

1 Single-cell mean rank gene set scoring method for between-dataset comparison of  
2 scRNA-seq data

3

4 Dazhou Li<sup>1</sup>, Guohao Meng<sup>2\*</sup>, Jieyu Wu<sup>1</sup>, Guihui Tong<sup>1</sup>

5 1 Department of Pathology, First Affiliated Hospital of Guangzhou Medical University, Guangzhou,  
6 510150, China

7 2 Department of Pathophysiology, Key Laboratory of Cell Differentiation and Apoptosis of Chinese  
8 Ministry of Education, Shanghai Jiaotong University School of Medicine, Shanghai, 200025, China

9 \* Correspondence author

10 Email: [mgh0525@outlook.com](mailto:mgh0525@outlook.com)

11

12 **Abstract**

13 The surge in single-cell RNA sequencing (scRNA-seq) data offers a unique chance for researchers  
14 to understand functional changes in biological processes and diseases through gene set scoring  
15 across diverse datasets. Despite this, current methods for comparing scRNA-seq data at the  
16 signaling pathway level across datasets remain untested. To bridge this gap, we introduce the single-  
17 cell mean rank gene set scoring (scMRGSS) method, which assesses gene set activity between  
18 different scRNA-seq datasets. Leveraging gene expression ranks within each dataset, scMRGSS  
19 calculates mean rank scores for gene sets, enabling the comparison of their relative enrichment or  
20 depletion across datasets. Demonstrating its efficacy through simulated and real datasets, scMRGSS  
21 proves to be a simple yet informative tool for comparing gene set activity between cell types across  
22 diverse datasets. Its robustness against sequencing depth and dropout rate variations underscores its  
23 value for integrative scRNA-seq data analysis. Applying the method, we uncover that abnormal  
24 activity in oxidative phosphorylation and NF- $\kappa$ B signaling pathways in glioblastoma cancer cells  
25 may not solely stem from neurodevelopmental programs. Notably, the highest activity of these

26 pathways is observed in the mesenchymal cancer cell type, emphasizing the need to target specific  
27 cell types in glioblastoma drug development.

28 **Keywords:** single-cell analysis, RNA sequencing, gene set analysis, glioblastoma, NF- $\kappa$ B pathway

29

## 30 **Introduction**

31 Single-cell RNA sequencing (scRNA-seq) has revolutionized various fields of biological research  
32 by enabling the comprehensive profiling of gene expression at the individual cell level. Over the  
33 past decade, a substantial amount of scRNA-seq data has been accumulated through collaborative  
34 efforts in studying diverse biological systems and diseases(1–5). While scRNA-seq data offers  
35 several advantages in terms of scale and resolution, it presents unique challenges for data analysis  
36 and interpretation due to its noisy and zero-inflated nature(6). These challenges arise from inherent  
37 biological phenomena and the limited capture efficiency of the technology and are therefore  
38 difficult to overcome(7). As a result, specialized analytical methods are necessary to effectively  
39 process scRNA-seq data and extract meaningful biological insights.

40 A particularly noteworthy approach is the aggregation of gene expression profiles into biologically  
41 functional representations, commonly referred to as gene set analysis. This method aims to better  
42 comprehend the biological relevance of scRNA-seq data. A plethora of tools have been developed  
43 to estimate the activity of *a priori* gene sets or pathways based on scRNA-seq data, including  
44 single-cell signature explorer (SCSE)(8), AUCell(9), single-cell gene set enrichment analysis  
45 (scGSEA)(7), variance-adjusted Mahalanobis (VAM)(10), Pagoda2(11), and Vision(12). These tools  
46 enable researchers to assess the enrichment or depletion of specific gene sets in different cell  
47 populations, thereby providing insight into their functional differences. Considering the rapid  
48 increase in scRNA-seq data, it has become appealing to compare gene set activity underlying  
49 essential biological processes across various scRNA-seq datasets to identify commonalities and  
50 differences in cellular states. Regrettably, existing single-cell gene set scoring methods have yet to  
51 address this issue.

52 In this work, we present the single-cell mean rank gene set scoring (scMRGSS), a straightforward  
53 method for comparing gene set activity in single-cell RNA sequencing (scRNA-seq) data across  
54 different datasets. The scMRGSS method calculates a rank-based score for each gene set in a single  
55 cell mainly by averaging the ranks of genes belonging to that gene set. This simple approach  
56 facilitates easy interpretation of results, and by utilizing rank-based scores instead of raw or  
57 transformed expression values, it accounts for the noisy nature of scRNA-seq data. Our study  
58 demonstrates that scMRGSS is a reliable and robust method for comparing gene set activity  
59 between scRNA-seq datasets, as evidenced by simulations and real scRNA-seq datasets of cell lines  
60 and peripheral blood. Additionally, we applied the method to cancer biology and found that the up-  
61 regulation of oxidative phosphorylation and NF- $\kappa$ B signaling pathways in glioblastoma did not  
62 solely reflect its resemblance to normal neurodevelopmental lineages, and there was notable  
63 heterogeneity across glioblastoma subtypes.

64

## 65 **Methods**

### 66 **Single-cell mean rank gene set scoring**

67 The scMRGSS method, which is adopted from the algorithm developed by Noureen et al.(13),  
68 calculates the normalized mean gene rank for each gene set among expressed genes in a cell,  
69 resulting in a score that ranges from 0 to 1. This method is robust to commonly used scRNA-seq  
70 normalization methods, such as counts per million (CPM), due to its reliance on gene ranks. In this  
71 study, scMRGSS was used to compare the scores between two datasets from different datasets in  
72 order to determine if there is a significant difference in the activity of a gene set between the  
73 datasets.

74 The method begins with basic data filtering, in which each dataset is subject to filtering of cells and  
75 genes in order to reduce technical and biological noise. The intersection of genes between the two  
76 datasets is selected as the input for the method. Gene sets are represented as a list of HUGO gene  
77 symbols, and those with less than the threshold proportion of expressed genes in a dataset (typically

78 0.6 in this study) are ignored for further procedure. Only expressed genes (with a count or  
79 expression value greater than 0) are considered for the calculation in order to minimize the  
80 influence of dropouts across cells in the same dataset. The genes are then ranked in each cell, and  
81 the rank-based score for each gene set is computed using the following formula:

$$82 \text{ score}^{(s,c)} = \frac{\sum_{i=1}^m \text{rank}_i^{(s,c)}}{m \cdot n}$$

83 where  $m$  is the number of expressed genes of gene set  $s$  in cell  $c$ , and  $n$  is the total number of  
84 expressed genes in this cell. The metric is essentially the effective mean rank normalized by the  
85 number of expressed genes, and the normalization procedure allows for more comparable scores  
86 across cells.

87

### 88 **Simulation of single-cell RNA-seq datasets and gene sets**

89 The zero-inflated negative binomial (ZINB) model was employed to generate simulated single-cell  
90 RNA sequencing (scRNA-seq) datasets comprised of 4000 cells divided into four groups, each with  
91 varying library sizes or dropout rates. The VGAM R package was used to execute the simulation  
92 processes(14). Each group of cells expressed background genes and group-specific genes from eight  
93 gene sets, the sizes of which ranged from 50 to 120, with increments of 10 for both simulation  
94 scenarios. All genes not specific to the group were considered as background genes. The datasets,  
95 which corresponded to individual parameterisations, thus consisted of 2720 genes in total. The  
96 negative binomial (NB) distribution can be viewed as a gamma-Poisson mixture where the lambda  
97 parameter of Poisson distribution is distributed as a gamma distribution parameterised by shape and  
98 scale. For the library size simulation, the shape (size) parameter of the rzinegbin function was  
99 adjusted from 4 to 8 in increments of 1, while the munb parameter equalled product of shape and  
100 scale, where the scale parameter remained fixed at 3. The pstr0 parameter, which represents the  
101 probability of structural zero in the ZINB distribution, was set at 0.5. In the scenario of the dropout  
102 simulation, the shape and scale parameters were maintained at fixed values of 4 and 3, respectively.

103 Meanwhile, the  $pstr0$  parameter varied from 0.2 to 0.8 in increments of 0.15. Similar to the  
104 distribution for group-specific genes in the dataset, background gene counts were also simulated  
105 using the ZINB distribution, of which the shape parameter was the same, the scale parameter was  
106 scaled by 0.5, and the  $pstr0$  parameter was increased by 0.3 and bounded above by 0.9.

107

## 108 **Performance evaluation**

109 Simulation of scRNA-seq datasets was performed to evaluate the capacity of the proposed method  
110 in producing distinct gene set scores for different biological groups. The Splatter R package's  
111 `splatSimulateGroups` function(15) was utilised to generate datasets with two biological groups. The  
112 simulated dataset encompassed 4000 cells and 5000 genes, with approximately 20% of the genes  
113 designated as differentially expressed genes (DEGs). The DEGs exhibited a 50% probability of  
114 being down-regulated. Subsequent to dataset generation, random sampling procedures were  
115 employed to construct five distinct collections, each comprising 80 gene sets. These gene sets  
116 varied in size, ranging from 50 to 145 genes, with incremental steps of 5 genes. Importantly, the  
117 collections were distinguished by varying percentage of DEGs, spanning from 20% to 100%. For  
118 each gene set size, four gene sets were systematically generated using a consistent procedure for  
119 both biological groups. Specifically, within each group and gene set size, a group-specific gene set  
120 and a parallel group-unspecific gene set were simulated. The group-specific gene set was composed  
121 of group-specific genes sampled from the group-specific gene pool, along with background genes  
122 sampled from the group-unspecific gene pool. In contrast, the group-unspecific gene set solely  
123 comprised group-unspecific genes.

124 To quantitatively assess the performance, a diagnostic metric denoted as  $d$  was introduced, defined  
125 by the following equation:

$$126 \quad d = \frac{((- \log_{10}(adj.p)/2) + \max\{fc, 1/fc\})}{2}$$

127 where  $adj.p$  is the adjusted p-value, and  $fc$  is the fold change (see below). Subsequently, this metric  
128  $d$  was employed to construct receiver operating characteristic (ROC) curve and precision-recall

129 (PR) curve for each DEG percentage using the EvalMetrics Julia package  
130 (<https://github.com/VaclavMacha/EvalMetrics.jl>). The objective was to meticulously evaluate the  
131 efficacy of scMRGSS in accurately classifying differences in gene set scores as either true or false  
132 under varying noise conditions, offering insights into the nuanced performance dynamics.

133

#### 134 **Public single-cell RNA sequencing datasets and gene sets**

135 We utilised nine public scRNA-seq datasets in this study. These datasets were obtained from a  
136 variety of sources, including previous research studies and public databases. Four separate datasets  
137 consisting of Jurkat cells, 293T cells, a 50/50 mixture of Jurkat cells and 293T cells, a 99/1 mixture  
138 of Jurkat cells and 293T cells were obtained from a previous study(16). PBMC 10k 3p and PBMC  
139 10k 5p datasets were downloaded from the dataset portal of 10x Genomics. Another human  
140 peripheral blood dataset was obtained from the Tabular Sapiens project(1). The glioblastoma and  
141 associated neural development datasets were obtained from previous study(17). Finally, another  
142 human fetal brain development dataset(18) was downloaded from the GEO database with accession  
143 number GSE162170. All dataset sources are listed in the Data Availability section. Furthermore, we  
144 obtained the Biocarta gene set collection, which consists of 292 gene sets, and the oxidative  
145 phosphorylation gene set of KEGG (v2023.1) from the MSigDB database(19).

146

#### 147 **Pre-processing of single-cell RNA sequencing datasets**

148 The raw counts of scRNA-seq datasets were utilized as input for the method. For the human-  
149 derived datasets, we applied filters to exclude genes with numbers of expressed cells less than 2,  
150 and cells with numbers of expressed genes less than 2000. The threshold for filtering cells was  
151 determined based on the distribution of the number of genes for the simulated datasets. These  
152 procedures were pre-implemented in the scMRGSS Julia package and easily executable. We utilized  
153 pre-defined cell type labels when possible, and manually curated the datasets otherwise.

154 For the PBMC 10k 3p and BMC 10k 5p datasets, we employed the Seurat R package to cluster cells  
155 separately(20). We first removed low-quality cells based on the number of expressed genes and  
156 mitochondrial gene expression percentage, resulting in 8,827 and 9,894 cells in the two datasets,  
157 respectively. The sctransform method was applied to normalize the data(21), followed by linear  
158 dimensional reduction and cell clustering using the Louvain algorithm(22). The identified cell  
159 clusters were then embedded in non-linear low-dimensional space using the uniform manifold  
160 approximation and projection (UMAP) algorithm to better display local relationships. Next, we  
161 conducted cluster biomarker identification through gene differential expression analysis and  
162 compared the results with reference cell types to curate the cell clusters. We employed the SingleR  
163 package for this comparison(23). Cell labels were consistent among three reference sources for  
164 most identified cell clusters in the PBMC 10k 3p dataset (Supplementary Fig. S2A-C) and the  
165 PBMC 10k 5p dataset (Supplementary Fig. S2E-G), and cell labels were assigned in both datasets  
166 accordingly (Supplementary Fig. S2D,H).

167 The glioblastoma dataset contained curated cell labels, although they were incomplete. Only cells  
168 derived from the whole tumor samples were included in the study. The cells without labels and  
169 duplicates were excluded. In total, 18,430 cells were selected. The fetal brain cells at 24 post-  
170 conceptional weeks (pcw) were excluded from the brain development dataset (GSE162170)(18) to  
171 match the donor age in the fetal brain dataset associated with the glioblastoma study(17).

172

### 173 **Comparison of gene set scores between two groups**

174 To compare gene set scores between two groups, whether from the same dataset or not, we  
175 leveraged the results of hypothesis testing and the fold change. By default, the two-tailed Mann-  
176 Whitney U test implemented by the HypothesisTests Julia package  
177 (<https://github.com/JuliaStats/HypothesisTests.jl>) was applied to determine if the observed gene set  
178 score difference was statistically significant. The Bonferroni multiple hypothesis correction  
179 implemented by the MultipleTesting Julia package

180 (<https://github.com/juliangehring/MultipleTesting.jl>) was used to adjust p-values pertaining to the  
181 tested gene sets. An adjusted p-value of less than 0.01 was considered indicative of statistical  
182 significance. The gene set score difference was deemed biologically meaningful if the adjusted p-  
183 value was less than 0.01 and the fold change or its reciprocal was greater than 1.2 in the context of  
184 simulation study and applications in Jurkat and 293T cells as well as human peripheral blood cells.  
185 For the application in glioblastoma, the Bonferroni procedure was utilised to correct p-values for  
186 multiple group pair-wise comparisons regarding the selected gene set and the fold change threshold  
187 was set at 1.1. We also calculated the proportion of meaningfully differential gene sets group pair-  
188 wise and displayed the results on a heatmap using the Makie Julia package(24) for comparisons  
189 related to the cell lines and peripheral blood cells. Volcano plots were created to identify differential  
190 gene sets between the two groups using the EnhancedVolcano R package  
191 (<https://bioconductor.org/packages/EnhancedVolcano>). Venn diagrams were created to illustrate the  
192 overlapping differential gene sets using the ggVennDiagram R package(25).

193

## 194 **Results**

### 195 **scMRGSS accurately identifies cell group-specific gene sets in simulation study**

196 Different scRNA-seq preparation kits and sequencing platforms often produce data with different  
197 library sizes and dropouts levels. To assess the capacity of the method to detect difference in gene  
198 set activity between different conditions, we sought to perform simulation studies on the effects of  
199 library size and dropout rates on the performance of scMRGSS. Firstly, five scRNA-seq count  
200 datasets with various library sizes (or total count numbers) were simulated based on the ZINB  
201 model. Each dataset was composed of 2720 genes and 4000 cells, which were divided into four  
202 groups. Eight gene sets with increasing number of genes were simulated at the same time  
203 exclusively for one group. Density plots verified that the simulated datasets varied in terms of  
204 library size, while numbers of genes were comparable (Supp. Fig. S1A,B). As expected, cells within  
205 the same group showed higher similarity in the low-dimensional embedding and were clustered



206 together (Fig. 1A). Library size had a minor effect on cell clustering compared to the group  
207 condition (Fig. 1B). Each group of cells specifically expressed genes from the corresponding  
208 simulated gene sets (Fig. 1C). We combined the adjusted p-value and fold change to determine  
209 whether the detected differential gene sets were meaningful in the sense of biological context (see  
210 Methods). Proportions of biologically meaningful differential gene sets were larger in the same-  
211 group pairs than in the different-group pairs independent of library size, indicating that scMRGSS  
212 accurately identifies cell group-specific gene sets despite the challenge of shrinking library size  
213 (Fig. 1D).

214 Additionally, we conducted the simulation study on the impact of dropout rates using the ZINB  
215 model likewise. As dropout rates increased in the simulated datasets, both the total number of  
216 counts and the number of genes decreased (Supp. Fig. S1C,D). Group condition dominated the cell  
217 clustering compared to the dropout rate (Fig. 1E,F). Each cell group showed unique gene set  
218 activity pattern, with group exclusive gene sets only expressed in their respective groups (Fig. 1G).  
219 Proportions of differential gene sets were close to zero between the same groups from two datasets,  
220 and approximated 50% between different groups, in line with the theoretical values (Fig. 1H). This  
221 pattern was consistent regardless of the dropout levels, highlighting the robustness of scMRGSS in  
222 accurately identifying cell group-specific gene sets even in the presence of varying dropout rates.  
223 Taken together, these results indicated that scMRGSS is a robust method for detecting difference in  
224 gene set activity in scRNA-seq data in the simulation settings.

225

### 226 **Simulation study demonstrates favourable performance of scMRGSS in yielding distinct** 227 **scores for distinct conditions**

228 We endeavoured to evaluate the method's capability to call true differential gene sets between two  
229 conditions using simulated datasets, taking into account the trade-off between sensitivity and  
230 specificity, as well as the trade-off between precision and recall. To assess the robustness of the  
231 method against noise, we incorporated varying proportions (0, 0.2, 0.4, 0.6, 0.8) of random

232 background genes into the group-specific gene sets. A diagnostic metric, integrating the adjusted p-  
233 value and fold change (see Methods), was developed to predict differential gene sets between the  
234 two conditions. We found that the 4000 simulated cells from the two groups displayed considerable  
235 overlap in a low-dimensional linear space, indicating a significant similarity between the two  
236 conditions (Fig. 2A). The scMRGSS demonstrated satisfactory sensitivity while maintaining high  
237 specificity even in the presence of 80% noise in the simulated gene sets (Fig. 2B). The PR curve is  
238 more informative than the ROC curve in imbalanced scenarios(26). The area under the PR curve  
239 surpassed 0.9 at all levels of noise (Fig. 2C), indicating that scMRGSS consistently exhibited high  
240 precision and recall in calling differential gene sets between the two conditions. Collectively, these  
241 findings demonstrate the robustness and favourable performance of scMRGSS in generating  
242 different gene set scores for two groups on simulated datasets.

243

#### 244 **scMRGSS differentiates biological conditions in human cell line and peripheral blood datasets**

245 In order to evaluate the utility of scMRGSS for comparing gene set activity between two datasets in  
246 real-world situations, we applied the method to both human cell line and peripheral blood datasets.  
247 We used scMRGSS to calculate Biocarta gene set scores for Jurkat cells and 293T cells from four  
248 separate datasets and then determined the proportions of meaningful differential gene sets between  
249 each group pair. Our results revealed that the proportions of differential gene sets were substantially  
250 lower between same-cell groups compared to different-cell groups, indicating that scMRGSS  
251 effectively distinguishes between biological conditions (Fig. 3A). Furthermore, we examined the  
252 overlap of detected differential gene sets and found that only two differential gene sets overlapped  
253 between Jurkat cells from the Jurkat datasets and Jurkat cells from the 50/50 mixture dataset out of  
254 all dataset pairs of Jurkat cells (Fig. 3B). On the other hand, nine differential gene sets were shared  
255 among the four Jurkat-293T dataset pairs, and the majority of differential gene sets for a dataset pair  
256 overlapped with at least one dataset pair, suggesting consistent and reproducible differentiation  
257 between cell lines using scMRGSS (Fig. 3C).

258 We utilized a comparable analytical approach to assess the efficacy of scMRGSS in human  
259 peripheral blood datasets. This analysis encompassed CD14<sup>+</sup> monocytes, CD16<sup>+</sup> monocytes, naive  
260 CD4<sup>+</sup> T cells, and B cells from three separate scRNA-seq sources generated by 10x Genomics  
261 microfluidic droplet and Smart-seq2 techniques(1). Similarly, we calculated the proportions of  
262 differential Biocarta gene sets between pairwise groups from the same or different datasets and  
263 observed similar patterns. In general, the proportions of differential gene sets were lower between  
264 same-cell groups compared to different-cell groups. Additionally, CD14<sup>+</sup> monocytes and CD16<sup>+</sup>  
265 monocytes shared a higher degree of similarity regarding this proportion compared to other cell  
266 group pairs (Fig. 3D), consistent with their intrinsic biological similarity and cell clustering results  
267 (Supplementary Fig. S2). The largest proportion of differential gene sets was observed between  
268 CD14<sup>+</sup> monocytes and naive CD4<sup>+</sup> T cells (Fig. 3D). We further investigated the detected  
269 differential gene sets in greater depth and found that the majority of differential gene sets were  
270 unique to each dataset pair in the cases of comparisons involving naive CD4<sup>+</sup> T cells (Fig. 3E-H).  
271 However, when comparing CD16<sup>+</sup> monocytes and naive CD4<sup>+</sup> T cells, there was a higher overlap  
272 of differential gene sets among the four dataset combinations (Fig. 3I-M), indicating that the  
273 identified gene set score difference was consistent across different sources of scRNA-seq data.  
274 Given the above, our results illustrate that scMRGSS is effective in identifying differential gene sets  
275 in various cell types between scRNA-seq datasets.

276

### 277 **scMRGSS reveals cancer cell heterogeneity of oxidative phosphorylation and NF- $\kappa$ B pathway** 278 **activity in glioblastoma**

279 Identifying differential gene sets and revealing potential heterogeneity in cellular activity across  
280 different cell types and sources can be of great utility to researchers. This is particularly true when  
281 comparing gene set activity in cells from various sources. One such application is the analysis of  
282 scRNA-seq datasets from glioblastoma samples. A previous study revealed that glioblastoma cancer  
283 cells resembled the human fetal brain lineages, including truncated radial glia (tRG), glial

284 progenitor cells (GPCs), oligo-lineage cells (OPCs), and interneurons, and hence the cancer cells  
285 were classified into mesenchymal, glial-progenitor, oligo-lineage, and neuronal subtypes(17).  
286 Another study uncovered four cell states and subtypes in glioblastoma along neurodevelopmental  
287 and metabolic axes, of which the mitochondrial cell state depends on the oxidative phosphorylation  
288 pathway for energy production(27).

289 NF- $\kappa$ B pathway is active in glioblastoma and plays an essential role in tumour progression and  
290 treatment resistance(28,29). We were interested in investigating if the up-regulated activity of  
291 oxidative phosphorylation and NF- $\kappa$ B pathways in glioblastoma could be explained by its  
292 resemblance to neurodevelopmental lineages and further if heterogeneity across cancer cell types  
293 exists regarding the activity of these pathways. In addition to the glioblastoma and neural  
294 development datasets included in the study of Couturier et al.(17), we incorporated another human  
295 brain development dataset(18) in the comparisons of gene set scores to demonstrate the consistency  
296 of our findings. The difference in gene set scores was viewed as significant in the biological sense if  
297 the related adjusted p-value was below 0.01 and the fold change exceeded 1.1. Our analysis  
298 revealed that gene set scores of the oxidative phosphorylation pathway were comparable between  
299 tRG from two neural development datasets (Fig. 4A). Notably, the gene set score of the  
300 mesenchymal cell type was higher when compared to tRG, whereas the scores of the other cancer  
301 cell types exhibited rising trends. When the scores of various cancer cell types were compared to  
302 those of GPCs and OPCs, respectively, this pattern remained (Supplementary Fig. S3A,B). This  
303 suggests that the up-regulated activity of the oxidative phosphorylation pathway in glioblastoma  
304 could not be merely attributed to its resemblance to neurodevelopmental lineages. Furthermore, the  
305 analysis revealed significant cellular heterogeneity within glioblastoma tumours with regard to the  
306 activity of the oxidative phosphorylation pathway.

307 The NF- $\kappa$ B pathway exhibited different patterns of activity across the various cancer cell types of  
308 glioblastoma. While tRG from two sources showed comparable scores of NF- $\kappa$ B pathway activity,  
309 the scores were significantly higher in the mesenchymal subtype compared to tRG (Fig. 4B).

310 However, the activity of the NF- $\kappa$ B pathway was lower in the neuronal subtype than in tRG, while  
311 the oligo-lineage and glial progenitor subtypes remained comparable to tRG. All other subtypes  
312 showed lower activity of the NF- $\kappa$ B pathway compared to the mesenchymal subtype, indicative of  
313 notable heterogeneity within glioblastoma with regard to the activity of the NF- $\kappa$ B pathway (Fig.  
314 4B). Using GPCs and OPCs as the reference groups, respectively, we found similar patterns  
315 (Supplementary Fig. S3C,D). Collectively, these findings suggest that the up-regulated activity of  
316 the oxidative phosphorylation and NF- $\kappa$ B pathways in glioblastoma are not solely due to their  
317 resemblance to neurodevelopmental lineages. Moreover, the results also demonstrate the  
318 consistency of the method in comparing gene set scores between datasets.

319

## 320 **Discussion**

321 In recent years, the volume of scRNA-seq data has increased significantly across various research  
322 domains. This surge in data has made it feasible to compare the activity of gene sets of interest  
323 between different cell types across datasets. However, existing methods have not yet been tested to  
324 address this issue. In this study, we introduced scMRGSS, a simple yet effective approach for  
325 performing between-dataset comparisons of gene set activity using scRNA-seq data. The reliance of  
326 scMRGSS on the gene rank rather than the expression value confers it with greater robustness  
327 against commonly adopted normalisation and transformation methods, such as CPM, log  
328 transformation and z-score normalisation. This simplifies the analysis pipeline and enhances the  
329 applicability of the method to diverse datasets. Additionally, the normalisation of mean rank  
330 facilitates more consistent comparisons of cells between datasets with different scales of gene  
331 expression. We demonstrated the potential of scMRGSS in cancer research by identifying  
332 differences in the activity of oxidative phosphorylation and NF- $\kappa$ B pathways between glioblastoma  
333 cancer cell types and tRG, as well as highlighting the cellular heterogeneity within glioblastoma  
334 concerning these pathway activity.

335 The efficacy of the method in identifying disparate gene sets between biological conditions while  
336 concurrently regulating false positive rates was demonstrated through both simulated and genuine  
337 data analysis. The performance of scMRGSS is contingent upon the premise that the technical bias  
338 of mRNA capture and amplification between the two data sources remains consistent or varies only  
339 marginally. To elucidate, when the biological variations between the two datasets are equivalent, the  
340 discrepancy in gene rank on average is negligible. This assumption is generally applicable to data  
341 derived from the 10x Genomics and Smart-seq2 platforms, particularly when obtained from the  
342 same platform. Nonetheless, when employing the method to compare pathway activity of two cell  
343 types from two datasets generated by distinct scRNA-seq protocols or platforms, it would be  
344 prudent to include one shared cell type between the two datasets to assess the validity of the  
345 assumption. In such cases, single-cell dataset integration approaches, such as Scanorama(16) and  
346 scGen(30), may be required to account for the technical bias. In light of the challenge in  
347 establishing strict rules for determining thresholds in adjusted p-values and fold changes, the  
348 interpretation of disparities in gene set scores demands meticulous attention. For instance, assuming  
349 the difference in gene set scores remains constant, larger group sizes often yield smaller p-values.  
350 Given that gene set scores range between 0 and 1, the upper limit for the fold change must be  
351 constrained by a value depending upon the reference score. Accumulating datasets from analogous  
352 sources will contribute valuable insights towards identifying pattern of gene set activity difference.  
353 scMRGSS is designed for gene set analysis following the clustering and labelling of cells within  
354 datasets. Although it is a non-parametric and rank-invariant method, it is essential to note that basic  
355 gene and cell filtering is necessary to ensure the reliability and accuracy of the analysis results. This  
356 is because noise in the data could potentially impact the gene rank(31) and the mean and dispersion  
357 of gene set scores within a cell type when performing hypothesis testing. As a general guideline, it  
358 is advisable to set the filtering threshold based on the distribution of the data and *a priori*  
359 knowledge about the cells of interest. For datasets generated by common scRNA-seq platforms such  
360 as 10x Genomics and Smart-seq2, we typically recommend removing genes with the number of

361 cells below 2 and removing cells with the number of expressed genes below 2000. Additionally, the  
362 percentage of mitochondrial genes may be integrated into the filtering scheme to account for  
363 potential cell stress or contamination. Furthermore, it is recommended to include gene sets whose  
364 percentage of expressed genes in both datasets exceeds some threshold, such as 0.6 in this study, to  
365 avoid identifying differential gene sets of little biological value depending on only a small subset of  
366 genes. The exact percentage may depend on both the gene set and the datasets.

367 In this study we also sought to investigate the relationships of the activity of oxidative  
368 phosphorylation and NF- $\kappa$ B pathways between neurodevelopmental lineages and cancer cell types  
369 in glioblastoma. Previous research has shown that glioblastoma cells epitomise the normal  
370 neurodevelopmental process and resemble several lineages during the process, including tRG,  
371 GPCs, OPCs and interneurons(17). It is established that most glioblastoma cancer cells produce  
372 energy through oxidative phosphorylation as opposed to glycolysis(32). The key role of oxidative  
373 phosphorylation in glioblastoma is also emphasized by the potent inhibitory effect of the oxidative  
374 phosphorylation inhibitor on glioblastoma cancer cells(33). Our findings suggest that the metabolic  
375 feature of the glioblastoma mesenchymal cell type is not associated with the neurodevelopmental  
376 programs, as we observed elevated oxidative phosphorylation activity in the mesenchymal cell  
377 type compared to brain development lineages. We speculate this cancer cell type may overlap with a  
378 recently identified mitochondrial subtype that depends exclusively on oxidative phosphorylation for  
379 energy production(27). Though not as strong as in mesenchymal cells, other glioblastoma cancer  
380 cell types also showed increasing trends in oxidative phosphorylation activity, indicating  
381 glioblastoma cellular heterogeneity in energy metabolism. Aberrant NF- $\kappa$ B activation in  
382 glioblastoma contributes to cancer cell proliferation, invasion, mesenchymal differentiation, and  
383 resistance to radiotherapy(34). We found that NF- $\kappa$ B activity were significantly elevated in the  
384 mesenchymal cell type of glioblastoma relative to neurodevelopmental lineages, but not other  
385 subtypes. This suggests that the neurodevelopmental programs hijacked by glioblastoma cannot  
386 fully explain NF- $\kappa$ B pathway activation in all cancer cell types. The heterogeneity of glioblastoma

387 also highlights the need for a more precise approach to target the mesenchymal cancer subtype  
388 when developing drugs aimed at the NF- $\kappa$ B pathway in glioblastoma treatment.  
389 In brief, we have presented a simple and efficient method to compare gene set activity amongst cell  
390 types from various scRNA-seq datasets. By applying this approach to glioblastoma, it has been  
391 discovered that the disease's aberrant oxidative phosphorylation and NF- $\kappa$ B pathway activity is not  
392 exclusively caused by neurodevelopmental programs. This finding has implications for precision  
393 medicine strategies aimed at addressing particular cancer subtypes.

394

### 395 **Availability of data**

396 All scRNA-seq datasets used in this study are available for download. We used the following public  
397 datasets:

398 pbmc\_10k\_3p dataset from 10x Genomics:

399 [https://cf.10xgenomics.com/samples/cell-exp/4.0.0/Parent\\_NGSC3\\_DI\\_PBMC/](https://cf.10xgenomics.com/samples/cell-exp/4.0.0/Parent_NGSC3_DI_PBMC/)

400 [Parent\\_NGSC3\\_DI\\_PBMC\\_filtered\\_feature\\_bc\\_matrix.h5](#);

401 pbmc\_10k\_5p dataset from 10x Genomics:

402 [https://cf.10xgenomics.com/samples/cell-vdj/5.0.0/sc5p\\_v2\\_hs\\_PBMC\\_10k/](https://cf.10xgenomics.com/samples/cell-vdj/5.0.0/sc5p_v2_hs_PBMC_10k/)

403 [sc5p\\_v2\\_hs\\_PBMC\\_10k\\_filtered\\_feature\\_bc\\_matrix.h5](#);

404 Tabular Sapiens dataset from figshare:

405 [https://figshare.com/articles/dataset/Tabula\\_Sapiens\\_release\\_1\\_0/14267219](https://figshare.com/articles/dataset/Tabula_Sapiens_release_1_0/14267219);

406 Glioblastoma dataset: [https://github.com/mbourgey/scRNA\\_GBM](https://github.com/mbourgey/scRNA_GBM);

407 Human neural development dataset: [https://github.com/mbourgey/scRNA\\_GBM](https://github.com/mbourgey/scRNA_GBM);

408 Human cortex development dataset: [https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE162170)

409 [acc=GSE162170](#).

410 All human gene sets in the study were downloaded from MSigDB

411 (<https://www.gsea-msigdb.org/gsea/msigdb>).



412 Simulated datasets have been deposited on figshare with the following DOI

413 <https://doi.org/10.6084/m9.figshare.24886047>.

414

### 415 **Availability of code**

416 scMRGSS is developed in Julia programming language v1.9.3 (<https://julialang.org/>) and available

417 at <https://github.com/giuseppedelnapalle/scmrgss>, while the scripts to reproduce results of the

418 manuscript can be accessed on zenodo at the DOI <https://doi.org/10.5281/zenodo.10418687>.

419

### 420 **Funding**

421 This work was carried out independent of funding.

422

### 423 **Acknowledgements**

424 We would like to thank Jiawei Li for helpful questions and comments about the manuscript.

425

### 426 **Author contributions**

427 G.M. conceived the problem and the algorithm. D.L. and G.M. developed the method and

428 performed the computational experiments. D.L., G.M., J.W. and G.T. wrote, reviewed and revised

429 the manuscript. All authors read and approved the final version of the manuscript.

430

### 431 **Competing interests**

432 The authors declare no competing in interests.

433

### 434 **References**

1. THE TABULA SAPIENS CONSORTIUM. The Tabula Sapiens: A multiple-organ, single-cell transcriptomic atlas of humans. Science. 2022 May 13;376(6594):eab14896.

2. Domínguez Conde C, Xu C, Jarvis LB, Rainbow DB, Wells SB, Gomes T, et al. Cross-tissue immune cell analysis reveals tissue-specific features in humans. *Science*. 2022 May 13;376(6594):eabl5197.
3. Eraslan G, Drokhyansky E, Anand S, Fiskin E, Subramanian A, Slyper M, et al. Single-nucleus cross-tissue molecular reference maps toward understanding disease gene function. *Science*. 2022 May 13;376(6594):eabl4290.
4. Yao Z, Liu H, Xie F, Fischer S, Adkins RS, Aldridge AI, et al. A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex. *Nature*. 2021 Oct;598(7879):103–10.
5. Jardine L, Webb S, Goh I, Quiroga Londoño M, Reynolds G, Mather M, et al. Blood and immune development in human fetal bone marrow and Down syndrome. *Nature*. 2021 Oct;598(7880):327–31.
6. Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert JP. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun*. 2018 Jan 18;9(1):284.
7. Franchini M, Pellicchia S, Viscido G, Gambardella G. Single-cell gene set enrichment analysis and transfer learning for functional annotation of scRNA-seq data. *NAR Genom Bioinform*. 2023 Mar;5(1):lqad024.
8. Pont F, Tosolini M, Fournié JJ. Single-Cell Signature Explorer for comprehensive visualization of single cell signatures across scRNA-seq datasets. *Nucleic Acids Res*. 2019 Dec 2;47(21):e133.
9. 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 Nov;14(11):1083–6.
10. Frost HR. Variance-adjusted Mahalanobis (VAM): a fast and accurate method for cell-specific gene set scoring. *Nucleic Acids Res*. 2020 Sep 18;48(16):e94.
11. 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 Jan;36(1):70–80.
12. DeTomaso D, Jones MG, Subramaniam M, Ashuach T, Ye CJ, Yosef N. Functional interpretation of single cell similarity maps. *Nat Commun*. 2019 Sep 26;10(1):4376.
13. Noureen N, Ye Z, Chen Y, Wang X, Zheng S. Signature-scoring methods developed for bulk samples are not adequate for cancer single-cell RNA sequencing data. *Elife*. 2022 Feb 25;11:e71994.
14. Yee TW. The VGAM Package for Categorical Data Analysis. *Journal of Statistical Software*. 2010 Jan 5;32:1–34.
15. Zappia L, Phipson B, Oshlack A. Splatter: simulation of single-cell RNA sequencing data. *Genome Biol*. 2017 Sep 12;18(1):174.
16. Hie B, Bryson B, Berger B. Efficient integration of heterogeneous single-cell transcriptomes using Scanorama. *Nat Biotechnol*. 2019 Jun;37(6):685–91.

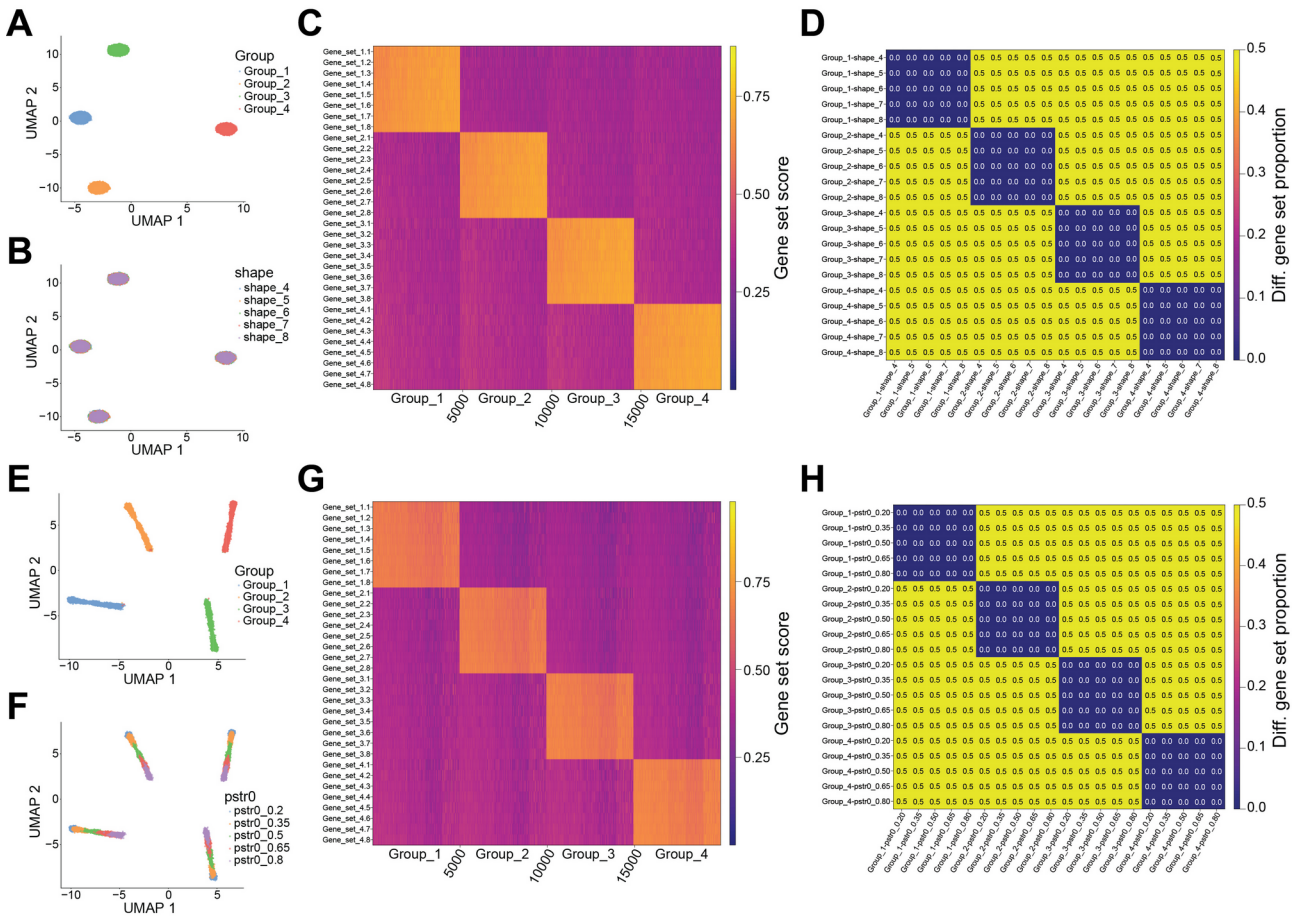
17. Couturier CP, Ayyadhury S, Le PU, Nadaf J, Monlong J, Riva G, et al. Single-cell RNA-seq reveals that glioblastoma recapitulates a normal neurodevelopmental hierarchy. *Nat Commun.* 2020 Jul 8;11(1):3406.
18. Trevino AE, Müller F, Andersen J, Sundaram L, Kathiria A, Shcherbina A, et al. Chromatin and gene-regulatory dynamics of the developing human cerebral cortex at single-cell resolution. *Cell.* 2021 Sep 16;184(19):5053-5069.e23.
19. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures database (MSigDB) 3.0. *Bioinformatics.* 2011 Jun 15;27(12):1739–40.
20. Hao Y, Hao S, Andersen-Nissen E, Mauck WM, Zheng S, Butler A, et al. Integrated analysis of multimodal single-cell data. *Cell.* 2021 Jun 24;184(13):3573-3587.e29.
21. Hafemeister C, Satija R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol.* 2019 Dec 23;20(1):296.
22. Blondel VD, Guillaume JL, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. *J Stat Mech.* 2008 Oct;2008(10):P10008.
23. Aran D, Looney AP, Liu L, Wu E, Fong V, Hsu A, et al. Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage. *Nat Immunol.* 2019 Feb;20(2):163–72.
24. Danisch S, Krumbiegel J. Makie.jl: Flexible high-performance data visualization for Julia. *Journal of Open Source Software.* 2021 Sep 1;6(65):3349.
25. Gao CH, Yu G, Cai P. ggVennDiagram: An Intuitive, Easy-to-Use, and Highly Customizable R Package to Generate Venn Diagram. *Front Genet.* 2021;12:706907.
26. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLoS One.* 2015;10(3):e0118432.
27. Garofano L, Migliozi S, Oh YT, D'Angelo F, Najac RD, Ko A, et al. Pathway-based classification of glioblastoma uncovers a mitochondrial subtype with therapeutic vulnerabilities. *Nat Cancer.* 2021 Feb;2(2):141–56.
28. Ji J, Ding K, Luo T, Zhang X, Chen A, Zhang D, et al. TRIM22 activates NF- $\kappa$ B signaling in glioblastoma by accelerating the degradation of I $\kappa$ B $\alpha$ . *Cell Death Differ.* 2021 Jan;28(1):367–81.
29. Xiang J, Alafate W, Wu W, Wang Y, Li X, Xie W, et al. NEK2 enhances malignancies of glioblastoma via NIK/NF- $\kappa$ B pathway. *Cell Death Dis.* 2022 Jan 14;13(1):58.
30. Lotfollahi M, Wolf FA, Theis FJ. scGen predicts single-cell perturbation responses. *Nat Methods.* 2019 Aug;16(8):715–21.
31. Foroutan M, Bhuvana DD, Lyu R, Horan K, Cursons J, Davis MJ. Single sample scoring of molecular phenotypes. *BMC Bioinformatics.* 2018 Nov 6;19(1):404.
32. Bonnay F, Veloso A, Steinmann V, Köcher T, Abdusselamoglu MD, Bajaj S, et al. Oxidative Metabolism Drives Immortalization of Neural Stem Cells during Tumorigenesis. *Cell.* 2020 Sep 17;182(6):1490-1507.e19.

33. Shi Y, Lim SK, Liang Q, Iyer SV, Wang HY, Wang Z, et al. Gboxin is an oxidative phosphorylation inhibitor that targets glioblastoma. *Nature*. 2019 Mar;567(7748):341–6.

34. Soubannier V, Stifani S. NF- $\kappa$ B Signalling in Glioblastoma. *Biomedicines*. 2017 Jun 9;5(2):29.

435 **Figures and Legends**

436



437 **Figure 1.** Simulation study of the single-cell mean rank gene set scoring (scMRGSS) method.

438 (A,B) UMAP representation of 20,000 simulated cells categorised into four groups across five

439 datasets with varying library sizes, labelled by group (A) and the shape (size) parameter (B) for the

440 zero-inflated negative binomial (ZINB) model, respectively. Each shape parameter value

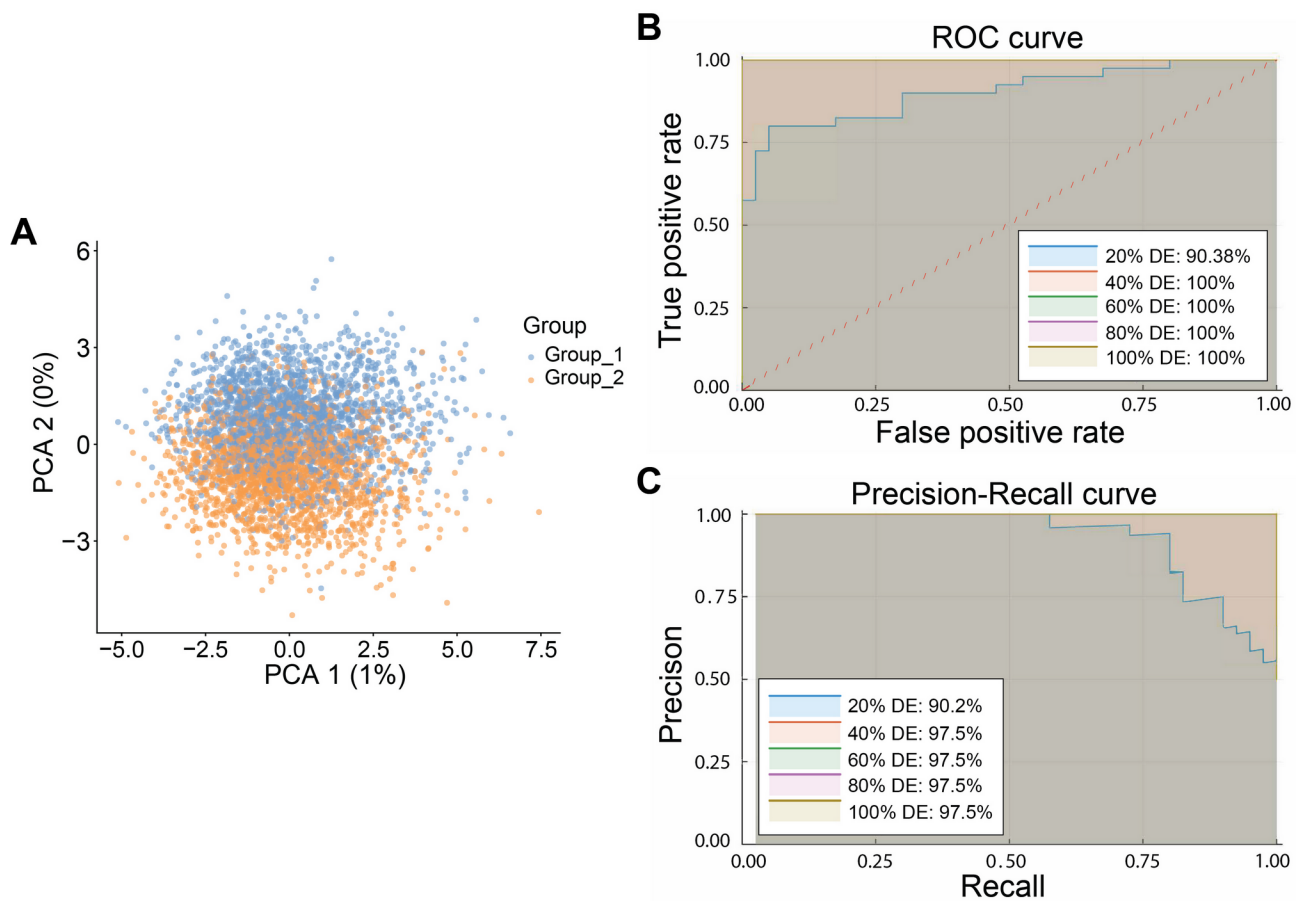
441 corresponds to a distinct dataset. (C) Gene set activity of 32 simulated gene sets computed by

442 scMRGSS across 20,000 cells from five distinct datasets as presented in (A,B). Rows represent

443 gene sets while columns represent cells in the heatmap. The hierarchical clustering algorithm

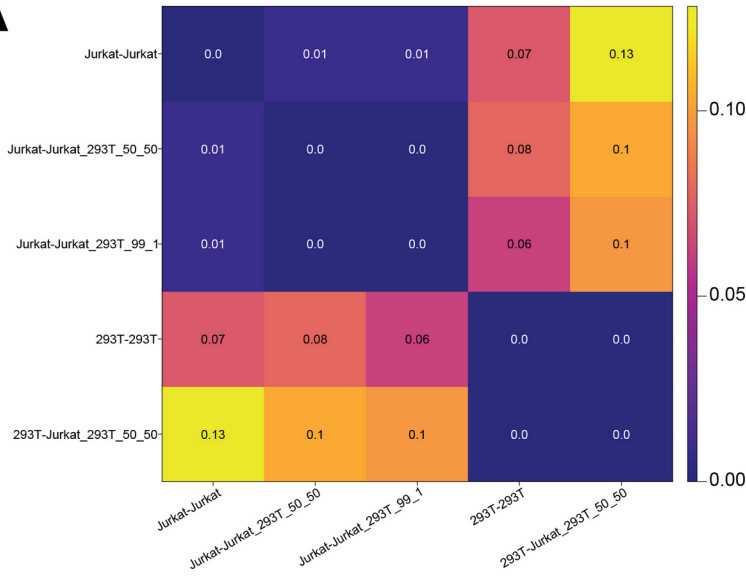
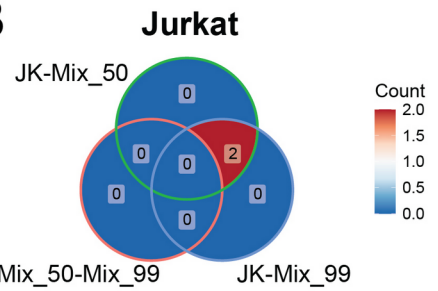
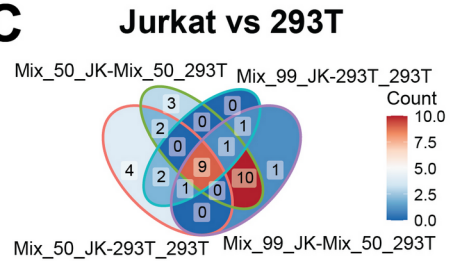
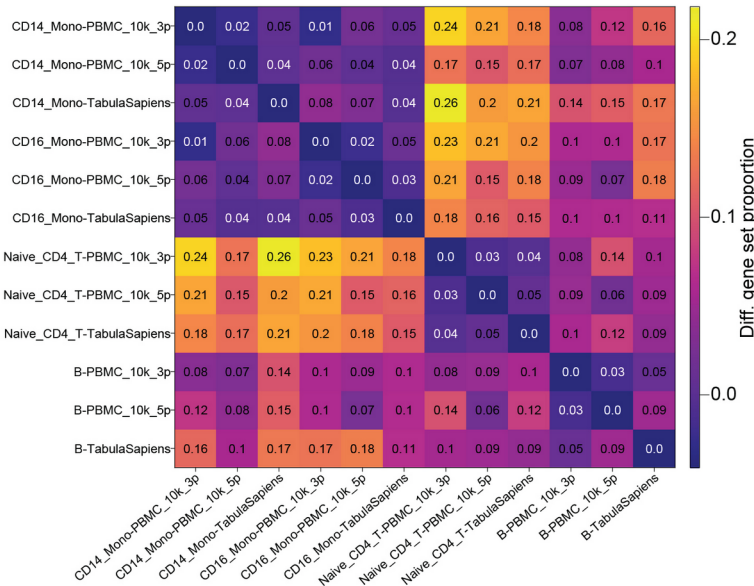
444 groups cells into four categories mirroring the simulated groups. (D) Heatmap illustrating the

445 proportions of biologically meaningful differential gene sets between groups from the same or  
446 distinct datasets as observed in **(A,B)**. The difference in gene set activity is considered valid if the  
447 associated adjusted p-value is below 0.01 and the fold change or its reciprocal is larger than 1.2.  
448 **(E,F)** UMAP plots depicting 20,000 simulated cells categorised into four groups across five datasets  
449 with varying dropout rates, labelled by group **(E)** and the pstr0 parameter **(F)** for the zero-inflated  
450 negative binomial (ZINB) model, respectively. pstr0 controls the probability of structural zero in the  
451 ZINB model, and each pstr0 value corresponds to a different dataset. **(G)** Heatmap presenting the  
452 activity of 32 simulated gene sets computed by scMRGSS across 20,000 cells from five distinct  
453 datasets as shown in **(E,F)**. Rows represent gene sets while columns represent cells. The  
454 hierarchical clustering algorithm categorises cells into four groups as expected. **(H)** Proportions of  
455 biologically meaningful differential gene sets between groups from same or distinct datasets as  
456 depicted in **(E,F)**.  
457



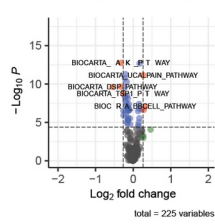
458 **Figure 2.** Performance evaluation of scMRGSS. **(A)** Principal component analysis (PCA) biplot of  
 459 4000 simulated cells from two groups. **(B)** Receiver operating characteristic (ROC) curve of five  
 460 classifiers associated with different differential expressed gene (DEG) percentage (ranging from  
 461 20% to 100%) of the group-specific gene sets. For each DEG percentage, 20 group 1-specific gene  
 462 sets, 20 group 2-specific gene sets, and 40 unspecific gene sets of varying sizes spanning from 50 to  
 463 145 genes were generated from the simulated single-cell RNA sequencing (scRNA-seq) dataset in  
 464 **(A)**. A diagnostic metric  $d$  incorporating the adjusted p-value and the fold change was introduced to  
 465 formulate classifying models for each DEG percentage to predict the veracity of the observed  
 466 differences in gene set scores between the two groups. The ROC curve serves as a visual  
 467 representation of the performance of the classifying models. **(C)** Precision-recall (PR) curve of the  
 468 five classifiers as described in **(B)**.

469

**A****B****C****D****E**

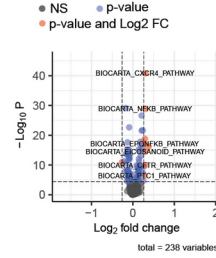
Naive CD4+ T cells: 5p vs 3p

- NS
- Log2 FC
- p-value
- p-value and Log2 FC

**F**

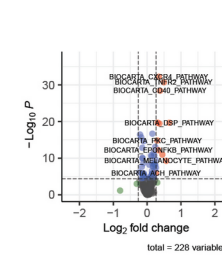
Naive CD4+ T cells: TS vs 3p

- NS
- p-value
- p-value and Log2 FC

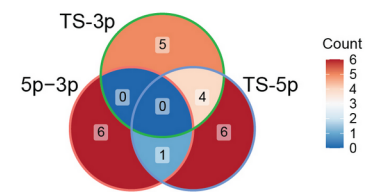
**G**

Naive CD4+ T cells: TS vs 5p

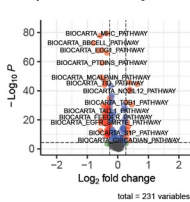
- NS
- p-value
- p-value and Log2 FC

**H**

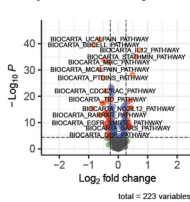
**Naive CD4+ T cells**

**I**

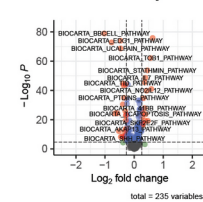
3p-Naive CD4+ T cells vs 3p-CD16+ Monocytes

**J**

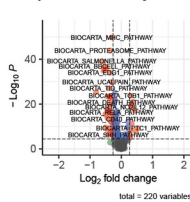
3p-Naive CD4+ T cells vs 5p-CD16+ Monocytes

**K**

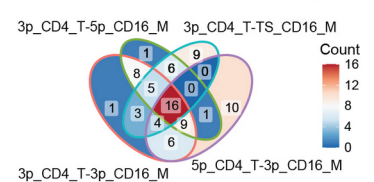
3p-Naive CD4+ T cells vs TS-CD16+ Monocytes

**L**

5p-Naive CD4+ T cells vs 3p-CD16+ Monocytes

**M**

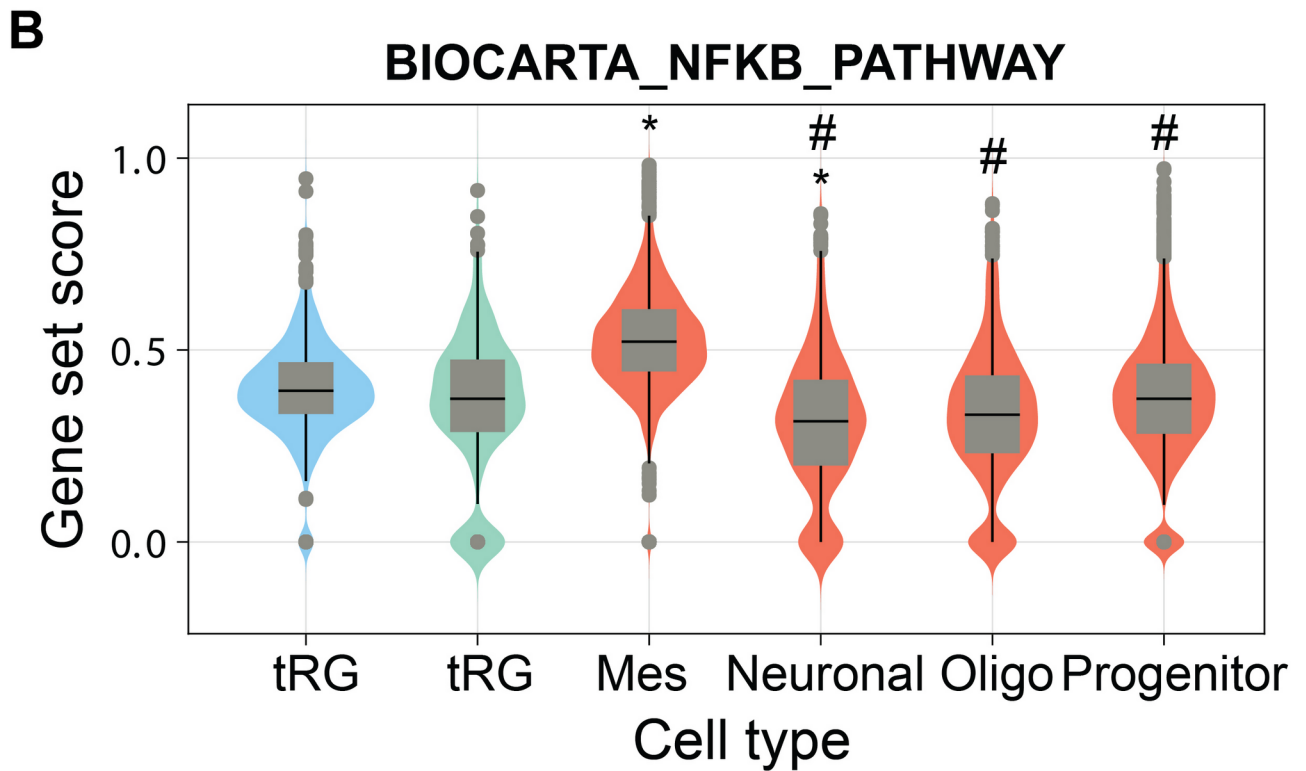
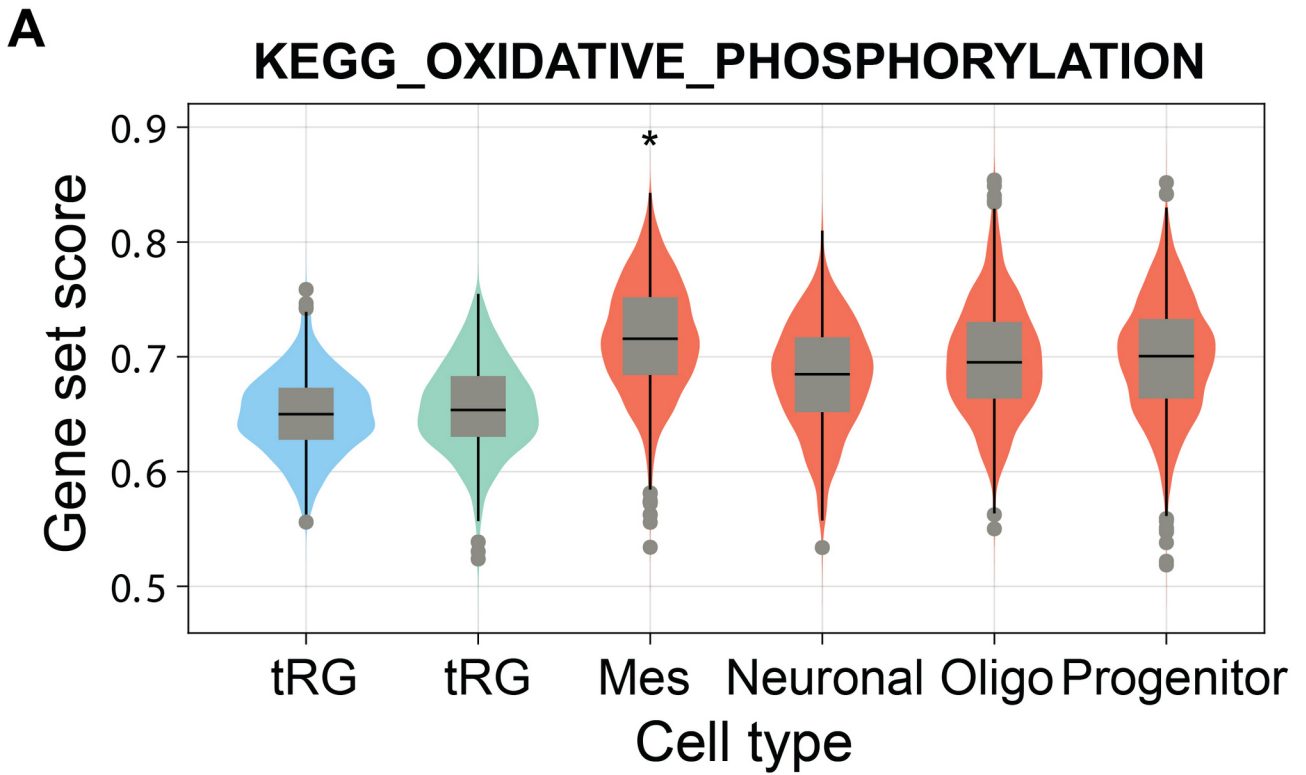
Naive CD4+ T cells vs CD16+ Monocytes



471 **Figure 3.** scMRGSS distinguishes biological groups in real datasets. **(A)** Heatmap displaying the  
472 proportions of biologically meaningful differential Biocarta gene sets between cell lines from four  
473 datasets comprising Jurkat and/or 293T cells. scMRGSS was applied to a collection of four datasets  
474 consisting of one entirely of Jurkat cells (Jurkat), one entirely of 293T cells (293T), a 50/50 mixture  
475 of Jurkat and 293T cells (Jurkat\_293T\_50\_50), and a 99/1 mixture of Jurkat and 293T cells  
476 (Jurkat\_293T\_99\_1). Cell line names and Dataset name are split by a hyphen (“-”) in each row or  
477 column. The difference in gene set scores is considered meaningful if the adjusted p-value is below  
478 0.01 and the fold change or its reciprocal is larger than 1.2. **(B)** Venn diagram illustrating how  
479 differential gene sets between Jurkat cells from three different datasets as described in **(A)** overlap.  
480 **(C)** Overlap of differential gene sets between Jurkat and 293T cells from four different datasets.  
481 Mix\_50: Jurkat\_293T\_50\_50; Mix\_99: Jurkat\_293T\_99\_1; JK: Jurkat cells; 293T: 293T dataset or  
482 293T cells. **(D)** Proportions of biologically meaningful differential Biocarta gene sets between  
483 human peripheral blood cells from the PBMC\_10k\_3p, PBMC\_10k\_5p and TabularSapienes  
484 datasets. The analysis takes into account four distinct cell groups: CD14+ monocytes, CD16+  
485 monocytes, naive CD4+ T cells and B cells. An identical process was carried out as in **(A)**. **(E-G)**  
486 Identification of differential gene sets of naive CD4+ T cells between PBMC\_10k\_5p and  
487 PBMC\_10k\_3p datasets **(E)**, between TabularSapienes and PBMC\_10k\_3p datasets **(F)**, and  
488 between TabularSapienes and PBMC\_10k\_5p datasets **(G)**, respectively. **(H)** Overlap of differential  
489 gene sets between naive CD4+ T cells from three different datasets depicted in **(D)**. **(I-L)**  
490 Identification of differential gene sets between naive CD4+ T cells and CD16+ monocytes from the  
491 three datasets. Pairwise comparisons of gene set activity were conducted between naive CD4+ T  
492 cells and CD16+ monocytes of the PBMC\_10k\_3p dataset **(I)**, between naive CD4+ T cells of the  
493 PBMC\_10k\_3p dataset and CD16+ monocytes of the PBMC\_10k\_5p dataset **(J)**, between naive  
494 CD4+ T cells of the PBMC\_10k\_3p dataset and CD16+ monocytes of the TabularSapienes dataset  
495 **(K)**, between naive CD4+ T cells of the PBMC\_10k\_5p dataset and CD16+ monocytes of the  
496 PBMC\_10k\_3p dataset **(L)**, respectively. **(M)** Overlap of differential gene sets between naive CD4+



497 T cells and CD16+ monocytes from the three datasets. 3p: PBMC\_10k\_3p; 5p: PBMC\_10k\_5p; TS:  
498 TabularSapienes; CD4\_T: naive CD4+ T cells; CD16\_M: CD16+ monocytes.  
499



501 **Figure 4.** Application of scMRGSS in glioblastoma. **(A-B)** Gene set activity of the KEGG oxidative  
502 phosphorylation **(A)** and the Biocarta NF- $\kappa$ B **(B)** pathways estimated by scMRGSS in distinct  
503 cancer cell types of glioblastoma compared to truncated radial glia (tRG) of two datasets. tRG of  
504 the cortex development dataset (first column) and the mesenchymal cancer cells are chosen as  
505 reference to calculate p-values. Bonferroni procedure was used to adjust the p-values. The  
506 difference in gene set scores is considered meaningful if the adjusted p-value is below 0.01 and the  
507 fold change or its reciprocal is larger than 1.1. \* denotes the difference is meaningful when  
508 compared to tRG, whereas # denotes the difference is meaningful when compared to the  
509 mesenchymal cancer cells.