

1 **QUBIC2: A novel biclustering algorithm for large-scale bulk RNA-sequencing and single-**
2 **cell RNA-sequencing data analysis**

3
4 Juan Xie^{1,2}, Anjun Ma^{1,2}, Yu Zhang³, Bingqiang Liu⁴, Changlin Wan⁵, Sha Cao⁵, Chi Zhang^{5,*},
5 Qin Ma^{1,2,*}

6 ¹Bioinformatics and Mathematical Biosciences Lab, Department of Agronomy, Horticulture
7 and Plant Science, South Dakota State University, BioSNTR, Brookings, SD, 57007, USA

8 ²Department of Mathematics and Statistics, South Dakota State University, Brookings, SD, 57006,
9 USA,

10 ³Colleges of Computer Science and Technology, Jilin University, Changchun, 130012, China,

11 ⁴School of Mathematics, Shandong University, Jinan, 250100, China, and

12 ⁵Center for Computational Biology and Bioinformatics, Department of Medical and Molecular
13 Genetics, Indiana University, School of Medicine, Indianapolis, IN, 46202, USA.

14 *To whom correspondence should be addressed. Tel: +1 605-688-6315; Email:
15 qin.ma@sdstate.edu. Correspondence is also addressed to Chi Zhang. Tel: +1 317-278-9625;
16 Email: c Zhang87@iu.edu.

17

18 **ABSTRACT**

19 The combination of biclustering and large-scale gene expression data holds a promising potential
20 for inference of the condition specific functional pathways/networks. However, existing
21 biclustering tools do not have satisfied performance on high-resolution RNA-sequencing (RNA-
22 Seq) data, majorly due to the lack of (i) a consideration of high sparsity of RNA-Seq data, e.g.,
23 the massive zeros or lowly expressed genes in the data, especially for single-cell RNA-Seq
24 (scRNA-Seq) data, and (ii) an understanding of the underlying transcriptional regulation signals
25 of the observed gene expression values. Here we presented a novel biclustering algorithm namely
26 QUBIC2, for the analysis of large-scale bulk RNA-Seq and scRNA-Seq data. Key novelties of the
27 algorithm include (i) used a truncated model to handle the unreliable quantification of genes with
28 low or moderate expression, (ii) adopted the mixture Gaussian distribution and an information-
29 divergency objective function to capture shared transcriptional regulation signals among a set of
30 genes, (iii) utilized a Core-Dual strategy to identify biclusters and optimize relevant parameters,
31 and (iv) developed a size-based *P*-value framework to evaluate the statistical significances of all
32 the identified biclusters. Our method validation on comprehensive data sets of bulk and single cell
33 RNA-seq data suggests that QUBIC2 had superior performance in functional modules detection
34 and cell type classification compared with the other five widely-used biclustering tools. In addition,

35 the applications of temporal and spatial data demonstrated that QUBIC2 can derive meaningful
36 biological information from scRNA-Seq data. The source code for QUBIC2 can be freely accessed
37 at <https://github.com/maqin2001/qubic2>.

38

39 INTRODUCTION

40 As next-generation sequencing technologies have become more affordable in these years (1,2),
41 it is possible to generate large-scale biological data with higher resolution, better accuracy, and
42 lower technical variation than the array-based counterparts (3,4). RNA-Seq measures the
43 abundance of RNA transcripts, giving rise to genome-scale gene expression data in a biological
44 sample at a given moment (5). Nowadays, researchers can isolate individual cells from complex
45 organisms and measure transcriptional activity using single-cell sequencing. Hundreds of RNA-
46 seq data sets with more than hundreds of sample were emerged in the public domain in the past
47 five years, and their tremendous values have been confirmed in many research areas, e.g.,
48 elucidation of cell-type-specific gene regulatory networks (6) and cancer & complex diseases (7-
49 9).

50

51 The abundance of gene expression datasets provides an opportunity to computationally identify
52 condition based functional gene modules (FGMs), each of which is defined by a similar expression
53 patterns over a certain gene set, which tend to be functionally related or co-regulated by the same
54 transcriptional regulatory signals (TRSs) under a specific condition. Thus, successfully derivation
55 of the FGMs may grant a higher-level interpretation of gene expression data, improve functional
56 annotation of genes, facilitate inference of gene regulatory relationships, and provide a better
57 mechanism level understanding of diseases such as cancer. The identification of FGMs can be
58 naturally modeled as a specific data pattern over unknown subset of genes and samples, and
59 solved with a bi-clustering approach (10), a two-dimensional data mining technique that can
60 simultaneously identify co-expressed genes under a subset of conditions (i.e., samples or cells).
61 This unique feature makes it more useful than clustering when applied to large-scale gene
62 expression data, as genes are usually co-expressed under certain instead of all conditions.

63

64 Besides the identification of FGMs in bulk tissue data, a similar formulation may also be applied
65 to scRNA-Seq data, to identify individual cells or cell types as well as their complex interactions
66 under specific stimuli, e.g., cell types classification and clustering. In multicellular organisms,
67 biological function emerges when heterogeneous cell types form complex organs (11).
68 Investigations into organ development, cell function, and disease microenvironment highly

69 depend upon an accurate identification and categorization of cell types, sometimes along with
70 their temporal and spatial features (12). Traditionally, a cell type was predicted based on
71 morphological properties or marker proteins, yet this method failed to characterize the full diversity
72 of cells. scRNA-Seq data provides the possibility to group cells based on their genome-wide
73 transcriptome profiles, and several studies have already been carried out using scRNA-Seq data
74 to identify novel cell types, proving its power to unravel the full diversity of cells in human and
75 mouse (13). Mathematically, the problem of scRNA-seq based cell types classification can be
76 naturally formulated as biclustering problems, since the essence is to find sub-populations of cells
77 sharing common expression patterns among subsets of genes.

78
79 Substantial efforts have been made in biclustering algorithm and tool development since 2000
80 (14-26), and a few review studies have provided considerable guidance in choosing suitable
81 algorithms in different contexts (27-29): Eren *et al.* compared 12 algorithms and concluded that
82 our previously developed method, QUBIC, is one of the top performed methods, as it has
83 achieved the highest performance in synthetic datasets and captured a high proportion of
84 enriched biclusters on real datasets, comparing to Plaid, FABIA, ISA and Bimax, which were also
85 recommended for capturing upregulated biclusters (27). In 2018, Saelens *et al.* ranked ISA,
86 FABIA and QUBIC as the top biclustering methods in terms of predicting gene modules from
87 human and/or synthetic data (30).

88
89 Although numerous bi-clustering methods have been developed for gene expression data
90 analysis, the most existing algorithm are designed and evaluated using microarray rather than
91 RNA-Seq data. One of the unique features of gene expression data derived from RNA-Seq,
92 especially the scRNA-Seq data, is the massive zeros (up to 60% of all the genes in a cell have
93 read counts being zeros) (31,32). The normalized read counts roughly follow lognormal
94 distributions; however, the raw zero counts of specific genes will lead to negative infinity after
95 logarithmic transformation (33-36), resulting in unquantifiable errors. Therefore, the biclustering
96 methods that are successful for microarray cannot be directly applicable to RNA-Seq data (37),
97 and novel methods taking full consideration of characteristics of RNA-Seq data are urgently
98 needed in the public domain.

99
100 In this paper, we developed a novel bi-clustering algorithm, namely QUBIC2, for large-scale RNA-
101 seq data analysis. We demonstrated the performance of QUBIC2 on capturing FGMs by applying
102 it to four datasets and benchmarking against five widely used biclustering algorithms. QUBIC2

103 turned out to be a superior player as it has identified a significant higher proportion of enriched
104 and diverse FGMs. Besides, QUBIC2 can also identify cell types with a higher accuracy,
105 comparing to the five biclustering tools and SC3, a state-of-the-art cell clustering method.
106 Furthermore, we also illustrated the application power of QUBIC2 on inferring time- and spatial-
107 related insights from two temporal and two spatial scRNA-seq datasets.

108

109 **RESULTS**

110 ***Overall design of QUBIC2***

111 Inheriting the qualitative representation and graph-theory based model from QUBIC (19), QUBIC2
112 has four unique features: (i) developed a rigorous truncated model to handle the unquantifiable
113 errors caused by zeros, and used a reliable qualitative representation of gene expression to reflect
114 expression states corresponding to various TRSs; (ii) integrated an information-divergence
115 objective function in the biclustering framework in support of functional gene modules
116 identification; (iii) employed a Core-Dual strategy to optimize consistency level of a to-be-identified
117 bicluster; and (iv) developed a robust P -value framework to support statistical evaluation of all the
118 identified biclusters. Details of these four features are showcased as follows (Figure 1).

119

120 A mixture of left-truncated Gaussian distributions (LTMG) model was designed to fit the RNA-Seq
121 data, rather than discarding zeros or adding a small constant to original counts (34,38). The basic
122 idea is to treat the large number of observed zeros and low expressions as left censored data in
123 the mixture Gaussian model of each gene (39,40), assuming that the observed frequency of
124 expressions on the left of the censoring point should be equal to the area of the cumulative
125 distribution function of the mixture Gaussian distribution left of the censoring point. Furthermore,
126 we assumed that a gene should receive K possible TRSs under all the conditions, and its
127 expression profile would follow a mixture of K left truncated Gaussian distributions. The LTMG
128 model was applied to fit the expression value of each gene and the gene expression value under
129 a specific condition was labeled to the most likely distribution. Accordingly, a row consisting of
130 discrete values (1,2, ..., K) for each gene was generated (Figure 1A). Then this qualitative row
131 was split into K new rows, such that in the i^{th} row those labeled initially as i are labeled as 1, while
132 the rest were labeled as 0. Finally, a binary representing matrix M_R was generated.

133

134 A weighted graph $G = (V, E)$ was constructed based on M_R , where nodes V correspond to genes,
135 edges E connecting every pair of genes (Figure 1B). The edge weight indicates the similarity
136 between the two corresponding genes, which is defined as the number of conditions in which the

137 two genes have 1s in M_R . Intuitively, two genes from a bicluster should have a heavy edge in G
138 innately while two random genes may have a heavy edge only accidentally. Hence, a bicluster
139 should correspond to a maximal subgraph of G , with edges typically heavier than the edges of an
140 arbitrary subgraph. Identifying all the biclusters equals to identifying all the heavy subgraphs in G ,
141 which is an NP-hard problem. Therefore, a heuristic strategy was designed as follows.

142

143 The algorithm would iterate a seed list (S), which is the sorted list of edges in G in the decreasing
144 order of their weights (i.e., $w(e_1) \geq w(e_2) \geq \dots, w(e_{|E|})$). An edge $e_{ij} = g_i g_j$ is selected as a seed
145 if and only if at least one of g_i and g_j is not in any previously identified biclusters, or g_i and g_j are
146 in two nonintersecting biclusters in terms of genes. QUBIC2 first built a core bicluster from a seed
147 and then expanded to recruit more genes and conditions into a to-be-identified bicluster, until the
148 Kullback-Leibler divergence score (KL score) was locally optimized. It was proposed based on
149 the assumption that the difference between a bicluster and its background should be larger than
150 the difference between an arbitrary same-size submatrix and its background. The KL score of a
151 bicluster was designed to quantify this difference as the larger of the difference was, the larger of
152 the score is (Figure 1C. See **Methods** for details).

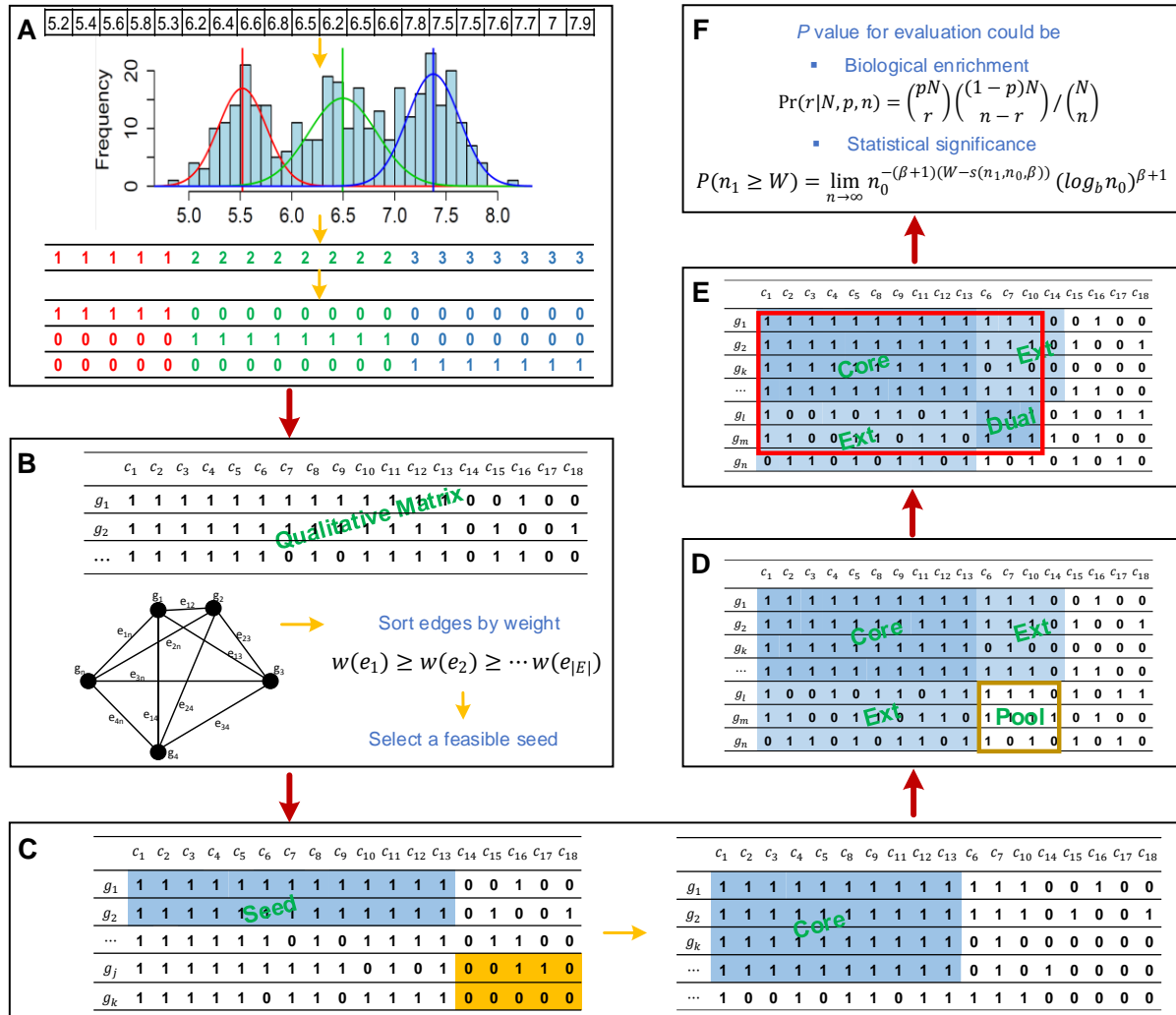
153

154 During bicluster expansion, the algorithm controlled the consistency level for a bicluster, which is
155 defined as the minimum ratio of the number of 1s in a column/row and the number of
156 rows/columns in the bicluster. In QUBIC, a pre-specified value c ($0 < c \leq 1.0$) was used to control
157 the overall consistency level of the bicluster. While this parameter was dynamically optimized by
158 a Dual searching method in QUBIC2 (Figure 1D-E), giving rise to a submatrix (I, J) of M_R (i.e., a
159 bicluster) with optimized consistency level and maximal KL score can be identified. Biclusters
160 expanded using Dual strategy tend to be more significant than those without Dual (See Example
161 S1 in Supplementary File 1).

162

163 Furthermore, for the first time, a statistical framework based on the size of the biclusters was
164 implemented to calculate a P -value for each of the identified biclusters. The problem of assessing
165 the significance of identified biclusters was formulated as calculating the probability of finding at
166 least one submatrix enriched by 1 from a binary matrix with given size, with a beta distribution
167 employed during the process. This P -value framework enables users systematically evaluate the
168 statistical significance of all the identified biclusters, especially for those from less-annotated
169 organisms (Figure 1F).

170



171
 172 **Figure 1.** QUBIC2 workflow. **A.** Discretization of gene expression data from RNA-Seq. The LTMG model
 173 will be applied to fit each gene's expression profile. A representing row for each gene will be generated with
 174 integers denote the most likely component distribution that each value belongs to. Then this representing
 175 row will be split into multiple rows. Finally, a binary representing matrix will be generated; **B.** Graph
 176 construction and seed selection. A weighted group will be constructed based on the representing matrix
 177 from A. By sorting the edges in decreasing order of their weight, and an initial seed list will be obtained.
 178 QUBIC2 will select a feasible seed from the list; **C.** Build an initial core based on the selected seed. During
 179 seed expand, QUBIC2 will search for genes with higher weight with the seed. In case of two genes have
 180 the same weight, the one with higher KL score will be selected. Thus, gene k (KL=0.1914) instead of gene
 181 j (KL=0.0622) will be added to the core first; **D.** Expand core and determine pool. QUBIC2 will expand the
 182 core vertically and horizontally to recruit more genes and conditions under a preset consistency level,
 183 respectively. The intersected zone created by extended genes and conditions as a Dual searching pool
 184 (brown box); **E.** Dual search in the pool and output the bicluster with genes and conditions that come from

185 Core and Dual as final bicluster (red box); **F**. Statistical evaluation of identified biclusters based on either
186 biological annotations or the size of the bicluster.

187

188 **Functional gene modules detection from RNA-Seq data**

189 Compared with five biclustering algorithms (Bimax(18), ISA(41), FABIA(20), Plaid(15), and
190 QUBIC(19), with more details in Table S1 of Supplementary File 1), the performance of QUBIC2
191 in identifying FGMs was systematically evaluated using four gene expression datasets: a
192 simulated RNA-Seq dataset based on an in-house method (22,846 rows \times 100 columns), a bulk
193 RNA-Seq dataset from *Escherichia coli* (*E. coli*, 4,497 rows \times 155 columns), a bulk RNA-Seq
194 dataset from TCGA (3,084 rows \times 8,555 columns), and a scRNA-Seq dataset from human
195 embryos (3,798 genes \times 90 cells). For the identified biclusters from a specific tool, *precision*
196 showcases the fraction of biclusters whose genes are significantly enriched with certain biological
197 pathways (i.e., relevance), and *recall* reflects the fraction of captured known modules/pathways
198 among all known modules in a functional annotation database, e.g., KEGG (42) and RegulonDB
199 (43) (i.e., diversity). The harmonic mean value of precision and recall, referred to as the *F* score,
200 was used as the integrated criteria in performance evaluation.

201

202 Evaluation studies usually used default parameters of the to-be-analyzed tools, which were
203 optimized for specific benchmark datasets. However, when applied to datasets coming from a
204 different organism (e.g., *E. coli* vs. human), or be acquired by other technologies (e.g., microarray
205 vs. RNA-Seq), the default parameters often fail to achieve satisfying performance and need
206 further optimization/adjustment. To minimize the biases in performance comparison among
207 multiple tools, for each of the four datasets, we run the six tools under more than 50 parameter
208 combinations by adjusting their critical parameters around default/recommended values (see
209 **Methods** and Table S2 in Supplementary File 1). Then the *F* score of identified biclusters under
210 each parameter combination was calculated. In this way, we can test a tool's robustness and infer
211 how sensitive of its performance is to parameter adjustment, besides the basic performance
212 comparison among different tools.

213

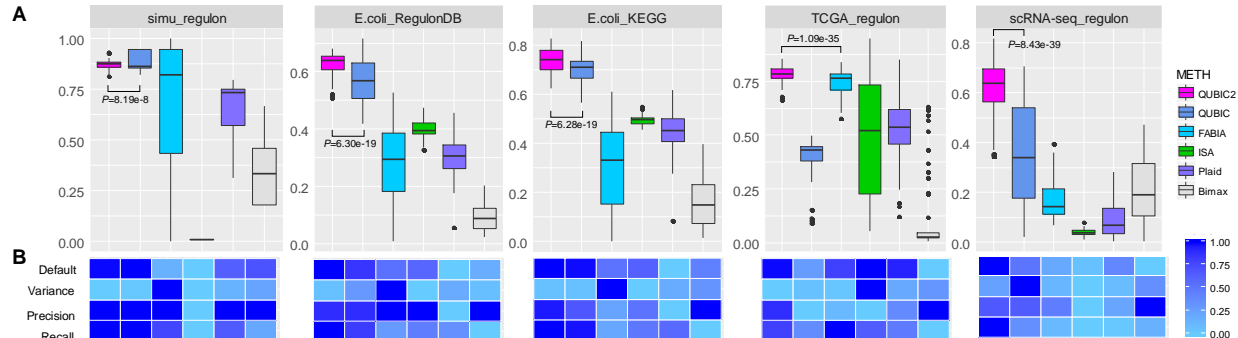
214 As showcased in Figure 2, QUBIC2 achieved the highest median *F* scores and the highest *F*
215 scores with the default parameter on all the four datasets, and its *F* scores were significantly
216 higher than the second-best algorithms in all the comparison circumstances (Wilcoxon test *P*-
217 value <0.01). QUBIC2 performed well in both precision and recall, indicating that the identified
218 FGMs are relevant and diverse; and it had relatively small variance, while the performance of

219 some algorithms on certain dataset was very sensitive to parameter change (e.g., FABIA on *E.*
 220 *coli*). Regarding median *F* scores, QUBIC was the second-best algorithm on simulated data, *E.*
 221 *coli* RNA-Seq data, and human scRNA-Seq data, while FABIA was the second-best one for TCGA
 222 data. As regards the default settings, QUBIC ranked as the top ones on simulated data and *E.*
 223 *coli* data, and ISA and Plaid had relative higher rank on TCGA data. ISA was generally very stable,
 224 and its variances were the smallest on three datasets. As for Bimax, although its recall was
 225 relatively low, it was characterized with high precision on the four datasets. It is noteworthy that
 226 QUBIC2 is the only program, among all the six biclustering algorithms, which did not encounter a
 227 dramatic performance drop on scRNA-Seq data compared to RNA-Seq data, suggesting the
 228 unique applicative power of QUBIC2 on FGMs detection from scRNA-Seq data (Figure S1 in
 229 Supplementary File1).

230

231 Furthermore, the performance of all the biclustering algorithms on *E. coli* data was better than on
 232 human data, with the possible reason that *E. coli* data has more completed functional annotation
 233 and affects the evaluation of module significance. Therefore, for less annotated organisms, we
 234 need a statistical evaluation framework for all the identified biclusters.

235



236

237 **Figure 2.** Overall performance comparison between QUBIC2 and five popular biclustering methods based
 238 on the agreement between identified biclusters and known modules. **A.** Distribution of *F* scores on each of
 239 the four datasets under multiple runs ($n > 40$). Black line in the box denote median value, whiskers denote
 240 10% and 90% percentiles, while the box denotes 25% and 75% percentiles; **B.** relative performance of six
 241 algorithms in terms of *F* score under default parameters, variance of *F* scores under multiple sets of
 242 parameters, median value for the precision and median value for the recall, respectively (normalized over
 243 six algorithms). Note that the variance of *F* scores depends on the increment of parameters, and therefore
 244 only indicative.

245

246

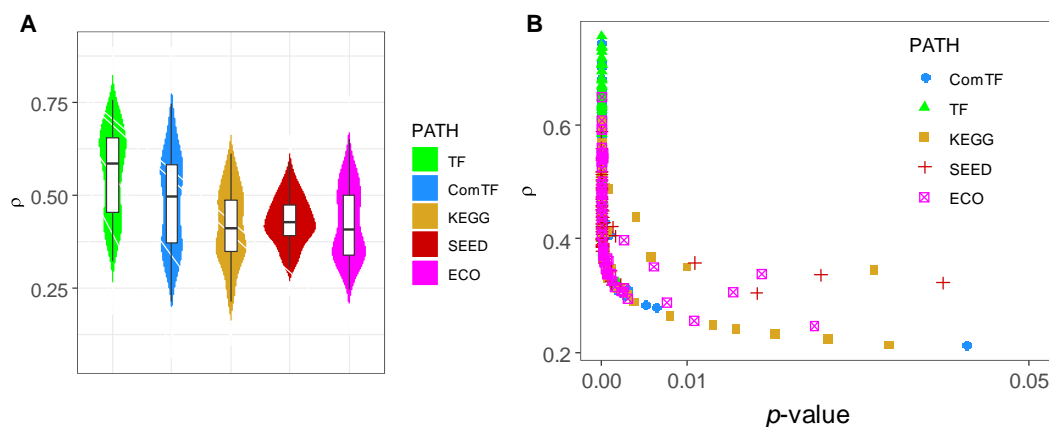
247 **A statistical evaluation framework for identified biclusters**

248 The significances of gene modules from the identified biclusters were usually evaluated by
249 pathway enrichment analysis. However, many organisms (including human) have limited
250 functional annotations supported by experimentally verifications, which makes a systematic
251 evaluation of all identified biclusters non-trivial. To fill this gap, a statistical method was proposed
252 in this study, which can calculate a P -value for a bicluster purely based on their size (number of
253 genes and conditions).

254

255 Interestingly, we found that there is a strong association between the P -values of biclusters
256 calculated via pathway enrichment analysis (named knowledge-based P -value) and the
257 corresponding size-based P -values. Specifically, spearman correlation tests were conducted
258 between size-based P -values and five groups of knowledge-based P -value (see **Methods**). The
259 average spearman correlation coefficients (ρ) were higher than 0.40 (ComTF_ ρ =0.48, TF_ ρ =0.56,
260 KEGG_ ρ =0.42, SEED_ ρ =0.43 and ECO_ ρ =0.42), and the average p -values for the correlation
261 test were smaller than 0.01. As showcased in Figure 3A, all the ρ s in the five groups are positive.
262 In addition, ρ s related with regulatory pathways (i.e., TF_ ρ and ComTF_ ρ) were generally larger
263 than ρ s those related to metabolic pathways (i.e., KEGG_ ρ and SEED_ ρ). This indicated that the
264 size-based P -value seemed to be more suitable for the evaluation of biclusters' regulatory
265 significance. Furthermore, all the corresponding p -values were less than 0.05 (Figure 3B),
266 suggesting that the correlations between knowledge-based P -values and size-based P -values
267 were statistically significant at the 0.05 level. In addition, the parameter f which controls the level
268 of overlaps between biclusters had a negative association with ρ (Figure S2 in Supplementary
269 File1), suggesting that the size-based P -values would have a stronger association with
270 knowledge-based P -values when the overlaps between biclusters are relatively low.

271



272

273 **Figure 3. A.** The distribution of correlation coefficients(ρ) between P -value obtained from enrichment
274 analysis and size-based P -value. We run QUBIC2 under 63 different parameter settings, and ρ was
275 calculated under each run; **B.** Scatter plot of ρ and p -value. The y-axis denotes ρ , the correlation coefficient
276 for the spearman association test, the x-axis denotes the p -value of the association test. Note that to
277 distinguish, italic lowercase p was used to denote the p -value of the Spearman correlation test, while italic
278 uppercase P was used to denote the significance of biclusters.

279

280 ***Cell type classification based on scRNA-Seq data***

281 The above sections demonstrated the outstanding performance of QUBIC2 on FGMs
282 identification and its unique feature of statistical evaluation for all the identified biclusters. In this
283 section, we showed the predictive power of biclustering methods on cell types identification from
284 scRNA-Seq data.

285

286 We developed a pipeline to group cells into different types with the assumption that two cells
287 belonging to the same bicluster have a higher likelihood of being the same cell type than two
288 randomly selected cells (see **Methods**). Briefly, a biclustering tool was first used to identify
289 biclusters from a scRNA-Seq expression data. Then, a weighted graph $G = (C, E)$ was
290 constructed to model the relationship between cell pairs, where nodes C represent cells, edges E
291 connect pairs of cells, and edge weight indicates the number of biclusters that the two
292 corresponding cells appear in simultaneously. Finally, cell types were predicted via the Markov
293 Cluster Algorithm (MCL) clustering on the weighted graph (Figure 4A).

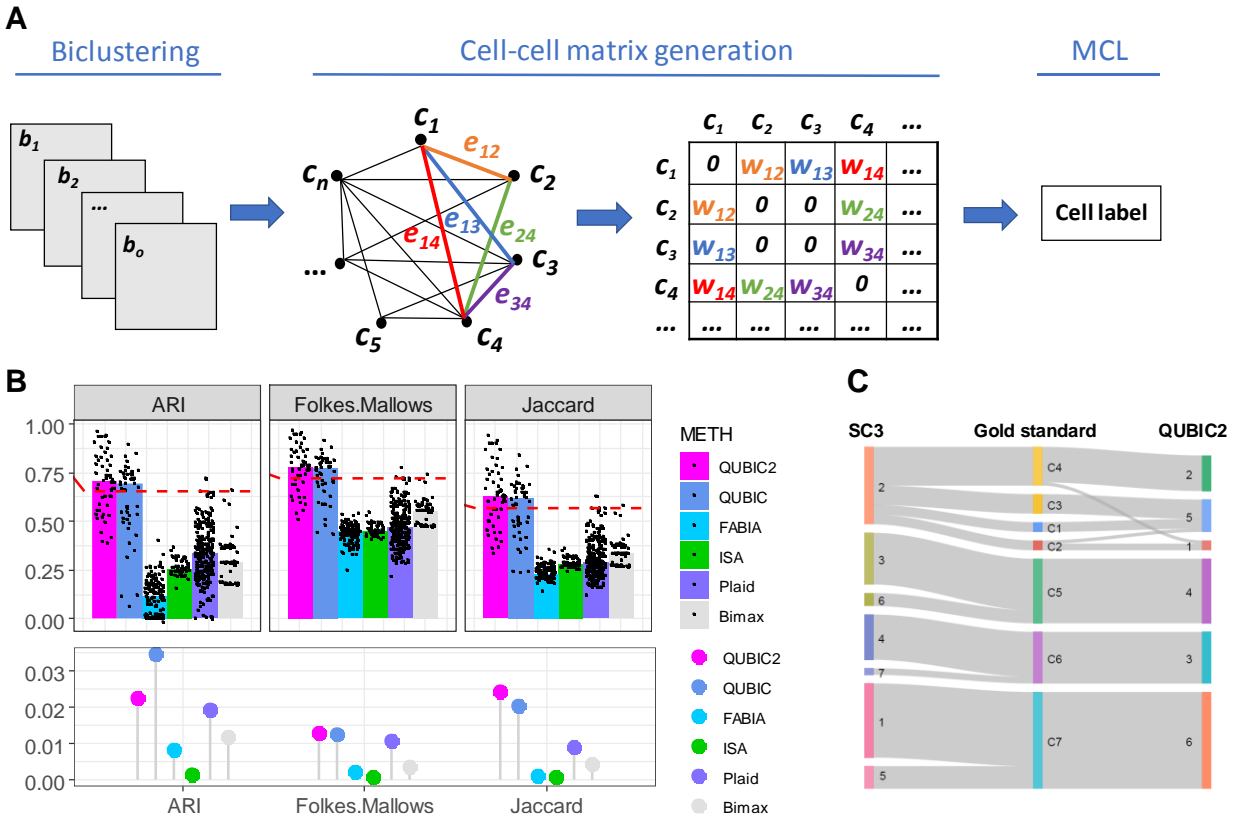
294

295 For each of the six biclustering methods in Figure 2, we applied this pipeline to a benchmark
296 dataset with 20,214 genes and 90 cells (41), which have been experimentally classified into seven
297 types (46). The Adjusted Rand Index (ARI) was adopted as the evaluation criteria to assess the
298 agreement between predicted cell types and these 'ground truth' (46). Two more external
299 validation criteria, namely Jaccard Index (JI) and Fowlkes Mallows Index (FW), were also used
300 here aiming to provide a comprehensive evaluation.

301

302 As Figure 4B showed, QUBIC2 and QUBIC were the top two biclustering tools, respectively, in
303 terms of median values on the three criteria. Both surpassed the performance of SC3 (41), which
304 was used as the benchmark (median value from 100 runs) and was denoted by the red dash line
305 in each panel of Figure 4B. In addition, ISA always demonstrated the smallest variance across
306 the three validation criteria. The FW values of each tool were more stable than other two values.

307 Figure 4C showcased one cell type classification result obtained by QUBIC2 (parameter f was set
 308 to 0.85, c set to 0.85, k set to 13, o set to 2000). The result was in good agreement with the
 309 reference cell labels and QUBIC2 correctly grouped the three major cell types (8_cell_embryo,
 310 Morulae, and late_blastoCyst).
 311



312
 313 **Figure 4. A.** Computational pipeline for cell type classification. This pipeline consists of three steps:
 314 biclustering, generation of weighted cell-cell matrix and clustering using MCL. The input is biclusters and
 315 output is cell type labels; **B.** Benchmark of QUBIC2 against five popular biclustering algorithms. Upper layer:
 316 each panel shows the similarity between the inferred labels and the reference labels quantified by the four
 317 indices, i.e., Adjusted Rand Index (ARI), Folkes Mallows's index and Jaccard Index, respectively. Each
 318 algorithm was applied >40 times to the same dataset to evaluate accuracy and stability. The three indices
 319 were calculated for each run of the respective methods (black dots). Bars represent the median of the
 320 distribution of black dots. The red dash lines correspond to the benchmark performance of SC3 (ARI:
 321 0.6549, FW: 0.7243, JI: 0.5671). Lower layer: the variance of each tool in terms of three validation criteria;
 322 **C.** Sankey diagram comparing the 7 clusters obtained with SC3 (left layer) and 6 clusters obtained with
 323 QUBIC2 (right layer). The middle layer corresponds to the seven reference clusters. The widths of the lines
 324 linking nodes from two layers correspond to the number of cells they have in common.
 325

326 ***QUBIC2 inferred the temporal and spatial organization of cells from scRNA-Seq data***

327 When spatial and temporal information is available, scRNA-Seq can reveal more biological
328 insights beyond cell types. In this section, QUBIC2 was applied on two temporal (and) and two
329 spatial scRNA-Seq datasets, respectively, to explore the temporal and spatial organization of cells.

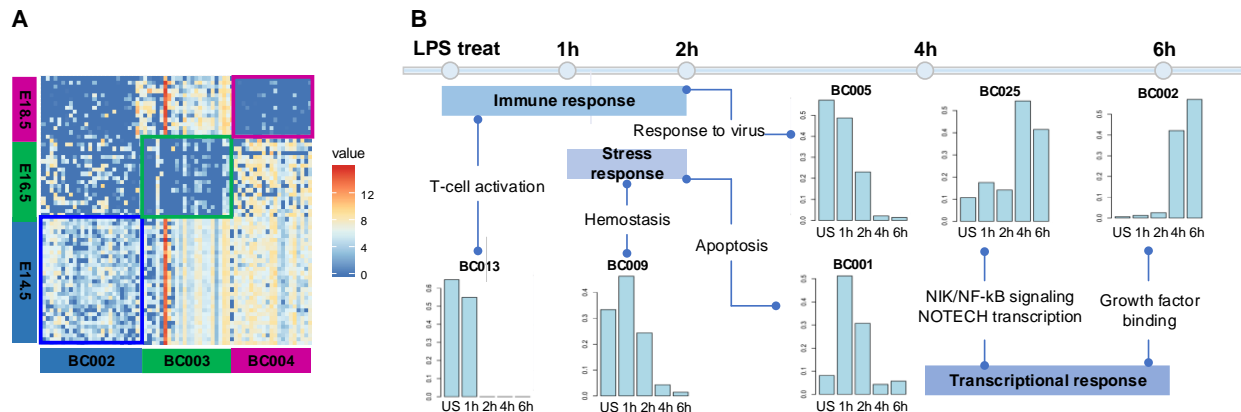
330

331 Five biclusters were identified by QUBIC2 from a time series lung scRNA-Seq data (*GSE52583*),
332 which consists of 152 cells collected at E14, E16 and E18, respectively (47). Three of the five
333 biclusters contain time-specific cells. In particular, bicluster BC002 consists of cells exclusively
334 from E14; bicluster BC003 contains cells that only from E16; and bicluster BC004 has cells coming
335 from E18 (Figure 5A). Functional enrichment analyses of the component genes from these three
336 biclusters were carried out based on DAVID (48) and the results showed that genes in BC002
337 mainly related to cell cycle, cell division, and mitosis; BC003 genes were enriched with ribosome,
338 translation, and structural constituent of ribosome; and spliceosome-related genes were grouped
339 in BC004 (see details in Supplementary File 2).

340

341 In addition to identifying biclusters corresponding to specific time point, QUBIC2 can also be used
342 to find biclusters with time-dependent patterns. Here QUBIC2 was used to analyze a scRNA-Seq
343 data with mouse dendritic cells (DCs) collected at 1h, 2h, 4h and 6h after treatment with
344 pathogenic agent lipopolysaccharide (LPS) and untreated controls (*GSE48968*) (49). In total, 51
345 biclusters were identified in the datasets treated with LPS. For each bicluster, the Fisher exact
346 test was conducted on its constituting samples to assess if significant over-representation by any
347 time points could be found within the bicluster. For those biclusters showing significant association
348 with the time-course, a pathway enrichment analysis was conducted to infer the biological
349 characteristics of the bicluster. In the end, 30 biclusters that are significantly over-represented by
350 one or several consecutive time points were identified in the LPS dataset ($\alpha=0.005$, $P<1e-22$),
351 and six of them showed distinct time dependence (Figure 5B). Specifically, bicluster BC013
352 consists of untreated samples and samples collected at 1h, which represents the earliest
353 response to LPS and enriches multiple immune response pathways. Bicluster BC005 consists
354 largely of untreated samples and samples collected at 1h and 2h, which also is enriched with
355 immune response pathways but with more responses to a virus, T cell chemotaxis and so on.
356 BC009 and BC001 are enriched by samples collected at 1h and 2h, covering a wider range of
357 stress-response pathways, suggesting that the activation of stress response pathways and
358 altered metabolisms as secondary responses after the early immune response. BC025 and
359 BC002 consist of samples collected at 4h and 6h, and their genes enrich pathways associated

360 with alterations in cell morphogenesis, migration, cell-cell junction and so on. Overall these
 361 observations suggest that our analysis can identify all the major responses to the LPS treatment
 362 in a time-dependent manner. Detailed pathways enriched by the six biclusters are given in Figure
 363 S3 in Supplementary File 1. The detailed information of these biclusters is given in Supplementary
 364 File 3.

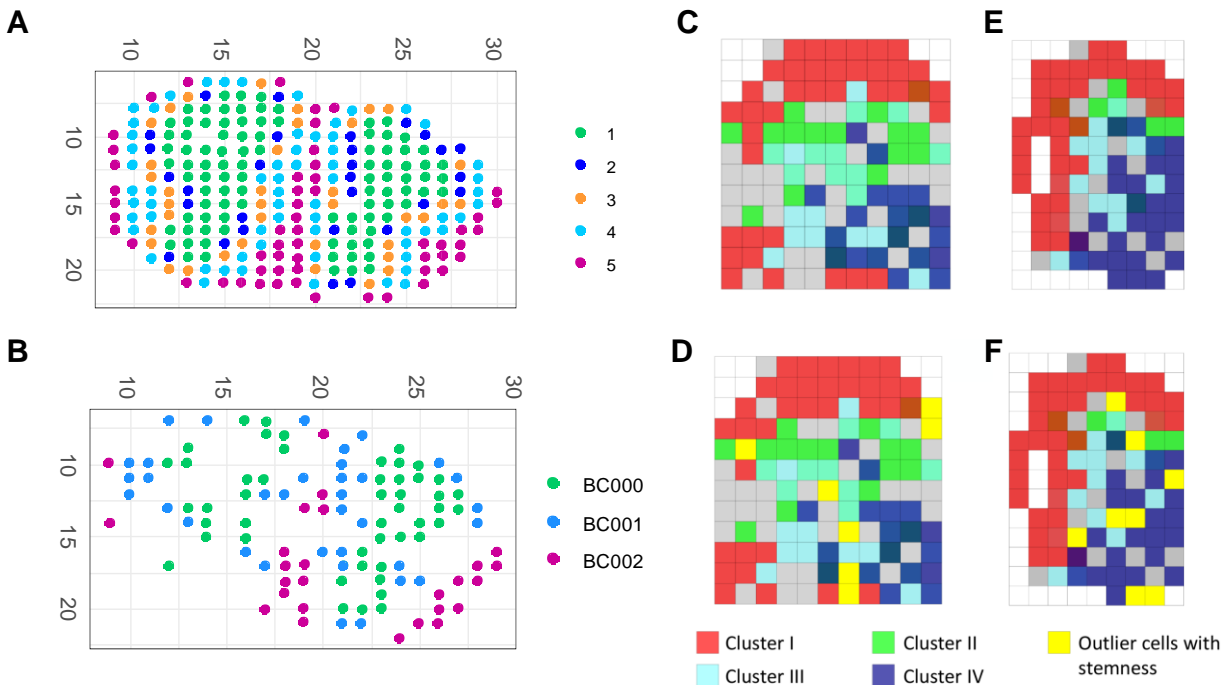


365 **Figure 5. A.** Visualization of three biclusters (BC002, BC003, and BC004) selected based on the specificity
 366 to time point; **B.** Time-dependent distribution of cells in six selected biclusters identified in the LPS data. In
 367 each histogram, the five bars from left to right show the proportion of the untreated samples and samples
 368 collected at 1h, 2h, 4h and 6h after the LPS treatment.

369
 370
 371 Then QUBIC2 was applied to a mouse spatial scRNA-Seq dataset with 280 cells. The cells were
 372 classified into five clusters that correspond to five well-defined morphological layers in (50) (Figure
 373 6A). Five biclusters were predicted. Among them, the bicluster BC000 consists of cells mainly
 374 from the granular layer; the bicluster BC001 contains cells from the mitral layer and glomerular
 375 layer; and the bicluster BC002 contains cells mainly from the olfactory nerve layer (Figure 6B).
 376 Functional annotation showed that BC000 mainly enriches plasma membrane, cell membrane,
 377 and cell projection; BC001 enriches synapse, neuron projection, and cell projection; and BC002
 378 enriches cell projection (Details in Supplementary File 4).

379
 380 Finally, another spatial scRNA-Seq dataset (*GSE60402*) with samples dissected from three
 381 mouse medial ganglionic eminence tissues and known spatial coordinates was analyzed.
 382 QUBIC2 was applied and 37, 40, and 120 biclusters were identified in the mutant, wild-type 1,
 383 and wild-type 2 datasets, respectively (Details in Supplementary File 5). Further investigation on
 384 the spatial distribution of cells in each bicluster showed that all the four spatial biclusters with
 385 distinct expression patterns by cell cycle, cell morphogenesis, and neuron development genes,
 386 as reported in the original study (51), were identified by QUBIC2. It is noteworthy that the outliers

387 with highly expressed stem cell markers tend to be located at the intermediate region between
 388 two adjacent (or overlapping) biclusters in the three datasets as shown in Figure 6D and 6F. Our
 389 interpretation is that these location-dependent expression patterns may be caused by parallel and
 390 independent differentiations from common stem cells.



391
 392 **Figure 6. A.** The coordinates of cells correspond to five morphological layers (1. Granular cell layer; 2.
 393 Mitral cell layer; 3. Outer plexiform layer; 4. Glomerular layer; 5. Olfactory nerve layer); **B.** The coordinates
 394 of cells from three selected biclusters; **C.** The spatial coordinates of samples in the four biclusters identified
 395 in wild-type 1 mouse; Colors red, green, cyan and dark blue represent samples in four different biclusters;
 396 **D.** In addition to the coordinates of bicluster samples, the yellow cubes represent significant outlier samples;
 397 **E.** The same information as in C except the samples are from wild-type 2 mouse; **F.** The same information
 398 as in D except the samples are from wild-type 2 mouse.

399 400 METHODS AND MATERIALS

401 *Data acquisition*

402 A total of four expression datasets were used in the **Functional gene modules detection from**
 403 **RNA-Seq data** section, that is, one synthetic RNA-Seq data, one *E. coli* RNA-Seq data and two
 404 human datasets (one RNA-Seq and one scRNA-Seq). The synthetic dataset was simulated
 405 using our in-house simulation method (details in **Simulation of co-regulated gene expression**
 406 **data** section). It contains 22,846 genes and 100 samples. A total of 10 co-regulated modules
 407 was embedded in this dataset, covering 2,240 up-regulated genes. The *E. coli* RNA-Seq data

408 consists of 4,497 genes and 155 samples, which was integrated and aggregated by our group.
409 In short, 155 fastq files were downloaded from <ftp://ftp.sra.ebi.ac.uk/vol1/fastq/> using the
410 sratoolkit (v2.8.1, <https://github.com/nbci/sra-tools/wiki/Downloads>), and they are processed
411 following quality check (FastQC), reads trimming (Btrim), reads mapping (HISAT2) and
412 transcript counting (HTseq). Then, raw read counts were RPKM normalized. The human RNA-
413 Seq data contains 3,084 genes and 8,555 samples, which was obtained from (30). The scRNA-
414 Seq data was downloaded from (13) as an RPKM expression matrix with 20,214 gene and 90
415 cells and then 3,798 genes were kept for the analysis in this study by removing the genes
416 without annotation.

417
418 Multiple sets of known modules/biological pathways were provided or collected to support the
419 enrichment analysis of the above four datasets. For synthetic data, the 10 groups of pre-defined
420 up-regulated genes were used as co-regulated modules. For *E. coli* data, we used five kinds of
421 biological pathways, which are complex regulons and regulons extracted from the RegulonDB
422 database (version 9.4, accessed on 05/08/2017), KEGG pathways collected from the KEGG
423 database (accessed on 08/08/2017), SEED subsystems from the SEED genomic database
424 (accessed on 08/08/2017) (44), and EcoCyc pathways from the EcoCyc database (version 21.1,
425 as of 08/08/2017) (45). Complex regulons were defined as a group of genes that are regulated
426 by the same transcription factor (TF) or the same set of TFs. In total, 457 complex regulons, 204
427 regulons, 123 KEGG pathways, 316 SEED subsystems, and 424 EcoCyc pathways were
428 retrieved, respectively. For the human TCGA and scRNA-Seq data, we used three sets of
429 modules provided by (30).

430
431 One golden-standard scRNA-Seq data (52) was downloaded from [https://scrnaseq-public-](https://scrnaseq-public-datasets.s3.amazonaws.com/manual-data/yan/nsmb.2660-S2.csv)
432 [datasets.s3.amazonaws.com/manual-data/yan/nsmb.2660-S2.csv](https://scrnaseq-public-datasets.s3.amazonaws.com/manual-data/yan/nsmb.2660-S2.csv) in the cell type classification
433 section. It consists of 20,214 genes and 90 cells, where the cells were assigned into seven
434 subgroups with the true cell subtypes information provided in (52).

435 The time series lung scRNA-Seq dataset (GSE52583) with 152 cells and 15,174 genes from was
436 downloaded from (47). The cells were collected at three time points: E14, E16, and E18. Another
437 time series scRNA-Seq data with 527 cells and 13991 genes (GSE48968) was downloaded from
438 the GEO database, in which the RPKM values are available.

439
440 The Mouse olfactory bulb spatial transcriptomic data was downloaded from (50), which contains
441 280 cells and 15,981 genes. Ståhl *et al.* (50) classified the cells into five clusters that correspond

442 to well-defined morphological layers. The cells use coordinates as IDs, and the cell layers
 443 information was manually extracted using the ST viewer
 444 (https://github.com/SpatialTranscriptomicsResearch/st_viewer), based on the coordinate
 445 information (see Supplementary File 6). The raw reads of mouse spatial scRNA-Seq data
 446 GSE60402 was retrieved from the SRA database (53,54), and the RPKM values for it were
 447 calculated using software packages TopHat (55) and Cufflink (56). GSE60402 was split into three
 448 subsets according to sample information. The detailed information of the selected and split
 449 datasets is listed in Table 1.

450

451 **Table 1.** Summary of GSE60402

GEO Accession ID	Data ID	Description	#Cells	#Genes
GSE60402	GSE60402-Mutant	From Gfra1 mutant sample	94	11094
GSE60402	GSE60402-Wildtype1	From wild type mouse 1	124	10037
GSE60402	GSE60402-Wildtype2	From wild type mouse 2	94	10714

452

453 ***Left Truncated Mixed Gaussian (LTMG) model and qualitative representation***

454 To accurately model the gene expression profile of RNA-Seq and scRNA-Seq data, we
 455 specifically developed a mixed Gaussian model with left truncation assumption. Denotes the log
 456 transformed FPKM, RPKM or CPM expression values of gene X over N conditions as $X =$
 457 $\{x_1, \dots, x_n\}$, we assumed that $x_j \in X$ follows a mixture of k Gaussian distributions, corresponding to
 458 k possible TRSs. The density function of x_j is:

459
$$p(x_j; \Theta) = \sum_{i=1}^k \alpha_i p(x_j; \theta_i) = \sum_{i=1}^k \alpha_i \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(x_j - \mu_i)^2}{2\sigma_i^2}}$$

460 And the density function of X is:

461
$$p(X; \Theta) = \prod_{j=1}^n p(x_j; \Theta) = \prod_{j=1}^n \sum_{i=1}^k \alpha_i p(x_j; \theta_i) = \prod_{j=1}^n \sum_{i=1}^k \alpha_i \frac{1}{\sqrt{2\pi}\sigma_i} e^{-\frac{(x_j - \mu_i)^2}{2\sigma_i^2}} = L(\Theta; X)$$

462 where α_i is the mixing weight, μ_i and σ_i are the mean and standard deviation of i^{th} Gaussian
 463 distribution, which can be estimated by an EM algorithm with given X :

464
$$\Theta^* = \arg \max_{\Theta} L(\Theta; X)$$

465

466 Parameters θ can be estimated by iteratively running the estimation (E) and maximization (M)
 467 steps. In this study, Z_{cut} is set for each gene as the logarithm of the minimal non-zero
 468 RPKM/FPKM/TPM value in the gene's expression profile. The EM algorithm is conducted for $K =$
 469 1, ..., 9 to fit the expression profile of each gene and the K that gives the best fit is selected
 470 according to the Bayesian Information Criterion (BIC):

$$471 \quad BIC = -2 \ln(\theta^*) + 3K \ln(N)$$

472 where K is the number of TRS, K is the number of conditions. K that minimizes the BIC will be
 473 selected.

474

475 Then the original gene expression values will be labeled to the most likely distribution under each
 476 cell. In detail, the probability that x_j belongs to distribution i is formulated by:

$$477 \quad p(x_j \in TRS\ i|K, \theta^*) \propto \frac{\alpha_i}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_i)^2}{2\sigma_i^2}}$$

478 And x_j is labeled by TRS i if $p(x_j \in TRS\ i|K, \theta^*) = \max_{i=1, \dots, K} (p(x_j \in TRS\ i|K, \theta^*))$. In such a
 479 way, a row consisting of discrete values (1, 2, ..., K) for each gene will be generated.

480

481 **KL score**

482 A Kullback-Leibler divergence score (**KL** score) is introduced in QUBIC 2 to guide candidate-
 483 selection and biclustering optimization. The KL score of a bicluster is defined as:

$$484 \quad KL_B = \frac{1}{N} \sum_{j=1}^N \sum_{i \in \{0,1\}} R(i, j) \times \log \frac{R(i, j)}{Q(i, j)} + \frac{1}{M} \sum_{k=1}^M \sum_{i \in \{0,1\}} C(i, k) \times \log \frac{C(i, k)}{P(i, k)}$$

485 where N and M are the numbers of rows and columns of a submatrix B in M_R , respectively. $R(i, j)$
 486 represents the proportion of element i in row j of B , $Q(i, j)$ is the proportion of i in the
 487 corresponding entire row, $C(i, k)$ is the proportion of i in column k of B , and $P(i, k)$ is the
 488 proportion of i in the entire corresponding column.

489

490 Meanwhile, the KL score for a gene quantify the similarity between a candidate gene j and a
 491 bicluster, which is defined as follows:

$$492 \quad KL_j = \sum_{i \in \{0,1\}} R(i, j) \times \log \frac{R(i, j)}{Q(i, j)}$$

493 where $R(i, j)$ represent the proportion of i under corresponding columns of the current bicluster.

494

495 **Simulation of co-regulated gene expression data**

496 We utilized a single cell RNA-Seq dataset of human melanoma (58) (with 22,846 genes and 4,645
497 cells) to simulate bulk tissue RNA-Seq data with known co-regulated modules. Specifically, a
498 single cell RNA-Seq pool consists counts data of 4,466 cells of six annotated cell types namely
499 B-, T-, endothelial, fibroblast, macrophage and cancer cells were constructed. The top 1,000 cell
500 type specifically expressed genes of each cell type were identified by using Z score of the mean
501 of each gene's expression level in each cell type.

502

503 For each round of simulation, the number of to be simulated bulk tissue samples and co-regulation
504 modules is first defined. Then the genes of each co-regulation module denoted as X_k , will be
505 specified by randomly selecting M_k genes from the top 1,000 cell type specifically expressed
506 genes of one cell type. A co-regulation strength matrix P is then simulated from a bimodal
507 distribution over $(0,1)$, with $P[i, k]$ denotes the proportion of cells with the transcriptional regulatory
508 signal of co-regulation module k in bulk sample i . A bulk tissue data is simulated by randomly
509 drawing cells from the cell pool by following a multinomial distribution, with predefined parameters
510 and the total number of cells. For co-regulation module k in bulk sample i , genes X_k in a
511 proportion $P[i, k]$ of the selected cells of the cell type corresponds to k are perturbed by an X -fold
512 increase of the gene expression. Then the bulk data i with simulated co-regulations are formed
513 by summing the perturbed gene expression profile the selected cells and normalized to RPKM
514 expression scale. The Pseudo code of the simulation approach is provided Method S1 in
515 Supplementary File 1.

516

517 The rationales of this simulation approach include (1) gene expression level and noise in the bulk
518 data are purely simulated by sum of real single-cell data, without using artificially assigned
519 expressions scale and noise; (2) co-regulation genes are modeled as a specific fold increase of
520 a number of cell-type-specific genes in a particular subset of the cells, which characterizes the
521 heterogeneity of transcriptional regulation among cells in a tissue; (3) multiple co-regulation
522 modules in specific to different cell types can be simultaneously simulated. Hence, we believe the
523 gene expression data simulated by this way can satisfactorily reflect genes co-regulated by a
524 perturbed transcriptional regulation signal in real bulk tissue data.

525

526 **Evaluation of the functional modules**

527 The capability of algorithms to recapitulate known functional modules are assessed using
528 precision and recall. First, for each identified bicluster, we use the P -value of its most enriched

529 functional class (biological pathway) as the P -value of the bicluster. Specifically, the probability of
530 having x genes of the same functional class in a bicluster of size n from a genome with a total of
531 N genes can be computed using the following hypergeometric function(59):

$$532 \quad P(X = x|N, p, n) = \frac{\binom{pN}{x} \binom{(1-p)N}{n-x}}{\binom{N}{n}}$$

533 where p is the percentage of that pathway among all pathways in the whole genome. The
534 P -value of getting such enriched or even more enriched bicluster is calculated as:

$$535 \quad P - \text{value} = P(X \geq x) = 1 - P(X < x) = 1 - \sum_{i=0}^{x-1} \frac{\binom{pN}{i} \binom{(1-p)N}{n-i}}{\binom{N}{n}}$$

536 The bicluster is deemed enriched with that function if its p -value is smaller than a specific
537 cutoff (e.g., 0.05).

538
539 Given a group of biclusters identified by a tool under a parameter combination, the precision is
540 defined as the fraction of observed biclusters significantly enriched with the one biological
541 pathway/known modules (Benjamini-Hochberg adjusted $p < 0.05$),

$$542 \quad \text{Precision} = \frac{\# \text{ of significant biclusters}}{\# \text{ of biclusters}}$$

543
544
545 For recall, we compute the fraction of known modules that were rediscovered by the algorithms,

$$546 \quad \text{Recall} = \frac{\# \text{ of significant modules}}{\# \text{ of modules}}$$

547 Finally, the harmonic mean of precision and recall were calculated to represent the performance
548 of an algorithm on a given dataset and parameter setting, denoted as F score:

$$549 \quad F = \frac{2}{\frac{1}{\text{Precision}} + \frac{1}{\text{Recall}}}$$

550 Note that the number of biclusters used to calculate precision and recall may affect the results.
551 To make sure the evaluation is as fair as possible, for each dataset, we select the first 30
552 biclusters.

553 554 **Parameter adjustment of biclustering tools**

555 To assess the robustness of selected algorithms' performance, each tool is run multiple times by
556 varying parameters that affect the size and number of biclusters. In general, parameters are
557 adjusted around their default or recommended (if available) value. The parameters that varied are

558 listed in Table2, and details about the range and increment of parameters can be found in
559 Supplementary File.

560

561

Table 2. Main parameters adjusted for each algorithm

Algorithm	Implementation	Parameters
Bimax	R package 'biclust'	minr, minc, number
ISA	R package 'isa2'	set.seed
FABIA	R package 'fabia'	alpha, spl, spz, cyc, p
Plaid	R package 'biclust'	row.release, col.release, max.layer
QUBIC	R package 'QUBIC'	f, c, k, o
QUBIC2	C++	f, c, k, o

562

563 ***Spearman correlation test***

564 QUBIC2 was run on the *E. coli* RNA-Seq data in Figure 2 under 63 parameter settings. For each
565 setting, around 100 biclusters were identified. Five sets regulatory or metabolic pathways were
566 extracted from four databases of *E. coli* (RegulonDB, KEGG, SEED (46) and EcoCyc (47)) to
567 support this association study. In specific, for each set of ~100 biclusters obtained under the
568 same settings, six groups of *P*-values for all these biclusters were calculated, with five knowledge-
569 based groups and one size-based group. Spearman correlation test was conducted to investigate
570 the rank-order correlation among the six groups of *P*-values. Five correlation coefficients (ρ),
571 which demonstrated the extent of correlation between size-based *P*-values and five biological
572 knowledge-based *P*-values, as well as five corresponding *p*-values, were recorded from the test.
573 Note that the *p*-value of correlation test denotes the probability of observing such a correlation or
574 even stronger correlation, under the null hypothesis that no correlation exists. For simplicity, the
575 correlation coefficient between the size-based *P*-value and biological knowledge-based *P*-value
576 was prefixed with the name of a pathway, e.g., TF_ ρ and KEGG_ ρ . In the end, a total of 5×63
577 ρ (63 parameter settings, each with five ρ s) and a same number of *p*-values were obtained.

578

579 **Cell type classification pipeline**

580 By using biclustering, we can group genes and cells simultaneously. However, since biclustering
581 aims to find sets of genes that are co-expressed across a subset of conditions, it is possible that
582 genes may co-expressed across multiple cell types. Therefore, one bicluster may consist of cells
583 from different types, and cells from the same types may appear in different biclusters. In a word,
584 it is not guaranteed that one bicluster corresponds to one cell type. However, it is assumed that
585 two cells from a bicluster are more likely to be of the same subtypes than the two cells that are
586 randomly selected. It is believed that biclusters can capture this feature to some extent. If there
587 are multiple biclusters and when we condense them together, we can distinguish sets of cells
588 belonging to the same type from sets of cells that are grouped by chance.

589
590 Based on the above idea, we developed a pipeline to obtain cell type classification based on
591 biclustering results (Figure 4A). First, a biclustering tool was applied to the expression data (rows
592 represent genes and columns represent cells) to identify a set of biclusters. Then a weighted
593 graph $G = (C, E)$ was constructed to model the relationship between cell pairs among biclusters.
594 A node c_i in G represented a cell, and $e_{i,j}$ represented the edge connecting c_i and c_j , where $i \neq j$.
595 We assigned weight $w_{i,j}$ to $e_{i,j}$ to represent the number of biclusters that contain both c_i and c_j .
596 Intuitively, a higher $w_{i,j}$ value indicates that c_i and c_j are simultaneously involved in more
597 biclusters, hence, are more likely to be the same cell type than cell pairs with lower weight. A
598 symmetrical cell-cell matrix with diagonal as 0 was then constructed to record $w_{i,j}$ and Markov
599 Cluster Algorithm (**MCL**) was performed to cluster cells into cell types and produce cell labels. In
600 specific, the MCL clustering was run 100 times by varying inflation factor, resulting 100 cell labels.
601 A binary similarity matrix was constructed for each cell label: if two cells belong to the same cluster,
602 their similarity is 1; otherwise, the similarity is 0. Then a consensus matrix was built by averaging
603 all similarity matrices. The resulting consensus matrix was clustered using hierarchical clustering
604 with complete agglomeration, and the clusters were inferred at the k level of the hierarchy.

605

606 **External cluster validity indices**

607 External validation measures the extent to which cluster labels match externally supplied class
608 labels. Generally, they are based on counting the pairs of points on which two classifiers
609 agree/disagree. Denote two partitions of the same data set as R and Q. The reference partition,
610 R, encode the class labels, i.e., it partitions the data into k known classes. Partition Q, in turn,
611 partitions the data into v categories, which is the one to be evaluated.

612

613 Adjusted Rand Index (ARI) is defined as

$$614 \quad ARI = \frac{a - \frac{(a+c)(a+b)}{d}}{\frac{(a+c) + (a+b)}{2} - \frac{(a+c)(a+b)}{d}}$$

615 *a*: Number of pairs of data objects belonging to the same class in R and the same cluster in Q.

616 *b*: Number of pairs of data objects belonging to the same class in R and different clusters in Q.

617 *c*: Number of pairs of data objects belonging to different classes in R and the same cluster in Q.

618 *d*: Number of pairs of data objects belonging to different classes in R and different clusters in Q.

619 Terms *a* and *d* are measures of consistent classifications (agreements), whereas terms *b*
620 and *c* are measures of inconsistent classifications (disagreements).

621

622 Jaccard Index is defined as:

$$623 \quad JI = \frac{a}{a + b + c}$$

624 The Jaccard Index can be seen as a proportion of good pairs with respect to the sum of non-
625 neutral (good plus bad) pairs.

626

627 Folkes-Mallow's index is defined as

$$628 \quad FI = \frac{a}{\sqrt{(a+b)(a+c)}}$$

629 Fowlkes–Mallow's index can be seen as a non-linear modification of the Jaccard coefficient that
630 also keeps normality.

631

632 **Pathway enrichment analysis**

633 Pathway enrichment analysis is conducted and the statistical significance of each enriched
634 pathway is assessed by using a hypergeometric test (statistical significance cutoff = 0.005)
635 against 4,725 curated gene sets in the MsigDB database, which includes 1,330 canonical KEGG,
636 Biocarta and Reactome pathways, and 3,395 gene sets representing expression signatures
637 derived from experiments with genetic and chemical perturbations, together with 6,215 Mouse
638 GO terms each containing at least 5 genes (62,63).

639

640 **CONCLUSION**

641 QUBIC2 is a novel biclustering algorithm developed for bulk RNA-Seq and scRNA-Seq data
642 analysis in this study. It has four unique characteristics: (*i*) used a left-truncated mixture model to
643 fit the log-transformed RPKM/CPM/TPM values of each gene and qualitatively represent gene

644 expression; (ii) integrates an information-divergence objective function in the biclustering
645 framework; (iii) applies a Dual strategy to optimize consistency level of a to-be-identified bicluster;
646 and (iv) develops a robust *P*-value framework to evaluate the significance of all the identified
647 biclusters. QUBIC2 proved to have significant advantages in the functional module detection area,
648 outperforming five widely-used biclustering methods based upon our test on four datasets. The
649 proposed *P*-value calculation method based on bicluster size did make sense, which may facilitate
650 the evaluation of all the identified biclusters, especially from less-annotated organisms. The cell
651 type classification pipeline, based on QUBIC2, worked well and outperformed the state-of-the-art
652 performance of SC3. By utilizing time-dependent data, QUBIC2 discovered biclusters specific to
653 time point and identified a cascade of immune responses to the external pathogenic treatment.
654 From the spatial transcriptomic data, QUBIC2 discovered that spatially adjacent single cells may
655 have high co-expression patterns, and particularly, two distinct spatially clustered cells may be
656 derived initially from the same stem cell. We believe that QUBIC2 can serve biologists as a useful
657 tool to extract novel biological insights from large-scale RNA-Seq data (The tutorial for QUBIC2
658 program is provided in Supplementary File 7).

659

660 **DISCUSSION**

661 Single-cell sequencing has enabled new transcriptome-based studies, including the study of
662 distinct responses by different cell types in the same population when encountered by the same
663 stimuli or stresses, and identification of the complex relationships among different cells in complex
664 biological environments such as tissues. However, to fully excavate the potential of scRNA-Seq
665 data, we must overcome several technical challenges.

666

667 As sequencing costs decrease, larger scRNA-Seq datasets will become increasingly common;
668 thus, the scalability to large dataset and efficiency of tools will become more and more important.
669 Currently, the discretization and Dual searching functions of QUBIC2 are time consuming on
670 large-scale datasets. Based on our test, it takes 17 minutes to discretize a dataset with 4,297
671 rows and 466 columns (a desktop with 48.0GB memory, Intel Core i7-6700 and 3.40GHz). Given
672 a dataset with 22,846 genes and 100 conditions, the running time while using Dual strategy are
673 generally 2 minutes longer than that without Dual. The openMP method will be implemented in
674 the EM steps for discretization and more efficient heuristics algorithm will be designed to optimize
675 the dual searching of biclustering.

676

677 Another challenge involves the interpretation of time-series and spatial data. For example, in the
678 GSE52583 data, QUBIC2 could only separate cells collected at different time points, yet the
679 further differentiation stage information was not captured. For the mouse olfactory bulb data,
680 QUBIC2 did not separate cells from adjacent layers. To deal with this drawback, we need to
681 combine biclustering with other statistical methods specifically designed for time series and spatial
682 gene expression data.

683
684 It is noteworthy that many other kinds of methods can be used for gene expression data analysis.
685 Forty-two module detection tools covering five main approaches were reviewed in (30) and the
686 authors concluded that decomposition methods outperformed all other strategies, including
687 biclustering methods. Meanwhile, they also observed that QUBIC and FABIA had higher
688 performance on human and synthetic data. We compared two top rated decomposition methods
689 and two top clustering methods with QUBIC2 and QUBIC on a human scRNA-Seq data; and the
690 results showed that QUBIC2 surpassed both decomposition and clustering methods (Figure S4
691 in Supplementary File 1). In the future, we will carry out more comprehensive comparison
692 between QUBIC2 and other decomposition and network-based methods, aiming to give a
693 systematical evaluation of the power of computational techniques on scRNA-seq data.

694

695 REFERENCES

- 696 1. Van Dijk, E.L., Auger, H., Jaszczyszyn, Y. and Thernes, C. (2014) Ten years of next-
697 generation sequencing technology. *Trends in genetics*, **30**, 418-426.
- 698 2. Goodwin, S., McPherson, J.D. and McCombie, W.R. (2016) Coming of age: ten years of
699 next-generation sequencing technologies. *Nature Reviews Genetics*, **17**, 333-351.
- 700 3. Marioni, J.C., Mason, C.E., Mane, S.M., Stephens, M. and Gilad, Y. (2008) RNA-seq: an
701 assessment of technical reproducibility and comparison with gene expression arrays.
702 *Genome research*, **18**, 1509-1517.
- 703 4. Nagalakshmi, U., Wang, Z., Waern, K., Shou, C., Raha, D., Gerstein, M. and Snyder, M.
704 (2008) The transcriptional landscape of the yeast genome defined by RNA sequencing.
705 *Science*, **320**, 1344-1349.
- 706 5. Wang, Z., Gerstein, M. and Snyder, M. (2009) RNA-Seq: a revolutionary tool for
707 transcriptomics. *Nature reviews genetics*, **10**, 57-63.
- 708 6. Aibar, S., Gonzalez-Blas, C.B., Moerman, T., Huynh-Thu, V.A., Imrichova, H.,
709 Hulselmans, G., Rambow, F., Marine, J.C., Geurts, P., Aerts, J. *et al.* (2017) SCENIC:
710 single-cell regulatory network inference and clustering. *Nat Methods*, **14**, 1083-1086.
- 711 7. Prince, M.E., Sivanandan, R., Kaczorowski, A., Wolf, G.T., Kaplan, M.J., Dalerba, P.,
712 Weissman, I.L., Clarke, M.F. and Ailles, L.E. (2007) Identification of a subpopulation of
713 cells with cancer stem cell properties in head and neck squamous cell carcinoma. *Proc*
714 *Natl Acad Sci U S A*, **104**, 973-978.
- 715 8. Navin, N., Kendall, J., Troge, J., Andrews, P., Rodgers, L., McIndoo, J., Cook, K.,
716 Stepansky, A., Levy, D., Esposito, D. *et al.* (2011) Tumour evolution inferred by single-
717 cell sequencing. *Nature*, **472**, 90-94.

- 718 9. Xu, X., Hou, Y., Yin, X., Bao, L., Tang, A., Song, L., Li, F., Tsang, S., Wu, K., Wu, H. *et al.* (2012) Single-cell exome sequencing reveals single-nucleotide mutation
719 characteristics of a kidney tumor. *Cell*, **148**, 886-895.
720
- 721 10. Ulitsky, I., Maron-Katz, A., Shavit, S., Sagir, D., Linhart, C., Elkon, R., Tanay, A., Sharan,
722 R., Shiloh, Y. and Shamir, R. (2010) Expander: from expression microarrays to networks
723 and functions. *Nat Protoc*, **5**, 303-322.
- 724 11. Jaitin, D.A., Kenigsberg, E., Keren-Shaul, H., Elefant, N., Paul, F., Zaretsky, I., Mildner,
725 A., Cohen, N., Jung, S., Tanay, A. *et al.* (2014) Massively Parallel Single-Cell RNA-Seq
726 for Marker-Free Decomposition of Tissues into Cell Types. *Science*, **343**, 776-779.
- 727 12. Shekhar, K., Lapan, S.W., Whitney, I.E., Tran, N.M., Macosko, E.Z., Kowalczyk, M.,
728 Adiconis, X., Levin, J.Z., Nemesh, J., Goldman, M. *et al.* (2016) Comprehensive
729 Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics. *Cell*, **166**, 1308-
730 1323.e1330.
- 731 13. Kiselev, V.Y., Kirschner, K., Schaub, M.T., Andrews, T., Yiu, A., Chandra, T., Natarajan,
732 K.N., Reik, W., Barahona, M., Green, A.R. *et al.* (2017) SC3: consensus clustering of
733 single-cell RNA-seq data. *Nat Methods*, **14**, 483-486.
- 734 14. Kluger, Y., Basri, R., Chang, J.T. and Gerstein, M. (2003) Spectral biclustering of
735 microarray data: coclustering genes and conditions. *Genome research*, **13**, 703-716.
- 736 15. Lazzeroni, L. and Owen, A. (2002) Plaid models for gene expression data. *Statistica*
737 *sinica*, 61-86.
- 738 16. Madeira, S.C. and Oliveira, A.L. (2009) A polynomial time biclustering algorithm for
739 finding approximate expression patterns in gene expression time series. *Algorithms Mol*
740 *Biol*, **4**, 8.
- 741 17. Bergmann, S., Ihmels, J. and Barkai, N. (2003) Iterative signature algorithm for the
742 analysis of large-scale gene expression data. *Physical review E*, **67**, 031902.
- 743 18. Prelić, A., Bleuler, S., Zimmermann, P., Wille, A., Bühlmann, P., Grissem, W., Hennig,
744 L., Thiele, L. and Zitzler, E. (2006) A systematic comparison and evaluation of
745 biclustering methods for gene expression data. *Bioinformatics*, **22**, 1122-1129.
- 746 19. Li, G., Ma, Q., Tang, H., Paterson, A.H. and Xu, Y. (2009) QUBIC: a qualitative
747 biclustering algorithm for analyses of gene expression data. *Nucleic Acids Res*, **37**,
748 e101.
- 749 20. Hochreiter, S., Bodenhofer, U., Heusel, M., Mayr, A., Mitterecker, A., Kasim, A.,
750 Khamiakova, T., Van Sanden, S., Lin, D., Talloen, W. *et al.* (2010) FABIA: factor
751 analysis for bicluster acquisition. *Bioinformatics*, **26**, 1520-1527.
- 752 21. Barkow, S., Bleuler, S., Prelić, A., Zimmermann, P. and Zitzler, E. (2006) BicAT: a
753 biclustering analysis toolbox. *Bioinformatics*, **22**, 1282-1283.
- 754 22. Cheng, K.O., Law, N.F., Siu, W.C. and Lau, T.H. (2007) BiVisu: software tool for
755 bicluster detection and visualization. *Bioinformatics*, **23**, 2342-2344.
- 756 23. Wu, C.J. and Kasif, S. (2005) GEMS: a web server for biclustering analysis of
757 expression data. *Nucleic Acids Res*, **33**, W596-599.
- 758 24. Zhou, F., Ma, Q., Li, G. and Xu, Y. (2012) QServer: a biclustering server for prediction
759 and assessment of co-expressed gene clusters. *PloS one*, **7**, e32660.
- 760 25. Kaiser, S., Santamaria, R., Theron, R., Quintales, L. and Leisch, F. (2009) biclust:
761 Bicluster algorithms. *R package version 0.7*, **2**.
- 762 26. Zhang, Y., Xie, J., Yang, J., Fennell, A., Zhang, C. and Ma, Q. (2016) QUBIC: a
763 bioconductor package for qualitative biclustering analysis of gene co-expression data.
764 *Bioinformatics*, btw635.
- 765 27. Eren, K., Deveci, M., Küçükünç, O. and Çatalyürek, Ü.V. (2013) A comparative analysis
766 of biclustering algorithms for gene expression data. *Briefings in bioinformatics*, **14**, 279-
767 292.

- 768 28. Chia, B.K.H. and Karuturi, R.K.M. (2010) Differential co-expression framework to
769 quantify goodness of biclusters and compare biclustering algorithms. *Algorithms for*
770 *molecular biology*, **5**, 23.
- 771 29. Padilha, V.A. and Campello, R.J. (2017) A systematic comparative evaluation of
772 biclustering techniques. *BMC Bioinformatics*, **18**, 55.
- 773 30. Saelens, W., Cannoodt, R. and Saeys, Y. (2018) A comprehensive evaluation of module
774 detection methods for gene expression data. *Nature Communications*, **9**, 1090.
- 775 31. Lun, A.T., Bach, K. and Marioni, J.C. (2016) Pooling across cells to normalize single-cell
776 RNA sequencing data with many zero counts. *Genome Biol*, **17**, 75.
- 777 32. Bacher, R. and Kendzierski, C. (2016) Design and computational analysis of single-cell
778 RNA-sequencing experiments. *Genome biology*, **17**, 63.
- 779 33. Bengtsson, M., Stahlberg, A., Rorsman, P. and Kubista, M. (2005) Gene expression
780 profiling in single cells from the pancreatic islets of Langerhans reveals lognormal
781 distribution of mRNA levels. *Genome Res*, **15**, 1388-1392.
- 782 34. Lu, C. and King, R.D. (2009) An investigation into the population abundance distribution
783 of mRNAs, proteins, and metabolites in biological systems. *Bioinformatics*, **25**, 2020-
784 2027.
- 785 35. Hebenstreit, D., Fang, M., Gu, M., Charoensawan, V., van Oudenaarden, A. and
786 Teichmann, S.A. (2011) RNA sequencing reveals two major classes of gene expression
787 levels in metazoan cells. *Mol Syst Biol*, **7**, 497.
- 788 36. Glaus, P., Honkela, A. and Rattray, M. (2012) Identifying differentially expressed
789 transcripts from RNA-seq data with biological variation. *Bioinformatics*, **28**, 1721-1728.
- 790 37. Rau, A. and Maugis-Rabusseau, C. (2017) Transformation and model choice for RNA-
791 seq co-expression analysis. *Brief Bioinform*.
- 792 38. Reuter, J.A., Spacek, D.V., Pai, R.K. and Snyder, M.P. (2016) Simul-seq: combined
793 DNA and RNA sequencing for whole-genome and transcriptome profiling. *Nature*
794 *Methods*.
- 795 39. Cohen, A.C. (1959) Simplified estimators for the normal distribution when samples are
796 singly censored or truncated. *Technometrics*, **1**, 217-237.
- 797 40. Stegle, O., Teichmann, S.A. and Marioni, J.C. (2015) Computational and analytical
798 challenges in single-cell transcriptomics. *Nat Rev Genet*, **16**, 133-145.
- 799 41. Bergmann, S., Ihmels, J. and Barkai, N. (2003) Iterative signature algorithm for the
800 analysis of large-scale gene expression data. *Phys Rev E Stat Nonlin Soft Matter Phys*,
801 **67**, 031902.
- 802 42. Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. and Morishima, K. (2017) KEGG: new
803 perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res*, **45**, D353-
804 d361.
- 805 43. Gama-Castro, S., Salgado, H., Santos-Zavaleta, A., Ledezma-Tejeida, D., Muniz-
806 Rascado, L., Garcia-Sotelo, J.S., Alquicira-Hernandez, K., Martinez-Flores, I., Pannier,
807 L., Castro-Mondragon, J.A. *et al.* (2016) RegulonDB version 9.0: high-level integration of
808 gene regulation, coexpression, motif clustering and beyond. *Nucleic Acids Res*, **44**,
809 D133-143.
- 810 44. Overbeek, R., Begley, T., Butler, R.M., Choudhuri, J.V., Chuang, H.-Y., Cohoon, M., de
811 Crécy-Lagard, V., Diaz, N., Disz, T., Edwards, R. *et al.* (2005) The Subsystems
812 Approach to Genome Annotation and its Use in the Project to Annotate 1000 Genomes.
813 *Nucleic Acids Research*, **33**, 5691-5702.
- 814 45. Keseler, I.M., Mackie, A., Santos-Zavaleta, A., Billington, R., Bonavides-Martínez, C.,
815 Caspi, R., Fulcher, C., Gama-Castro, S., Kothari, A., Krummenacker, M. *et al.* (2017)
816 The EcoCyc database: reflecting new knowledge about *Escherichia coli* K-12. *Nucleic*
817 *Acids Research*, **45**, D543-D550.

- 818 46. Yan, L., Yang, M., Guo, H., Yang, L., Wu, J., Li, R., Liu, P., Lian, Y., Zheng, X., Yan, J.
819 *et al.* (2013) Single-cell RNA-Seq profiling of human preimplantation embryos and
820 embryonic stem cells. *Nature structural & molecular biology*, **20**, 1131-1139.
- 821 47. Treutlein, B., Brownfield, D.G., Wu, A.R., Neff, N.F., Mantalas, G.L., Espinoza, F.H.,
822 Desai, T.J., Krasnow, M.A. and Quake, S.R. (2014) Reconstructing lineage hierarchies
823 of the distal lung epithelium using single-cell RNA-seq. *Nature*, **509**, 371.
- 824 48. Dennis, G., Sherman, B.T., Hosack, D.A., Yang, J., Gao, W., Lane, H.C. and Lempicki,
825 R.A. (2003) DAVID: database for annotation, visualization, and integrated discovery.
826 *Genome biology*, **4**, R60.
- 827 49. Shalek, A.K., Satija, R., Shuga, J., Trombetta, J.J., Gennert, D., Lu, D., Chen, P.,
828 Gertner, R.S., Gaublomme, J.T., Yosef, N. *et al.* (2014) Single-cell RNA-seq reveals
829 dynamic paracrine control of cellular variation. *Nature*, **510**, 363-369.
- 830 50. Ståhl, P.L., Salmén, F., Vickovic, S., Lundmark, A., Navarro, J.F., Magnusson, J.,
831 Giacomello, S., Asp, M., Westholm, J.O., Huss, M. *et al.* (2016) Visualization and
832 analysis of gene expression in tissue sections by spatial transcriptomics. *Science*, **353**,
833 78-82.
- 834 51. Zechel, S., Zajac, P., Lonnerberg, P., Ibanez, C.F. and Linnarsson, S. (2014)
835 Topographical transcriptome mapping of the mouse medial ganglionic eminence by
836 spatially resolved RNA-seq. *Genome Biol*, **15**, 486.
- 837 52. Kiselev, V.Y., Kirschner, K., Schaub, M.T., Andrews, T., Yiu, A., Chandra, T., Natarajan,
838 K.N., Reik, W., Barahona, M. and Green, A.R. (2017) SC3: consensus clustering of
839 single-cell RNA-seq data. *Nature methods*.
- 840 53. Barrett, T., Wilhite, S.E., Ledoux, P., Evangelista, C., Kim, I.F., Tomashevsky, M.,
841 Marshall, K.A., Phillippy, K.H., Sherman, P.M. and Holko, M. (2013) NCBI GEO: archive
842 for functional genomics data sets—update. *Nucleic acids research*, **41**, D991-D995.
- 843 54. Barrett, T., Troup, D.B., Wilhite, S.E., Ledoux, P., Rudnev, D., Evangelista, C., Kim, I.F.,
844 Soboleva, A., Tomashevsky, M. and Edgar, R. (2007) NCBI GEO: mining tens of millions
845 of expression profiles—database and tools update. *Nucleic acids research*, **35**, D760-
846 D765.
- 847 55. Trapnell, C., Pachter, L. and Salzberg, S.L. (2009) TopHat: discovering splice junctions
848 with RNA-Seq. *Bioinformatics*, **25**, 1105-1111.
- 849 56. Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., Van Baren, M.J.,
850 Salzberg, S.L., Wold, B.J. and Pachter, L. (2010) Transcript assembly and quantification
851 by RNA-Seq reveals unannotated transcripts and isoform switching during cell
852 differentiation. *Nature biotechnology*, **28**, 511-515.
- 853 57. Sha Cao, T.S., Xin Chen, Qin Ma, Chi Zhang. (2017) A probabilistic model-based bi-
854 clustering method for single-cell transcriptomic data analysis. *bioRxiv*.
- 855 58. Tirosh, I., Izar, B., Prakadan, S.M., Wadsworth, M.H., 2nd, Treacy, D., Trombetta, J.J.,
856 Rotem, A., Rodman, C., Lian, C., Murphy, G. *et al.* (2016) Dissecting the multicellular
857 ecosystem of metastatic melanoma by single-cell RNA-seq. *Science*, **352**, 189-196.
- 858 59. Castillo-Davis, C.I. and Hartl, D.L. (2003) GeneMerge—post-genomic analysis, data
859 mining, and hypothesis testing. *Bioinformatics*, **19**, 891-892.
- 860 60. Monk, J., Nogales, J. and Palsson, B.O. (2014) Optimizing genome-scale network
861 reconstructions. *Nat Biotechnol*, **32**, 447-452.
- 862 61. Sun, X. and Nobel, A.B. (2008) On the size and recovery of submatrices of ones in a
863 random binary matrix. *Journal of Machine Learning Research*, **9**, 2431-2453.
- 864 62. Liberzon, A., Subramanian, A., Pinchback, R., Thorvaldsdóttir, H., Tamayo, P. and
865 Mesirov, J.P. (2011) Molecular signatures database (MSigDB) 3.0. *Bioinformatics*, **27**,
866 1739-1740.

867 63. Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P.,
868 Dolinski, K., Dwight, S.S. and Eppig, J.T. (2000) Gene Ontology: tool for the unification
869 of biology. *Nature genetics*, **25**, 25-29.

870

871