene sets for scATAC-seq

384    data. Taking the GSS results of the matched scRNA-seq datasets and those of UniPathATAC

385    as the reference, we confirmed that RNA-seq GSS tools are applicable to scATAC-seq. First,

386    we performed GSS for the matched scATAC-seq and scRNA-seq data from PBMCs and Brain,

387    and found that the performance of GSS tools on scATAC-seq for clustering cells or

388    distinguishing cell types was comparable to that on scRNA-seq (Figure 2). Second, by the

389    enrichment analysis of marker gene sets for cell types using PBMC10K scATAC-seq data, we

390    found that the top few (1-10) gene sets with high scores can be used to determine the cell types

391    of most cells (Figure 5). Third, the comprehensive evaluation of various scATAC-seq datasets

392    shows that several RNA-seq GSS tools, e.g., Pagoda2, PLAGE, and Vision, even have much

393    better results under different evaluation scenarios than the GSS tool specially designed for

394    scATAC-seq -- UniPathATAC (Figures 2-6). After demonstrating the applicability of RNA-seq

GSS tools on scATAC-seq, we systematically evaluated 10 GSS tools and found that Pagoda2 and PLAGE have the best overall performance for the raw peak-cell profile, which is similar to the previous benchmark results of GSS tools on scRNA-seq data [15]. In particular, Pagoda2 is developed for scRNA-seq and PLAGE is for bulk RNA-seq, both of which are PCA-based RNA-seq methods but also provide good performance on scATAC-seq. Several previous studies have shown that GSS tools developed for bulk RNA-seq are applicable to scRNA-seq data [14, 15], and tools for scRNA-seq imputation is also widely used in recovering scATAC-seq dropouts [25]. Our benchmark further confirmed that GSS tools designed for RNA-seq is also suitable for scATAC-seq data.

We also comprehensively evaluated the impact of data preprocessing of scATAC-seq on GSS, including dropout imputation and GA transformation. We found that GSS results using data after imputation are significantly better than those using raw data, regardless of GSS tools or imputation tools (Figure 4). Among the four imputation tools, SCALE performs generally better than other three scRNA-seq tools, while the scRNA-seq tool DrImpute provides comparable performance to SCALE. Previously, Liu et al. [25] benchmarked multiple scRNA-seq imputation tools on scATAC-seq including MAGIC and SAVER, and found that MAGIC provides much better performance than SAVER. This is consistent to our observation that SAVER shows the worst performance on scATAC-seq data. Moreover, the two tools included in our benchmark that have overall high performance, SCALE and DrImpute, were not involved in the previous benchmark [25]. Particularly, the performance of Pagoda2 and PLAGE, which provide the best performance on raw data, is not significantly improved after imputation, while the performance of several other tools, including GSVA, Vision, Pagoda2, ssGSEA, and AUCell, is greatly improved after imputation, surpassing Pagoda2 and PLAGE (Figure 4). Compared to the positive impact of dropout imputation on GSS, the impact of different GA methods or gene set collections on GSS is uncertain and limited (Figures 6 & 7). Therefore, we recommend users to try different GA tools and different gene sets for GSS in practical applications. Moreover, we found that although the open-chromatin profile obtained from scATAC-seq data can be preprocessed using different imputation tools and different GA tools, GSS results are highly dependent on scATAC-seq datasets. Some datasets, such as Hematopoiesis and Leukemia, have extremely poor results regardless of the evaluation scenarios (dimensionality reduction, clustering or classification) or the representation of the data (peak profile, gene-level activity score or gene set score) (Figures 3,4, and S1). The low quality of the raw scATAC-seq data could be alleviated to some extent by dropout imputation

428    rather than choosing a different GA tool. However, no matter how the raw data is preprocessed,

429    GSS results on data with very poor quality of raw data often cannot reach the ideal level.

430         In our benchmark study, the performance of GSS tools was quantitively evaluated under

431    four scenarios of dimensionality reduction, clustering, classification, and cell type

432    determination. These scenarios, especially clustering, are critical steps of single-cell analysis

433    in most scRNA-seq and scATAC-seq studies. We acknowledged that the ARI score that

434    represents the consistency between the predicted cell type labels from clustering and the true

435    reference is not high throughout our benchmarking of GSS tools (< 0.5 in most cases), which

436    means that the clustering results solely based on gene set scores may be poor. However, for

437    scATAC-seq data, which is even sparser than the already sparse scRNA-seq data, the ARI value

438    is normally very low. For example, the ARI value in these pioneering scATAC-seq studies [25,

439    37-39] is also < 0.5 in most cases. Nevertheless, clustering is a routine step in most single-cell

440    analysis pipelines and the outputs of different tools or methods are frequently used as the input

441    for clustering algorithms to produce clustering results. Therefore, evaluating the clustering

442    ability would be a useful measure for assessing the performance of different GSS tools. We

443    estimated that the value of ARI can reflect the performance of different GSS tools under the

444    clustering scenario. At the same time, the low ARI value indicates that the clustering results

445    should be used in caution. Moreover, we also speculated that the low ARI value may be also

446    due to the poor annotation or high similarity of some cell types, and/or the inability to

447    completely restore the true cell types only through the scATAC-seq data. As such, integrating

448    information of additional modalities with gene set scores, such as the gene expression profile

449    from scRNA-seq and the peak-level profile from scATAC-seq, would help to obtain better

450    clustering results for better cell type distinguishing.

451         Currently, matched scRNA-seq and scATAC-seq data on dynamic processes (e.g.

452    differentiation of induced pluripotent stem cells) are increasingly available [40-44]. It would

453    be interesting to examine whether and how well the cell transition trajectory could be inferred

454    based on gene set scores obtained by different GSS tools. However, trajectory analysis is a

455    more complex procedure that requires more biological interpretation than clustering analysis,

456    and its results are difficult to quantify using performance indicators like ARI in clustering

457    analysis. Nevertheless, evaluating GSS tools under the scenario of trajectory analysis could be

458    a future direction upon the availability of appropriate quantification methods for evaluation the

459    accuracy of trajectory inference.

460    **Material and methods**

**Datasets**

We used ten publicly available scATAC-seq datasets (Table S2), including InSilico [2], GM12878HEK [3], GM12878HL [3], Leukemia [45], Hematopoiesis [46], Forebrain [47], SNAREmix [48], and three matched datasets from 10X Genomics (Brain, PBMC3K, and PBMC10K) [8]. The InSilico dataset is an *in silico* mixture of four independent scATAC-seq experiments performed on different cell lines [2]. The GM12878HEK and GM12878HL datasets are mixtures of two commonly-used cell lines, respectively [3]. The Leukemia dataset includes mononuclear cells and lymphoid-primed pluripotent progenitor cells isolated from a healthy human donor, and leukemia stem cells and blast cells isolated from two patients with acute myeloid leukemia [45]. The Forebrain dataset is derived from P56 mouse forebrain cells [47]. The Hematopoiesis dataset was used to characterize the epigenome pattern and heterogeneity of human hematopoiesis [46]. The Brain, PBMC3K, and PBMC10K datasets are publicly available datasets generated by 10x Genomics [8], which jointly profiled mRNA abundance and DNA accessibility in human peripheral blood mononuclear cells (PBMCs) and human healthy brain tissue of cerebellum, respectively. The SNAREmix dataset is a mixture of cultured human BJ, H1, K562, and GM12878 cells [48]. These diverse datasets were generated from both microfluidics-based and cellular indexing platforms with substantially different number of cells and peaks, which were widely used in previous studies for benchmarking [25] or validating computational tools for scATAC-seq, such as scMVP [38], scABC [49], SCALE [50], and Signac [8]. We used Azimuth [51] to annotate cell types in the PBMC3K and PBMC10K datasets by label transfer from a publicly available multimodal PBMC reference dataset [51] and in Brain dataset by label transfer from the human cerebellum dataset [52]. Cell types of other datasets were obtained from relevant studies.

**Preprocessing of scATAC-seq data**

For scATAC-seq datasets without publicly available peak-cell matrix, the raw FASTQ files downloaded from NCBI were aligned to the reference genome (human: hg19; mouse: mm10) using Bowtie 2 [53], resulting in alignment files of BAM format. Then these BAM files were used as inputs for MACS2 [54] for peak calling and then SnapTools (https://github.com/r3fang/SnapTools) was adopted to generate the peak-cell matrix. Similar to the previous study [55], we filtered peaks with read counts >=2 and present in at least 10 cells for InSilico, GM12878HEK and GM12878HL data. We filtered peaks with read counts >=2 and present in at least 50 cells for Forebrain. For Hematopoiesis, Leukemia, SNAREmix, Brain, PBMC3K and PBMC10K, we followed the routine preprocessing following the tutorial of Signac to filter peaks and cells.

495    We chose four tools for dropout imputation of scATAC-seq data, including SCALE [35]

496    which is currently the only method specifically designed for scATAC-seq and three widely used

497    scRNA-seq tools – MAGIC [26], DrImpute [34] and SAVER [27]. The peak-cell matrix was

498    used as the input for these tools with default parameters for recovering dropout peaks. Of note,

499    because Signac, ArchR, and snapATAC require a fragment file of the raw scRNA-seq data to

500    calculate gene-level activity, we can only use MAESTRO [5] to obtain GA matrix directly from

501    the recovered peak-cell matrix. We used liftOver [56] to convert coordinates between different

502    genome versions, if necessary.

503    **GA conversion**

504    The peak-level profile of scATAC-seq data needs to be converted into the gene-level activity

505    before using RNA-seq GSS tools. We chose four GA tools, including MAESTRO [5], Signac

506    [8], ArchR [6], and snapATAC [7], to transform the open-chromatin profile obtained from

507    scATAC-seq into the gene-level activity scores. MAESTRO obtains a regulatory weight based

508    on the distance from the peak center to the gene transcription start site, and associates it with

509    the peak-cell matrix to get the gene activity score. Signac is used in the Seurat package [22] for

510    GA conversion, which simply sums the gene body with the peaks that intersect in the 2-kbp

511    upstream region in each cell. SnapATAC obtains a score for each gene by normalizing the

512    number of fragments overlapping genes in cells. ArchR infers gene expression from chromatin

513    accessibility by using a custom distance-weighted accessibility model. Among these tools,

514    MAESTRO use the peak-cell matrix for GA conversion, while other three tools use the

515    fragment file. The fragment file [8] is a coordinate-sorted file for storing scATAC-seq data,

516    which contains five columns: chromosome, start coordinate, end coordinate, cell barcode, and

517    duplicate count. This file can be generated from a BAM file using Cellranger or the Sinto

518    package (https://pypi.org/project/sinto/). It should be noted that, only the peak-cell matrix

519    rather than the fragment file can be imputed by imputation tools, therefore, only MAESTRO

520    can be used for GA conversion on the peak-cell data after imputation.

521    We used three matched scRNA-seq and scATAC-seq datasets (Brain, PBMC3K, and

522    PBMC10K) to evaluate the performance of different GA tools in predicting the gene expression

523    level from scATAC-seq data. First, we used each GA tool to convert the raw peak-cell matrix

524    into the GA matrix for each dataset. As MAESTRO is applicable to the imputed peak-cell

525    profile, we also used MAESTRO to obtain the GA matrix based on the SCALE-imputed peak-

526    cell matrix. Then we calculated the Pearson's correlation between the raw or imputed GA

527    profile from scATAC-seq and the gene expression profile from scRNA-seq for each cell. The

528 correlation profiles of all cells obtained from the four GA tools for each matched scRNA-seq

529 and scATAC-seq dataset were compared.

530 **GSS tools**

531 Ten GSS tools were evaluated in our benchmark. We run these tools with default parameters

532 according to the tutorials provided in the respective studies.

533 PLAGE (Pathway Level Analysis of Gene Expression) [29] scores gene sets for RNA-seq

534 by singular value decomposition (SVD). The gene expression matrix is normalized, and the

535 first coefficient of the right singular vector obtained by SVD is considered as the gene set score.

536 Combined z-score (z-score) [30] is a classic strategy to aggregate the expression of

537 multiple genes. Gene expression is scaled by the mean and standard deviation of the cells. Then,

538 gene expression levels of all genes within each gene set are averaged to score the gene set of

539 each cell.

540 ssGSEA (Single Sample Gene Set Enrichment Analysis) [31] is an extension of GSEA.

541 ssGSEA ranks genes by expression levels within each cell individually, then scores gene sets

542 by enrichment analysis using random walk statistics such as Kolmogorov-Smirnov (K-S)

543 statistic.

544 GSVA (Gene Set Variation Analysis) [32] utilizes the K-S statistic to assess gene set

545 variation. GSVA first estimates the cumulative density function for each gene, using the classic

546 maximum deviation method by default. The score matrix is obtained by calculating the score

547 of the gene set from the gene density profile using the K-S statistic.

548 AUCell [18] employs the area under the curve (AUC) to calculate the enrichment of a

549 pathway (i.e., gene set) in the expressed genes of each cell. AUCell first ranks genes based on

550 their expression levels in each cell, resulting in a ranking matrix. The AUC of the recovery

551 curve is then used to determine whether the gene set is enriched at top genes in each cell. To

552 calculate AUC, only the top 5% of genes are used by default, which means to examine how

553 many genes in the gene set are within the top 5% genes in the respective cell.

554 Pagoda2 (Pathway and Gene Set Overdispersion Analysis) [16] is a computational

555 framework to detect cellular heterogeneity from scRNA-seq data. The method fits an error

556 model to each cell to characterize its properties, and then renormalizes the residual variance

557 for each gene in the cell. Then, the scoring matrix for each gene set is quantified by its first

558 weighted principal component.

559 Vision [17] uses autocorrelation statistics to identify biological variation across cells,

560 which performs directly on the manifold of cell-cell similarity. It first identifies the K-Nearest

561     Neighbors (KNNs) of each cell to generate a KNN map of the cell, then the GSS matrix is

562     calculated based on the average gene expression of each gene set.

563     VAM (Variance-Adjusted Mahalanobis) [33] is a fast and accurate method for cell-specific

564     gene set evaluation, which is integrated with the Seurat framework to accommodate the

565     characteristics of high technical noise, sparsity and large sample size of scRNA-seq data. It

566     calculates cell-specific pathway scores to convert a gene-by-gene matrix into a pathway-by-

567     pathway matrix that can be used for data visualization and statistical enrichment analysis.

568     UniPath [28] is a uniform approach for pathway and gene-set based analysis for both

569     scRNA-seq and scATAC-seq. For scRNA-seq, it first converts gene expression profiles to p-

570     values assuming a Gaussian distribution, according to the mean and variance of each cell. Then

571     p-values of genes in each gene set are combined using Brown's method and then an adjusted

572     p-value is obtained for each gene set. For scATAC-seq, UniPath first highlights enhancers by

573     normalizing read counts of scATAC-seq peaks using their global accessibility scores and

574     performs a hypergeometric or binomial test using proximal genes of peaks, which then converts

575     the open-chromatin profile to pathway enrichment scores for gene sets. UniPath provides

576     functions for scoring gene sets in scRNA-seq and scATAC-seq, respectively. In this study, we

577     referred to the method for scRNA-seq as UniPath and the method for scATAC-seq as

578     UniPathATAC.

579     **Benchmarking gene set scoring tools**

580     *Cell type clustering*

581     We evaluated the performance of different GSS tools in the context of unsupervised clustering,

582     using Louvain which is imbedded in the Seurat package. Given a GSS-ATAC matrix obtained

583     by a GSS tool, we employed PCA for dimensionality reduction and then performed Louvain

584     clustering on the first 10 PCs. Louvain clustering provides a tuneable parameter 'resolution'

585     for determining the number of clusters based on a binary search algorithm, which was set to

586     0.5 in our benchmark. We used ARI (Adjust Random Index), a widely-used indicator, to

587     measure the consistency between two clustering results. The ARI is the adjusted value of the

588     raw RI (Random Index) score; the RI computes a similarity metric between two clustering

589     results by considering all sample pairs and counting pairs assigned in the same or different

590     clusters in the predicted and true clusters (Eq. 1). An ARI close to 0 means random labelling

591     and ARI = 1 means perfect matching of the two clustering results. ARI is calculated with the

592     'adjustedRandIndex' function in the mclust [57] package.

593     $$ARI = \frac{RI - Exp(RI)}{\max(RI) - Exp(RI)} \qquad (1)$$

*Dimensionality reduction*

594 We first performed dimensionality reduction by PCA on the GSS-ATAC matrix obtained by a GSS tool with Seurat (PCs = 10). Then UMAP (Uniform Manifold Approximation and Projection) [58] was performed with the first 10 PCs and the average Silhouette width of all cells was calculated using the 'silhouette' function provided in the R package cluster. The Silhouette score was used to evaluate the performance of dimensionality reduction for each GSS-ATAC matrix. Silhouette score ranges from -1 to 1, with a high value indicating that cells of the same cell type group together and are far from cells of a different type. The silhouette score for cell $i$ is defined as:

$$
s(i) = \begin{cases} 1 - \frac{x(i)}{y(i)} & if\ x(i) < y(i) \\ 0 & if\ x(i) = y(i) \\ \frac{y(i)}{x(i)} - 1 & if\ x(i) > y(i) \end{cases} \tag{2}
$$

Here, $x(i)$ and $y(i)$ is the average distance from cell $i$ to all other cells in cell $i$'s cluster and cell $i$'s nearest cluster, respectively.

*Classification*

To evaluate the performance of GSS tools in the context of classification, we implemented a multi-normal logistic regression model with k-fold cross-validation using the Python scikit-learn package. The inverse of the regularization strength of the multinormal logistic regression model was set to 1. The parameter $k$ of the k-fold cross-validation was set to 5. Gene set scores in the GSS-ATAC matrix were scaled between 0 and 1 before model training and testing. The classification accuracy of the test dataset is calculated.

*Enrichment analysis of marker gene sets*

Similar to the previous study [28], we used marker genes of known cell types as the reference to examine whether gene sets scored by different GSS tools are enriched on known cell types. We obtained human marker genes from CellMarker [36] to make a collection of gene sets for 467 cell types (Table S3) and then organized these gene sets as the form of the gene set representation in MSigDB. Each GSS tool was used to score these marker gene sets for each scATAC-seq dataset to obtain a GSS-ATAC matrix. Based on the GSS-ATAC matrix, for each cell the top $N$ gene sets ranking by the gene set score can be obtained. If a cell's cell type falls within cell types of the top $N$ gene sets, then the cell is considered as correctly recognized. Finally, given a scATAC-seq dataset, the percentage of cells annotated with correct cell type was calculated for each GSS tools.

*Running time evaluation*

625 We used scATAC-seq datasets and gene sets with different sizes to test the running time of
626 GSS tools. Three datasets with different orders of magnitude were used for evaluation,
627 including InSilico, Hematopoiesis and PBMC10K, which contain approximately 500, 2000 and
628 10K cells, respectively. Four sources of gene sets with different sizes were selected from
629 MSigDB, including KEGG (186 pathways), TFT (1133 pathways), REACTOME (1797
630 pathways) and GO:BP (7350 pathways). The computer processor for evaluation is
631 intel@Xeon(R) CPU E5-2680 v4 @ 2.40GHz × 56. One CPU core is allocated to each task of
632 running a GSS tool on a dataset with given gene sets. Only the running time of the GSS tool is
633 counted, excluding the time consumption of data and package loading, preprocessing, data
634 imputation and gene activity conversion.

## CRediT author statement

636 **Xi Wang**: Investigation, Methodology, Data curation, Formal analysis. **Qiwei Lian**:
637 Investigation, Methodology, Data curation, Formal analysis. **Haoyu Dong**: Data curation.
638 **Shuo Xu**: Data curation. **Yaru Su**: Formal analysis. **Xiaohui Wu**: Conceptualization,
639 Writing - original draft, Writing - review & editing, Supervision, Project administration,
640 Funding acquisition. All authors read and approved the final manuscript.

## Competing interests

642 The authors have declared no competing interests.

## Acknowledgements

## ORCID

648 0000-0002-4515-5749 (Xi Wang)
649 0000-0003-3366-6127 (Qiwei Lian)
650 0000-0003-1163-3623 (Haoyu Dong)
651 0000-0001-8406-8958 (Shuo Xu)
652 0000-0002-9539-8511 (Yaru Su)
653 0000-0003-0356-7785 (Xiaohui Wu)

## References

655 [1] Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. Transposition of native chromatin
656 for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and nucleosome
657 position. Nat Methods 2013;10:1213-8.

[2] Buenrostro JD, Wu B, Litzenburger UM, Ruff D, Gonzales ML, Snyder MP, et al. Single-cell chromatin accessibility reveals principles of regulatory variation. Nature 2015;523:486-90.

[3] Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, et al. Multiplex single cell profiling of chromatin accessibility by combinatorial cellular indexing. Science 2015;348:910-4.

[4] Pliner HA, Packer JS, McFaline-Figueroa JL, Cusanovich DA, Daza RM, Aghamirzaie D, et al. Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin Accessibility Data. Mol Cell 2018;71:858-71.e8.

[5] Wang C, Sun D, Huang X, Wan C, Li Z, Han Y, et al. Integrative analyses of single-cell transcriptome and regulome using MAESTRO. Genome Biol 2020;21:198.

[6] Granja JM, Corces MR, Pierce SE, Bagdatli ST, Choudhry H, Chang HY, et al. ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. Nat Genet 2021;53:403-11.

[7] Fang R, Preissl S, Li Y, Hou X, Lucero J, Wang X, et al. Comprehensive analysis of single cell ATAC-seq data with SnapATAC. Nat Commun 2021;12:1337.

[8] Stuart T, Srivastava A, Madad S, Lareau CA, Satija R. Single-cell chromatin state analysis with Signac. Nat Methods 2021;18:1333-41.

[9] Cusanovich DA, Hill AJ, Aghamirzaie D, Daza RM, Pliner HA, Berletch JB, et al. A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. Cell 2018;174:1309-24.e18.

[10] Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, et al. A survey of best practices for RNA-seq data analysis. Genome Biology 2016;17:1-19.

[11] Rahmatallah Y, Emmert-Streib F, Glazko G. Gene set analysis approaches for RNA-seq data: performance evaluation and application guideline. Brief Bioinform 2016;17:393-407.

[12] Das S, McClain CJ, Rai SN. Fifteen Years of Gene Set Analysis for High-Throughput Genomic Data: A Review of Statistical Approaches and Future Challenges. Entropy 2020;22:427.

[13] Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The Molecular Signatures Database (MSigDB) hallmark gene set collection. Cell Syst 2015;1:417-25.

[14] Holland CH, Tanevski J, Perales-Patón J, Gleixner J, Kumar MP, Mereu E, et al. Robustness and applicability of transcription factor and pathway analysis tools on single-cell RNA-seq data. Genome Biol 2020;21:36.

[15] Zhang Y, Ma Y, Huang Y, Zhang Y, Jiang Q, Zhou M, et al. Benchmarking algorithms for pathway activity transformation of single-cell RNA-seq data. Comput Struct Biotechnol J 2020;18:2953-61.

[16] Lake BB, Chen S, Sos BC, Fan J, Kaeser GE, Yung YC, et al. Integrative single-cell analysis of transcriptional and epigenetic states in the human adult brain. Nat Biotechnol 2018;36:70-80.

[17] DeTomaso D, Jones MG, Subramaniam M, Ashuach T, Ye CJ, Yosef N. Functional interpretation of single cell similarity maps. Nat Commun 2019;10:4376.

[18] Aibar S, González-Blas CB, Moerman T, Huynh-Thu VA, Imrichova H, Hulselmans G, et al. SCENIC: single-cell regulatory network inference and clustering. Nat Methods 2017;14:1083.

[19] Zappia L, Theis FJ. Over 1000 tools reveal trends in the single-cell RNA-seq analysis landscape. Genome Biol 2021;22:301.

[20] Vieth B, Parekh S, Ziegenhain C, Enard W, Hellmann I. A systematic evaluation of single cell RNA-seq analysis pipelines. Nat Commun 2019;10.

[21] Stuart T, Satija R. Integrative single-cell analysis. Nat Rev Genet 2019;20:257-72.

[22] Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, et al. Comprehensive integration of single-cell data. Cell 2019;177:1888-902.e21.

[23] Yan F, Powell DR, Curtis DJ, Wong NC. From reads to insight: a hitchhiker's guide to ATAC-seq data analysis. Genome Biol 2020;21:22.

[24] Baek S, Lee I. Single-cell ATAC sequencing analysis: From data preprocessing to hypothesis generation. Comput Struct Biotechnol J 2020;18:1429-39.

[25] Liu Y, Zhang J, Wang S, Zeng X, Zhang W. Are dropout imputation methods for scRNA-seq effective for scATAC-seq data? Brief Bioinform 2022;23.

[26] van Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, et al. Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. Cell 2018;174:716-29.e27.

712 [27] Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, et al. SAVER: gene expression
713 recovery for single-cell RNA sequencing. Nat Methods 2018;15:539-42.
714 [28] Chawla S, Samydurai S, Kong SL, Wu Z, Wang Z, TAM WL, et al. UniPath: a uniform approach
715 for pathway and gene-set based analysis of heterogeneity in single-cell epigenome and transcriptome
716 profiles. Nucleic Acids Res 2020.
717 [29] Tomfohr J, Lu J, Kepler TB. Pathway level analysis of gene expression using singular value
718 decomposition. BMC Bioinf 2005;6:225.
719 [30] Lee E, Chuang HY, Kim JW, Ideker T, Lee D. Inferring pathway activity toward precise disease
720 classification. PLoS Comput Biol 2008;4:e1000217.
721 [31] Barbie DA, Tamayo P, Boehm JS, Kim SY, Moody SE, Dunn IF, et al. Systematic RNA
722 interference reveals that oncogenic KRAS-driven cancers require TBK1. Nature 2009;462:108-12.
723 [32] Sonja Hänzelmann1, RobertCastelo1,2* and Justin Guinney3*. GSVA: gene set variation
724 analysis for microarray and RNA-Seq data. BMC Bioinf 2013.
725 [33] Frost HR. Variance-adjusted Mahalanobis (VAM): a fast and accurate method for cell-specific
726 gene set scoring. Nucleic Acids Res 2020;48:e94.
727 [34] Gong W, Kwak I-Y, Pota P, Koyano-Nakagawa N, Garry DJ. DrImpute: imputing dropout events
728 in single cell RNA sequencing data. BMC Bioinf 2018;19.
729 [35] Xiong L, Xu K, Tian K, Shao Y, Tang L, Gao G, et al. SCALE method for single-cell ATAC-seq
730 analysis via latent feature extraction. Nat Commun 2019;10:4576.
731 [36] Zhang X, Lan Y, Xu J, Quan F, Zhao E, Deng C, et al. CellMarker: a manually curated resource
732 of cell markers in human and mouse. Nucleic Acids Res 2019;47:D721-D8.
733 [37] Chen H, Lareau C, Andreani T, Vinyard ME, Garcia SP, Clement K, et al. Assessment of
734 computational methods for the analysis of single-cell ATAC-seq data. Genome Biol 2019;20:241.
735 [38] Li G, Fu S, Wang S, Zhu C, Duan B, Tang C, et al. A deep generative model for multi-view
736 profiling of single-cell RNA-seq and ATAC-seq data. Genome Biol 2022;23:20.
737 [39] Liu Q, Chen S, Jiang R, Wong WH. Simultaneous deep generative modeling and clustering of
738 single cell genomic data. Nat Mach Intell 2021;3:536-44.
739 [40] Ranzoni AM, Tangherloni A, Berest I, Riva SG, Myers B, Strzelecka PM, et al. Integrative
740 Single-Cell RNA-Seq and ATAC-Seq Analysis of Human Developmental Hematopoiesis. Cell Stem
741 Cell 2021;28:472-87.e7.
742 [41] Tedesco M, Giannese F, Lazarević D, Giansanti V, Rosano D, Monzani S, et al. Chromatin
743 Velocity reveals epigenetic dynamics by single-cell profiling of heterochromatin and euchromatin.
744 Nat Biotechnol 2022;40:235-44.
745 [42] Giles JR, Ngiow SF, Manne S, Baxter AE, Khan O, Wang P, et al. Shared and distinct biological
746 circuits in effector, memory and exhausted CD8(+) T cells revealed by temporal single-cell
747 transcriptomics and epigenetics. Nat Immunol 2022;23:1600-13.
748 [43] Bielecki P, Riesenfeld SJ, Hütter JC, Torlai Triglia E, Kowalczyk MS, Ricardo-Gonzalez RR, et
749 al. Skin-resident innate lymphoid cells converge on a pathogenic effector state. Nature 2021;592:128-
750 32.
751 [44] Ameen M, Sundaram L, Shen M, Banerjee A, Kundu S, Nair S, et al. Integrative single-cell
752 analysis of cardiogenesis identifies developmental trajectories and non-coding mutations in congenital
753 heart disease. Cell 2022;185:4937-53.e23.
754 [45] Corces MR, Buenrostro JD, Wu B, Greenside PG, Chan SM, Koenig JL, et al. Lineage-specific
755 and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. Nat
756 Genet 2016;48:1193-203.
757 [46] Buenrostro JD, Corces MR, Lareau CA, Wu B, Schep AN, Aryee MJ, et al. Integrated Single-
758 Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation.
759 Cell 2018;173:1535-48 e16.
760 [47] Preissl S, Fang R, Huang H, Zhao Y, Raviram R, Gorkin DU, et al. Single-nucleus analysis of
761 accessible chromatin in developing mouse forebrain reveals cell-type-specific transcriptional
762 regulation. Nat Neurosci 2018;21:432-9.
763 [48] Chen S, Lake BB, Zhang K. High-throughput sequencing of the transcriptome and chromatin
764 accessibility in the same cell. Nat Biotechnol 2019;37:1452-7.
765 [49] Zamanighomi M, Lin Z, Daley T, Chen X, Duren Z, Schep A, et al. Unsupervised clustering and
766 epigenetic classification of single cells. Nat Commun 2018;9:2410.

767 [50] Xiong L, Xu K, Tian K, Shao Y, Tang L, Gao G, et al. SCALE method for single-cell ATAC-seq
768 analysis via latent feature extraction. Nat Commun 2019;10:4576.
769 [51] Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, et al. Integrated analysis of
770 multimodal single-cell data. Cell 2021;184:3573-87.e29.
771 [52] Aldinger KA, Thomson Z, Phelps IG, Haldipur P, Deng M, Timms AE, et al. Spatial and cell type
772 transcriptional landscape of human cerebellar development. Nat Neurosci 2021;24:1163-75.
773 [53] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods 2012;9:357-
774 9.
775 [54] Feng J, Liu T, Qin B, Zhang Y, Liu XS. Identifying ChIP-seq enrichment using MACS. Nat
776 Protoc 2012;7:1728-40.
777 [55] Zamanighomi M, Lin Z, Daley T, Chen X, Duren Z, Schep A, et al. Unsupervised clustering and
778 epigenetic classification of single cells. Nat Commun 2018;9:2410.
779 [56] Nassar LR, Barber GP, Benet-Pagès A, Casper J, Clawson H, Diekhans M, et al. The UCSC
780 Genome Browser database: 2023 update. Nucleic Acids Res 2023;51:D1188-d95.
781 [57] Browne RP, McNicholas PD, Sparling MD. Model-based learning using a mixture of mixtures of
782 Gaussian and uniform distributions. IEEE Trans Pattern Anal Mach Intell 2012;34:814-7.
783 [58] McInnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation and
784 Projection. Journal of Open Source Software 2018;3:861.
785

786 **Figure legends**

787 **Figure 1  Overview of the benchmark workflow**

788 Before applying GSS tools, scATAC-seq dropout peaks can be recovered by imputation tools

789 and then the peak-level open-chromatin profile is converted into gene-level activity scores

790 using GA transforming tools. Using gene sets from MSigDB as prior information, ten GSS

791 tools are benchmarked in the context of diverse evaluation scenarios of dimensionality

792 reduction, clustering, classification and cell type determination based on a variety of

793 performance indicators. Tools marked with solid borders, including SCALE, the four GA tools

794 and UniPathATAC, are specifically designed for scATAC-seq. MAESTRO can be used for

795 GA transformation from both raw peaks and recovered peaks, while other three GA tools can

796 be only applied to raw peaks as they require a fragment file which is not available for the

797 imputed peak data. GSS, gene set scoring; scATAC-seq, single-cell assay for transposase

798 accessible chromatin using sequencing; GA, gene activity; MSigDB, the molecular signatures

799 database.

800 **Figure 2  GSS results using matched datasets of scATAC-seq and scRNA-seq**

801 **A.** Comparison of the performance of GSS tools on scRNA-seq (RNA) and scATAC-seq

802 (ATAC) data in the context of dimensionality reduction measured by Silhouette, clustering

803 measured by ARI, and classification measured by accuracy. Signac was employed to convert

804 the peak-cell matrix into the gene-cell activity matrix, and KEGG gene sets were used as prior

805 information. Three datasets including Brain, PBMC3K, and PBMC10K were used and the

806 average performance was calculated. **B.** UMAP visualization of cell types using gene set scores

807 obtained by applying different GSS tools on scRNA-seq and scATAC-seq PBMC10K data,

808 respectively. The plot was created using the DimPlot function provided in the Seurat package.

809 **C.** Comparison of the impact of different GA transformation tools on GSS of the PBMC10K

810 data. Signac, MAESTRO, ArchR, and snapATAC were used for transformation and then ten

811 GSS tools were applied on the GA matrix for scoring gene sets. Each violin plot summarizes

812 ARI scores of the ten GSS tools, with each dot representing one tool. *P* values of Wilcoxon

813 Rank Sum test used to compare ARI values between the scRNA-seq group and the other four

814 groups of Signac, MAESTRO, ArchR, and snapATAC are 0.60, 0.60, 0.22, and 0.86,

815 respectively. ARI, adjust random index; UMAP, uniform manifold approximation and

816 projection.

817 **Figure 3. Comparison of the performance of GSS tools**

818 The comparison was performed in the context of dimensionality reduction measured by

819 Silhouette, clustering measured by ARI and classification measured by accuracy. In each

820 column, the index values of the top performer for the respective dataset are displayed in red.

821 The 'Average' column is the average score of each row.

822 **Figure 4. Comparison of the impact of different dropout imputation tools on GSS**

823 **A.** Average performance of GSS tools on nine scATAC-seq datasets before or after imputation

824 in the context of dimensionality reduction measured by Silhouette, clustering measured by ARI

825 and classification measured by accuracy. **B.** The change of ARI scores of different GSS tools

826 before and after imputation by SCALE. In each column, the index value of the best performer

827 for the respective dataset is coloured in red. The 'Avg.' column is the average score of each

828 GSS tools on the nine datasets before (RAW) or after imputation (SCALE). **C.** UMAP

829 visualization of cell types using gene set scores obtained from the raw or imputed peak profile

830 of the InSilico data by each GSS tool. Datasets: Leuke., Leukemia; Hemat., Hematopoiesis;

831 HL., GM12878HL; HEK., GM12878HEK; Fore., Forebrain; SNAR., SNAREmix; InSil.,

832 InSilico; PBMC., PBMC3K.

833 **Figure 5. Evaluation of the enrichment and relevance of gene sets in single cells quantified**

834 **by different GSS tools**

835 PBMC3K and PBMC10K datasets were used, with six main cell types and 25 sub-types.

836 Marker genes of 467 known cell types from the CellMarker database were used as the reference.

837 Each GSS tool was used to score the 467 marker gene sets for each PBMC dataset, and the top

838 *N* gene sets ranking by the gene set score can be obtained for each cell. If a cell's cell type falls

839 within cell types of the top *N* gene sets, then the cell is considered as correctly recognized. The

840 Y-axis denotes the average percentage of cells annotated with correct cell type of the two

841   PBMC datasets based on the results of each GSS tool. The X-axis denotes the number of top

842   gene sets used for cell type recognition.

843   **Figure 6. Comparison of the impact of different GA transformation tools on GSS**

844   **A.** The correlation between the GA profile obtained by four GA transformation tools from

845   scATAC-seq and the gene expression profile from scRNA-seq for three matched scRNA-seq

846   and scATAC-seq datasets. Labels with 'raw' means the GA tools was performed on the raw

847   scATAC-seq profile, while 'SCALE' means that MAESTRO was used on the SCALE-imputed

848   scATAC-seq profile. B. Average performance on seven datasets in the context of

849   dimensionality reduction measured by Silhouette, clustering measured by ARI and

850   classification measured by accuracy. Results of UniPathATAC that is designed for scATAC-

851   seq without needing GA transformation are displayed as horizontal dotted lines for comparison.

852   For three of the ten scATAC-seq datasets used in this study (GM12878HEK, GM12878HL,

853   and SNAREmix), the fragment file that is needed for GA conversion of Signac, snapATAC,

854   and ArchR was not available, therefore, the remaining seven datasets were used here for

855   evaluation, including Leukemia, Hematopoiesis, Forebrain, InSilico, PBMC3K, PBMC10K,

856   and Brain.

857   **Figure 7. Comparison of the impact of different gene sets on GSS**

858   **A.** Average ARI score of ten GSS tools on nine scATAC-seq datasets using six gene set

859   collections from MSigDB. Dots in the 'Average' column represent the average ARI score of all

860   GSS tools using the respective gene set collection. Average ARI scores: GO:CC = 0.419;

861   GO:MF = 0.412; REACTOME = 0.401; KEGG = 0.4; GO:BP = 0.388; TFT = 0.368. Dropout

862   peaks in each scATAC-seq dataset were recovered by SCALE, followed by MAESTRO for

863   GA transformation. **B.** Each boxplot summarizes the ARI scores by applying a GSS tool on the

864   six gene set collections. KEGG, Kyoto encyclopedia of genes and genomes; GO, gene ontology;

865   GO:BP, GO biological process; GO:MF, GO molecular function; GO:CC, GO cellular

866   component; TFT, transcription factor target. **C.** Standard deviation (SD) of ARI scores on

867   different datasets (left) or GSS tools (right). To obtain the SD for each dataset, the average of

868   the SD of ARI scores of all GSS tools using different gene set collections was calculated. To

869   obtain the SD for each GSS tool, SD of ARI scores of the GSS tool on each dataset using

870   different gene set collections was calculated. Then the average of SD on different datasets for

871   each GSS tool was calculated.

872   **Figure 8. Evaluation of running time (in minute) of different GSS tools**

873   Three datasets were tested, including InSilico, Hematopoiesis and PBMC10K, which contain

874   approximately 500, 2000, and 10,000 cells, respectively. Four gene set collections were used,

875 including KEGG, TFT, REACTOME, and GO:BP, which contain approximately 200, 1000,

876 2000, and 7000 pathways, respectively. Cases where Pagoda2 failed to complete the calculation

877 are marked with '-'.

878 **Figure 9. Summarization of the performance of different GSS tools in various evaluation**

879 **pipelines measured by ARI**

880 **A.** ARI scores of different scATAC-seq datasets were averaged. Cases guiding the tool

881 recommendation are coloured in red. Each column denotes an evaluation task, which involves

882 GA transformation with each of the four tools, dropout imputation (no imputation or imputation

883 with each of the four tools), and selection of six gene set collections. Of note, when the dropout

884 imputation is performed for the peak-cell matrix, only MAESTRO can be used for GA

885 transformation because the other three GA tools are only applicable to the fragment file. **B.**

886 Practical guidelines for choosing appropriate tools for GSS. The GSS tool with border is the

887 most recommended tool with the best overall performance in the respective group.

888 **Supplementary materials**

889 **Figure S1  UMAP plots showing 2D-embeddings of the raw peak-cell matrix of eight**

890 **scATAC-seq datasets**

891 **Table S1 Size of gene set collections used in this study**

892 **Table S2  Detailed information of scATAC-seq and scRNA-seq datasets used in this**

893 **study**

894 **Table S3.  Human marker gene sets collected from the CellMarker database**

895

896

**A**

| | Silhouette | | ARI | | Accuracy | |
|---|---|---|---|---|---|---|
| | RNA | ATAC | RNA | ATAC | RNA | ATAC |
| AUCell | 0.16 | -0.04 | 0.50 | 0.39 | 0.94 | 0.85 |
| GSVA | 0.10 | 0.04 | 0.47 | 0.51 | 0.94 | 0.87 |
| Pagoda2 | 0.23 | 0.04 | 0.62 | 0.60 | 0.96 | 0.91 |
| PLAGE | 0.21 | 0.08 | 0.64 | 0.57 | 0.96 | 0.90 |
| ssGSEA | 0.09 | -0.10 | 0.44 | 0.34 | 0.94 | 0.87 |
| UniPath | 0.15 | -0.08 | 0.47 | 0.41 | 0.92 | 0.83 |
| VAM | 0.01 | -0.02 | 0.36 | 0.38 | 0.88 | 0.86 |
| Vision | 0.16 | 0.01 | 0.49 | 0.50 | 0.95 | 0.88 |
| z-score | 0.13 | -0.04 | 0.46 | 0.43 | 0.94 | 0.88 |
| **UniPath ATAC** | - | -0.01 | - | 0.55 | - | 0.83 |

High

Low

**C**



**B**

scRNA-seq

scATAC-seq

AUCell  GSVA  Pagoda2  PLAGE  ssGSEA  UniPath  VAM  Vision  z-score  **UniPath ATAC**

B
DC
Mono
NK
other
T

ARI

scRNA-seq  Signac  MAESTRO  ArchR  snapATAC

**Silhouette**

| Method | Hematopoiesis | Leukemia | SNAREmix | Forebrain | InSilico | Brain | PBMC3K | PBMC10K | Average |
|---|---|---|---|---|---|---|---|---|---|
| AUCell | -0.13 | -0.14 | -0.05 | -0.10 | 0.06 | -0.04 | -0.08 | 0.02 | -0.06 |
| GSVA | -0.24 | -0.23 | -0.10 | -0.09 | -0.02 | 0.06 | -0.05 | 0.12 | -0.07 |
| Pagoda2 | -0.13 | -0.08 | -0.05 | 0.02 | 0.20 | -0.05 | -0.03 | 0.21 | 0.01 |
| PLAGE | -0.28 | -0.23 | -0.07 | 0.03 | 0.42 | 0.02 | 0.03 | 0.19 | 0.01 |
| ssGSEA | -0.24 | -0.27 | -0.09 | -0.10 | -0.05 | -0.15 | -0.19 | 0.05 | -0.13 |
| UniPath | -0.14 | -0.16 | -0.09 | -0.12 | -0.04 | -0.10 | -0.11 | -0.03 | -0.10 |
| VAM | -0.09 | -0.11 | -0.06 | -0.08 | -0.06 | -0.05 | -0.06 | 0.05 | -0.06 |
| Vision | -0.18 | -0.16 | -0.05 | -0.06 | 0.00 | -0.01 | -0.04 | 0.08 | -0.05 |
| z-score | -0.24 | -0.24 | -0.10 | -0.09 | -0.06 | -0.11 | -0.08 | 0.06 | -0.11 |
| **UniPathATAC** | -0.18 | -0.14 | -0.08 | -0.09 | -0.00 | -0.01 | -0.07 | 0.06 | -0.06 |

**ARI**

| Method | Hematopoiesis | Leukemia | SNAREmix | Forebrain | InSilico | Brain | PBMC3K | PBMC10K | Average |
|---|---|---|---|---|---|---|---|---|---|
| AUCell | -0.00 | 0.03 | 0.00 | 0.07 | 0.23 | 0.38 | 0.29 | 0.49 | 0.19 |
| GSVA | 0.01 | 0.04 | 0.06 | 0.07 | 0.02 | 0.57 | 0.48 | 0.48 | 0.22 |
| Pagoda2 | 0.00 | 0.01 | -0.00 | 0.19 | 0.53 | 0.40 | 0.80 | 0.60 | 0.32 |
| PLAGE | 0.00 | 0.04 | 0.06 | 0.20 | 0.44 | 0.53 | 0.51 | 0.65 | 0.30 |
| ssGSEA | 0.00 | 0.04 | 0.06 | 0.05 | 0.02 | 0.41 | 0.29 | 0.33 | 0.15 |
| UniPath | 0.00 | 0.03 | 0.08 | 0.02 | -0.00 | 0.59 | 0.19 | 0.44 | 0.17 |
| VAM | 0.01 | 0.05 | 0.05 | 0.05 | 0.03 | 0.45 | 0.27 | 0.43 | 0.17 |
| Vision | 0.00 | 0.03 | 0.02 | 0.08 | 0.16 | 0.49 | 0.45 | 0.57 | 0.22 |
| z-score | 0.01 | 0.03 | 0.06 | 0.03 | 0.01 | 0.40 | 0.35 | 0.52 | 0.18 |
| **UniPathATAC** | 0.02 | 0.06 | 0.07 | 0.04 | 0.11 | 0.59 | 0.54 | 0.51 | 0.24 |

**Accuracy**

| Method | Hematopoiesis | Leukemia | SNAREmix | Forebrain | InSilico | Brain | PBMC3K | PBMC10K | Average |
|---|---|---|---|---|---|---|---|---|---|
| AUCell | 0.28 | 0.38 | 0.55 | 0.55 | 0.81 | 0.76 | 0.84 | 0.94 | 0.64 |
| GSVA | 0.33 | 0.43 | 0.59 | 0.56 | 0.84 | 0.77 | 0.89 | 0.95 | 0.67 |
| Pagoda2 | 0.29 | 0.30 | 0.50 | 0.67 | 0.91 | 0.83 | 0.93 | 0.97 | 0.67 |
| PLAGE | 0.37 | 0.42 | 0.65 | 0.67 | 0.94 | 0.82 | 0.92 | 0.96 | 0.72 |
| ssGSEA | 0.33 | 0.43 | 0.60 | 0.55 | 0.84 | 0.78 | 0.89 | 0.95 | 0.67 |
| UniPath | 0.23 | 0.27 | 0.48 | 0.35 | 0.73 | 0.76 | 0.80 | 0.92 | 0.57 |
| VAM | 0.27 | 0.34 | 0.55 | 0.48 | 0.68 | 0.78 | 0.87 | 0.95 | 0.61 |
| Vision | 0.33 | 0.45 | 0.55 | 0.59 | 0.85 | 0.79 | 0.90 | 0.95 | 0.68 |
| z-score | 0.30 | 0.37 | 0.57 | 0.55 | 0.81 | 0.79 | 0.89 | 0.95 | 0.65 |
| **UniPathATAC** | 0.37 | 0.34 | 0.55 | 0.34 | 0.86 | 0.73 | 0.83 | 0.94 | 0.62 |

Low — High

**A**

Silhouette / ARI / Classification plots with legend: SCALE, DrImpute, MAGIC, Saver, Raw

X-axis categories: AUCell, GSVA, Pagoda2, PLAGE, ssGSEA, UniPath, VAM, Vision, z-score, UniPath ATAC

**B**

| | Leuke. | | Hemat. | | HL. | | HEK. | | Fore. | | SNAR. | | InSil. | | PBMC | | Brain | | Avg. | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE | RAW | SCALE |
| AUCell | 0.04 | 0.16 | 0.00 | 0.23 | -0.00 | 0.58 | 0.03 | 0.60 | 0.01 | 0.40 | 0.00 | 0.43 | 0.02 | 0.56 | 0.05 | 0.45 | 0.05 | 0.28 | 0.02 | 0.42 |
| GSVA | 0.04 | 0.16 | 0.01 | 0.24 | 0.01 | 0.46 | 0.04 | 0.62 | 0.02 | 0.51 | 0.06 | 0.63 | 0.12 | 0.85 | 0.23 | 0.39 | 0.49 | 0.38 | 0.11 | 0.47 |
| Pagoda2 | 0.00 | 0.15 | 0.01 | 0.23 | -0.02 | 0.79 | 0.16 | 0.51 | 0.10 | 0.39 | -0.00 | 0.57 | 0.49 | 0.62 | 0.42 | 0.37 | 0.37 | 0.30 | 0.17 | 0.44 |
| PLAGE | 0.04 | 0.15 | 0.01 | 0.23 | 0.01 | 0.36 | 0.06 | 0.57 | 0.09 | 0.36 | 0.06 | 0.57 | 0.29 | 0.47 | 0.23 | 0.30 | 0.45 | 0.31 | 0.14 | 0.37 |
| ssGSEA | 0.04 | 0.15 | 0.00 | 0.29 | -0.01 | 0.30 | 0.04 | 0.68 | 0.03 | 0.48 | 0.06 | 0.52 | 0.04 | 0.63 | 0.12 | 0.42 | 0.33 | 0.36 | 0.07 | 0.43 |
| UniPath | 0.06 | 0.09 | 0.01 | 0.22 | -0.00 | 0.52 | 0.01 | 0.69 | 0.03 | 0.48 | 0.08 | 0.55 | 0.04 | 0.54 | 0.09 | 0.40 | 0.24 | 0.39 | 0.06 | 0.43 |
| VAM | 0.01 | 0.03 | 0.01 | 0.16 | 0.00 | 0.26 | 0.02 | 0.26 | 0.01 | 0.42 | 0.05 | 0.42 | 0.16 | 0.60 | 0.31 | 0.36 | 0.50 | 0.34 | 0.12 | 0.29 |
| Vision | 0.02 | 0.10 | 0.00 | 0.27 | -0.00 | 0.49 | 0.02 | 0.68 | 0.03 | 0.43 | 0.02 | 0.53 | 0.13 | 0.88 | 0.25 | 0.43 | 0.33 | 0.36 | 0.09 | 0.46 |
| z-score | 0.03 | 0.08 | 0.00 | 0.26 | 0.00 | 0.32 | 0.00 | 0.42 | 0.04 | 0.49 | 0.06 | 0.51 | 0.15 | 0.61 | 0.17 | 0.38 | 0.33 | 0.36 | 0.09 | 0.38 |
| UniPath ATAC | 0.06 | 0.17 | 0.02 | 0.14 | 0.01 | 0.37 | -0.00 | 0.37 | 0.04 | 0.11 | 0.07 | 0.58 | 0.11 | 0.41 | 0.31 | 0.36 | 0.55 | 0.33 | 0.13 | 0.32 |

ARI  0  0.2  0.4  0.6  0.8

**C**

| Imputation | Raw |
|---|---|
| AUCell | |
| GSVA | |
| Pagoda2 | |
| PLAGE | |
| ssGSEA | |
| UniPath | |
| VAM | |
| Vision | |
| z-score | |
| UniPath ATAC | |

BJ   H1ESC   GM   TF1

| | 500~200 | 500~1000 | 500~2000 | 500~7000 | 2000~200 | 2000~1000 | 2000~2000 | 2000~7000 | 10k~200 | 10k~1000 | 10k~2000 | 10k~7000 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AUCell | 0.28 | 1.06 | 1.08 | 4.15 | 1.13 | 4.25 | 4.37 | 6.96 | 2.96 | 19.73 | 20.72 | 31.80 |
| GSVA | 2.09 | 5.32 | 9.74 | 20.29 | 15.24 | 39.98 | 54.12 | 57.56 | 57.94 | 60.57 | 67.45 | 65.75 |
| Pagoda2 | 0.43 | 2.01 | - | - | 3.22 | 18.81 | 25.10 | 43.52 | 3.92 | 11.21 | 21.75 | 31.92 |
| PLAGE | 0.37 | 3.18 | 2.39 | 11.90 | 0.78 | 10.84 | 5.24 | 8.75 | 2.95 | 23.91 | 9.57 | 11.87 |
| ssGSEA | 5.08 | 24.72 | 38.99 | 174.94 | 19.47 | 109.53 | 170.54 | 286.13 | 22.68 | 117.00 | 201.19 | 251.72 |
| UniPath | 0.93 | 131.16 | 19.73 | 156.20 | 2.14 | 328.84 | 39.54 | 38.27 | 5.91 | 267.83 | 90.49 | 91.31 |
| VAM | 0.33 | 14.61 | 2.86 | 23.12 | 1.47 | 78.30 | 12.36 | 13.73 | 5.67 | 344.93 | 60.96 | 55.99 |
| Vision | 0.71 | 1.17 | 0.96 | 5.16 | 1.77 | 2.72 | 2.41 | 2.57 | 4.71 | 6.17 | 7.18 | 5.51 |
| z-score | 0.40 | 0.87 | 0.37 | 0.65 | 1.82 | 2.29 | 2.00 | 2.24 | 3.80 | 4.64 | 3.63 | 3.50 |
| UniPathATAC | 0.64 | 8.30 | 18.38 | 61.55 | 2.94 | 12.79 | 20.67 | 20.15 | 8.22 | 29.60 | 40.63 | 43.95 |

gene sets →
cells →

**A**

| | ArchR | MAESTRO | Signac | snapATAC | DrImpute | MAGIC | SAVER | SCALE | GO.BP | GO.CC | GO.MF | KEGG | REACTOME | TFT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| AUCell | 0.01 | 0.06 | 0.17 | 0.20 | 0.23 | 0.24 | 0.04 | 0.42 | 0.43 | 0.44 | 0.49 | 0.42 | 0.40 | 0.49 |
| GSVA | 0.22 | 0.18 | 0.18 | 0.14 | 0.39 | 0.23 | 0.14 | 0.47 | 0.43 | 0.43 | 0.43 | 0.47 | 0.43 | 0.35 |
| Pagoda2 | 0.19 | 0.25 | 0.29 | 0.31 | 0.43 | 0.18 | 0.21 | 0.44 | 0.23 | 0.36 | 0.25 | 0.44 | 0.39 | 0.35 |
| PLAGE | 0.19 | 0.21 | 0.27 | 0.22 | 0.45 | 0.22 | 0.22 | 0.37 | 0.37 | 0.40 | 0.41 | 0.37 | 0.38 | 0.32 |
| ssGSEA | 0.17 | 0.13 | 0.12 | 0.09 | 0.39 | 0.27 | 0.19 | 0.43 | 0.47 | 0.43 | 0.43 | 0.43 | 0.39 | 0.34 |
| UniPath | 0.02 | 0.11 | 0.16 | 0.16 | 0.48 | 0.20 | 0.14 | 0.43 | 0.47 | 0.47 | 0.47 | 0.43 | 0.48 | 0.41 |
| VAM | 0.10 | 0.20 | 0.16 | 0.15 | 0.28 | 0.30 | 0.14 | 0.29 | 0.21 | 0.33 | 0.33 | 0.29 | 0.32 | 0.24 |
| Vision | 0.21 | 0.15 | 0.17 | 0.17 | 0.39 | 0.26 | 0.16 | 0.46 | 0.47 | 0.50 | 0.50 | 0.46 | 0.49 | 0.44 |
| z-score | 0.14 | 0.14 | 0.12 | 0.14 | 0.31 | 0.25 | 0.14 | 0.38 | 0.35 | 0.42 | 0.41 | 0.38 | 0.42 | 0.30 |
| UniPath ATAC | 0.22 | 0.22 | 0.22 | 0.22 | 0.24 | 0.21 | 0.16 | 0.32 | 0.45 | 0.42 | 0.40 | 0.32 | 0.31 | 0.44 |
| **Average** | **0.15** | **0.17** | **0.19** | **0.18** | **0.36** | **0.24** | **0.15** | **0.40** | **0.39** | **0.42** | **0.41** | **0.40** | **0.40** | **0.37** |

ARI High / Low

**Preprocess**

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Imputation | | Raw | | | | | | | | | SCALE | | | |
| GA tools | | MAESTRO | | | | MAESTRO | | | | | MAESTRO | | | |
| Gene sets | | KEGG | | | | KEGG | | | | | | KEGG | | |

**B**

GSS

scATAC-seq data → Imputation **SCALE** → GA conversion **MAESTRO** → Vision, Pagoda2, GSVA, ssGSEA, AUCell

**Any gene sets** →

Raw peak data → GA conversion **Signac / snapATAC** → Pagoda2, PLAGE

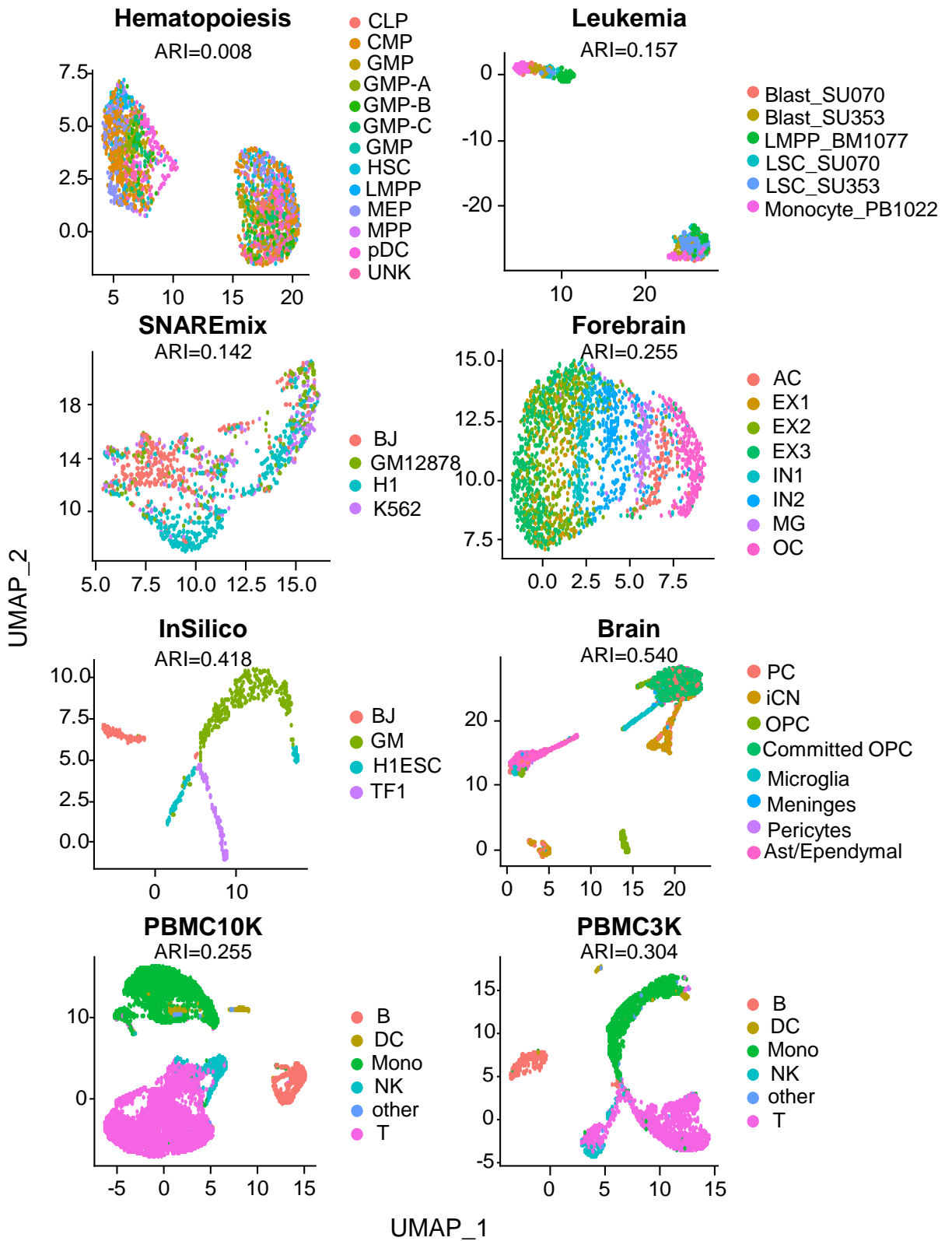Raw peak data **No GA conversion** → UniPathATAC

**Figure S1. UMAP plots showing 2D-embeddings of the raw peak-cell matrix of eight scATAC-seq datasets.**