

Single-cell biclustering for cell-specific transcriptomic perturbation detection in AD progression

Yuqiao Gong¹, Jingsi Xu¹, Ruitian Gao¹, Jianle Sun¹, Zhangsheng

Yu^{1,2,3,4,*} and Yue Zhang^{1,*}

¹Department of Bioinformatics and Biostatistics, School of Life Sciences and Biotechnology, Shanghai Jiao Tong University, 800 Dongchuan Road, Minhang District, 200240, Shanghai, China

²SJTU-Yale Joint Center for Biostatistics and Data Science Organization, Shanghai Jiao Tong University, Shanghai, China

³Clinical Research Institute, Shanghai Jiao Tong University School of Medicine, Shanghai, China

⁴Center for Biomedical Data Science, Translational Science Institute, Shanghai Jiao Tong University School of Medicine, Shanghai, China

*Corresponding author. yuzhangsheng@sjtu.edu.cn, yue.zhang@sjtu.edu.cn

Abstract

Alzheimer's disease (AD) is a complex and debilitating neurodegenerative disorder that has received widespread attention and research in recent years. The pathogenesis of AD involves complex regulation between genes and intricate changes in different cell types. Searching for a single pathogenic gene or investigating changes in a single cell type has limitations in understanding the disease. Therefore, there is a need for a more comprehensive approach to analyze AD. Here, we propose a single-cell Bayesian biclustering (scBC) framework for the cell-specific detection of network gene biomarkers in scRNA-seq data, utilizing the biclustering method to analyze the perturbations in functional gene modules of complex diseases at single-cell level. By applying our framework to AD scRNA-seq data, we uncover the perturbations of functional gene modules under different cell groups and shed light on gene-cell correlations during AD progression. Our method can handle the challenges of single-cell data such as batch effects and the high proportion of dropout, which are common issues in single-cell data analysis. Moreover, by incorporating existing biological information into our model, we obtain biologically meaningful results. We conducted comparisons on simulated datasets and several highly heterogeneous real-world datasets at different levels, demonstrating that our method is more precise and robust than other state-of-the-art biclustering methods. Our study highlights the great potential of scBC in uncovering the mechanism of polygenic diseases with complex gene co-expression patterns and providing potential treatment options.

Keywords: Functional gene modules, Biclustering, scRNA-seq, scBC, Alzheimer's disease

40 **Introduction**

41 In recent years, the advancement of single-cell sequencing technology has enabled
42 the analysis of single-cell data to reveal meaningful biological information at the
43 cellular level. Specifically, single-cell RNA sequencing (scRNA-seq) enables the
44 sequencing of cells that are hard-to-retrieve or challenging to isolate[1]. This
45 unprecedented resolution into cell states provides us with new insights into the
46 function and dysfunction of cells[2], which is particularly necessary for complex
47 diseases like Alzheimer's disease (AD), as changes in gene expression are related to
48 cell type[3, 4]. Recently, there has been a surge of single-cell studies aimed at
49 understanding the mechanism of AD based on transcriptional profiles[3, 5, 6], which
50 have provided valuable insights into cellular diversity. However, these studies often
51 lack an integrative analysis of functional gene modules, which can reveal how genes
52 work together to regulate biological processes. A recent study used a network-based
53 approach to identify functional gene modules involved in the selective vulnerability of
54 neurons in Alzheimer's disease, demonstrating the importance of analyzing functional
55 gene modules (FGMs) to gain insights into the underlying mechanisms of complex
56 diseases like AD[7]. FGMs are groups of genes that work together to perform a specific
57 biological function and can exhibit complex co-expression or co-regulation patterns,
58 rather than solely comprising differentially expressed genes[8, 9]. Moreover, these
59 local patterns are often cell-specific and may change with disease progression[10-12].
60 Therefore, it is crucial to identify functional gene modules and their corresponding
61 functional cell groups simultaneously in studies of complex diseases. In this study, we
62 focus on FGMs, the network gene biomarkers, to investigate their potential role in AD.

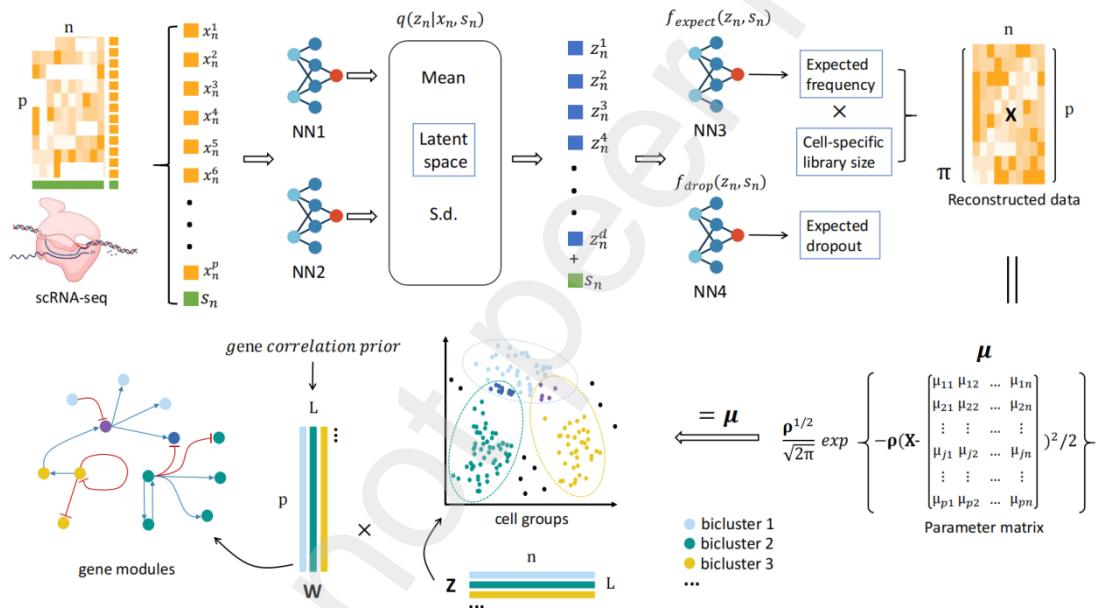
63

64 Unlike clustering methods, which can only conduct clustering in either cell space or
65 gene space, biclustering can identify functional gene modules and their corresponding
66 functional cell groups simultaneously. A cells-genes pair is called a bicluster, the genes
67 in a specific bicluster can be deemed as a FGM shared across related cells. Meanwhile,
68 these cells behave similarly with respect to these relative genes. Therefore, through
69 biclustering, we can easily identify functional gene modules and find the cell
70 populations in which they are active at the same time. It is worth noting that in the
71 complex cell machinery, multiple FGMs are active in a cell group, and different cell
72 groups may share a common FGM (Fig. S1). Fortunately, biclustering with overlapped
73 biclusters can easily capture such complex features[13]. Through biclustering, cell
74 population-specific network gene biomarker and potential gene-cell connections can
75 be identified in a single pass.

76

77 Although biclustering is an exquisite tool, it encounters some intractable problems
78 when applied to scRNA-seq data. First, batch effects, due to laboratory conditions,
79 reagent lots, and personnel differences, are widespread and critical to address[14]. If
80 batch effects are not properly accounted for, biclustering algorithms may falsely
81 identify batch-specific co-expression patterns instead of true biological patterns,
82 leading to incorrect conclusions. Second, due to low mRNA content per cell and

83 molecule losses during the experiment (known as "dropout"), the gene expression
 84 matrix has a substantial amount of zero read counts which can cause problems for
 85 biclustering algorithms that assume continuous expression values[15]. Biclustering
 86 algorithms that are not designed to handle dropout may either ignore the zero read
 87 counts, leading to incomplete biclusters, or consider them as low-expressed genes,
 88 leading to spurious biclusters. Although algorithms have been designed to address
 89 these inherent problems pervasive in scRNA-seq data, they can only improve the
 90 performance of the biclustering algorithm in one particular aspect – cell clustering,
 91 FGM finding, or the simultaneity of co-clustering[16-18]. However, in complex
 92 polygenic diseases, functionally related potential cell groups are finely divided, cell-
 93 type conditional gene co-expression patterns are complicated, and the cell-gene
 94 correlation changes throughout the progression. Therefore, an algorithm with better
 95 performance in functionally related cell group discovery, FGM finding, and cell-gene
 96 correlation pattern detection is urgently needed to assist the research of these
 97 diseases.



98
 99 **Fig. 1 | flowchart of scBC procedure.** The scRNA-seq data with high proportion of dropout and
 100 batch annotation (if available) is first fed into the variational autoencoder (VAE). We use x_n^g to
 101 denote the g_{th} gene in cell n . s_n is an extra dimension added for each cell to denote the batch
 102 annotation. Through training process, we can get a low dimensional approximate posterior
 103 distribution $q(z_n | x_n, s_n)$ conditional on s_n . At inference stage, the low dimensional
 104 representation of each cell z_n is taken to reconstruct the expression data through non-linear
 105 mapping. The like likelihood function of gene g from cell n is $\pi(x_{ji} | \mu_j, \rho_j) = \frac{\rho_j^{1/2}}{\sqrt{2\pi}} \exp\{-\rho_j(x_{ji} - \mu_j)^2/2\}$. To
 106 reduce randomness, we decompose the parameter matrix μ of the reconstructed data matrix X
 107 rather than directly on X . Gene correlation prior is used to guide the variable selection at gene
 108 level, and results are presented as matrix W with each column denoting a module. Matrix Z is the
 109 result at cell level with each row representing a functionally-related cell group.

110
 111 Here we propose a novel single-cell biclustering method (scBC) that can handle the
 112 problems mentioned above and provide a more reliable result. We utilize a variational

113 autoencoder (VAE) to model gene expression in single cells, enabling us to gracefully
114 remove batch effects and impute missing data[19]. By estimating the variational
115 posterior distribution, we can obtain a low-dimensional representation of each cell
116 that is conditioned on the batch annotation, enabling us to obtain batch-corrected
117 expression through the generating process. Additionally, we can manually control the
118 procedure of "dropout" during the generating process, leading to the pre-dropout
119 imputed expression. By reconstructing the original data matrix in this way, we can
120 obtain more precise results when conducting biclustering. Furthermore, we
121 incorporate existing biological information (e.g., gene interaction and regulation) into
122 the biclustering procedure through the Bayesian framework, which guides variable
123 selection to more likely capture pathway information and true biological signals[20].
124 The flowchart of our procedure is depicted in Fig. 1.
125

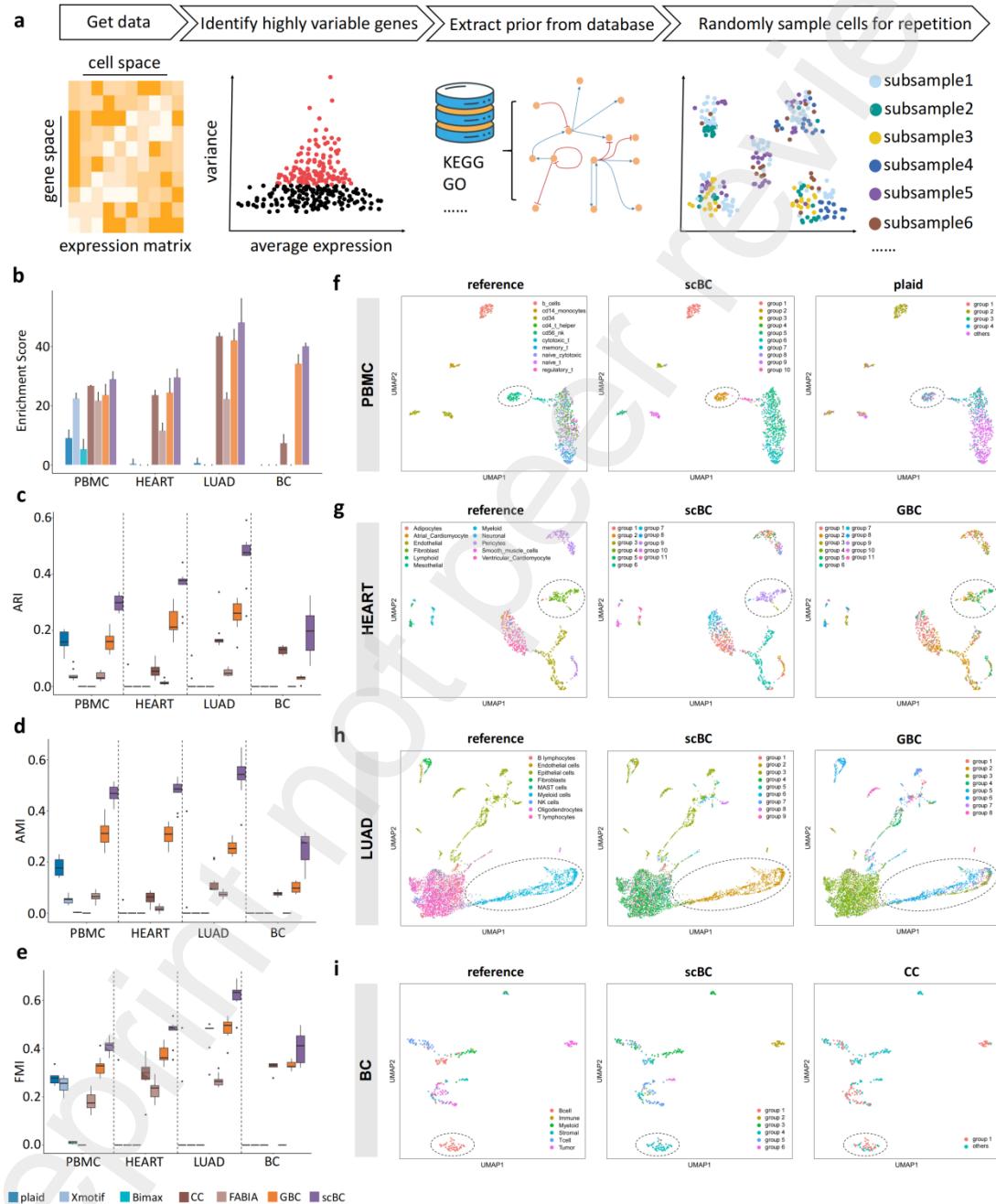
126 Results

127 scBC outperforms other methods on FGM detecting

128 To investigate whether scBC can detect biologically meaningful functional gene
129 modules, we analyzed four highly heterogeneous single-cell datasets obtained from
130 different parts and tissues of the human body under different pathological conditions
131 (purified Peripheral Blood Mononuclear Cells dataset, PBMC; cardiac cells with
132 annotation from Heart Cell Atlas, HEART; Single cell RNA sequencing of lung
133 adenocarcinoma, LUAD and single cell RNA sequencing of primary breast cancer, BC).
134 Figure 2a illustrates the preprocessing procedure, and the STAR Methods section
135 provides further details. To verify the feasibility of our method, we compared scBC
136 with six state-of-the-art biclustering algorithms, namely CC[21], xMotifs[22],
137 FABIA[13], Bimax[23], plaid[24] and GBC[20] to see how scBC outperforms alternative
138 methods. For methods that need to set the number of biclusters in advance (eg.
139 Xmotif, CC, Bimax, FABIA, GBC and scBC), we just set the max number of biclusters as
140 the number of cell types (based on cell lable). We then conducted GO enrichment
141 analysis for each FGM detected by each method (see STAR Methods), and used -
142 log10(p) (BH adjusted) as the enrich score. Methods that failed to detect any bicluster
143 were assigned a score of zero (Fig. 2b-d).

144 We found that the FGMs detected by scBC were consistently more significant than
145 those detected by other algorithms, even in highly heterogeneous settings (Fig. 2b,
146 Table S1). In the PBMC dataset, all methods were able to capture the specific FGM,
147 indicating a relatively simple data structure. Among these methods, scBC performed
148 the best, followed by CC. GBC, Xmotif, and FABIA also performed well, but plaid and
149 Bimax had unsatisfactory results. In the HEART dataset, Xmotif and Bimax failed to
150 detect any biologically meaningful gene modules, and plaid was almost invalid (Fig.
151 2b), suggesting that biclustering on cardiac tissue data may be more challenging. GBC
152 and CC had passable performance, second only to scBC. Similar results were observed
153 in the LUAD dataset, which is related to the tumor-associated immune
154 microenvironment. In the BC dataset, where sample size was minimal (only up to three
155

156 hundred cells after negative sampling) and tumor cells were mixed with normal cells,
 157 most methods failed to detect meaningful functional gene modules (Fig. 2b). Even
 158 under such challenging conditions, scBC was able to identify more biologically
 159 meaningful functional gene modules. Despite both CC and GBC had similar excellent
 160 performance in the first few conditions, CC was far inferior to GBC in BC dataset.
 161 Overall, scBC demonstrated robust superior performance in FGM detection across
 162 highly heterogeneous conditions.



163
 164 **Fig. 2 | scBC outperforms other methods on FGM detecting and cell clustering in 4 real-world**
 165 **datasets. a, The preprocess procedure of the four datasets. In each dataset, highly variable genes**
 166 **are filtered out and used to extract prior co-expression information in database like GO or KEGG.**
 167 **Cells are sampled along with the highly variable genes to generate 10 subsampled datasets for**

168 repetition. **b**, Enrichment score of different methods in four highly heterogenous datasets. The x-
169 axis represents different datasets and y-axis represents the enrichment score (-log10(p), BH
170 adjusted) of different methods. Error bar stands for standard deviation of 10 subsample's results.
171 **c-e**, Benchmarking clustering results at cell level with different criteria in the four datasets. **f-i**, Cell
172 representation of UMAP dimentionality reduction. In addition to the reference labels, shown here
173 are the methods with highest ARI (scBC) and second-highest ARI in the last subsample dataset. The
174 whole comparison of cell clustering can be found in Fig. S2-5. Some methods output too many
175 categories so that we merge some into "others" whose number of cells is less than one percent of
176 the total sample. Highlighted with black dotted lines are cell populations that are correctly
177 identified by scBC but not identified by the second-best method.

178

179 **Benchmarking clustering results at cell level**

180 Intuitively, cell groups identified by biclustering are more functionally related since
181 each group corresponds to a similar functional gene module. Functionally related cells
182 are also naturally more likely to belong to the same cell type since they have similar
183 functions, although sometimes it is not necessarily correct. In this study, we
184 investigated the clustering performance at the cell level to determine if scBC provides
185 more meaningful clustering results, even when focusing solely on the cell-level
186 clustering results. As mentioned earlier, biclustering results may have overlap
187 between each bicluster, which can result in single cells belonging to different groups.
188 However, the results from scBC and GBC enable us to assign each cell to its most
189 involved groups. For the remaining methods with less well-defined cell-level clustering
190 results, we used the Markov clustering algorithm (MCL)[25] to transform the
191 biclustering results, fully utilizing the information from the biclustering results (see
192 STAR Methods). We used the Adjusted Rand Index (ARI)[26, 27], the Fowlkes Mallows
193 score (FMI)[28], and the Adjusted Mutual Information (AMI)[29] as recommended
194 metrics to quantify the agreement between clusters (see STAR Methods). Their values
195 range from -1 to 1, with higher values indicating better performance. We evaluated
196 the clustering performance at the cell level using the four real-world datasets.

197

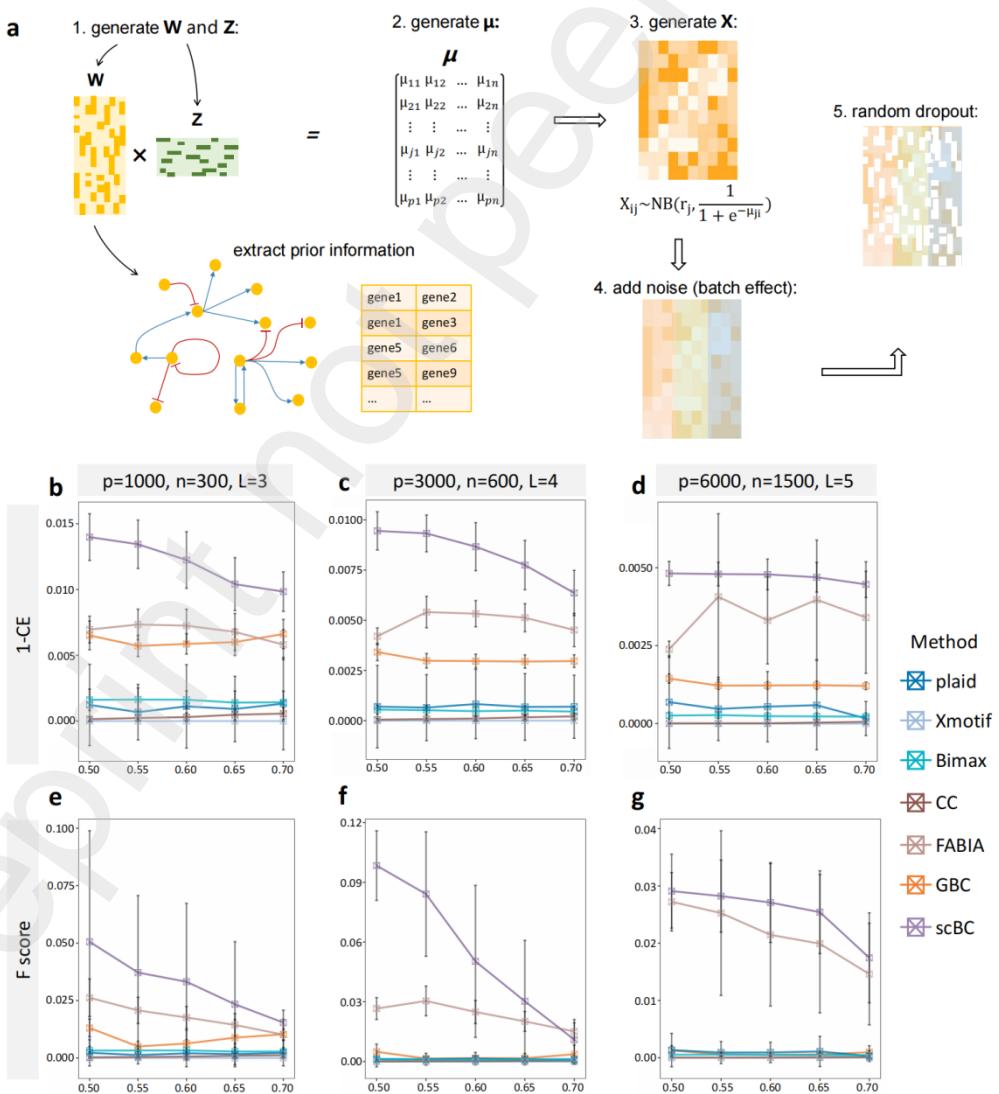
198 The cell clustering results obtained by scBC are consistently more precise than those
199 of other methods across different heterogeneous conditions (Fig. 2c-e, Table S2-S4).
200 We found that Bimax performed the worst not only in FGM detection, but also cell
201 clustering in all datasets. Although CC performed well in FGM detection in the PBMC
202 dataset, it was unable to perform cell clustering tasks simultaneously. Plaid had the
203 second-best performance on ARI, but was not as good as GBC on AMI and FMI (Fig.
204 2c-e). In PBMC dataset, multiple T cell subtypes are mixed together and difficult to
205 distinguish, while several other cell types are highly differentiated. This pattern was
206 successfully captured by scBC but not by CC (Fig. 2f). In the HEART dataset, scBC still
207 outperformed other methods, with GBC coming in second. Plaid, Xmotif, and Bimax
208 were invalid in this dataset (Fig. 2c-e). Similar results were found in the LUAD dataset.
209 We observed that scBC captures certain features of these cell populations that GBC
210 does not (Fig. 2h). In the BC dataset, CC ranks second only to scBC on ARI and FMI but
211 is inferior to GBC on AMI (Fig. 2c-e). Plaid, Xmotif, Bimax, and FABIA fail to cluster cells

212 into functionally related groups as they cannot detect FGMs in this dataset. These
 213 results demonstrate that scBC can capture the complex patterns involved in clustering
 214 functional cell groups and is more robust and precise than other methods across
 215 heterogeneous datasets.

216

217 scBC performs best on a bicluster level

218 Gene co-expression patterns differ across different cell types. These complex gene-
 219 cell correlations are of particular interest to us. When we compare the performance
 220 of different biclustering methods at the bicluster level, we pay more attention to
 221 whether the method can detect cell subgroups with similar FGMs and present these
 222 cells and FGMs at the same time. Once such a biosignal is found, we can make
 223 guidelines for downstream analysis. In this study, we introduced two evaluation
 224 methods, 1-CE and F score, to compare the performance of our scBC with six other
 225 state-of-the-art biclustering methods (see STAR Methods for simulation detail). We
 226 used simulated datasets with different dropouts under varying sc ales to elucidate
 227 how the performance of these methods varies along with the conditions (Fig. 3, Table
 228 S5-10).



229

230 **Fig. 3 | scBC performs best on a bicluster level.** **a**, Data simulation process. The parameter μ is
231 computed by the multiplicative model $\mu = \mathbf{W}\mathbf{Z}$. The prior edge information is generated along with \mathbf{W} .
232 When generating \mathbf{X} , each element is generated from $\text{NB}\left(r_j, \frac{1}{1 + e^{-\mu_{ji}}}\right)$. To simulate different batches, we
233 divided the dataset into three parts, with different intensities of noise. The implementation of dropout
234 is to perform Bernoulli censoring. See STAR Methods for detail. **b-g**, Performance of different methods
235 under different conditions. For each plot, x axis represents the dropout rate and y axis represents the
236 quantified performance (1-CE or F score). We ran 100 independent simulations for each setting, the
237 data points represent mean value and the error bars represent the standard deviation calculated across
238 repeated simulations.
239

240 Since the single-cell data itself is highly sparse, our simulation started directly from
241 dropout=0.5 and explored with a step size of 0.05. It can be observed that scBC was
242 significantly ahead of other methods across different scales with respect to 1-CE in
243 most cases, with FABIA coming in second (Fig. 3b-c), indicating that FABIA may be
244 more concerned about the overall effect of biclustering. However, when the scale of
245 datasets became larger, the performance of FABIA became extremely unstable (Fig.
246 3d). Methods like CC and Xmotif were invalid mostly (Fig. 3b-d). As expected, the
247 ability of scBC to detect biclusters decreases as dropout rate and data scale increase
248 (Fig. 3b-d). The F score of scBC was quite unstable when the data scale was too small,
249 although it still had the best performance in terms of mean performance (Fig. 3e).
250 When the dropout rate was extremely high, the performance of scBC may be
251 exceeded by FABIA with respect to F score (Fig. 3f). The performance with respect to
252 F score of FABIA also became unstable when the data scale became large, as was
253 observed in 1-CE. Overall, the results at the bicluster level indicate that scBC is more
254 reliable in uncovering complex gene-cell correlation patterns than other methods.
255

256 **scBC uncover the pathway perturbation in AD progression**

257 Neuropsychiatric disorders involve complex polygenic determinants as well as brain
258 alterations[30]. Biclustering methods can reveal cell population-specific gene co-
259 expression patterns and discover potential gene-cell connections, making them
260 inherently more suitable for the analysis and mining of complex polygenic disorders
261 such as neurodegenerative diseases. At the same time, single-cell-level resolution is
262 critical for neurodegenerative diseases like Alzheimer's disease because changes in
263 gene expression are related to specific cell types[3]. Therefore, scBC is more reliable
264 for analyzing the single-cell data of diseases with complex traits due to its excellent
265 performance.

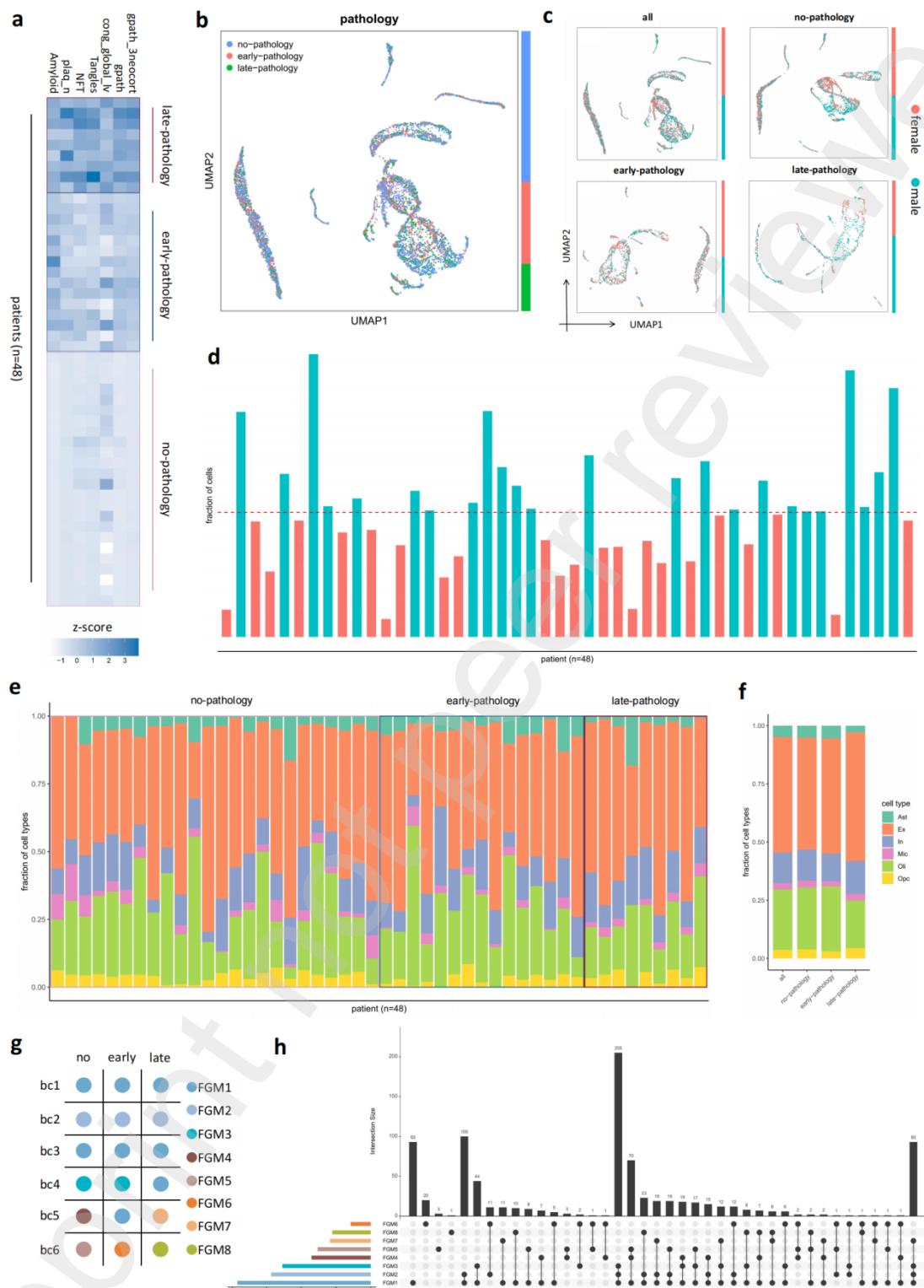
266 Alzheimer's disease (AD) is a neurodegenerative disorder associated with aging,
267 characterized by the accumulation of amyloid plaques and neurofibrillary tangles in
268 the brain parenchyma. Recent research, using a single-nucleus RNA sequencing
269 (snRNA-seq) dataset from Alzheimer's disease patients, has shown that AD is a
270 complex disease involving multiple brain cell types, as evidenced by marker gene
271 expression[3]. In this study, we aim to investigate further transcriptomic perturbations
272

273 during AD progression using network gene biomarkers identified by our scBC model.
274 This dataset includes 48 post-mortem human brain samples, with or without
275 Alzheimer's disease. The pathology groups are defined based on several pathological
276 traits (Table S11): 'no-pathology' (no amyloid burden, no neurofibrillary tangles, and
277 no cognitive impairment), 'early-pathology' (amyloid burden, but modest
278 neurofibrillary tangles and modest cognitive impairment), and 'late-pathology' (higher
279 amyloid burden, increased neurofibrillary tangles, global pathology, and cognitive
280 impairment) (Fig. 4a). After subsampling, we ensured that cells from different donors
281 were well-blended and not dominated by any one donor or biased by sex (Fig. 4b-d).
282 We also ensured that cells of the same type across individuals were consistent (Fig.
283 4e-f). To make sure that the results of multiple biclustering analyses corresponded
284 with each other, we matched biclusters in different pathological progression stages
285 and merged some FGM according to the degree of overlap in gene sets (Fig. 4g, STAR
286 Methods). We found that the overlap of each FGM, as expected, is considerable (Fig.
287 4h).

288
289 It is commonly believed that multiple FGMs can be simultaneously active in one cell
290 type, and a single FGM can be shared across different cell types, but the composition
291 percentage in different cells will vary. Our method, single-cell biclustering (scBC),
292 captured this structure perfectly (Fig. 5a-f), indicating that it is very suitable for such
293 analysis. In this study, we focused on the perturbation of FGMs for each cell type
294 during the progression of AD to gain a better understanding of the mechanisms
295 underlying AD and to provide potential recommendations for therapy. To clarify the
296 functional changes represented by specifically altered FGMs, we performed
297 enrichment analysis of specific gene sets before and after a progression stage (see
298 STAR Methods for details) to identify associated pathways that are disrupted during
299 the progression (Fig. 5g-k). The complete enrichment analysis results can be found in
300 Table S12.

301
302 For astrocytes and oligodendrocyte precursor cells, FGM perturbation occurs almost
303 only at the early pathology (Fig. 5a, f), indicating transcriptional patterns has largely
304 changed in these two cell types before individual develop severe pathological features,
305 which is consistent with previous research[3]. Inhibitory neurons' change in FGM
306 throughout the disease progression is minimal (Fig. 5c), indicating that this cell type
307 does not have many alterations in transcriptional patterns during AD progression.

308
309 Astrocytes are involved in neuronal trophic support, extracellular ion homeostasis and
310 brain fluid balance[31]. Energy metabolism is largely altered in AD astrocyte (Fig. 5g),
311 indicating the inflammatory state of the brain following injury and neurodegeneration
312 since astrocytes are a central driver of energy homeostasis in the brain, which is also
313 mentioned in previous studies[31, 32]. Consistent with previous studies, we found ion
314 transporters are dysregulated in AD astrocytes (Fig. 5g). At the same time, we also
315 found pathways related to myelination and neuron ensheathment are altered with
316 the progression of AD (Fig. 5g).



319 **Fig. 4 | overview of the subsampled AD dataset.** **a**, Clinic-pathological variables (columns) of 48
320 individuals (rows). Amyloid, overall amyloid level; plaq_n, neuritic plaque burden; NFT, neurofibrillary
321 tangle burden; Tangles, neuronal neurofibrillary tangle density; cogn_global_lv, global cognitive
322 function (last valid score). Since the lower the value, the more serious the disease, here we use its
323 opposite number, in order to be consistent with other indicators, so as to more intuitively show the

324 differences between different pathology groups. gpath, global AD-pathology burden; gpath_3neocort,
325 global measure of neocortical pathology; **b**, UMAP visualization of all cells (n=7063) indicates cells from
326 different donors of different pathological states are well blended. Color bar at the right side
327 represents the fraction composition of cells under different pathology. **c**, UMAP visualization of cells
328 from all sample, no-pathology group, early-pathology group and late-pathology group, but colored by
329 sex. Color bar at the right side represents the fraction composition of cells of different sex. **d**, The
330 proportion of cells provided across individuals (columns). Bars represent the fraction of cells
331 corresponding to each individual. Bar color indicates whether the corresponding value exceeds
332 (blue-green) or does not exceed (rose red) the average value measured across all the donors in the
333 row. Red dashed line indicates the average. **e**, Fraction of cells of each type isolated from each
334 individual (columns; n = 48). **f**, Fraction of cells of each type isolated across all (n = 48), no-
335 pathology (n = 24), early-pathology (n = 15) and late-pathology (n = 9) individuals. **g**, Merge result
336 between different biclusters. Gene sets from different biclusters in different pathology groups
337 labeled with the same color are combined as a new FGM. **h**, The overlap of FGMs. The FGM marked
338 with a solid black dot below the bar graph indicates that it is included in the comparison, and the FGM
339 marked with a black transparent dot indicates that it is not included. For example, the first bar chart
340 indicates the number of genes appearing in FGM1 but not appearing in any other FGMs is 93. This result
341 shows the overlap between different FGMs is considerable.

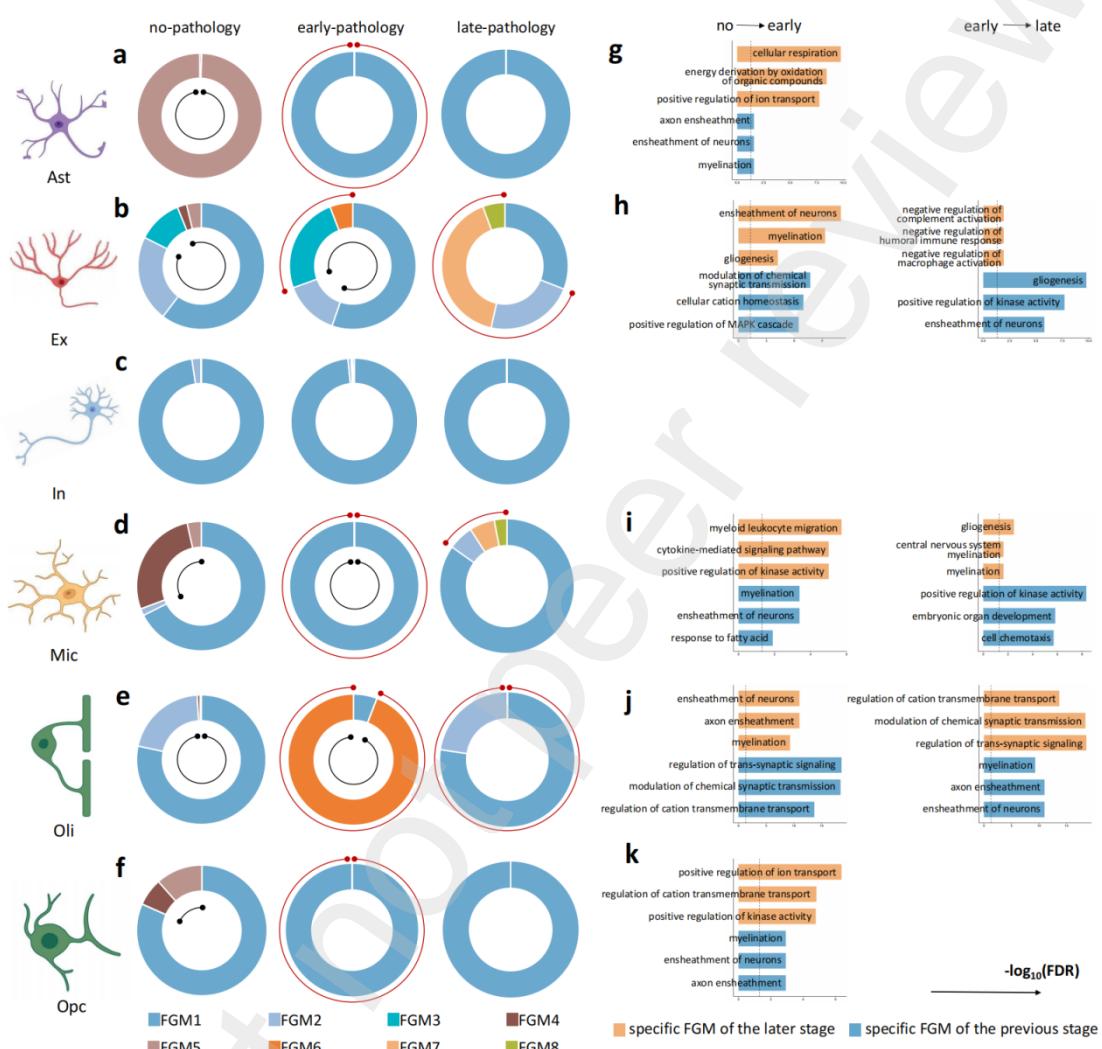
342

343 It has been found that gliogenesis and neuron ensheathment-related pathways are
344 largely impaired in AD progression[33]. We found the same conclusion in the
345 progression of AD pathology in excitatory neurons (Fig. 5h). However, the
346 perturbation of FGM composition in excitatory neurons from normal to early
347 pathology is very subtle compared to the change from the early stage to the late stage
348 (Fig. 5b), indicating that the rate at which the cells become abnormal may be slow at
349 first and then fast, at least from the point of view of gene module. Another continuous
350 change in excitatory neurons in AD progression is a general dysregulation in kinase
351 activity (Fig. 5h), which is closely related to neuronal DNA damage, well known to
352 occur in AD neurons[34]. Previous studies mentioned that the immune response is
353 also affected in the progression of AD[31]. Here we found the only gene specific in
354 latter one in the change of AD from early to late pathology is VSIG4 (see Table S12),
355 the encoded protein by which may be a negative regulator of T-cell responses and
356 largely have to do with impaired immune response (Fig. 5h). We also found that
357 cellular cation homeostasis pathway and synapse function are altered in the
358 progression from normal to early pathology (Fig. 5h).

359

360 Similar to excitatory neurons, pathways associated with kinase activity are also
361 continuously altered throughout AD progression in microglia (Fig. 5i). Cytokine-
362 mediated signaling pathway is altered in early pathology (Fig. 5i), which may be related
363 to changes in the immune response in AD progression and is also in accordance with
364 previous research[31]. Pathways related to gliogenesis and myelination also altered
365 throughout the disease progression in microglia (Fig. 5i), which is similar to astrocytes
366 and excitatory neurons. Cell migration-related pathways are dysregulated in the early
367 AD microglia (Fig. 5i), which is also consistent with several studies[3, 5, 35, 36] and

368 largely related to microglial plaque clustering phenotypes, a phenomenon of
 369 inappropriate interactions with amyloid. Response to fatty acid becomes odd in early
 370 AD microglia (Fig. 5i), which is also an indicator of lipid metabolism dysfunction. We
 371 also find cell chemotaxis becomes abnormal in late AD microglia (Fig. 5i), indicating an
 372 inflammatory state in AD microglia.
 373



374
 375 **Fig. 5 | scBC uncover the pathway perturbation in AD progression.** a-f, perturbation of FGM
 376 composition in each cell type during AD progression. Each pie chart quantifies the FGM composition of
 377 a cell under a specific progression condition. The outer red circles indicate FGMs whose composition is
 378 increasing compared to the previous stage; The inner black circle represents FGMs whose composition is
 379 decreasing compared to the latter stage. g-k, results of enrichment analysis of FGMs altered in two
 380 phases. Orange represents the specific FGM of the later stage, and blue represents the specific FGM of
 381 the previous stage, both representing the set of genes that are perturbed during the progression.

382
 383 We observed that in oligodendrocytes, the main changes in functional gene modules
 384 (FGMs) during AD progression occurred in myelination-related and synaptic signalling-
 385 related pathways (Fig. 5j). Since memory preservation is thought to require new
 386 myelin formation, the impaired capacity of oligodendrocytes to adaptively monitor

387 neural activity and facilitate myelin remodeling may govern cognitive decline in AD[37].
388 Moreover, synaptic signaling and axon development are critical for the transmission
389 of excitement in the nervous system, and dysregulation of these processes can result
390 in slower propagation of neural excitation. The changes in FGMs in oligodendrocytes
391 are directly related to the reduction of nervous system excitability. Previous research
392 also suggests that changes in oligodendrocytes may affect the function of other cells
393 in the central nervous system[38-41]. Thus, targeting oligodendrocytes may be a
394 promising strategy for the treatment of AD and other neurological disorders.

395

396 Oligodendrocyte precursor cells (OPCs), which are distributed throughout grey and
397 white matter, are thought to dynamically sense and modulate neural activity[40] like
398 oligodendrocytes do. So not surprisingly, pathways related to myelination and the
399 ensheathment of neurons become abnormal along with the progression of AD (Fig.
400 5k). Pathways related to ion transportation are also dysregulated, providing supports
401 for previous findings that genes related to ion channels are dysregulated in AD OPCs[3,
402 5]. Besides, pathways related to kinase activity and cellular cation homeostasis are
403 altered at the early stage in AD progression (Fig. 5k).

404

405 Except for inhibitory neurons, which did not change significantly throughout disease
406 progression, several other cell types exhibited specific functional gene module (FGM)
407 perturbations, highlighting the importance of single-cell analysis. Notably, we
408 observed that pathways related to myelination and gliogenesis were more or less
409 altered across all these cell types, indicating similar alterations among AD-associated
410 cells and suggesting that AD progression is largely related to dysregulation of this
411 pathway, which is further confirmed in a recent study[42].

412

413 Discussion

414 Molecular biomarkers have been widely used in clinical practice to identify diseases,
415 but they often suffer from low coverage and high false-positive or false-negative rates,
416 limiting their further application[43]. Network biomarkers, also known as module
417 biomarkers, have attracted attention as a more robust form of biomarker than
418 individual molecules for characterizing diseases[44, 45]. This is particularly important
419 for analyzing single-cell data, which is inherently more complex than bulk tissue data
420 due to the heterogeneity of individual cells within a sample. However, network
421 biomarkers are usually cell-specific and may change during the disease progression.
422 To detect cell-specific network biomarkers, we developed scBC, a single-cell Bayesian
423 biclustering method that combines variational autoencoder (VAE) for batch removal
424 and data imputation with matrix factorization-based Bayesian biclustering for utilizing
425 known biological information. Our method outperforms other state-of-the-art
426 methods in finding functional gene modules (FGMs), discovering functionally-related
427 cell groups, and detecting cell-gene correlation patterns in highly heterogeneous
428 single-cell RNA sequencing (scRNA-seq) datasets and simulated data. This makes it
429 well-suited for analyzing diseases with multifactorial etiologies whose functionally-

430 related potential cell groups are finely divided, cell-type conditioned gene co-
431 expression patterns are complicated, and cell-gene correlation changes throughout
432 the disease progression.

433
434 In this study, we applied scBC to a snRNA AD dataset to explore how the
435 transcriptional functional modules of each cell type change as the disease progresses.
436 Our results further confirmed the complex interplay of virtually every major brain cell
437 type in Alzheimer's disease[3, 33]. We found that FGM composition largely changed
438 in astrocytes and oligodendrocyte precursor cells before individuals developed severe
439 pathological features. However, inhibitory neurons showed minimal changes in FGM
440 throughout the disease progression, indicating that this cell type does not have many
441 alterations in transcriptional patterns during AD progression. A consistent FGM
442 perturbation across all other cell types, except inhibitory neurons, was the alteration
443 in pathways related to myelination and gliogenesis, suggesting that this pathway may
444 play a decisive role in the progression of AD.

445
446 Specific to each cell type, energy metabolism and ion transporters are dysregulated in
447 AD astrocyte, indicating the inflammatory state of the brain following injury and
448 neurodegeneration. The perturbation of FGM composition in excitatory neurons from
449 normal to early pathology is very subtle compared to the change from the early stage
450 to the late stage, indicating that the rate at which the cells become abnormal may be
451 slow at first and then fast. Another continuous change in excitatory neurons in AD
452 progression is the general dysregulation in kinase activity, which is closely related to
453 neuronal DNA damage. Besides, immune response, cellular cation homeostasis and
454 synapse function are also altered in AD excitatory neurons. Microglia shares a similar
455 alteration in kinase activity with excitatory neurons throughout AD progression.
456 Pathways like immune response-related cytokine-mediated signaling, amyloid
457 interaction-related cell migration and lipid metabolism-related fatty acid response are
458 dysregulated in the early AD microglia. Cell chemotaxis becomes abnormal in late AD
459 microglia, indicating an inflammatory state in AD microglia. Oligodendrocyte is a cell
460 that needs to be more focused on for disease treatment since FGM perturbations in
461 such a cell type are mainly concentrated in myelination-related and synaptic signaling-
462 related pathways, directly related to the reduction of the excitability of the nervous
463 system. At last, pathways related to ion transportation, kinase activity and cellular
464 cation homeostasis are dysregulated in oligodendrocyte precursor cells.

465
466 In the context of high-throughput sequencing data, network biomarker-based
467 analytical methods preserve the complex co-expression or co-regulation patterns in
468 the gene module and is more robust to the analysis of complex diseases. We believe
469 that scBC, as a novel technique for cell-specific network biomarkers detection, creates
470 an unprecedented opportunity towards effectively delineating mechanisms of
471 complex diseases at single-cell resolution, providing advice on the treatment of such
472 diseases. However, although the network biomarker may contain complex co-
473 expression or co-regulation patterns, its internal precise and quantitative regulatory

474 relationship has not been clarified. Future research can focus on the explanation of
475 the regulatory relationship within the cell-specific network structure, so as to have a
476 more accurate inference on the principle of FGM perturbations during disease
477 progression.

478

479 Conclusion

480 Our method, scBC, is a powerful tool for detecting cell-specific network biomarkers in
481 highly heterogeneous single-cell RNA sequencing datasets. It outperforms other state-
482 of-the-art methods in finding functional gene modules (FGMs) and discovering
483 functionally-related cell groups, making it a top contender for analyzing diseases with
484 multifactorial etiologies, where cell-type conditioned gene co-expression patterns are
485 complicated and cell-gene correlation changes throughout disease progression. To
486 demonstrate its potential, we applied scBC to a snRNA AD dataset, revealing FGM
487 perturbations in each cell type as the disease progresses. Our findings support
488 previous studies and offer new insights, such as minimal changes in FGM for inhibitory
489 neurons and a consistent FGM perturbation in pathways related to myelination and
490 gliogenesis across all other cell types. These results suggest that scBC provides
491 unprecedented opportunities for effectively delineating the mechanisms of complex
492 diseases at single-cell resolution on a gene-module-based view, providing valuable
493 insights for the treatment of such diseases.

494

495 STAR Methods

496 Details for Data Reconstruction Using Variational Inference

497 Taking advantage of recent work by Romain et.al[19], here we also adopt the idea of
498 using variational inference to estimate the posterior distribution for the low-
499 dimensional, latent variables z_n for each cell n which should reflect biological
500 differences among cells. To remove the nuisance variation due to technique factors
501 such as batch effects, it's reasonable to model the sampling distribution conditioned
502 on the batch annotations s_n [46, 47]. That is, the observed expression x_n^g of each
503 gene g in each cell n is drawn from $p(x_n^g|z_n, s_n)$. Considering the substantiate
504 amount of dropouts in the scRNA-seq data, here we choose zero-inflated negative
505 binomial(ZINB) distribution as the proper sampling distribution which is
506 recommended by Grün et.al[47].

507

508 Now we can present in-detail steps of our pre-defined generating process of how we
509 can get the gene expression of each cell. The generating process can be concluded as
510 follows:

511

$$512 z_n \sim N(0, I)$$

$$\begin{aligned}
513 \quad \rho_n &\sim f_{expect}(z_n, s_n) \\
514 \quad w_n^g &\sim \text{Gamma}(\rho_n^g, \theta^g) \\
515 \quad y_n^g &\sim \text{Poisson}(l_n w_n^g) \\
516 \quad h_n^g &\sim \text{Bernoulli}(f_{drop}(z_n, s_n))
\end{aligned}$$

518 Where z_n is the low-dimensional, latent variable. Here we use a standard
519 multivariate normal prior for z because it can be reparametrized in a differentiable
520 way into any arbitrary multivariate Gaussian random variable, which is extremely
521 helpful in the inference process. f_{expect} is a neuron network which maps the latent
522 space and batch annotations of each cell back to the full dimension of the gene
523 expression: $\mathbb{R}^{d+1} \rightarrow \mathbb{R}^G$. At the generating stage, f_{expect} is constrained by a softmax
524 activation function at the last layer so that each element of ρ_n sum up to 1 when we
525 decode the gene expression data. So it is reasonable to deem that ρ_n denotes the
526 mean proportion of transcripts expressed across all genes[19]. This also makes it easy
527 to reconstruct the expression data at an arbitrary scale only to specify the cell-specific
528 library size l_n . $f_{drop}(z_n, s_n)$ is also a neuron network which maps the latent space
529 and batch annotations of each cell to their respective dropout probabilities. w_n^g and
530 y_n^g are two intermediate variable and it can be shown that through this process, h_n^g
531 is a r.v. following ZINB distribution[47] with mean $l_n \rho_n^g$, gene-specific dispersion θ^g
532 and zero-inflation probability $f_{drop}(z_n, s_n)$.

533
534 When we conduct data reconstruction to get the batch-removal, imputed gene
535 expression data, we only take advantage of the intermediate variable ρ_n and scale it
536 to our expected library size. That is, multiplying it by a given parameter l_n , which we
537 just use the empirical library size (total number of transcripts per cell) of each cell
538 throughout our experiments. But one should notice we can re-scale it to any expected
539 library size if additional information is given.

540 **Model Training at Learning Stage**

541 Here we introduce a recognition model $q_\varphi(z_n|x_n, s_n)$: an approximation to the
542 intractable true posterior $p_\theta(z_n|x_n, s_n)$. The marginal likelihood can be written as:
543
$$\log p_\theta(x_n|s_n) = D_{KL}(q_\varphi(z_n|x_n, s_n)||p_\theta(z_n|x_n, s_n)) + L(\boldsymbol{\theta}, \boldsymbol{\varphi}; x_n)$$

545 Where $L(\theta, \varphi; x_n) = E_{q_\varphi(z_n|x_n, s_n)} [-\log q_\varphi(z_n|x_n, s_n) + \log p_\theta(z_n|x_n, s_n)]$. Since the KL-
546 divergence is always non-negative. We have:

$$\log p_\theta(x_n|s_n) \geq E_{q_\varphi(z_n|x_n, s_n)} [-\log q_\varphi(z_n|x_n, s_n) + \log p_\theta(z_n|x_n, s_n)]$$

548 The evidence lower bound (ELBO) $L(\theta, \varphi; x_n)$ can also be written as:

$$L(\theta, \varphi; x_n) = E_{q_\varphi(z_n|x_n, s_n)} [\log p_\theta(x_n|z_n, s_n)] - D_{KL}(q_\varphi(z_n|x_n, s_n)||p_\theta(z_n|s_n))$$

549 Optimizing the ELBO means optimizing both the variational parameters φ and
550 generative parameters θ at the same time. Assuming the true latent variable z_n is

551 batch-free (independent with batch annotation s_n) and the prior follows standard
552 multivariate Gaussian distribution, we can get the closed-form expression of the
553 derivative of $D_{KL}(q_\varphi(z_n|x_n, s_n)||p_\theta(z_n|s_n))$. To get the low-variance Monte Carlo
554 estimation of the gradient of term $E_{q_\varphi(z_n|x_n, s_n)} [\log p_\theta(x_n|z_n, s_n)]$, we use the
555 reparameterization trick in the learning stage[48]:

$$557 \tilde{E}_{q_\varphi(z_n|x_n, s_n)} [\log p_\theta(x_n|z_n, s_n)] \cong \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_n|g_\varphi(\epsilon^l, x_n), s_n)$$

558 Where $g_\varphi(\epsilon^l, x_n)$ is a differentiable transformation to reparameterize the random
559 variable $z_n \sim q_\varphi(z_n|x_n, s_n)$ and $\epsilon \sim p(\epsilon)$ is an auxiliary noise variable.

560 For a single data point $x_n^{(i)}$ (cell i) we have:

$$561 \tilde{L}(\theta, \varphi; x_n^{(i)}) = \frac{1}{L} \sum_{l=1}^L \log p_\theta(x_n^{(i)}|z^{(i,l)}, s_n) - D_{KL}(q_\varphi(z_n|x_n^{(i)}, s_n)||p_\theta(z_n))$$

562 Where $z^{(i,l)} = g_\varphi(\epsilon^{(l)}, x_n^{(i)})$ and $\epsilon^{(l)} \sim p(\epsilon)$.

563 At learning stage, we use mini-batch stochastic optimization to optimize the ELBO,
564 suppose our dataset contains N cells and the size of each mini-batch is M, we can get
565 the estimator of marginal likelihood lower bound of the stochastic mini-batch:

$$566 L(\theta, \varphi; x_n^{(i)}) \cong L^M(\theta, \varphi; x_n^{(i)}) = \frac{1}{M} \sum_{i=1}^M \tilde{L}(\theta, \varphi; x_n^{(i)})$$

567 When M is large enough, Diederik et al[48] found that the number of samples L per
568 datapoint can even be set to 1, hence decrease the time consumption when conduct
569 expectation estimation for $E_{q_\varphi(z_n|x_n, s_n)} [\log p_\theta(x_n|z_n, s_n)]$. Throughout our experiment we
570 set M=128 data points to guarantee the large-sample requirement. We use Adam
571 optimizer with learning rate=0.01. We also use deterministic warm-up and batch
572 normalization during learning to learn an expressive model which is recommended by
573 Sonderby et al[49].

574

575 Bayesian Biclustering Incorporate biological information

576 After reconstructing the original expression matrix, we can conduct biclustering
577 procedure to detect condition-specific FGMs and identify cell subpopulations with
578 distinct functions. Relevant studies have shown that if we can introduce existing
579 biological information (such as the metabolic pathways from the KEGG database) into
580 the process of biclustering, then the accuracy of the biclustering results will be

581 improved[20, 50-54]. Therefore, we adopt a Bayesian analysis framework, which can
 582 introduce prior information to guide variable selection.

583
 584 Suppose our reconstructed data matrix is \mathbf{X} of size $p \times n$, where p represents the
 585 number of genes and n is the number of cells. In order to reduce randomness, here
 586 we do not directly decompose the data matrix \mathbf{X} , but decompose its parameter matrix.
 587 We denote the parameter matrix of \mathbf{X} as $\boldsymbol{\mu}$ (eg. mean) and decompose it: $\boldsymbol{\mu} = \mathbf{m}\mathbf{1}^T + \mathbf{W}\mathbf{Z}$, where
 588 \mathbf{m} is a $p \times 1$ bias vector, $\mathbf{1}$ is a $n \times 1$ vector of 1, \mathbf{W} is a $p \times L$
 589 matrix, \mathbf{Z} is a $L \times n$ matrix. Since the observation of gene j from cell i x_{ji} is
 590 generated independently, the likelihood function of \mathbf{X} is the product of the likelihood
 591 functions of each independent observation. Since we have assumed the generating
 592 process of the reconstructed data, we set x_j to be a random variable that follows
 593 Gaussian distribution with a likelihood function π_j in the discussion following on.

594

595 The likelihood function of an individual observation is:

$$596 \quad \pi_j(x_{ji}|\mu_{ji},\rho_j) = \frac{\rho_j^{1/2}}{\sqrt{2\pi}} e^{-\rho_j(x_{ji}-\mu_{ji})^2/2}, x_{ji} = 0, 1, \dots \quad (1)$$

597 Now we discuss how to introduce prior information. To obtain a sparse estimate of \mathbf{W} ,
 598 first use the Laplace prior on the matrix \mathbf{W} :

$$599 \quad \log \pi(\mathbf{W} | \boldsymbol{\lambda}) = C + \sum_{j,l} \log \lambda_{jl} - \sum_{j,l} \lambda_{jl} |w_{jl}|$$

600 Here the prior parameter λ controls the degree of shrinkage of w , so the prior of the
 601 multivariate normal is applied on λ next:

$$602 \quad \log \pi(\boldsymbol{\alpha} | \boldsymbol{\Omega}) = C_{v_2} + \frac{L}{2} \log |\boldsymbol{\Omega}| - \frac{1}{2v_2} \sum_l (\boldsymbol{\alpha}_l - v_1 \mathbf{1}) \boldsymbol{\Omega} (\boldsymbol{\alpha}_l - v_1 \mathbf{1})^T \quad (2)$$

603 Where $\alpha_{jl} = \log \lambda_{jl}$, $\boldsymbol{\alpha}_l = (\alpha_{1l}, \dots, \alpha_{pl})^T$, v_1 and v_2 are hyperparameters. Finally, we
 604 apply a suitable prior to the precision matrix $\boldsymbol{\Omega}$ to connect the correlated λ , which
 605 will be defined as:

$$606 \quad \boldsymbol{\Omega} = \begin{bmatrix} 1 + \sum_{j \neq 1} \omega_{1j} & -\omega_{12} & \cdots & -\omega_{1p} \\ -\omega_{21} & 1 + \sum_{j \neq 2} \omega_{2j} & \ddots & -\omega_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ -\omega_{p1} & -\omega_{p2} & \cdots & 1 + \sum_{j \neq p} \omega_{pj} \end{bmatrix}$$

607 This ensures the symmetry and positive definiteness of the precision matrix, and the
 608 prior of the ω is:

$$609 \quad \pi(\boldsymbol{\omega}) \propto |\boldsymbol{\Omega}|^{-\frac{L}{2}} \prod_{(j,k) \in E} \omega_{jk}^{a_\omega - 1} \exp(-b_\omega \omega_{jk}) \mathbb{1}(\omega_{jk} > 0) \prod_{(j,k) \notin E} \delta_0(\omega_{jk}) \quad (3)$$

610 $\delta_0(\cdot)$ is the Dirac function centered at 0, $\mathbb{1}(\cdot)$ is an indicative function. Therefore, if
 611 x_j and x_k are directly connected in the graph G , then (3) will try to make the

612 precision matrix components ω_{jk} to be non-zero, and make the contraction term λ_{jj}
 613 and λ_{kl} related through (2). Since w_{jl} and w_{kl} are subject to a similar degree of
 614 contraction under this condition, they tend to be both zero or non-zero at the same
 615 time. In other words, if genes j and k are directly connected in a pathway, they are
 616 encouraged to be selected together (or not selected together) in bicluster l .
 617 Therefore, a standout feature of this approach is that the selected feature set in each
 618 bicluster tends to include functional gene module rather than individual genes,
 619 resulting in more biologically meaningful results.
 620

621 Since the \mathbf{Z} matrix represents the results on the cell set, there is no special pathway
 622 information between the samples, so it is sufficient to perform Laplace sparse prior
 623 on it:
 624

$$\log \pi(\mathbf{Z} | \xi) = C + \sum_{lj} \log \xi_{li} - \sum_{lj} \xi_{li} |z_{li}|$$

625 Where ξ is the contraction factor, on which a conjugate prior is applied, i.e. Gamma
 626 prior:
 627

$$\log \pi(\xi) = C_{v_3, v_4} + (v_3 - 1) \sum_{li} \log \xi_{li} - \frac{1}{v_4} \sum_{li} \xi_{li}$$

628 This is the general framework of the Bayesian biclustering model we use.
 629

630 **MAP estimation for biclustering result**

631 In the optimization stage, we adopt the Pólya-Gamma latent variable proposed by
 632 Polson et al.[55]. We use the identity formula provided in Polson et al.[55]:
 633

$$\frac{e^{\mu_{ji}x_{ji}}}{(1 + e^{\mu_{ji}x_{ji}})^{b_{ji}}} = 2^{-b_{ji}} e^{\kappa_{ji}\mu_{ji}} \int_0^\infty e^{-\rho_{ji}\mu_{ji}^2/2} \pi_{ji}(\rho_{ji}) d\rho_{ji}$$

634 Where $\kappa_{ji} = x_{ji} - b_{ji}/2$, $\pi_{ji}(\rho_{ji})$ is of the Pólya-Gamma class $\mathcal{PG}(b_{ji}, 0)$. So (1) can be
 635 written as:
 636

$$\pi_j(\mathbf{x}_j | \boldsymbol{\mu}_j) \propto e^{-\frac{1}{2} \sum_i \rho_{ji} (\mu_{ji} - x_{ji})^2} \pi_j^*(\rho_j)$$

637 Where $\rho_j \sim \mathcal{G}\left(\frac{\zeta_j + n}{2}, \frac{\zeta_j}{2}\right)$, ζ_j is the prior parameter for variance. After the introduction of
 638 latent variable ρ , LASSO can be efficiently solved in the M step of the EM algorithm.
 639 Here we use dynamic weighted LASSO algorithm to speed up the calculation[56].
 640 Additionally, we utilize maximum a posteriori estimation (MAP) to estimate the
 641 parameters, which is defined as:
 642

$$(\hat{\mathbf{W}}, \hat{\mathbf{Z}}, \hat{\mathbf{m}}, \hat{\boldsymbol{\alpha}}, \hat{\boldsymbol{\xi}}) = \operatorname{argmax}_{\mathbf{W}, \mathbf{Z}, \mathbf{m}, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\rho}, \boldsymbol{\Omega}} \int \int \pi(\mathbf{W}, \mathbf{Z}, \mathbf{m}, \boldsymbol{\alpha}, \boldsymbol{\xi}, \boldsymbol{\rho}, \boldsymbol{\Omega} | \mathbf{X}) d\boldsymbol{\rho} d\boldsymbol{\Omega}$$

643 This can be efficiently solved using the EM algorithm, and the objective function at t
 644

646 iterations is:

$$647 \quad Q_t(\mathbf{Z}, \mathbf{W}, \mathbf{m}, \boldsymbol{\alpha}, \boldsymbol{\xi}) = -\frac{1}{2} \sum_{i,j} \rho_j^{(t)} (\mu_{ji} - x_{ji})^2 + \sum_{j,l} \alpha_{jl} - \sum_{j,l} \lambda_{jl} |w_{jl}| + v_3 \sum_{l,i} \log \xi_{l,i} - \\ 648 \quad \sum_{l,i} \xi_{l,i} \left(|z_{li}| + \frac{1}{4} \right) - \frac{1}{2v_2} \sum_l (\boldsymbol{\alpha}_l - v_1 \mathbf{1})^T \boldsymbol{\Omega}^{(t)} (\boldsymbol{\alpha}_l - v_1 \mathbf{1})$$

649 Where $\boldsymbol{\mu} = \mathbf{m}^{(t-1)} + \mathbf{W}^{(t-1)} \mathbf{Z}^{(t-1)}$, $\rho_{ij}^{(t)} = E(\rho_{ij} | \mathbf{X}, \mathbf{W}^{(t-1)}, \mathbf{Z}^{(t-1)}, \mathbf{m}^{(t-1)}, \boldsymbol{\alpha}^{(t-1)}, \boldsymbol{\xi}^{(t-1)})$ and
650 $\boldsymbol{\Omega}^{(t)} = E(\omega_{ij} | \mathbf{X}, \mathbf{W}^{(t-1)}, \mathbf{Z}^{(t-1)}, \mathbf{m}^{(t-1)}, \boldsymbol{\alpha}^{(t-1)}, \boldsymbol{\xi}^{(t-1)})$.

651

652

653 **Strong classification of cell group for different biclustering methods**

654 For scBC and GBC, we can directly observe the contribution of each bicluster to the
655 parameter matrix in each cell from the result of the Z matrix. We assign a cell to the
656 most involved cluster, which is determined by the row with the largest absolute value.
657 For the remaining methods, we aim to achieve optimal cell classification results
658 without losing any information from the biclustering results. Due to the high degree
659 of cell overlap in the biclustering results, we convert the cell-level biclustering results
660 into a graph, where cells in the same bicluster are connected by edges. If the
661 occurrences of a pair of cells increase, the weights of the edges between them also
662 increase accordingly. We then apply the Markov clustering algorithm (MCL)[25] to
663 convert the graph into cell-level clustering results. For each method, we set the
664 number of iterations to a value between 1 and 10 that allows the method to achieve
665 the highest adjusted Rand index (ARI). In fact, we observed that the number of
666 iterations required for the best results usually does not exceed 7. After the
667 transformation, each cell is exclusively assigned to a cluster, and we can evaluate the
668 cell-level clustering results using any clustering evaluation criterion.

669

670

671 **Datasets and preprocessing**

672 Here we describe all of the datasets and the preprocessing steps used in the current
673 work as follows. The prior information for all the real-world datasets is extracted by
674 *biomaRt* using the highly variable genes.

675

676 **heart**

677 This is a combined single cell and single nuclei RNA-Seq data of 485K cardiac cells with
678 annotation from [Heart Cell Atlas](#). Here we use a subsampled version provided at
679 https://github.com/YosefLab/scVI-data/blob/master/hca_subsampled_20k.h5ad, which has
680 been filtered down randomly to 20k cells. In our study, we further filtered 1000 highly
681 variable genes using *scanpy* and generate 10 subsampled datasets with each
682 containing 1000 randomly selected cells.

683

684 **PBMC**

685 This actually is a purified PBMC dataset from[57]. An organized version can be
686 accessed from <https://github.com/YosefLab/scVI-data/raw/master/PurifiedPBMCDataSet.h5ad>.
687 We also conducted a subsampling procedure here: first screen out 853 highly variable

688 genes using *scipy*, then generate 10 subsampled datasets with each containing 1000
689 random selected cells.

690

691 **LUAD**

692 Single cell RNA sequencing of lung adenocarcinoma from[58], which can be accessed
693 from the NCBI Expression Omnibus database (accession code GSE131907). This is
694 single cell RNA sequencing (scRNA-seq) for 208,506 cells derived from 58 lung
695 adenocarcinomas from 44 patients, which covers primary tumour, lymph node and
696 brain metastases, and pleural effusion in addition to normal lung tissues and lymph
697 nodes. Here we use *Seurat* to conduct preprocessing: we first randomly select 10000
698 cells to filter 2000 highly variable genes, then generate 10 subsampled datasets with
699 each containing 5000 cells.

700

701 **BC**

702 Single cell RNA sequencing of primary breast cancer from[59], which can be accessed
703 from the NCBI Expression Omnibus database (accession code GSE75688). This dataset
704 contains 515 cells from 11 patients and most of the cell type is tumor. We first screen
705 out 2000 highly variable genes using *Seurat* then conduct subsampling. Due to the
706 serious category imbalance problem in this dataset (326 cells are labelled as “Tumor”),
707 we only sample 86 cells with the tumor label each time and all cells with other labels
708 are retained so that the results will not be unreliable due to class imbalance during
709 evaluation.

710

711 **AD**

712 A total of 80660 droplet-based single-nucleus RNA-seq (snRNA-seq) profiles for
713 Alzheimer’s disease from[3]. The post-mortem human brain samples came from 48
714 participants in the Religious Order Study (ROS) or the Rush Memory and Aging Project
715 (MAP), collectively known as ROSMAP with 24 individuals with high levels of β-amyloid
716 and other pathological hallmarks of AD (‘AD-pathology’), and 24 individuals with no or
717 very low β-amyloid burden or other pathologies (‘no-pathology’). The original study
718 clustered individuals based on nine clinico-pathological traits to further define the
719 pathology groups as ‘early-pathology’ and ‘late-pathology’. And that division is totally
720 adopted in our study. The snRNA-seq data are available on The Rush Alzheimer’s
721 Disease Center (RADC) Research Resource Sharing Hub at <https://www.radc.rush.edu/docs/omics.html> (snRNA-seq PFC) or at Synapse (<https://www.synapse.org/#!Synapse:syn18485175>) under the doi 10.7303/syn18485175. The data are available under
722 controlled use conditions set by human privacy regulations. To access the data, a data
723 use agreement is needed. Since we are not going to use this data set to conduct
724 benchmarking here, there is no need to repeatedly generate subsamples. When
725 preprocessing the dataset, we first use stratified sampling to draw one out of ten cells,
726 then 2000 highly variable genes are refined by Seurat. This sample is then used to be
727 explored later.

728

729 **Simulated data**

732 In each simulation setting, we generate 100 simulated datasets. For convenience, we
733 denote p as the number of genes, n as the number of cells. The scale of the FGM
734 increases adaptively with the size of the simulated dataset (actually the size of p). The
735 parameter μ is computed by the multiplicative model $\mu = WZ$, where W is a $p \times L$ matrix
736 and Z is an $L \times n$ matrix. The number of non-zero elements in each column of W is set
737 as $p/20$, and the number of non-zero elements in each row of Z is randomly drawn
738 from a Poisson distribution with a parameter of 30. The row indices of non-zero
739 elements in W and the column indices of Z with non-zero elements are randomly
740 drawn from 1 to p and 1 to n . The nonzero element values for both W and Z are
741 generated from a normal distribution with mean 1.5 and standard deviation 0.1, and
742 are randomly assigned to be positive or negative. The prior edge is generated along
743 with W .

744 When generating X , each element is generated from $NB\left(r_j, \frac{1}{1 + e^{-\mu_{ji}}}\right)$, and the parameter

745 r_j is randomly drawn from 5 to 20. Finally, in order to simulate different batches, we
746 divided the dataset into three parts, each of $n/3$ samples, with different intensities of
747 noise. The implementation of dropout is to perform Bernoulli censoring at each data
748 point according to the given dropout rate parameter.

749 The simulation data generation process is shown in Fig. 3a.

750

751 **Matching of biclusters when analysing AD dataset**

752 To avoid confusion, we first explain the difference between biclustering and bicluster,
753 two concepts we've been using throughout the paper. A biclustering refers to execute
754 one biclustering algorithm once (eg. scBC). After biclustering is conducted, we can get
755 several columns-rows pairs, each is called a bicluster. In our study, we conduct three
756 independent biclusterings on the three pathologically seperated AD datasets, each
757 with L biclusters. L is the number of biclusters we set beforehand.

758

759 When conducting scBC on the AD dataset, genes that are widely present in all
760 biclusters represent the commonality among all cells. We subtract these genes from
761 each bicluster to reduce the homogeneity of different biclusters, since we are not
762 interested in them in this case. Here, we set the number of biclusters in each round of
763 biclustering to 6. However, this is an empirical hyperparameter, and some biclusters
764 may have a high degree of similarity and be more reasonable to merge into a single
765 bicluster. This applies not only to different biclusters in a whole biclustering but also
766 to different biclusters in independent biclusterings. However, the methods for aligning
767 biclusters in a single biclustering and for biclusters in different biclusterings should be
768 different, since biclusters from the latter are somewhat more independent.

769

770 Due to their own homogeneity or correlation, more attention should be paid to the
771 exclusivity when merging biclusters from a single biclustering. We denote the number
772 of genes only appear in biclusters i and j as e , which means all the other biclusters

773 don't have these genes. And genes present in bicluster k is denoted as g_k , the overlap
774 score is defined as:

775

$$os_{ij} = \frac{e}{\max\{|g_i|, |g_j|\}}$$

776 In our study, a pair of biclusters with overlap score > 0.03 are combined as a FGM. The
777 biclusters correspondences in different biclusterings are independent, so more
778 attention should be paid to the degree of overlap. Here we use the "overlap over
779 union"(IoU) criterion to combine different biclusters:

780

$$IoU_{ij} = \frac{\text{intersect}\{|g_i|, |g_j|\}}{\text{union}\{|g_i|, |g_j|\}}$$

781 each pair of biclusters with IoU > 0.3 are combined as a FGM. The alignment results
782 are in Table S13.

783

784 **FGM perturbation during AD progression and enrichment analysis**

785 Before merging functional gene modules (FGMs) from different biclusters, we first
786 assign each cell exclusively to a bicluster using the strong classification method as
787 before. When similar FGMs are combined, functionally related cells are also merged
788 as a whole. Next, we obtain the FGM composition contained in each cell type. As
789 observed, multiple FGMs can be simultaneously active in one cell type. To illustrate
790 FGM perturbation for each cell type during AD progression, we first observe the
791 changes in the proportion of each FGM in each cell type. FGMs with an elevated ratio
792 are candidates for "increased activity," while those with a reduced ratio are candidates
793 for "decreased activity." We then examine the differences in the gene set makeup of
794 these two types of FGMs. The overlap between the two represents commonalities
795 exhibited in certain cell types, which are not of interest to us. We focus on the
796 exclusive genes of "increased activity" and "decreased activity," which may uncover
797 pathway perturbations in different pathological states. The exclusive genes in
798 "increased activity" and "decreased activity" are used for functional enrichment using
799 *clusterProfiler*. The results are used to reveal the pathway perturbation during AD
800 progression.

801

802 **Evaluation**

803 Here we will describe the evaluation metrics used in our study as follows.

804

805 **Comparison of FGM detection**

806 To quantify the performance of each method in detecting functional gene modules
807 (FGMs), we conduct gene ontology (GO) enrichment using *clusterProfiler* for each
808 gene set of each bicluster and record the most significant p value (BH adjusted). Since
809 the number of biclusters detected by each method differs, we take the most
810 significant p value of all the biclusters detected by a single method and transform it
811 using $-\log_{10}(p)$ to denote the performance of this method. Methods that fail to detect
812 any bicluster are labeled as 0. We use 10 subsamples from each dataset for repeated
813 evaluation.

814

815 **Criterions for clustering performance**816 There are three metrices we used to benchmark the clustering performance at cell
817 level: ARI, FMI and AMI. Here we will briefly describe how to compute these metrices:

818

819 *ARI*820 The Rand Index computes a similarity measure between two clusterings by
821 considering all pairs of samples and counting pairs that are assigned in the same or
822 different clusters in the predicted and true clusterings. The raw RI score is then
823 “adjusted for chance” into the ARI score using the following scheme:

824
$$ARI = \frac{RI - E(RI)}{\max(RI) - E(RI)}$$

825 To calculate this value, first calculate the contingency table like that:

	Y_1	Y_2	...	Y_s	Sums
X_1	n_{11}	n_{12}	...	n_{1s}	a_1
X_2	n_{21}	n_{22}	...	n_{2s}	a_2
:	:	:	:	:	:
X_r	n_{r1}	n_{r2}	...	n_{rs}	a_r
Sums	b_1	b_2	...	b_s	

826 each value in the table represents the number of data point located in both cluster (Y)
827 and true class (X), and then calculate the ARI value through this table:

828

$$\widehat{ARI} = \frac{\underbrace{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}_{\text{Index}}}{\underbrace{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right]}_{\text{Max Index}} - \underbrace{\left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}_{\text{Expected Index}}}$$

829 The adjusted Rand index is thus ensured to have a value close to 0.0 for random
830 labeling independently of the number of clusters and samples and exactly 1.0 when
831 the clusterings are identical (up to a permutation). The adjusted Rand index is
832 bounded below by -0.5 for especially discordant clusterings.833 *FMI*834 The Fowlkes-Mallows index (FMI) is defined as the geometric mean between of the
835 precision and recall:

836

$$FMI = \frac{TP}{\sqrt{(TP + FP)(TP + FN)}}$$

837 Where TP is the number of **True Positive** (i.e. the number of pair of points that
 838 belongs in the same clusters in both true labels and predicted labels), FP is the number
 839 of **False Positive** (i.e. the number of pair of points that belongs in the same clusters
 840 in true labels but not in predicted labels) and FN is the number of **False Negative** (i.e.
 841 the number of pair of points that belongs in the same clusters in predicted labels but
 842 not in true labels). The score ranges from 0 to 1. A high value indicates a good similarity
 843 between two clusters.

844 *AMI*

845 The Mutual Information is a measure of the similarity between two labels of the same
 846 data. Where $|U_i|$ is the number of the samples in cluster U_i and $|V_j|$ is the number of
 847 the samples in cluster V_j , the Mutual Information between clusterings U and V is
 848 given as:

849

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

850 Adjusted Mutual Information (AMI) is an adjustment of the Mutual Information (MI)
 851 score to account for chance. It accounts for the fact that the MI is generally higher for
 852 two clusterings with a larger number of clusters, regardless of whether there is
 853 actually more information shared. For two clusterings U and V , the AMI is given as:

854

$$AMI(U,V) = \frac{MI(U,V) - E(MI(U,V))}{\text{avg}(H(U), H(V)) - E(MI(U,V))}$$

855 Where $H(*)$ is the information entropy for a label's distribution (eg. $H(U) =$
 856 $\sum_{i=1}^{|U|} P(i) \log(P(i))$). This metric is independent of the absolute values of the labels:
 857 a permutation of the class or cluster label values won't change the score value in any
 858 way.

859

860 **Metrics for biclustering comparison**

861 Suppose $M : \{1\dots L\} \rightarrow \{1\dots L\}$ maps the ground true bicluster index to the index of the
 862 bicluster detected by an algorithm, T_i denote the i_{th} ground true bicluster and B_i denote

863 the i_{th} detected bicluster. The Cluster Error (CE) proposed by Anne et.al[60] is defined
864 as:

$$865 \quad 1 - CE(M) = \frac{\sum_{i=1}^L |T_i \cap B_{M(i)}|}{|\cup_{i=1}^L T_i \cup B_{M(i)}|}$$

866 This is a distance measure of subspace clustering with lower CE indicating better
867 consistency with ground truth. When we evaluate the performance, we choose a M
868 minimizing the CE as the optimal match and is used by other measurements. The
869 corresponding 1-CE is output with the higher the value, the better.

870
871 We also use F-score (F) to evaluate the performance. F-score is the harmonic mean of
872 precision (PRE) and recall (REC). Here we use the calculation way proposed by Zhong
873 et.al[18]:

$$874 \quad PRE_i = \frac{|T_i \cap B_{M(i)}|}{|B_{M(i)}|}$$

$$875 \quad REC_i = \frac{|T_i \cap B_M|}{|T_i|}$$

876 Where A denote all the elements of the expression data. PRE_i and REC_i are computed
877 for bicluster pair i , we finally output the average for each criterion as PRE and REC,
878 along with their harmonic mean as F-score(F), which is a combination of the two, and
879 we also pay more attention to it. The higher this indicator is, the better.

880

881 Supplementary Information

882 **Additional file 1:** **Fig. S1.** overlap of biclusters. **Fig. S2.** visualization of cell clustering
883 results on PBMC dataset. **Fig. S3.** visualization of cell clustering results on HEART
884 dataset. **Fig. S4.** visualization of cell clustering results on LUAD dataset. **Fig. S5.**
885 visualization of cell clustering results on BC dataset. **Table S1.** Enrichment score
886 comparison of different methods. **Table S2.** ARI of different methods in different
887 datasets. **Table S3.** FMI of different methods in different datasets. **Table S4.** AMI of
888 different methods in different datasets. **Table S5-S7.** 1-CE of different methods in
889 different datasets. **Table S8-S10.** F score of different methods in different datasets.

890 **Additional file 2.** **Table S11.** Clinico-pathological variables of 48 AD patients.

891 **Additional file 3.** **Table S12.** enrichment results of perturbated gene sets.

892 **Additional file 4.** **Table S13.** matching results of biclusters in different pathological
893 stage.

894

895 Declarations

896 Data and code availability

897 All data used in this research can be found in **Datasets and preprocessing** in **STAR**
898 **Methods** section. Our scBC method is available as a Python package on PyPI at

899 <https://pypi.org/project/scBC>, free for academic use, and the source code is openly
900 available from our GitHub repository at <https://github.com/GYQ-form/scBC>.

901

902 Competing interests

903 The authors declare no competing interests.

904

905 Funding

906 The research is supported partly by the National Natural Science Foundation of China
907 (11901387 for Y.Z. and 12171318 for Z.Y.).

908

909 Authors' contributions

910 YG performed the research, analyzed data and wrote the original manuscript, JX
911 participated in data collection, ZY and YZ supervised the research. YG, JX, RG, JS, ZY,
912 and YZ discussed and revised the manuscript.

913

914 Acknowledgements

915 Not applicable.

916

917

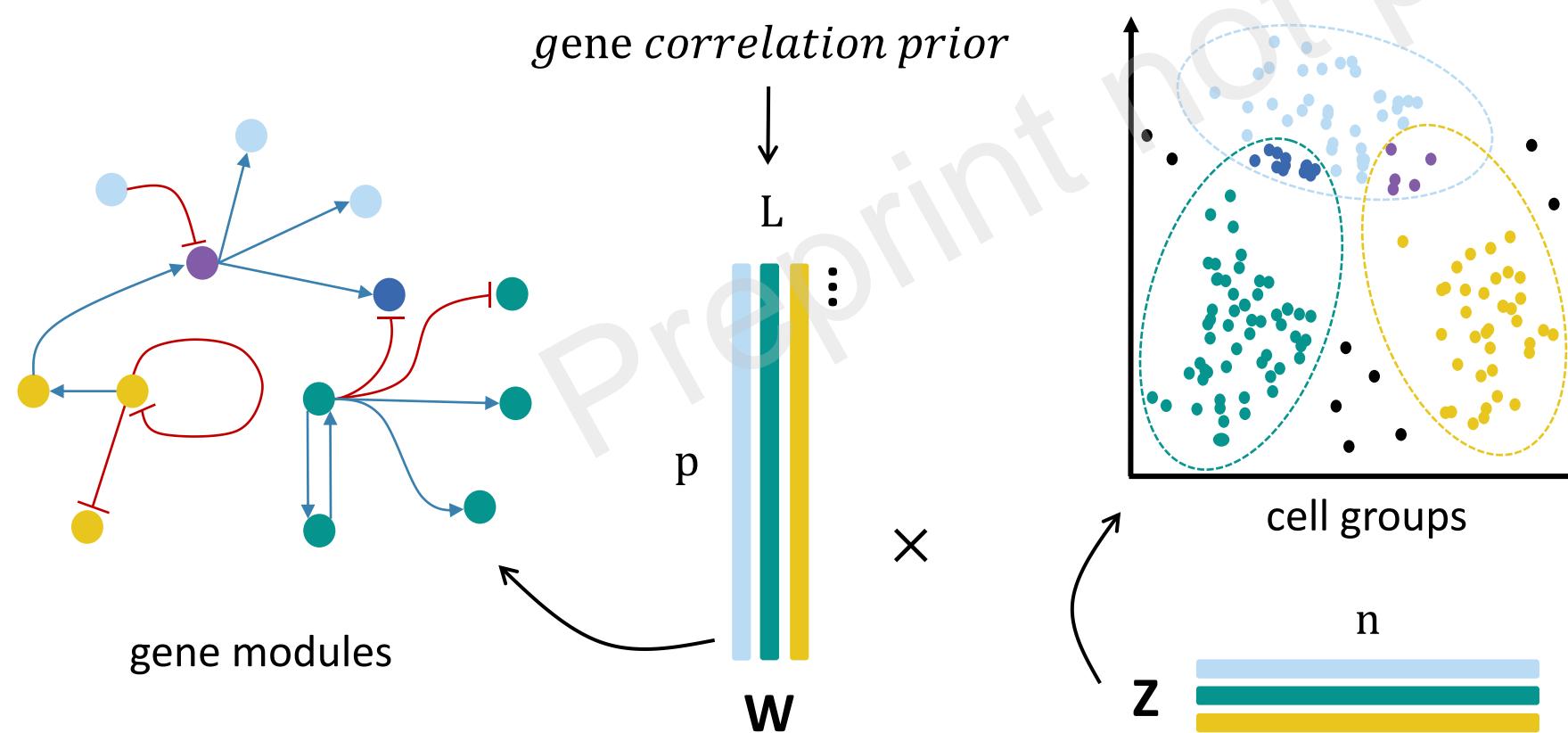
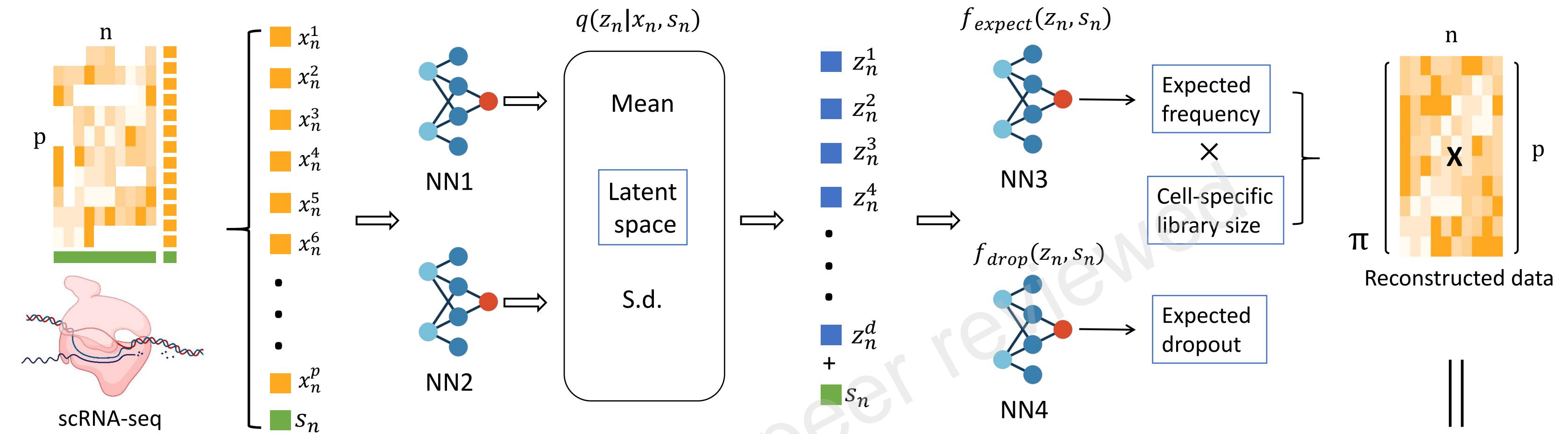
918 References

- 919 1. Shi, F. and H. Huang, *Identifying Cell Subpopulations and Their Genetic Drivers from Single-*
920 *Cell RNA-Seq Data Using a Bioclustering Approach*. *J Comput Biol*, 2017. **24**(7): p. 663-674.
- 921 2. Eisenberg, E. and E.Y. Levanon, *Human housekeeping genes, revisited*. *Trends Genet*, 2013.
922 **29**(10): p. 569-74.
- 923 3. Mathys, H., et al., *Single-cell transcriptomic analysis of Alzheimer's disease*. *Nature*, 2019.
924 **570**(7761): p. 332-337.
- 925 4. Lake, B.B., et al., *Integrative single-cell analysis of transcriptional and epigenetic states in the*
926 *human adult brain*. *Nat Biotechnol*, 2018. **36**(1): p. 70-80.
- 927 5. Grubman, A., et al., *A single-cell atlas of entorhinal cortex from individuals with Alzheimer's*
928 *disease reveals cell-type-specific gene expression regulation*. *Nat Neurosci*, 2019. **22**(12): p.
929 2087-2097.
- 930 6. Habib, N., et al., *Disease-associated astrocytes in Alzheimer's disease and aging*. *Nat*
931 *Neurosci*, 2020. **23**(6): p. 701-706.
- 932 7. Roussarie, J.P., et al., *Selective Neuronal Vulnerability in Alzheimer's Disease: A Network-*
933 *Based Analysis*. *Neuron*, 2020. **107**(5): p. 821-+.
- 934 8. Langfelder, P. and S. Horvath, *WGCNA: an R package for weighted correlation network*
935 *analysis*. *BMC Bioinformatics*, 2008. **9**: p. 559.
- 936 9. Zhang, B. and S. Horvath, *A general framework for weighted gene co-expression network*
937 *analysis*. *Stat Appl Genet Mol Biol*, 2005. **4**: p. Article17.
- 938 10. Ghazalpour, A., et al., *Comparative analysis of proteome and transcriptome variation in*
939 *mouse*. *PLoS Genet*, 2011. **7**(6): p. e1001393.
- 940 11. Zhang, B., et al., *Integrated systems approach identifies genetic nodes and networks in late-*
941 *onset Alzheimer's disease*. *Cell*, 2013. **153**(3): p. 707-20.

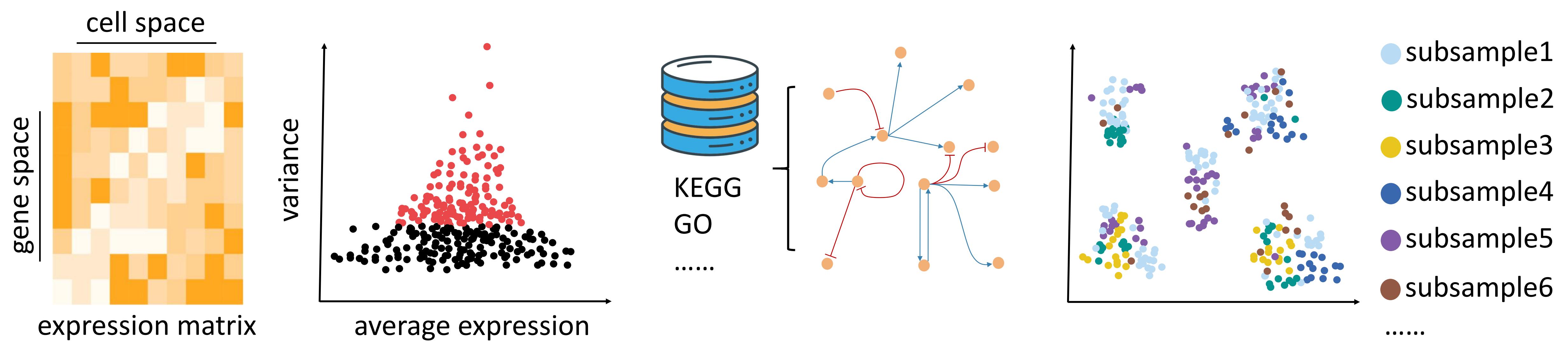
- 942 12. Risso, D., et al., *Normalization of RNA-seq data using factor analysis of control genes or*
943 *samples*. Nat Biotechnol, 2014. **32**(9): p. 896-902.
- 944 13. Hochreiter, S., et al., *FABIA: factor analysis for bicluster acquisition*. Bioinformatics, 2010.
945 **26**(12): p. 1520-7.
- 946 14. Leek, J.T., et al., *Tackling the widespread and critical impact of batch effects in high-*
947 *throughput data*. Nat Rev Genet, 2010. **11**(10): p. 733-9.
- 948 15. Hu, Z., S. Zu, and J.S. Liu, *SIMPLEs: a single-cell RNA sequencing imputation strategy*
949 *preserving gene modules and cell clusters variation*. NAR Genom Bioinform, 2020. **2**(4): p.
950 Iqaa077.
- 951 16. Fang, Q., et al., *An Effective Biclustering-Based Framework for Identifying Cell Subpopulations*
952 *From scRNA-seq Data*. IEEE/ACM Trans Comput Biol Bioinform, 2021. **18**(6): p. 2249-2260.
- 953 17. Xie, J., et al., *QUBIC2: a novel and robust biclustering algorithm for analyses and*
954 *interpretation of large-scale RNA-Seq data*. Bioinformatics, 2020. **36**(4): p. 1143-1149.
- 955 18. Zhong, Y. and J.Z. Huang, *Biclustering via structured regularized matrix decomposition*.
956 Statistics and Computing, 2022. **32**(3).
- 957 19. Lopez, R., et al., *Deep generative modeling for single-cell transcriptomics*. Nat Methods, 2018.
958 **15**(12): p. 1053-1058.
- 959 20. Li, Z., et al., *Bayesian generalized biclustering analysis via adaptive structured shrinkage*.
960 Biostatistics, 2020. **21**(3): p. 610-624.
- 961 21. Cheng, Y. and G.M. Church, *Biclustering of expression data*. Proceedings. International
962 Conference on Intelligent Systems for Molecular Biology, 2000. **8**: p. 93-103.
- 963 22. Murali, T.M. and S. Kasif, *Extracting conserved gene expression motifs from gene expression*
964 *data*. Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing, 2003: p. 77-
965 88.
- 966 23. Prelic, A., et al., *A systematic comparison and evaluation of biclustering methods for gene*
967 *expression data*. Bioinformatics, 2006. **22**(9): p. 1122-9.
- 968 24. Caldas, J. and S. Kaski, *Bayesian Biclustering with the Plaid Model*. 2008 Ieee Workshop on
969 Machine Learning for Signal Processing, 2008: p. 291-296.
- 970 25. Dongen, S.V., *Graph clustering by flow simulation*. PhD thesis, University of Utrecht, 2000.
- 971 26. Milligan, G.W. and M.C. Cooper, *A STUDY OF THE COMPARABILITY OF EXTERNAL CRITERIA*
972 *FOR HIERARCHICAL CLUSTER-ANALYSIS*. Multivariate Behavioral Research, 1986. **21**(4): p.
973 441-458.
- 974 27. Santos, J.M. and M. Embrechts, *On the Use of the Adjusted Rand Index as a Metric for*
975 *Evaluating Supervised Classification*. in *19th International Conference on Artificial Neural*
976 *Networks (ICANN 2009)*. 2009. Limassol, CYPRUS.
- 977 28. Fowlkes, E.B. and C.L. Mallows, *A method for comparing two hierarchical clusterings*. Journal
978 of the American statistical association, 1983. **78**(383): p. 553-569.
- 979 29. Strehl, A. and J. Ghosh, *Cluster ensembles- a knowledge reuse framework for combining*
980 *multiple partitions*. Journal of Machine Learning Research, 2003. **3**(3): p. 583-617.
- 981 30. Rahaman, M.A., et al., *Shared sets of correlated polygenic risk scores and voxel-wise grey*
982 *matter across multiple traits identified via bi-clustering*. Annu Int Conf IEEE Eng Med Biol Soc,
983 2021. **2021**: p. 2201-2206.
- 984 31. Murdock, M.H. and L.H. Tsai, *Insights into Alzheimer's disease from single-cell genomic*
985 *approaches*. Nat Neurosci, 2023. **26**(2): p. 181-195.

- 986 32. Hasel, P., et al., *Neuroinflammatory astrocyte subtypes in the mouse brain*. Nat Neurosci, 987 2021. **24**(10): p. 1475-1487.
- 988 33. Blanchard, J.W., et al., *APOE4 impairs myelination via cholesterol dysregulation in* 989 *oligodendrocytes*. Nature, 2022. **611**(7937): p. 769-779.
- 990 34. Welch, G. and L.H. Tsai, *Mechanisms of DNA damage-mediated neurotoxicity in* 991 *neurodegenerative disease*. EMBO Rep, 2022. **23**(6): p. e54217.
- 992 35. Zhou, Y., et al., *Author Correction: Human and mouse single-nucleus transcriptomics reveal* 993 *TREM2-dependent and TREM2-independent cellular responses in Alzheimer's disease*. Nat 994 Med, 2020. **26**(6): p. 981.
- 995 36. Lau, S.F., et al., *Single-nucleus transcriptome analysis reveals dysregulation of angiogenic* 996 *endothelial cells and neuroprotective glia in Alzheimer's disease*. Proc Natl Acad Sci U S A, 997 2020. **117**(41): p. 25800-25809.
- 998 37. Pan, S., et al., *Preservation of a remote fear memory requires new myelin formation*. Nat 999 Neurosci, 2020. **23**(4): p. 487-499.
- 1000 38. Fancy, S.P., et al., *Myelin regeneration: a recapitulation of development?* Annu Rev Neurosci, 1001 2011. **34**: p. 21-43.
- 1002 39. Franklin, R.J. and S.A. Goldman, *Glia Disease and Repair-Remyelination*. Cold Spring Harb 1003 Perspect Biol, 2015. **7**(7): p. a020594.
- 1004 40. Karadottir, R., et al., *Spiking and nonspiking classes of oligodendrocyte precursor glia in CNS* 1005 *white matter*. Nat Neurosci, 2008. **11**(4): p. 450-6.
- 1006 41. Mitew, S., et al., *Mechanisms regulating the development of oligodendrocytes and central* 1007 *nervous system myelin*. Neuroscience, 2014. **276**: p. 29-47.
- 1008 42. Depp, C., et al., *Myelin dysfunction drives amyloid-beta deposition in models of Alzheimer's* 1009 *disease*. Nature, 2023. **618**(7964): p. 349-357.
- 1010 43. Liu, R., et al., *Early diagnosis of complex diseases by molecular biomarkers, network* 1011 *biomarkers, and dynamical network biomarkers*. Med Res Rev, 2014. **34**(3): p. 455-78.
- 1012 44. Jin, G.X., et al., *The knowledge-integrated network biomarkers discovery for Major Adverse* 1013 *Cardiac Events*. Journal of Proteome Research, 2008. **7**(9): p. 4013-4021.
- 1014 45. Ideker, T. and R. Sharan, *Protein networks in disease*. Genome Res, 2008. **18**(4): p. 644-52.
- 1015 46. Risso, D., et al., *A general and flexible method for signal extraction from single-cell RNA-seq* 1016 *data*. Nat Commun, 2018. **9**(1): p. 284.
- 1017 47. Grun, D., L. Kester, and A. van Oudenaarden, *Validation of noise models for single-cell* 1018 *transcriptomics*. Nature Methods, 2014. **11**(6): p. 637-+.
- 1019 48. Kingma, D.P. and M. Welling *Auto-Encoding Variational Bayes*. 2013. arXiv:1312.6114 DOI: 10.48550/arXiv.1312.6114.
- 1020 49. Sonderby, C.K., et al., *Ladder Variational Autoencoders*. Advances in Neural Information 1021 Processing Systems 29 (Nips 2016), 2016. **29**.
- 1022 50. Li, C. and H. Li, *Network-constrained regularization and variable selection for analysis of* 1023 *genomic data*. Bioinformatics, 2008. **24**(9): p. 1175-82.
- 1024 51. Zhao, Y., et al., *Hierarchical Feature Selection Incorporating Known and Novel Biological* 1025 *Information: Identifying Genomic Features Related to Prostate Cancer Recurrence*. J Am Stat 1026 Assoc, 2016. **111**(516): p. 1427-1439.
- 1027 52. Li, Z., S.E. Safo, and Q. Long, *Incorporating biological information in sparse principal* 1028 *component analysis with application to genomic data*. BMC Bioinformatics, 2017. **18**(1): p.

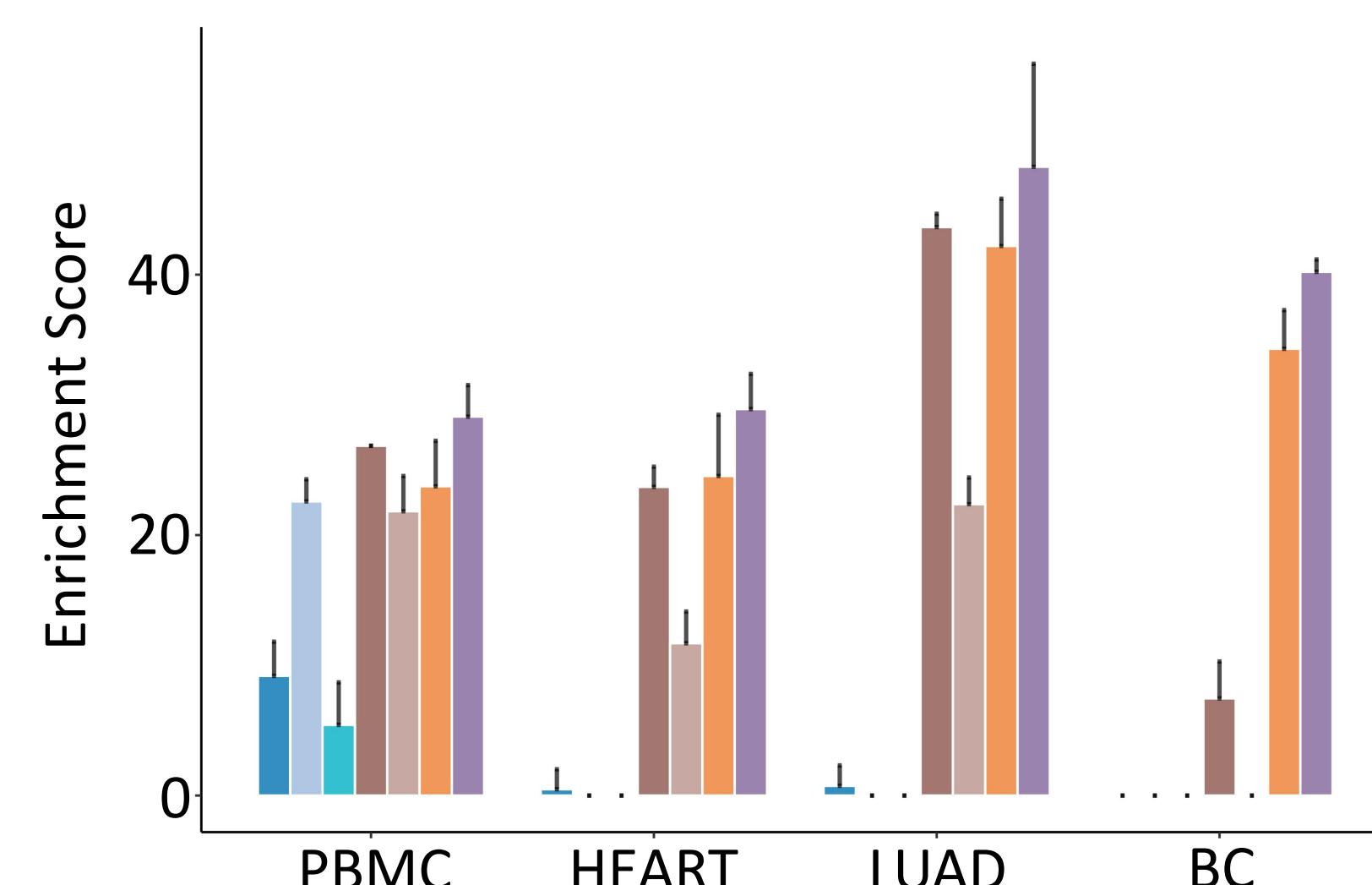
- 1030 332.
- 1031 53. Safo, S.E., S. Li, and Q. Long, *Integrative analysis of transcriptomic and metabolomic data via*
1032 *sparse canonical correlation analysis with incorporation of biological information*. Biometrics,
1033 2018. **74**(1): p. 300-312.
- 1034 54. Chang, C., S. Kundu, and Q. Long, *Scalable Bayesian variable selection for structured high-*
1035 *dimensional data*. Biometrics, 2018. **74**(4): p. 1372-1382.
- 1036 55. Polson, N.G., J.G. Scott, and J. Windle, *Bayesian Inference for Logistic Models Using Pólya–*
1037 *Gamma Latent Variables*. Journal of the American Statistical Association, 2013. **108**(504): p.
1038 1339-1349.
- 1039 56. Chang, C. and R.S. Tsay, *Estimation of covariance matrix via the sparse Cholesky factor with*
1040 *lasso*. Journal of Statistical Planning and Inference, 2010. **140**(12): p. 3858-3873.
- 1041 57. Zheng, G.X., et al., *Massively parallel digital transcriptional profiling of single cells*. Nat
1042 Commun, 2017. **8**: p. 14049.
- 1043 58. Kim, N., et al., *Single-cell RNA sequencing demonstrates the molecular and cellular*
1044 *reprogramming of metastatic lung adenocarcinoma*. Nat Commun, 2020. **11**(1): p. 2285.
- 1045 59. Chung, W., et al., *Single-cell RNA-seq enables comprehensive tumour and immune cell*
1046 *profiling in primary breast cancer*. Nat Commun, 2017. **8**: p. 15081.
- 1047 60. Patrikainen, A. and M. Meila, *Comparing subspace clusterings*. Ieee Transactions on
1048 Knowledge and Data Engineering, 2006. **18**(7): p. 902-916.
- 1049
- 1050



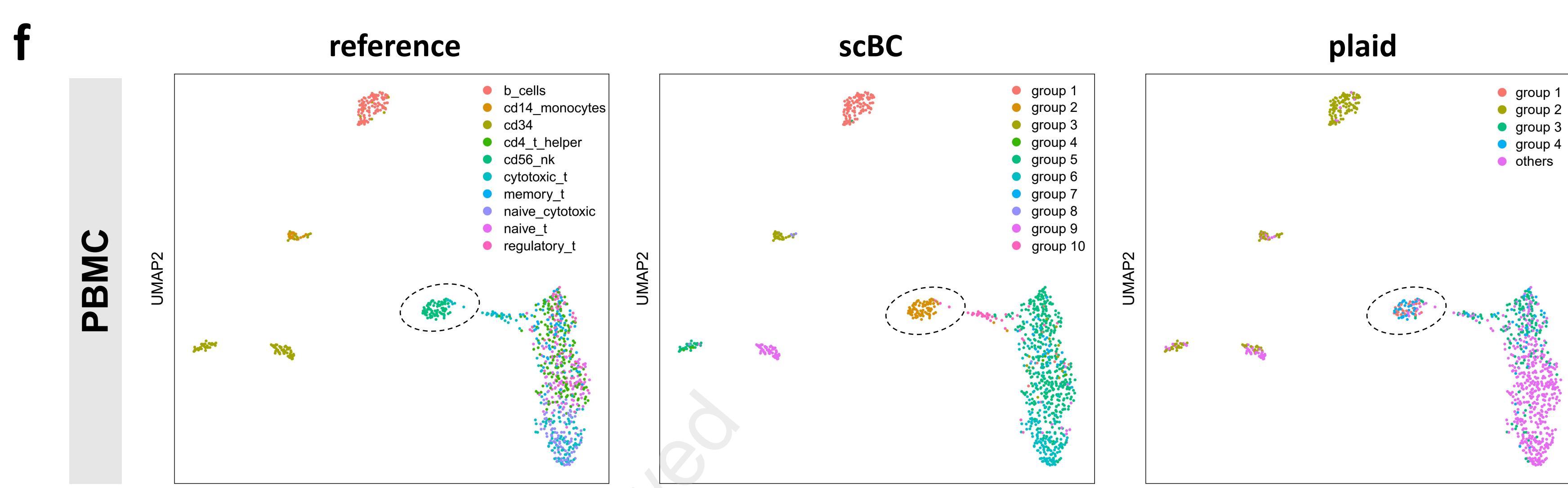
a Get data → Identify highly variable genes → Extract prior from database → Randomly sample cells for repetition



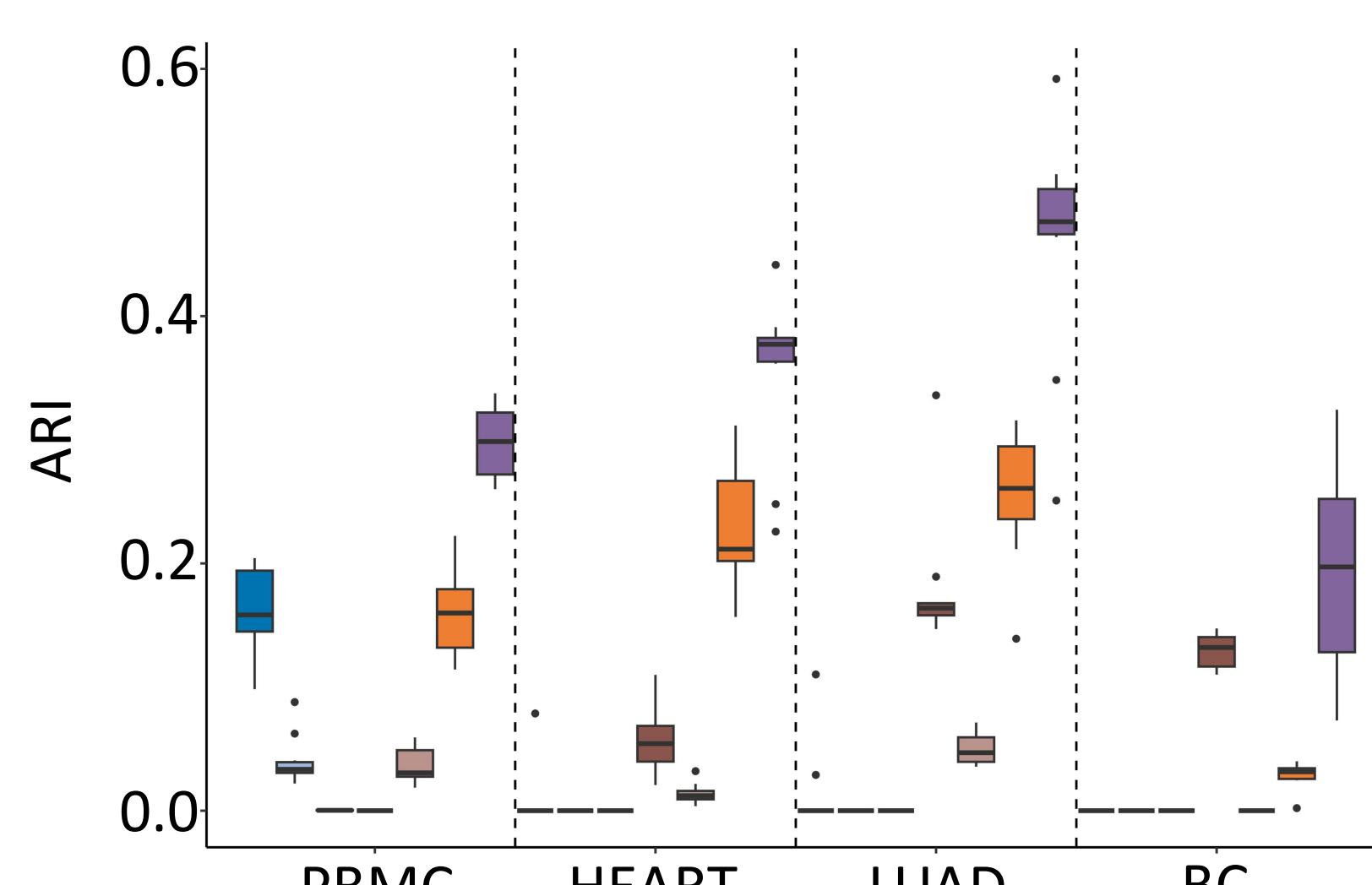
b



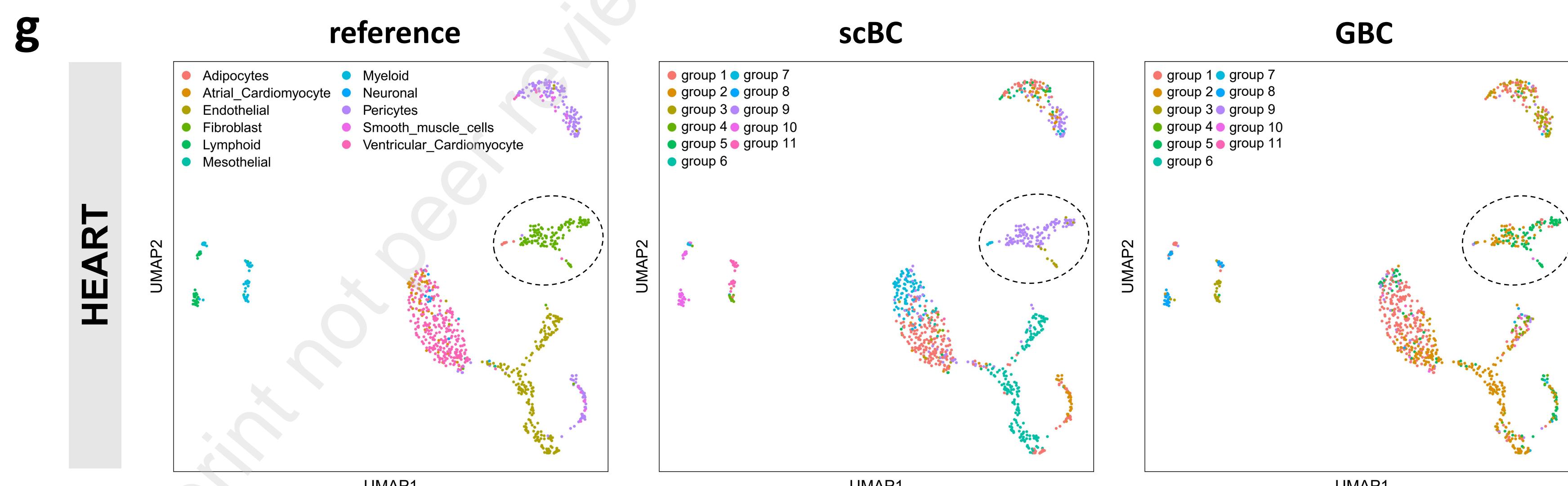
f



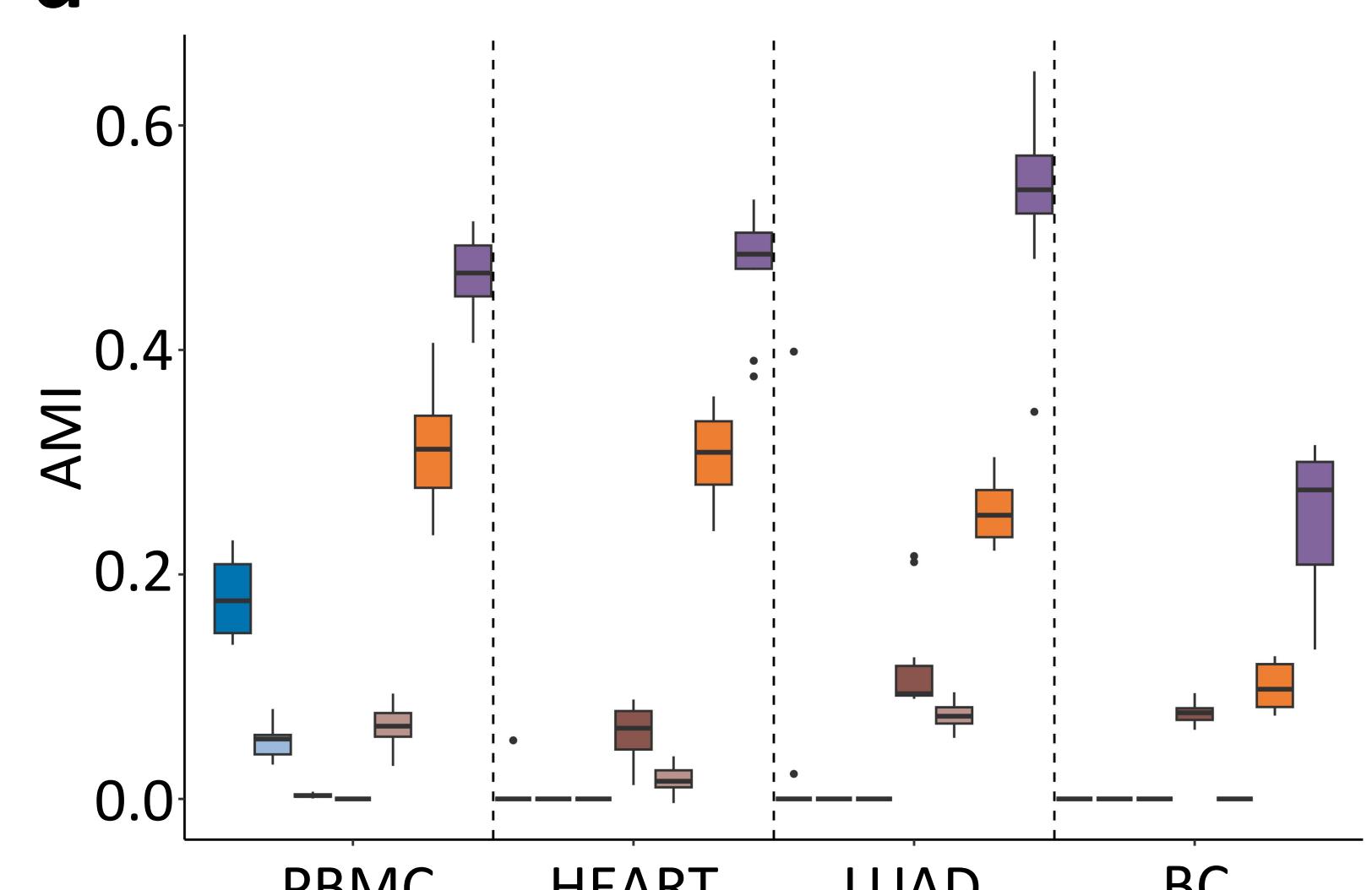
c



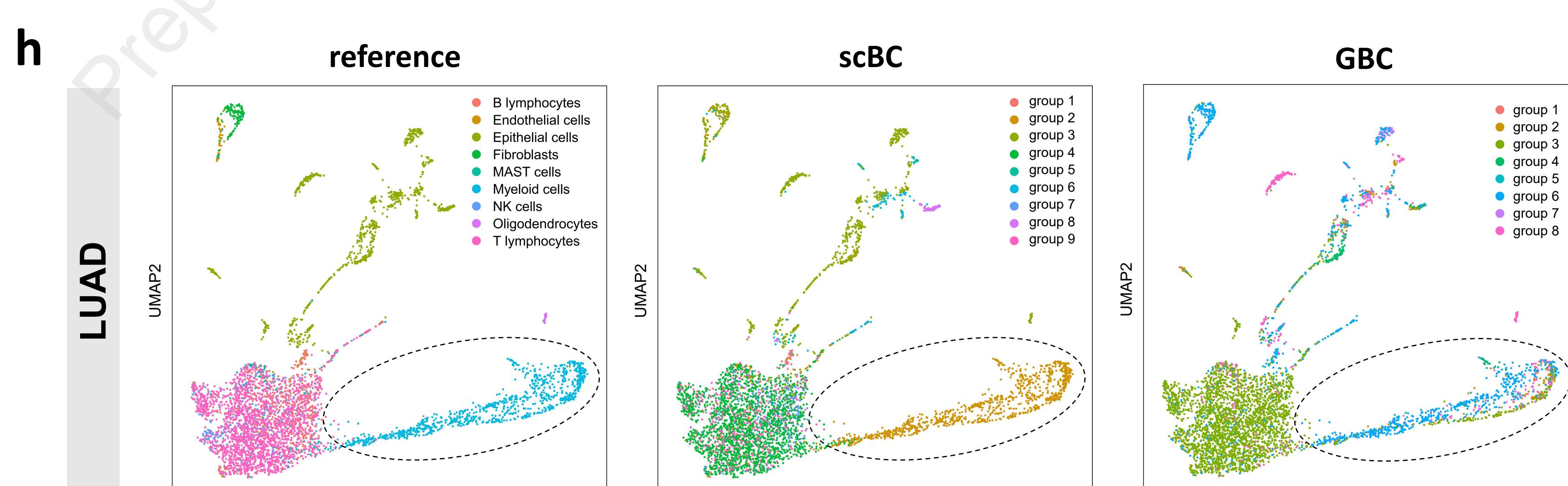
g



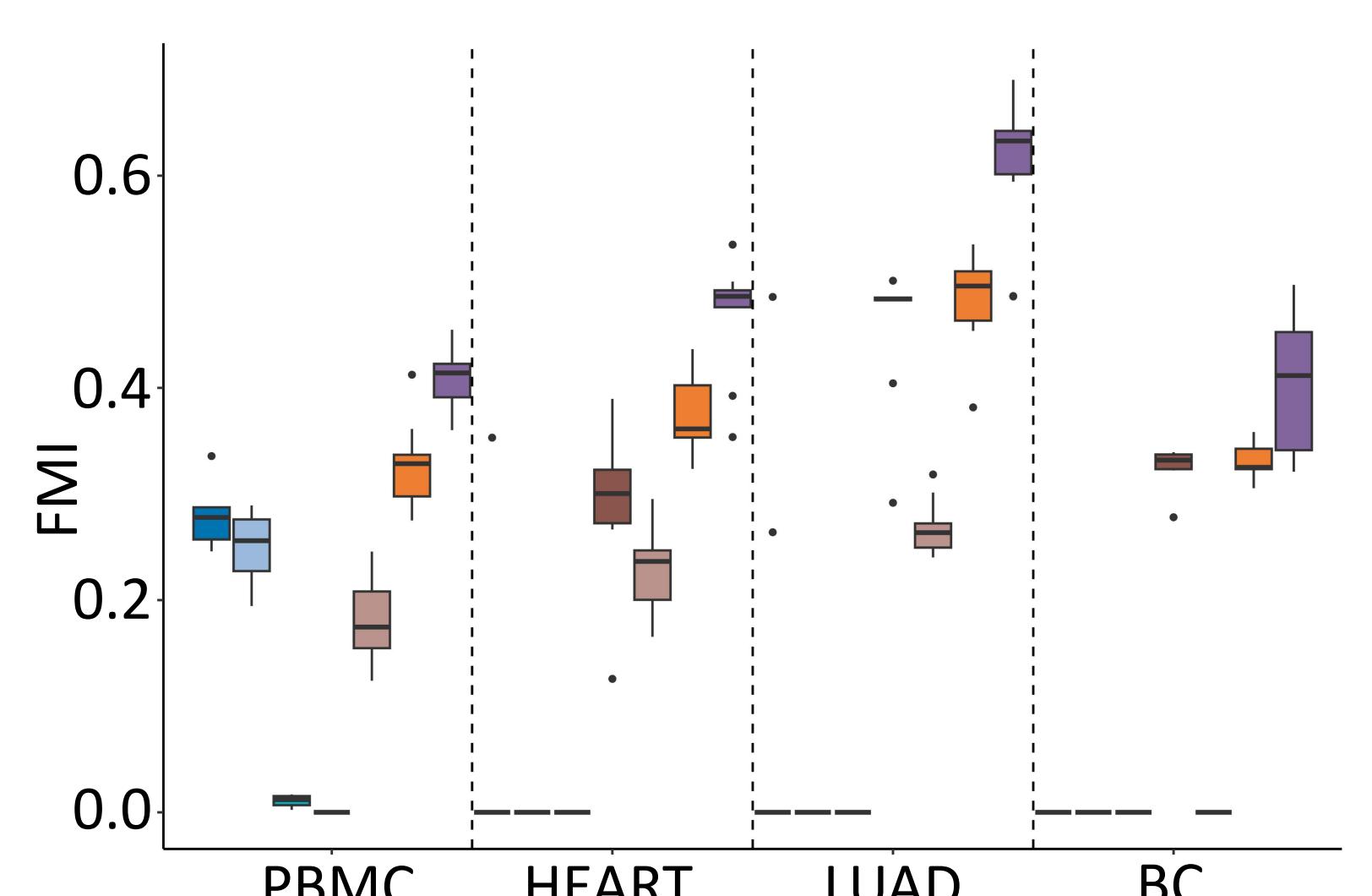
d



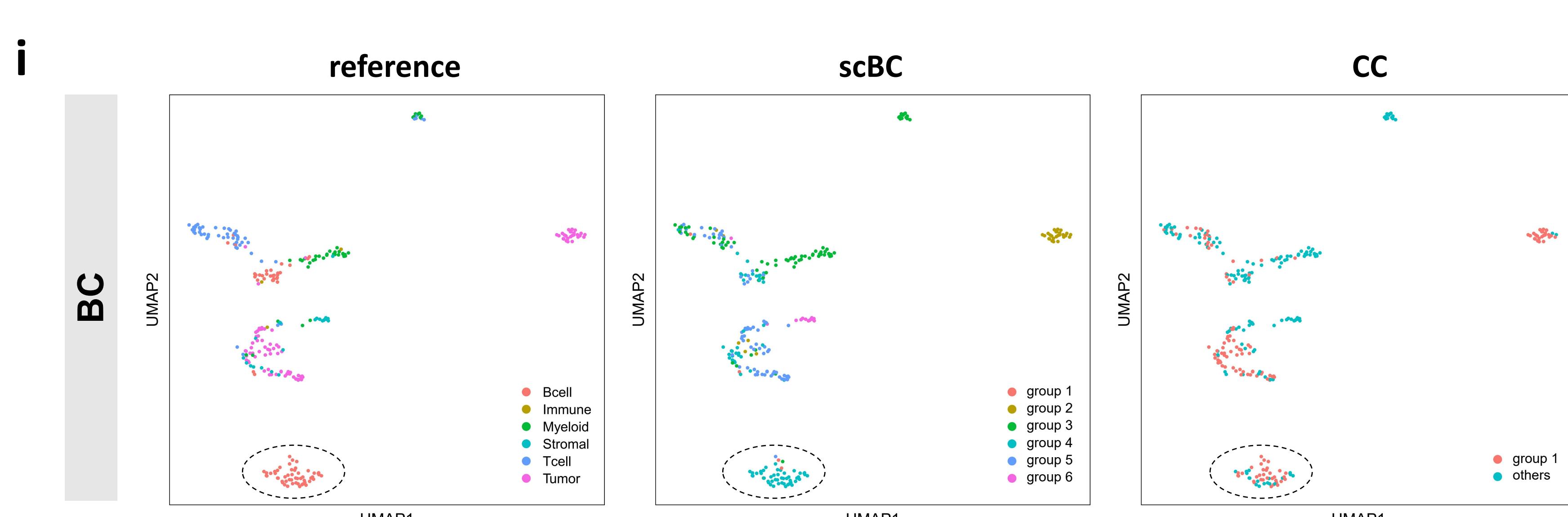
h

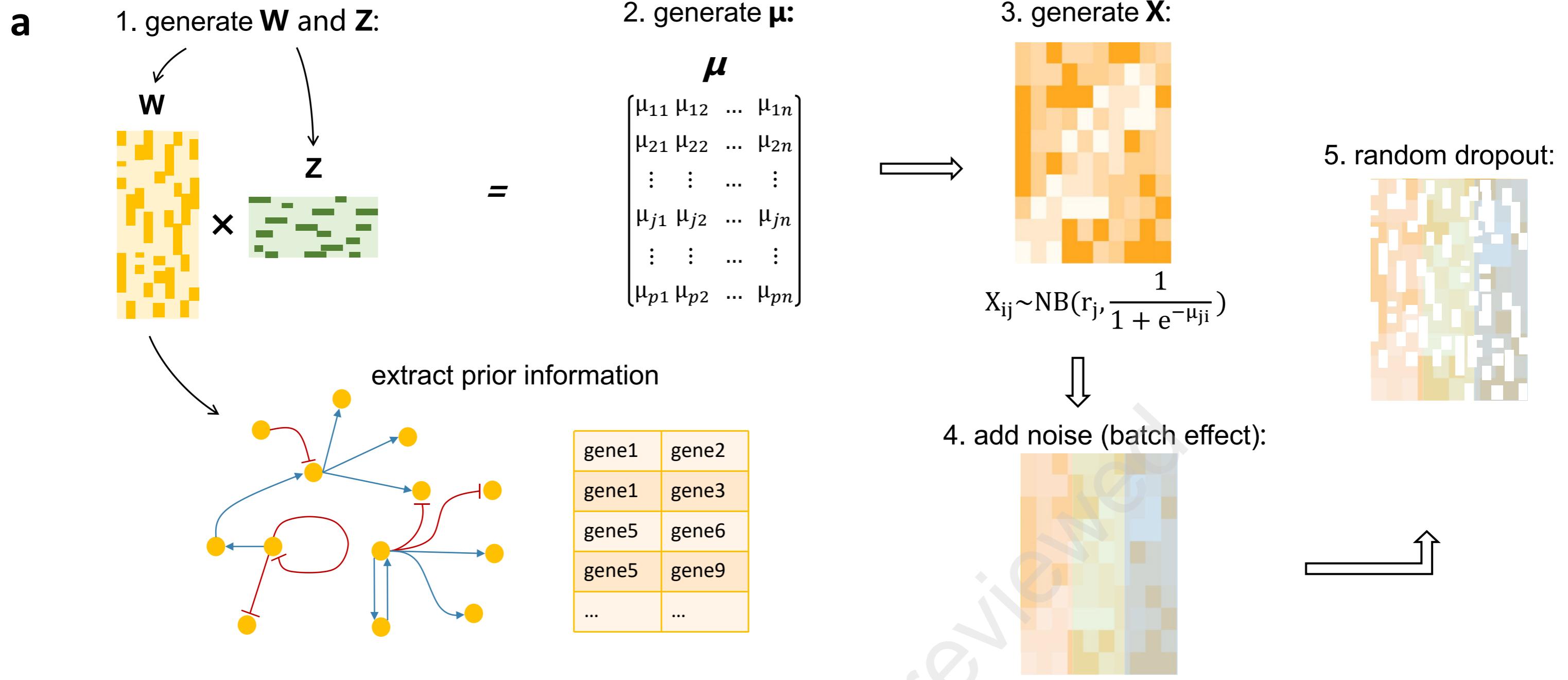


e

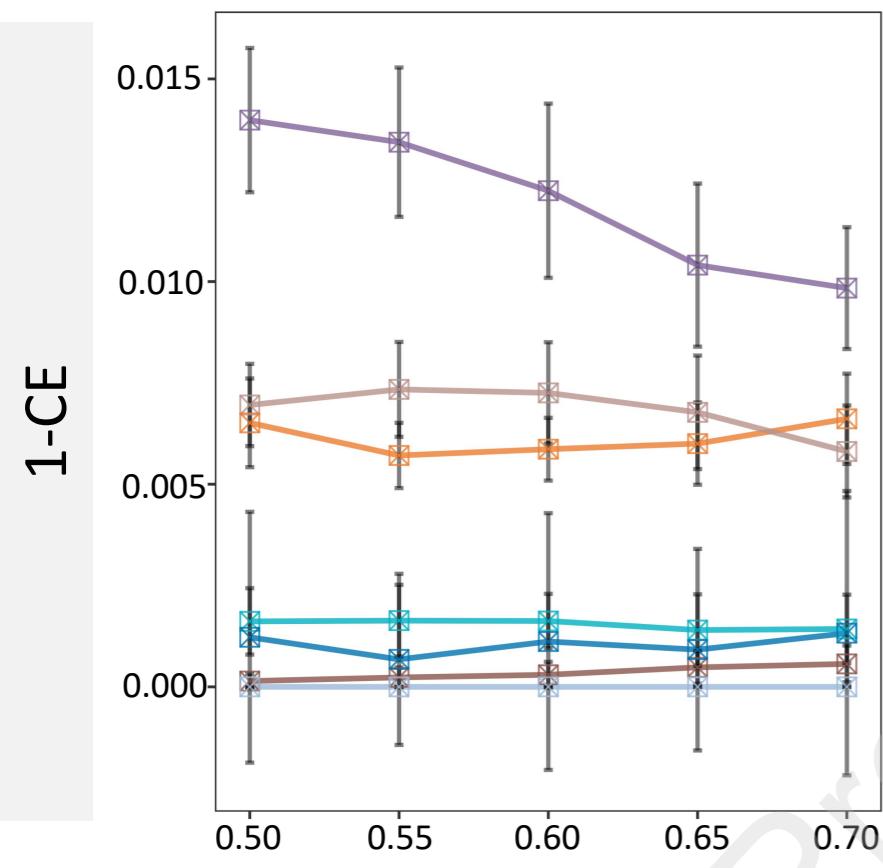


i

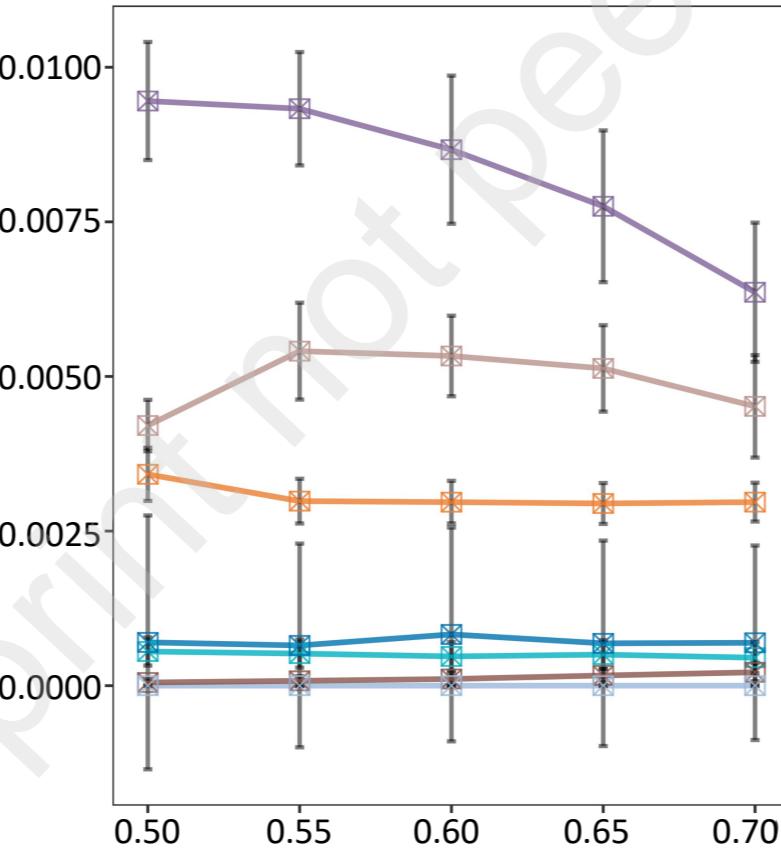




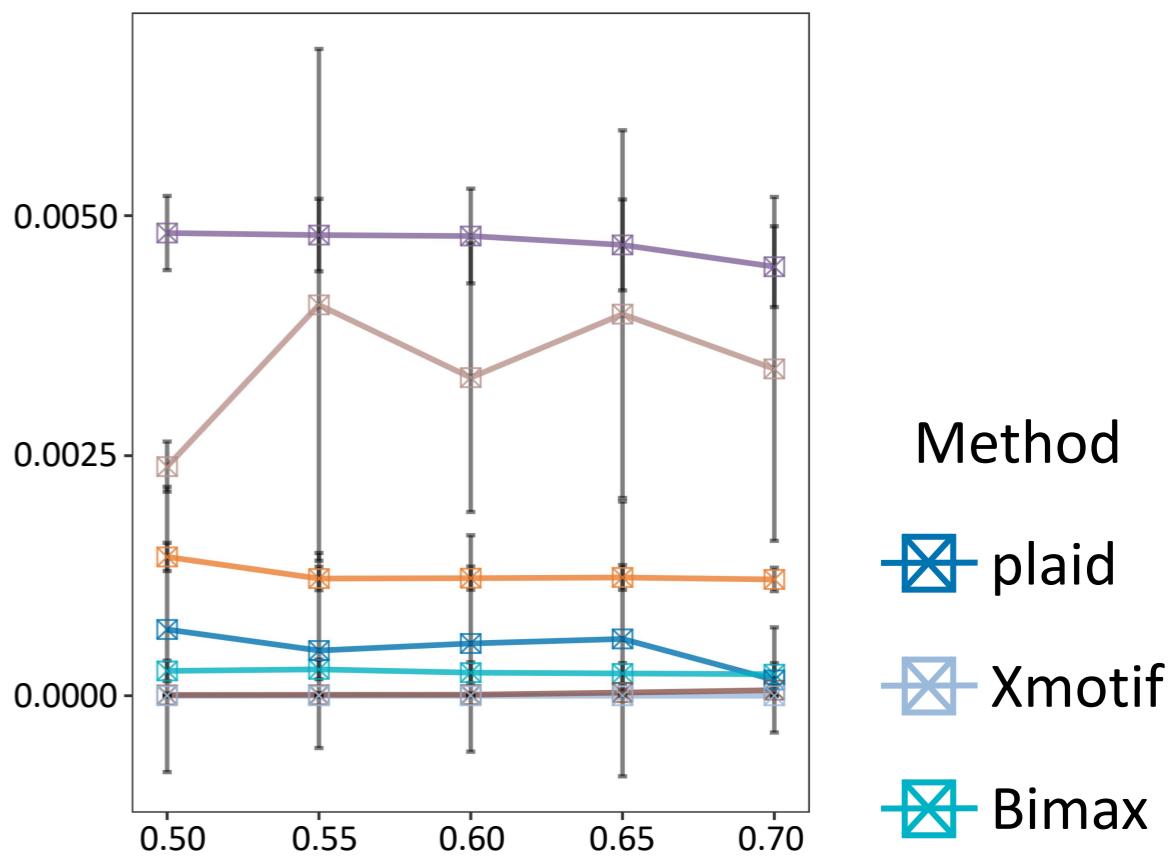
b $p=1000, n=300, L=3$



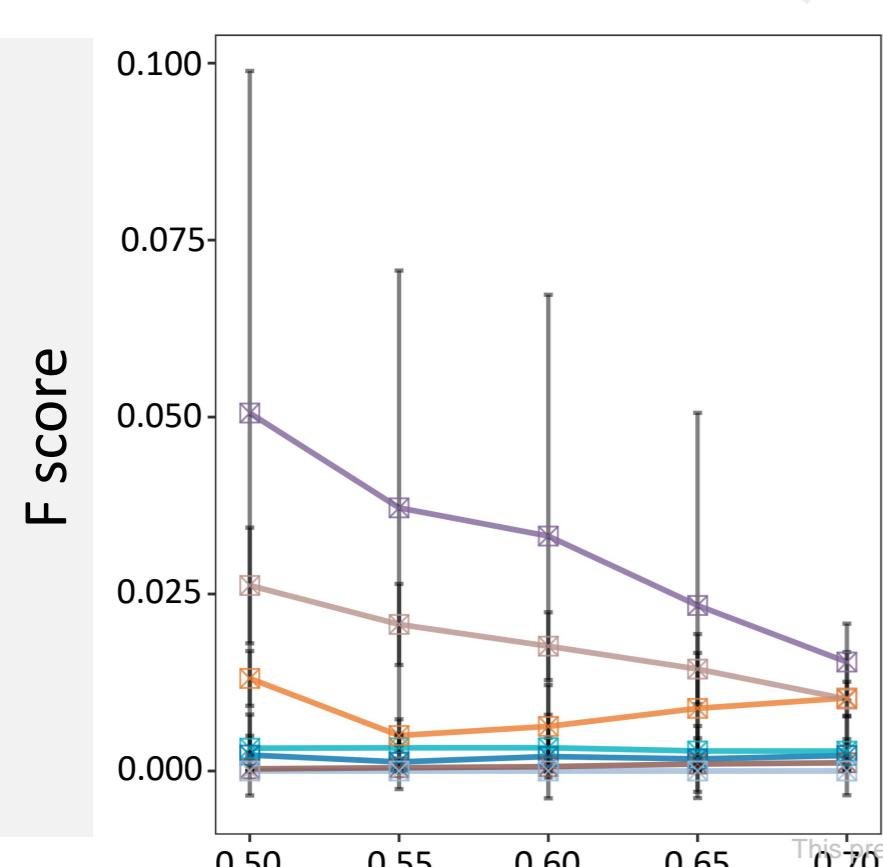
c $p=3000, n=600, L=4$



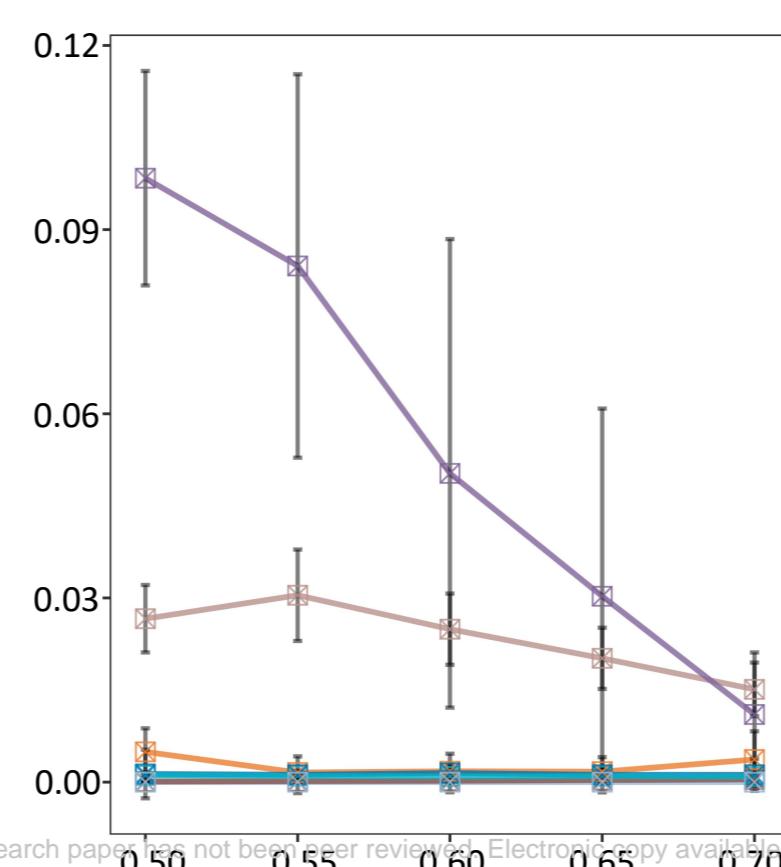
d $p=6000, n=1500, L=5$



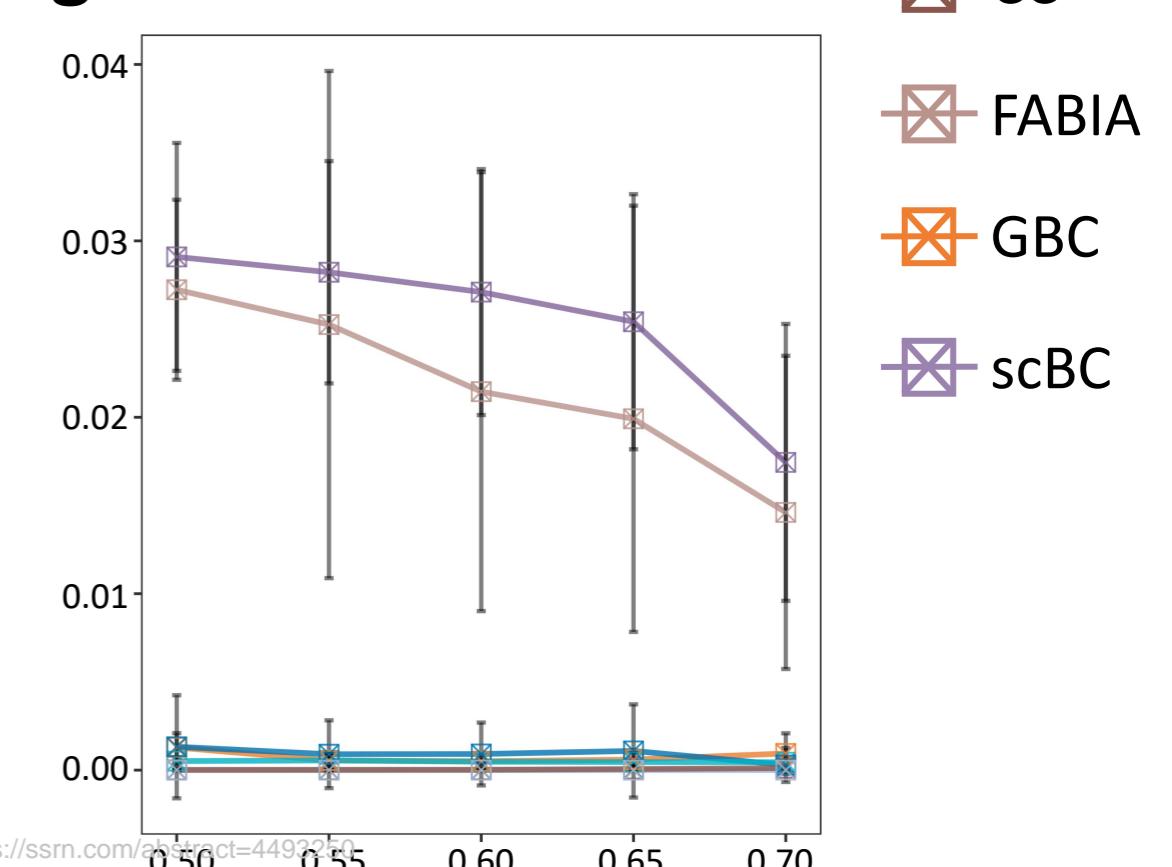
e

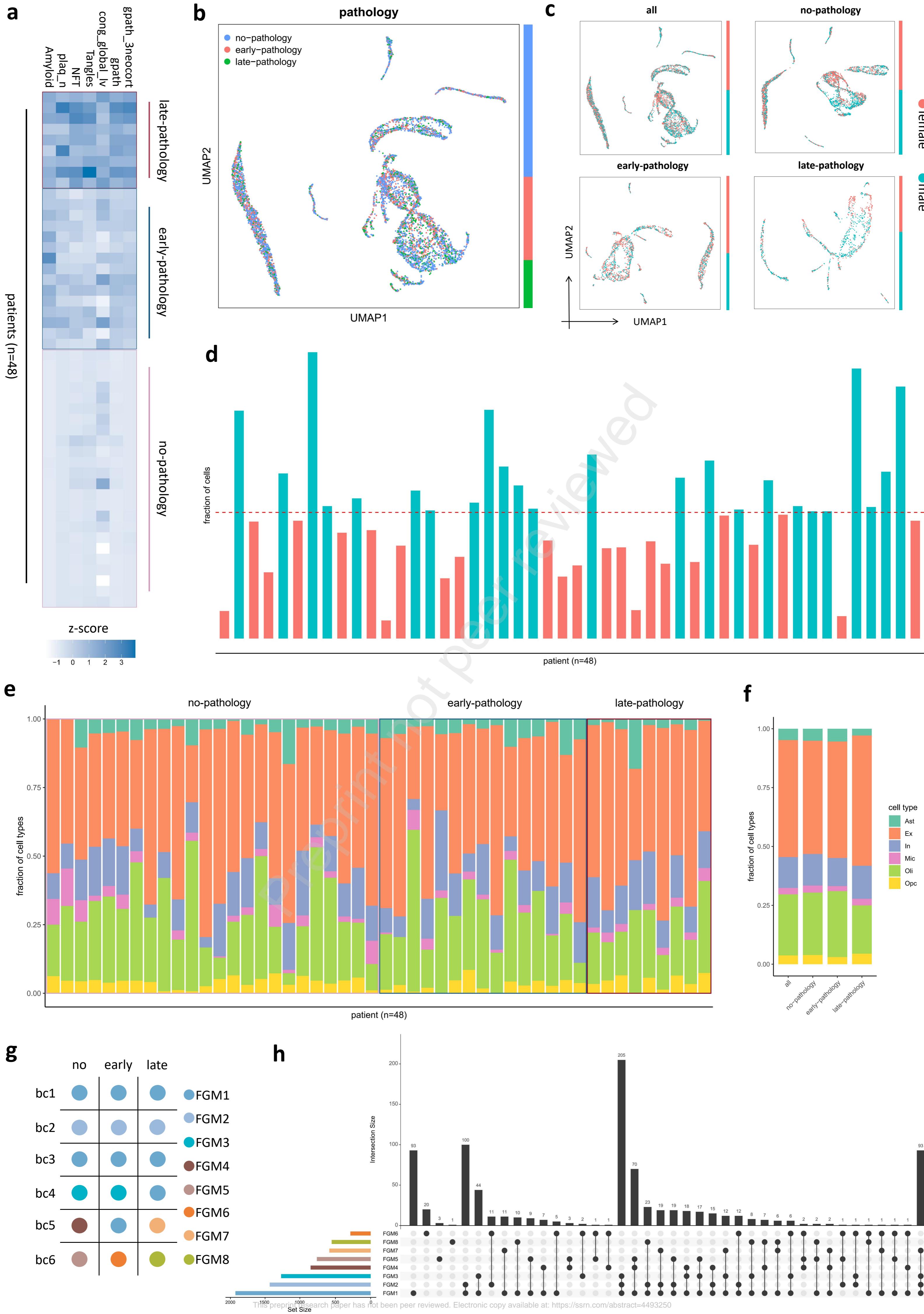


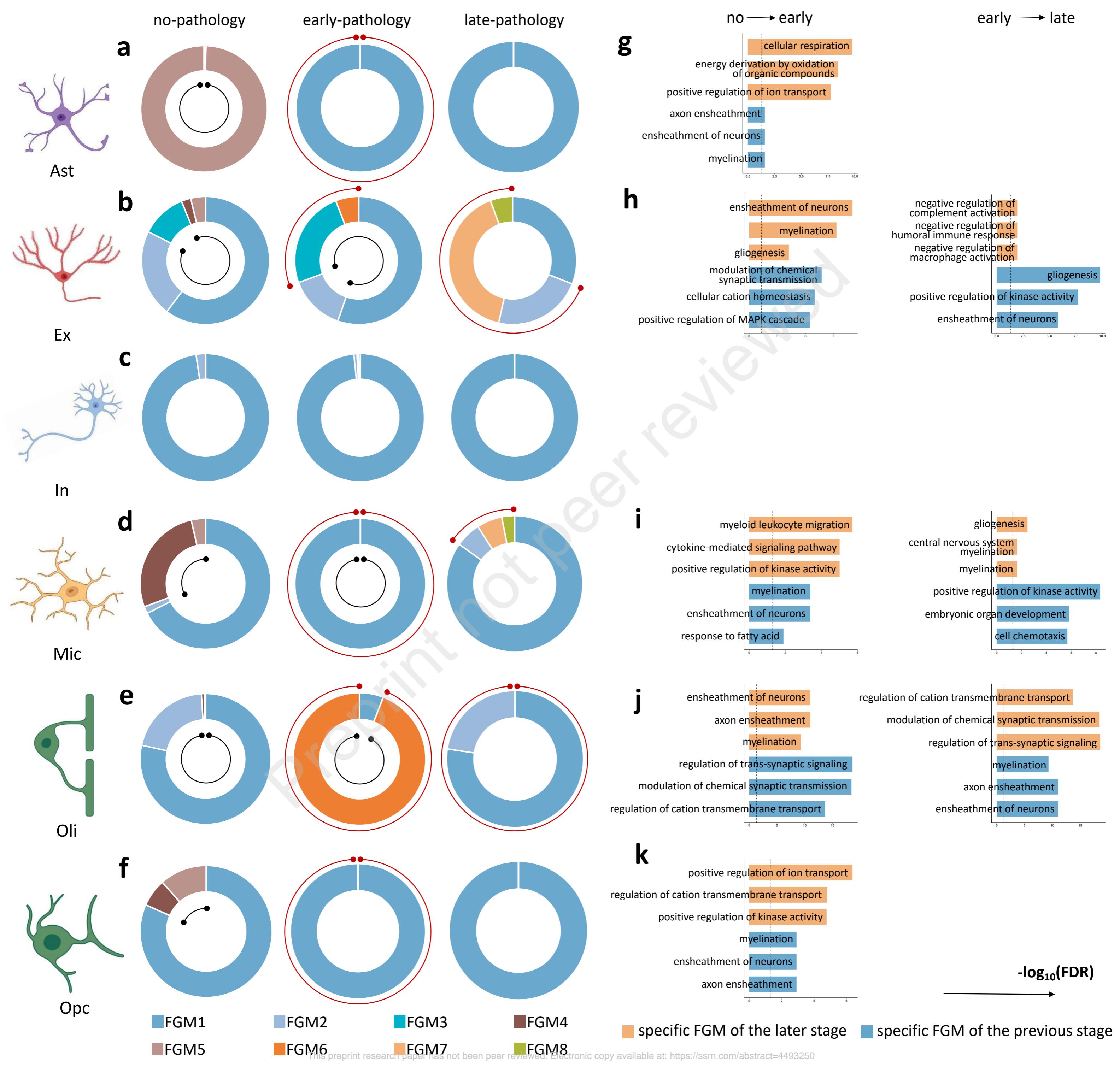
f



g







Supplementary Information

Supplementary Figures

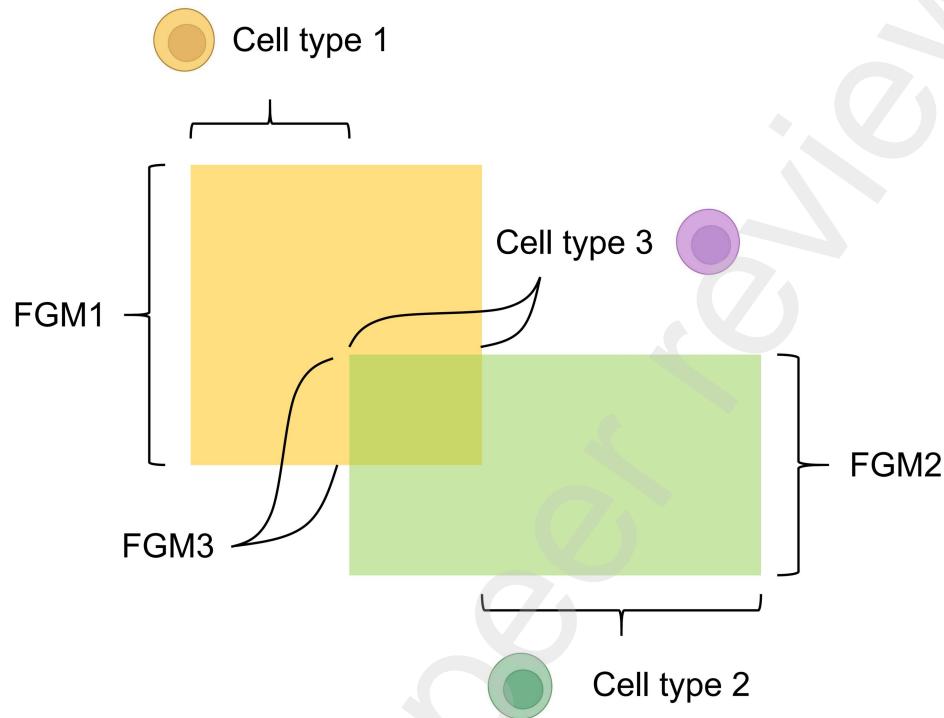


Fig. S1 | overlap of biclusters. This schematic diagram illustrate how biclustering with overlap can uncover the complex patterns in single-cell data. In the complex cell machinery, multiple FGMs are active in a cell group, which is represented by the figure that FGM1 and FGM2 are simultaneously active in cell type 3. Meanwhile, different cell groups may share a common FGM. This can be seen from the fact that cell type 1 and cell type 2 share a common FGM3, since FGM3 is contained in both FGM1 and FGM2.

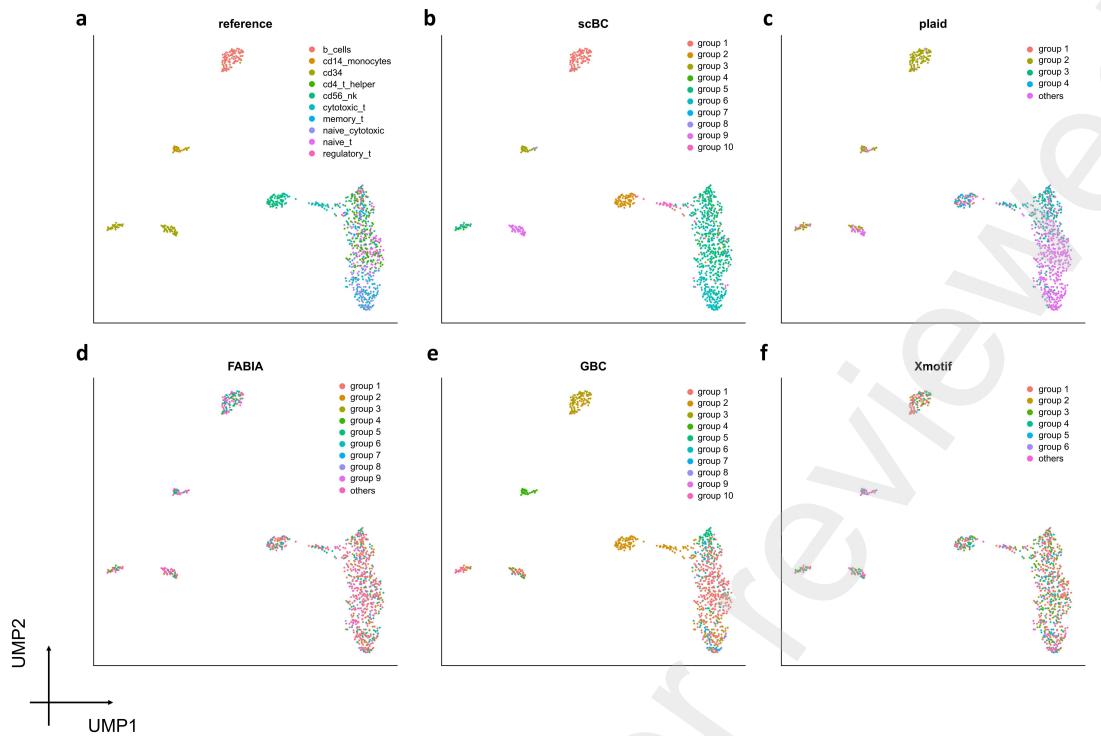


Fig. S2 | visualization of cell clustering results on PBMC dataset. These are the cell representation of the subsample from the last iteration of PBMC dataset. Some methods output too many categories so that we merge some into "others" whose number of cells is less than one percent of the total sample. Methods not shown here fail to detect functional cell groups.

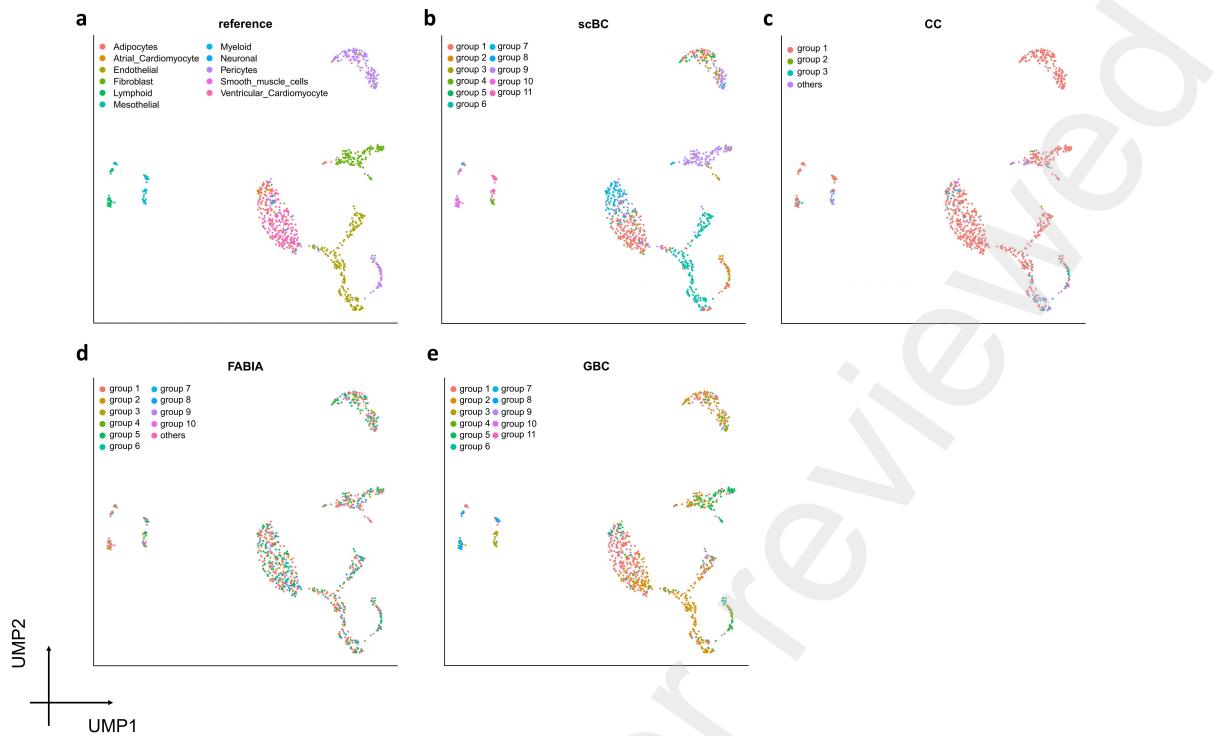


Fig. S3 | visualization of cell clustering results on HEART dataset. These are the cell representation of the subsample from the last iteration of HEART dataset. Some methods output too many categories so that we merge some into "others" whose number of cells is less than one percent of the total sample. Methods not shown here fail to detect functional cell groups.

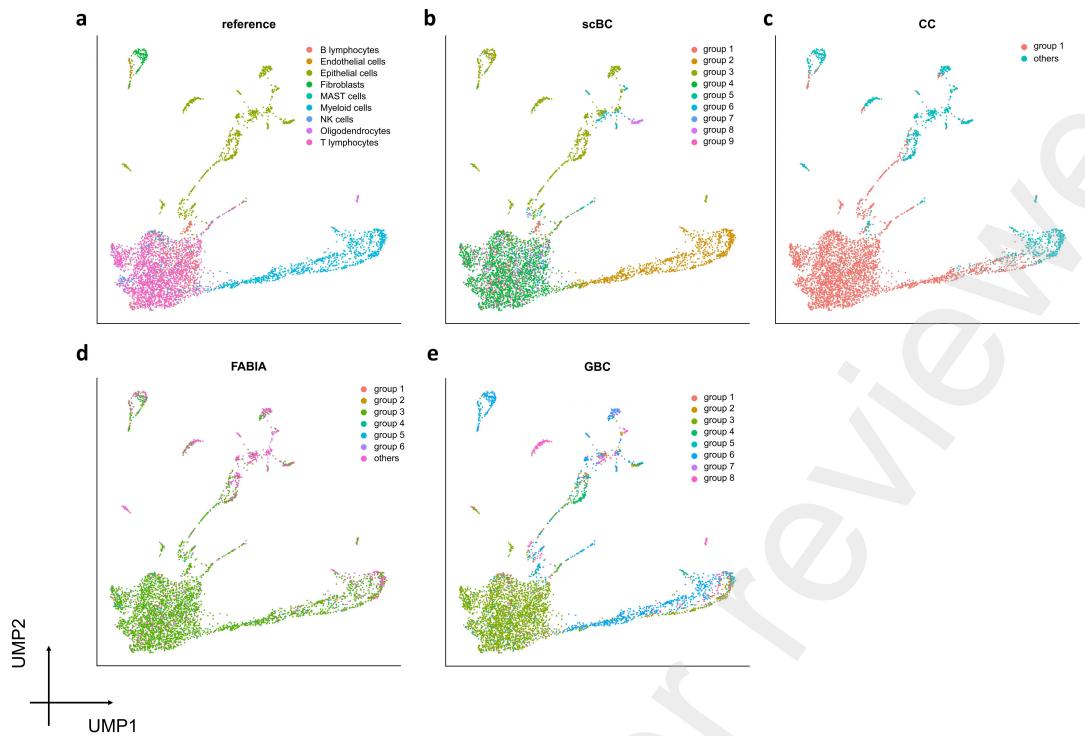


Fig. S4 | visualization of cell clustering results on LUAD dataset. These are the cell representation of the subsample from the last iteration of LUAD dataset. Some methods output too many categories so that we merge some into "others" whose number of cells is less than one percent of the total sample. Methods not shown here fail to detect functional cell groups.

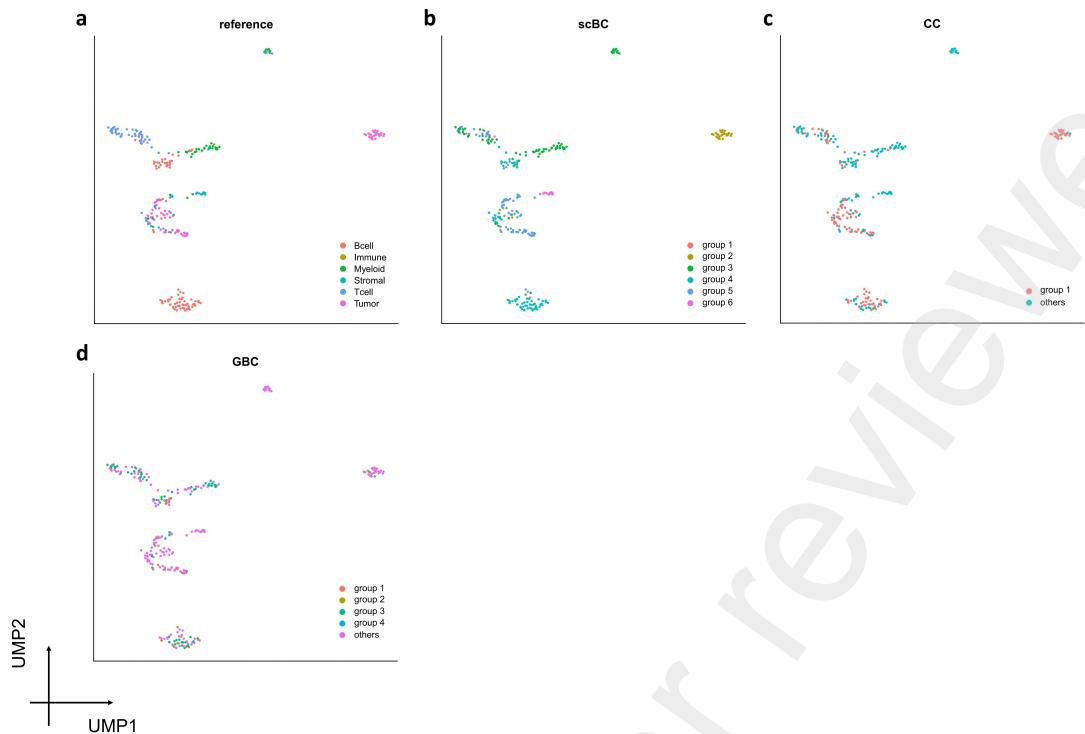


Fig. S5 | visualization of cell clustering results on BC dataset. These are the cell representation of the subsample from the last iteration of BC dataset. Some methods output too many categories so that we merge some into "others" whose number of cells is less than one percent of the total sample. Methods not shown here fail to detect functional cell groups.

Supplementary Tables

Table S1 shows the enrichment scores obtained by comparing different methods through repeated experiments (10 repetitions) on different datasets, corresponding to the performance of different methods in discovering functional gene modules (as shown in Fig. 2b of the main text).

Table S1 Enrichment score comparison of different methods

datasets \ methods	PBMC		HEART		LUAD		BC	
	mean	sd	mean	sd	mean	sd	mean	sd
GBC	23.75	3.49	24.54	4.70	42.20	3.63	34.31	2.97
plaid	9.19	2.63	0.48	1.53	0.74	1.57	0.00	0.00
Xmotif	22.59	1.70	0.00	0.00	0.00	0.00	0.00	0.00
CC	26.86	0.00	23.70	1.55	43.65	1.02	7.46	2.84
Bimax	5.43	3.266	0.00	0.00	0.00	0.00	0.00	0.00
FABIA	21.83	2.72	11.70	2.44	22.37	2.05	0.00	0.00
scBC	29.10	2.42	29.68	2.70	48.28	7.92	40.21	0.96

Tables S2 to S4 show the results of ARI, FMI, and AMI obtained by repeating experiments (10 repetitions) of different methods on different datasets, corresponding to the performance of different methods in cell clustering (as shown in Fig. 2c-d of the main text).

Table S2 ARI of different methods in different datasets

datasets \ methods	PBMC		HEART		LUAD		BC	
	mean	sd	mean	sd	mean	sd	mean	sd
GBC	0.16	0.03	0.23	0.05	0.26	0.05	0.03	0.01
plaid	0.16	0.04	0.01	0.02	0.01	0.04	0.00	0.00
Xmotif	0.04	0.02	0.00	0.00	0.00	0.00	0.00	0.00
CC	0.00	0.00	0.06	0.03	0.18	0.06	0.13	0.02
Bimax	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
FABIA	0.04	0.01	0.01	0.01	0.05	0.01	0.00	0.00
scBC	0.30	0.03	0.36	0.07	0.46	0.09	0.19	0.09

Table S3 FMI of different methods in different datasets

datasets		PBMC		HEART		LUAD		BC	
methods		mean	sd	mean	sd	mean	sd	mean	sd
GBC		0.33	0.04	0.37	0.04	0.48	0.05	0.33	0.02
plaid		0.28	0.03	0.04	0.11	0.07	0.17	0.00	0.00
Xmotif		0.25	0.04	0.00	0.00	0.00	0.00	0.00	0.00
CC		0.00	0.00	0.30	0.07	0.46	0.06	0.32	0.02
Bimax		0.01	0.01	0.00	0.00	0.00	0.00	0.00	0.00
FABIA		0.18	0.04	0.23	0.04	0.27	0.02	0.00	0.00
scBC		0.41	0.03	0.47	0.05	0.61	0.07	0.40	0.07

Table S4 AMI of different methods in different datasets

datasets		PBMC		HEART		LUAD		BC	
methods		mean	sd	mean	sd	mean	sd	mean	sd
GBC		0.31	0.06	0.31	0.04	0.26	0.03	0.10	0.02
plaid		0.18	0.04	0.01	0.02	0.04	0.13	0.00	0.00
Xmotif		0.05	0.02	0.00	0.00	0.00	0.00	0.00	0.00
CC		0.00	0.00	0.06	0.03	0.12	0.05	0.08	0.01
Bimax		0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
FABIA		0.06	0.02	0.02	0.01	0.07	0.01	0.00	0.00
scBC		0.47	0.03	0.47	0.05	0.53	0.08	0.25	0.07

Tables S5 to S10 show the results of 1-CE and F scores obtained by repeating experiments (100 repetitions for each setting) of different methods on simulated datasets under different settings. The values in parentheses represent standard deviations. They are corresponding to the comparison of the overall performance of biclustering (as shown in Fig. 3b-g of the main text).

Table S5 1-CE of different methods in different datasets ($p=1000$, $n=300$, $L=3$, $\times 10^{-3}$)

dropout methods \	0.5	0.55	0.6	0.65	0.7
GBC	6.51(1.09)	5.71(0.80)	5.86(0.77)	6.00(1.01)	6.61(1.11)
plaid	1.22(3.09)	0.68(2.11)	1.12(3.17)	0.92(2.49)	1.32(3.51)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.14(0.18)	0.23(0.24)	0.30(0.32)	0.48(0.40)	0.57(0.45)
Bimax	1.62(0.82)	1.63(0.88)	1.62(0.67)	1.40(0.88)	1.43(0.85)
FABIA	6.95(1.02)	7.34(1.17)	7.25(1.25)	6.77(1.40)	5.81(1.13)
scBC	13.98(1.78)	13.44(1.84)	12.24(2.15)	10.41(2.01)	9.84(1.50)

Table S6 1-CE of different methods in different datasets ($p=3000$, $n=600$, $L=4$, $\times 10^{-3}$)

dropout methods \	0.5	0.55	0.6	0.65	0.7
GBC	3.42(0.43)	2.98(0.36)	2.97(0.34)	2.95(0.33)	2.97(0.31)
plaid	0.70(2.05)	0.65(1.65)	0.83(1.73)	0.68(1.66)	0.69(1.57)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.05(0.05)	0.08(0.08)	0.11(0.10)	0.16(0.11)	0.22(0.15)
Bimax	0.55(0.23)	0.52(0.23)	0.47(0.23)	0.50(0.23)	0.45(0.25)
FABIA	4.21(0.42)	5.41(0.78)	5.33(0.65)	5.13(0.70)	4.52(0.83)
scBC	9.45(0.95)	9.33(0.92)	8.67(1.20)	7.75(1.23)	6.36(1.13)

Table S7 1-CE of different methods in different datasets ($p=6000$, $n=1500$, $L=5$, $\times 10^{-3}$)

dropout methods \	0.5	0.55	0.6	0.65	0.7
GBC	1.45(0.15)	1.22(0.12)	1.22(0.12)	1.23(0.13)	1.21(0.12)
plaid	0.69(1.49)	0.47(1.02)	0.54(1.13)	0.59(1.43)	0.16(0.55)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.01(0.02)	0.01(0.02)	0.01(0.01)	0.03(0.03)	0.06(0.04)
Bimax	0.26(0.11)	0.27(0.10)	0.24(0.11)	0.23(0.11)	0.22(0.11)
FABIA	2.39(0.26)	4.07(2.67)	3.31(1.40)	3.97(1.91)	3.40(1.79)
scBC	4.82(0.39)	4.80(0.38)	4.79(0.49)	4.70(0.47)	4.47(0.42)

Table S8 F score of different methods in different datasets ($p=1000$, $n=300$, $L=3$, $\times 10^{-2}$)

dropout methods	0.5	0.55	0.6	0.65	0.7
GBC	1.31(0.39)	0.50(0.23)	0.63(0.58)	0.88(0.78)	1.03(0.65)
plaid	0.23(0.57)	0.13(0.38)	0.20(0.59)	0.17(0.47)	0.22(0.56)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.03(0.04)	0.05(0.05)	0.06(0.07)	0.10(0.09)	0.12(0.09)
Bimax	0.32(0.18)	0.33(0.18)	0.33(0.14)	0.28(0.18)	0.28(0.17)
FABIA	2.62(0.82)	2.07(0.57)	1.76(0.48)	1.44(0.49)	1.02(0.25)
scBC	5.06(4.83)	3.72(3.35)	3.32(3.41)	2.34(2.72)	1.54(0.54)

Table S9 F score of different methods in different datasets (p=3000, n=600, L=4, $\times 10^{-2}$)

dropout methods	0.5	0.55	0.6	0.65	0.7
GBC	0.49(0.38)	0.16(0.03)	0.18(0.11)	0.17(0.07)	0.37(0.46)
plaid	0.13(0.40)	0.12(0.30)	0.15(0.31)	0.12(0.29)	0.12(0.26)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.01(0.01)	0.02(0.02)	0.02(0.02)	0.03(0.02)	0.04(0.03)
Bimax	0.11(0.05)	0.10(0.05)	0.09(0.05)	0.10(0.05)	0.09(0.05)
FABIA	2.66(0.55)	3.04(0.74)	2.49(0.58)	2.02(0.50)	1.51(0.44)
scBC	9.84(1.75)	8.41(3.12)	5.03(3.82)	3.03(3.06)	1.10(1.01)

Table S10 F score of different methods in different datasets (p=6000, n=1500, L=5, $\times 10^{-2}$)

dropout methods	0.5	0.55	0.6	0.65	0.7
GBC	0.13(0.08)	0.05(0.04)	0.05(0.01)	0.06(0.04)	0.09(0.11)
plaid	0.13(0.29)	0.09(0.19)	0.09(0.18)	0.11(0.26)	0.03(0.10)
Xmotif	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.00(0.00)
CC	0.00(0.00)	0.00(0.00)	0.00(0.00)	0.01(0.01)	0.01(0.01)
Bimax	0.05(0.02)	0.05(0.02)	0.05(0.02)	0.05(0.02)	0.04(0.02)
FABIA	2.72(0.51)	2.53(1.44)	2.15(1.24)	1.99(1.21)	1.46(0.89)
scBC	2.91(0.65)	2.82(0.63)	2.71(0.70)	2.54(0.72)	1.74(0.78)