

1 **Enhancing single-cell cellular state inference by incorporating molecular**
2 **network features**

3

4 Ji Dong^{1,3,*}, Peijie Zhou^{2,*}, Yichong Wu^{2,*}, Wendong Wang^{1,*}, Yidong Chen^{1,3,4*}, Xin
5 Zhou^{1,*}, Haoling Xie^{1,3}, Yuan Gao^{1,3}, Jiansen Lu^{1,3}, Jingwei Yang¹, Xiannian Zhang^{1,3},
6 Lu Wen^{1,3}, Wei Fu^{1,#}, Tiejun Li^{2,#}, Fuchou Tang^{1,3,4,#}

7

8 ¹ Beijing Advanced Innovation Center for Genomics (ICG), Department of General
9 Surgery, College of Life Sciences, Third Hospital, Peking University, Beijing
10 100871, P. R. China

11 ² LMAM and School of Mathematical Sciences, Peking University, 100871 Beijing,
12 China

13 ³ Biomedical Institute for Pioneering Investigation via Convergence and Center for
14 Reproductive Medicine, Ministry of Education Key Laboratory of Cell Proliferation
15 and Differentiation, Beijing 100871, China

16 ⁴ Peking-Tsinghua Center for Life Sciences, Academy for Advanced Interdisciplinary
17 Studies, Peking University, Beijing 100871, China

18 * These authors contributed equally

19 # Correspondence: fuwei@bjmu.edu.cn (W.F.), tieli@pku.edu.cn (T.L.),
20 tangfuchou@pku.edu.cn (F.T.),

21

22

1 **Abstract**

2 In biological systems, genes function in conjunction rather than in isolation. However,
3 traditional single-cell RNA-seq (scRNA-seq) analyses heavily rely on the
4 transcriptional similarity of individual genes, ignoring the inherent gene-gene
5 interactions. Here, we present SCORE, a network-based method, which incorporates
6 the validated molecular network features to infer cellular states. Using real
7 scRNA-seq datasets, SCORE outperforms existing methods in accuracy, robustness,
8 scalability, data integration and removal of batch effect. When applying SCORE to a
9 newly generated human ileal scRNA-seq dataset, we identified several novel
10 stem/progenitor clusters, including a Cripto-1+ cluster. Moreover, two distinct groups
11 of goblet cells were identified and only one of them tended to secrete mucus. Besides,
12 we found that the recently identified *BEST4+OTOP2+* microfold cells also highly
13 expressed *CFTR*, which is different from their colonic counterparts. In summary,
14 SCORE enhances cellular state inference by simulating the dynamic changes of
15 molecular networks, providing more biological insights beyond statistical
16 interpretations.

17

1 **Introduction**

2 Instead of executing function in isolation, genes tend to form complex molecular
3 networks and function in conjunction to determine the cellular or organismal
4 phenotypes^{1,2}. Although it has been recognized for a long time that molecular
5 networks are dynamic during organismal development and differentiation, how to
6 simulate this process is still challenging.

7

8 With the rapid development of single-cell RNA sequencing (scRNA-seq) techniques,
9 nowadays, researchers can easily profile the transcriptome of large number of cells at
10 single-cell resolution, and identify new cell types and intermediate cellular states
11 within a certain organismal system³⁻⁵. To fully utilize these rich datasets, many
12 efficient computational methods have been developed, such as Seurat, SCANPY, and
13 SINCERA⁶⁻⁸. However, these methods usually calculate the transcriptional similarity
14 inferred from individual genes to do cell clustering analysis and detect the marker
15 genes that can uniquely distinguish the identified cell types, while ignore the
16 integrated and synergistic nature of molecular interaction networks, which in fact
17 formulate the specific cell fates and provide more biological insights. In addition, the
18 performance of individual gene-based methods can be influenced by the prevalent
19 dropout events and other experimental artifacts in single-cell sequencing data, and the
20 performance strongly depends on the subjective selection of highly variable genes in
21 the pre-processing steps.

22

23 Recently, several methods such as SCENIC, PAGODA, and CSN are proposed to
24 identify cell types by incorporating some biological knowledge-based information,
25 which utilizes transcription factor-based regulatory networks, functional modules, and
26 cell type-specific networks to facilitate the biological interpretation, respectively⁹⁻¹¹.
27 However, due to strong technical noises there are still no optimal methods to
28 accurately infer the gene-gene or cell-cell relationship from the sparse scRNA-seq
29 datasets¹². In the meantime, accumulating public data on molecular interactions

1 derived from sound experimental or computational evidences can provide rich prior
2 biological knowledge to reduce the false-positive rate of molecular network inference
3 ^{13,14}. For example, with the aid of protein-protein interaction (PPI) databases, SCENT
4 is able to infer the single-cell entropy, while netSmooth can temper the noisy
5 scRNA-seq data^{15,16}.

6

7 We hypothesize that during the organismal development and differentiation, the
8 transition from a cellular state to another is accompanied by the destruction of critical
9 molecular networks of the former cellular state and the reconstruction of novel ones
10 (Fig. 1a). In view of this, we introduced a new method, **SCORE** (**S**ingle-**C**ell
11 **m**olecular **R**etwork), to simulate the dynamic changes of molecular networks from
12 scRNA-seq datasets by incorporating the experimentally validated and
13 high-confidence molecular interaction information from public databases. We
14 validated the accuracy, robustness and scalability of SCORE to uncover cell states
15 using gold-standard scRNA-seq datasets. The performance demonstrated the
16 superiority of the proposed method over previous proposals. Finally, with SCORE, we
17 accurately integrated five human fetal datasets and analyzed a newly generated human
18 adult ileal epithelium dataset. SCORE is freely available in
19 <https://github.com/wycwycpku/RSCORE>.

20

21

1 **Results**

2 **1. The workflow of SCORE**

3 In brief, as we reason that genes execute functions through interacting with other
4 genes in molecular networks, we need to extract the molecular networks that undergo
5 dynamic changes among different cellular states (Fig. 1b). Thus, we first generate a
6 weighted gene-gene Pearson correlation network. For scRNA-seq datasets, Pearson
7 correlation coefficient can cover more authentic gene-gene relationships compared
8 with other algorithms, although it will also introduce many false positives. To reduce
9 the false positive rates, we use a curated PPI network to trim the gene-gene
10 relationships; likewise, data-irrelevant interactions of PPI network are pruned by the
11 weighted correlation network. The obtained weighted network is then decomposed
12 into numbers of small molecular networks, termed as modules, by random walk
13 algorithm. Finally, each module is scored within each cell using the AUCell
14 algorithm⁹, and a cell-module matrix is obtained, which represents the activity of
15 individual module within each cell.

16

17 Next, this cell-module matrix can be utilized to perform downstream analyses such as
18 visualization, clustering, and cell lineage analysis. Importantly, inspired by the
19 concept of Steiner Tree in graph theory, SCORE also constructs the characteristic
20 molecular interaction network (CMIN) to annotate a certain cellular state (see
21 Methods for more details). For the convenience of users, SCORE is seamlessly
22 compatible with the popular R package Seurat⁶.

23

24 **2. Evaluation of the accuracy, robustness, and efficiency of SCORE**

25 To verify the rationale and assess the performance of the algorithm, we applied
26 **SCORE** to a gold-standard scRNA-seq dataset for cell clustering analysis, with 561
27 cells derived from seven human cell lines (GSE81861) after quality control¹⁷. Two
28 experimental batches exist in the cell lines for GM12878 lymphoblastoid cells and H1
29 embryonic stem cells. In the previous literature, several clustering methods were

1 benchmarked without providing the additional information on batch identity, and the
2 reference component analysis (RCA) achieved the highest adjusted rand index (ARI)
3 of 0.91, while other methods showed inferior performance (All=HC: 0.66;
4 HiLoadG-HC: 0.53; BackSPIN: 0.64; RaceID2: 0.15; Seurat: 0.70) (Fig. 2a). Here,
5 we intend to continue the benchmarking in the same set-up, by comparing the
6 performance of SCORE with other state-of-art feature extraction methods based on
7 gene regulatory network, and highlight the strength of SCORE in the accurate
8 identification of the cell lines and the effective alleviation of undesired experimental
9 noises.

10

11 As shown in Fig. 2b, SCORE yielded the highest accuracy, with the ARI amounting to
12 1.00 (clustering by SNN¹⁸) and 0.99 (clustering by SIMLR¹⁹), respectively. Both the
13 t-Distributed Stochastic Neighbor Embedding (t-SNE) plot and SC3²⁰ similarity
14 matrix demonstrated that the batch discrepancy was removed for the same cell lines,
15 while the sharp distinction among different cell lines was still retained. In comparison,
16 the direct clustering based on expression matrix within Seurat pipeline (Raw)
17 confronted with limitations, by blurring the distinct cell identities (e.g. K352 and
18 H1437 or IMR90 cells) and discriminating the unwanted experimental batches in both
19 GM12878 and H1 cell lines (indicated by the separate clusters in t-SNE plot and
20 independent blocks in SC3 matrix in Fig. 2b). Despite that the cell-specific network
21 (CSN) method separated different cell lines, clustering variations induced by technical
22 artifacts seemed to be strengthened. Clustering guided by SCENIC continuous
23 features generated relatively satisfactory results comparable to (but still inferior than)
24 the analysis by SCORE, but with considerable more computational costs.

25

26 In addition, the performance of SCORE was robust under the selection of input
27 variable genes for analysis. As shown in Fig. 2c and Supplementary Figs. 1-4,
28 SCORE always outperformed other methods for all the selected gene numbers
29 ranging from 2,000 to 14,000, both in terms of clustering ARI and the removal of

1 experimental batch effects, as indicated by the SC3 similarity matrices and t-SNE
2 plots. As long as the gene number exceeded 2,000, SCORE could robustly recover the
3 cell line identities with the optimal resolution. In comparison, while SCENIC
4 achieved satisfactory clustering results with 5,000-11,000 genes, the clustering ARI
5 was significantly reduced for the case of 14,000 genes due to the severe batch effects.

6
7 The implementation of SCORE is also efficient, mainly due to the fast module
8 decomposition of PPI by random walk distance, and the highly parallelizable
9 calculation of module activity matrix. With 14,000 top variable genes, SCORE
10 finalized the analysis within 2 minutes, while the implementation of SCENIC lasted
11 for 6 hours in R. In fact, from the running time tests over various groups of genes (Fig.
12 2c), the implementation of SCORE was typically 10-20 times faster than CSN and
13 100-400 times faster than SCENIC. Notably, SCORE could handle the task of 15,000
14 cells and 14,000 genes in less than 10 minutes (Fig. 2d).

15
16 Overall, application of SCORE to the gold standard dataset validates our rationale-
17 the incorporation of molecular interaction network in scRNA-seq analysis can
18 facilitate the accurate cell identity dissection and reduce the technical artifacts in
19 experiments. The comparisons with other methods also suggest the capability and
20 potential of SCORE for the unbiased and efficient analysis, and integration of
21 large-scale transcriptome datasets, which will be explored below.

22

23 **3. Integration and multi-scale comparison of five human fetal datasets**

24 Using scRNA-seq techniques, many studies have uncovered unprecedented findings
25 in terms of human fetal development. However, most of these studies focused on just
26 one organ or tissue and failed to consider the human fetal development in a holistic
27 view. In addition, it is also a challenge for existing methods to integrate the datasets
28 sampled from different organs and different time points because of batch effects
29 introduced by varied dissociation protocols, organ-specific differences, etc. Given this,

1 we used SCORE to integrate five high-quality human fetal datasets previously
2 generated by our group, including fetal gonads (ovary and testis)²¹, heart²², kidney²³,
3 prefrontal cortex (PFC)²⁴, and cerebral cortex²⁵, spanning from 4 to 26 weeks of fetal
4 development. Importantly, we re-organized the five datasets using uniform pipeline
5 and format, which provided a rich and convenient resource for studying human fetal
6 development (<https://github.com/zorrodong/HECA>). To evaluate the integration result,
7 we utilized the mixability of common cell types shared by these datasets, such as
8 immune cells, erythroid cells, and endothelial cells across different organs. Notably,
9 SCORE obtained a reasonable result with cell types grouped by their own identities,
10 while Seurat pipeline on individual gene expression matrix (Raw) failed to cluster
11 these common cell types (Fig. 3a,b and Supplementary Fig. 5). We also used two
12 batch effect correction methods, Harmony²⁶ and CCA²⁷, to integrate these datasets;
13 however, they all exhibited defective results. Even worse, both Harmony and CCA
14 introduced serious artificial results and unrelated cells were wrongly grouped
15 together.

16

17 In addition, we also analyzed each of these five datasets separately (Supplementary
18 Figs. 6-10). Take the human fetal kidney dataset as an example. Standard clustering
19 using Seurat pipeline showed a chaos state, especially in renal interstitium (RI)
20 (Supplementary Fig. 6a,b). Specifically, there was a distinctly isolated cell population
21 of 19-week fetal renal cells, which might be caused by batch effects. In contrast,
22 SCORE presented a much better result. All fetal cells were first separated into renal
23 cell populations and non-renal cell populations, including immune cells (IMMs),
24 erythrocytes (ERs), and endothelial cells (EDs). Renal cells comprising glomerular
25 cells, renal capsule cells, and renal tubular cells were then classified into 9 clusters
26 according to the anatomic structure of nephron, namely, cap mesenchyme (CM),
27 podocytes (PDs), proximal tubule (PT), loop of Henle (LH), distal convoluted tubule
28 (DT), collecting duct (CD), intraglomerular mesangium (MG), extraglomerular
29 mesangium (EM), and renal interstitium (RI). Importantly, the UMAP plot using

1 SCORE displayed an accurate developmental trajectory of nephrons from the cap
2 mesenchyme to the formation of epithelial tubules (Supplementary Fig. 6c-e). For
3 other four datasets, SCORE also performed better than Seurat. In heart dataset,
4 SCORE firstly divided cardiomyocyte cells (CMCs) into compact CMCs and
5 trabecular CMs, and then portioned compact CMCs into atrial CMCs (CMCs-A) and
6 ventricular CMCs (CMCs-V), which was consistent with the anatomic structures of
7 heart (Supplementary Fig. 7). In two cortex datasets and gonad dataset, slight batch
8 effects were detected in Seurat results but not in SCORE results (Supplementary Figs.
9 8-10).

10

11 **4. Expression landscape of human adult ileal epithelium**

12 Although small intestine epithelium has been studied widely in murine, a
13 comprehensive expression landscape in human is still lacking. To explore the cellular
14 diversity of human small intestine epithelium, we sampled the ileal crypts from two
15 patients suffered from right-sided colon cancer while their ilea were relatively normal
16 (Supplementary Fig. 11a). scRNA-seq libraries were generated using the 10X v3 Kit,
17 which could guarantee both throughput and quality (Supplementary Fig. 11b).
18 Compared with fetal datasets, adult ones tend to be more complex due to more
19 individual differences. As expected, Seurat pipeline grouped cells according to the
20 patient individuals, which indicated the existence of batch effect and would result in
21 confusion and inaccuracy for downstream analyses (Fig. 4a). In contrast, SCORE
22 successfully eliminated these unwanted effects and identified 19 clusters, which
23 covered all the known cell types, namely, stem/progenitor cells, paneth cells, tuft cells,
24 goblet cells, enteroendocrine cells, microfold (M) cells, enterocytes (Fig. 4b and
25 Supplementary Table1). The differentially expressed genes (DEGs), differentially
26 activated modules (DAMs) and the inferred differentiation potency all supported the
27 clustering accuracy of SCORE (Fig. 4c,d and Supplementary Fig. 12a).

28

29 Importantly, with SCORE, we could construct the CMIN of a certain cluster to further

1 explore the relationships and interactions of crucial genes (Fig. 4e and Supplementary
2 Fig. 12b). Unlike traditional methods that focus on the differences in expression levels,
3 CMIN ranks genes based on their topological importance in the optimal steiner tree.
4 Therefore, the CMIN provides a simplified representation of original PPI network,
5 and highlights the significance of non-marker interacting genes surpassing traditional
6 marker gene analysis. Taking the cluster 1_Stem as an example, *MYC*, *CD81*, *OLFM4*,
7 *JUN*, *REGIA*, *SMOC2*, *NPM1*, *NAPILI* and several ribosome genes were inferred to
8 be crucial for maintaining the cellular state of stem/progenitor cells (Fig. 4e).
9 Moreover, the CMIN could be further divided into 5 gene groups based on their
10 topological relationship, and the enriched gene ontology (GO) terms also supported
11 the stem/progenitor identity of this cluster and improved the readability of the CMIN.

12

13 **5. The heterogeneity within stem/progenitor cells, goblet cells and M cells**

14 Despite the achievements in exploring the intestinal epithelium, a key question
15 remains puzzling: do all intestinal stem cells have the equal differentiation potency?
16 Using SCORE, we found that the stem/progenitor cells were heterogeneous and we
17 identified 6 subgroups with the expression of specific marker genes, such as *TDGF1*
18 (also known as Cripto-1), *GEM*, *FNBPI*, *ICOSLG*, and *DURAS3* (Fig. 5a,b). Cripto-1
19 plays an important role in early embryonic development, the formation and
20 progression of several human tumor types²⁸. Moreover, Cripto-1 can also interact with
21 Wnt and Notch signaling pathways, which are crucial for the maintenance of
22 intestinal stem cells²⁹. Thus, the Cripto-1+ cells were important components of ileal
23 epithelium, though their actual function remained to be explored.

24

25 As known, goblet cells secrete mucus to create a protective layer to the intestinal
26 lumen²⁹. We identified 4 subgroups of goblet cells with SCORE, namely, Goblet1-4,
27 and they exhibited two varied differentiation routes: one is from Goblet1 to Goblet4
28 via Goblet3, while the other one is from Goblet1 to Goblet2 (Fig. 5c and
29 Supplementary Fig. 13a). Of note, only one subgroup of goblet cells (8_Goblet4) is

1 responsible for the secretion of mucus as they overrepresented *TFF1*, *MUC3A*,
2 *MUC13*; while the other differentiation route, Goblet2, highly expressed genes related
3 to respiratory electron transport, such as *ATP5G1*, *COX6C*, *COX7C*. (Fig. 5d,e).

4

5 Recently, a new cell type that distinctively expresses *BEST4* and *OTOP2* was found at
6 the top of the colonic crypts, which can sense pH and transport salt, ions and metals³⁰.

7 In our dataset, we also detected this cell population and found that they were
8 differentiated from the SPIB+ M cells (Fig. 5f and Supplementary Fig. 13b).

9 Surprisingly, this *BEST4+*/*OTOP2+* M cells also highly expressed functional cystic
10 fibrosis transmembrane conductance regulator (*CFTR*) (Fig. 5g). In human colon and

11 mouse small intestine, *CFTR* is mainly expressed in the crypt cells which helps these
12 cells secrete fluid to flush the crypt lumen and remove contaminants²⁹. In human ileal

13 epithelium, the highest expression level of *CFTR* is restricted to the *BEST4+*/*OTOP2+*
14 M cells, while this is not the situation in human colonic epithelium. Hence,

15 *BEST4+*/*OTOP2+* M cells might play a different role in ileum compared with their
16 colon counterparts (Fig. 5h). Interestingly, the other M cell population highly

17 expressed *LEFTY1*, which is also expressed in colonic M cells (Fig. 5g,h). *LEFTY1*
18 encodes a secreted ligand that binds to Cripto-1 to antagonize Nodal signaling. Thus,

19 these M cells might interact with the identified Cripto-1+ cells to regulate their
20 differentiation or other cellular behaviors.

21

1 **Discussion**

2 The close cooperation of different genes forms gene modules to fulfill specific
3 cellular functions, and a certain cellular state or cell type can be well depicted by the
4 activities of various gene modules¹. As cellular states transit rapidly during
5 organismal development and differentiation, how to simulate these dynamic processes
6 and uncover the corresponding cell fates becomes a fundamental biology issue. In this
7 study, we present a new computational method, SCORE, to infer this dynamic change
8 and reveal cell development trajectory from the molecular network point of view.
9 There are several advantages of SCORE compared with currently widely used
10 methods.

11

12 Firstly, the hypothesis of SCORE is more biologically reasonable. Genes function in
13 conjunction through molecular networks rather than in isolation. However, most
14 published methods ignored the importance of molecular networks and just calculate
15 the transcriptional similarity inferred from individual genes. A novel method called
16 SCENIC utilizes transcription factor (TF) based networks to retrieve regulatory
17 activity patterns, which can be used to annotate cellular states. However, SCENIC
18 only measures the activity patterns of approximately 1500 TFs, and a recent
19 evaluation showed that the sensitivity of SCENIC is only about 5% (about 75 TFs)³¹.
20 Thus the resolution of SCENIC is limited, especially when analyzing highly
21 heterogeneous single-cell data. In contrast, SCORE infers the molecular network from
22 most of the expressed genes, and uses the network signatures rather than individual
23 genes to define a certain cellular state, which highly improves the resolution and
24 accuracy.

25

26 Secondly, SCORE can significantly reduce the false positive rates and yield more
27 accurate results. SCORE uses the curated PPI network to correct the inferred
28 gene-gene relationship, and all the interactions are literature-proved. As shown in Fig.
29 2, compared with other popular methods, SCORE possesses higher accuracy and

1 robustness.

2

3 Thirdly, SCORE can overcome dropout effects and other technical variations to
4 successfully integrate different datasets. One of the major drawbacks in scRNA-seq
5 techniques is their high dropout rates. Since SCORE performs the analyses based on
6 gene modules rather than individual genes, it can effectively reduce the unwanted
7 dropout effects. In addition, PPI corrected correlation and AUC based module score
8 are both able to eliminate the batch effects. As shown in Figs. 3 and 4, SCORE
9 successfully integrated five human fetal datasets and one human adult ileal epithelium
10 dataset, while other methods at least partially failed. Compared with current
11 main-stream data integration methods, SCORE does not introduce artificial
12 alternations to the raw gene expression matrix, nor require the input of exact batch ID
13 for different datasets. The latter feature makes SCORE particularly useful to analyze
14 time-series scRNA-seq datasets during development process, since the boundary
15 between technical variations and biologically meaningful difference in various data
16 collection time points is often blurred in such case. The subjective choice of batch ID
17 for other data integration methods may omit the true temporal heterogeneity in gene
18 expression.

19

20 Fourthly, in addition to the DEGs, SCORE can also identify the differentially
21 activated gene modules (DAMs) and infer the gene relationship through the
22 constructed CMIN for each cell fate (Fig. 4e). Moreover, due to the flexible
23 framework of SCORE, it can be easily applied to deal with other networks, such as
24 pathway network, metabolism network, etc.

25

26 Finally, the molecular network decomposition step in SCORE workflow is realized
27 with high efficiency and independent of data size. The major computational cost of
28 SCORE lies in the quantification of module activities by AUCell, which scales
29 linearly with cell numbers and has been easily paralleled in the R implementation.

1 Thus, SCORE is highly scalable and is able to cater to the increasing size of
2 scRNA-seq datasets nowadays.

3

4 However, as SCORE relies on the curated PPI network, it may not be applicable for
5 the datasets of organisms without high-confidence molecular interaction information.

6 Besides, when the available gene number is too low, the accuracy of SCORE may be
7 impaired. Thus, we highly recommend researchers apply SCORE to analyze
8 high-quality scRNA-seq datasets with greater sequencing depth.

9

10 In summary, with the high accuracy, robustness, and scalability, SCORE can help to
11 explore the scRNA-seq datasets in a more biologically reasonable manner, and gain
12 more insights into the complex biological systems.

1 **Figure legends:**

2 **Fig. 1. The workflow of SCORE**

3 **a**, SCORE assumes that the cellular state transition is associated with the
4 activation/inactivation of functional modules of molecular interaction network, which
5 can be inferred from single-cell transcriptome data.

6 **b**, In the SCORE workflow, the input protein-protein interaction (PPI) network from
7 public database and gene correlation inferred from single-cell dataset are trimmed to
8 construct a weighted molecular interaction network (WMIN). The random walk
9 approach is then applied to WMIN to decompose molecular interaction modules via a
10 consensus strategy, and the activation score of modules for each cell is calculated
11 from AUCell. Downstream analysis are performed based on the obtained cell-module
12 activity matrix to cluster and visualize the cells against technical variations, and
13 construct the characteristic molecular interaction network (CMIN) for each cellular
14 state.

15

16 **Fig. 2. The performance assessment of SCORE.**

17 **a**, Accuracy of clustering methods evaluated on the gold-standard cell-line
18 benchmarking dataset in Li, et al. ¹⁷.

19 **b**, The performance of SCORE on the cell-line benchmarking dataset was compared
20 with the direct analysis on raw expression matrix within Seurat pipeline (denoted as
21 Raw) and other gene regulatory network (GRN) based methods (i.e. CSN and
22 SCENIC with continuous/binary features). SCORE outperforms other methods in
23 terms of both clustering accuracy (indicated by Adjusted Rand Index: ARI of different
24 clustering approach SNN and SIMLR, as well as the batch removal effect revealed by
25 the SC3 similarity matrix and t-SNE plot) and running time. The number of input
26 highly variable genes is set as 5,000, and the colors in t-SNE plot denote cell line
27 identities collected in the experiments.

28 **c**, The assessment of SCORE robustness to feature selection in gold-standard dataset
29 and comparison with other methods. SCORE always achieves the highest ARI and

1 shortest running time regardless of the number of selected HVGs.

2 **d**, The scalability of SCORE with the increase of numbers of input genes and cells,
3 tested by down sampling of the integrated human fetal datasets. The running time of
4 SCORE scales almost linearly with the number of cells, and does not witness the
5 significant increases as the gene number exceeds 11,000. The implementation of
6 SCORE completes within 10 minutes with 15,000 cells and 14,000 genes.

7

8 **Fig. 3. Unbiased integration of five human fetal datasets.**

9 **a**, The UMAP visualization of different data integration results implemented by
10 SCORE, direct merge of datasets (Raw), Harmony and CCA (aligned from left to
11 right in each row). The cells are colored by organ information (the top row),
12 developmental week information (middle), and clustering results based on integrated
13 data (bottom), respectively.

14 **b**, The UMAP visualization of data integration results by SCORE and other methods,
15 with cells colored by the expression level of marker genes for immune (top), erythroid
16 (middle), and endothelial (bottom) cells. The cells that highly express the marker
17 genes are denoted by the red circles.

18

19 **Fig. 4. Expression landscape of human adult ileal dataset by SCORE.**

20 **a**, UMAP visualization of human adult ileal epithelium based on raw expression
21 matrix within Seurat pipeline. Cells are colored by patient information.

22 **b**, UMAP visualization based on SCORE. Cells are colored by patient information
23 (left) and SCORE clusters (right). TA: transit-amplifying; Mcell: microfold cell.

24 **c**, Heatmaps displaying DEGs (left) and DAMs (right) within each cluster. The color
25 key from purple to yellow denotes low to high expression levels, respectively.

26 **d**, Dotplot displaying the expression levels of representative marker genes of each
27 cluster. Spot size denotes the percentage of cells expressing the gene within each
28 cluster and colour intensity denotes the expression levels of the gene.

29 **e**, CMIN of the epithelial stem/progenitor cells using the DEGs (left) and the related

1 gene ontology terms of CMIN (right). The node sizes in CMIN represent the
2 PageRank score of each gene and the genes of top 40 PageRank scores are displayed
3 with their names in CMIN. The node colors denote the classification of genes as
4 DEGs (red) and the connecting Steiner genes (gray) in CMIN. Spot size of gene
5 ontology terms denotes the percentage of the related gene number and the color key
6 denotes the adjusted p-values.

7

8 **Fig. 5. Subgroups of stem/progenitor cells, goblet cells and M cells.**

9 **a**, UMAP visualization of subgroups of cluster 1_Stem based on SCORE. Cells are
10 colored by subgroups.

11 **b**, Violinplots displaying the expression levels of representative marker genes of
12 stem/progenitor subgroups.

13 **c**, PCA plot displaying the subgroups of goblet cells. cells are colored by subgroups.

14 **d**, Enriched GO terms using the DEGs of each goblet cell subgroup. Spot size denotes
15 the percentage of the related gene number and the color key denotes the adjusted
16 p-values.

17 **e**, Dotplot displaying the expression levels of representative marker genes of goblet
18 cell subgroups. Spot size denotes the percentage of cells expressing the gene within
19 each cluster and colour intensity denotes their expression level.

20 **f**, Violinplots displaying the expression levels of representative marker genes of two
21 M cell subgroups.

22 **g**, UMAP plots displaying the expression levels of representative marker genes across
23 all the ileal cells. The color key denotes the expression levels.

24 **h**, Violinplots displaying the expression levels of *CFTR* and *LEFTY1* in the human
25 colon dataset published recently³².

26

1 **Supplementary figure and table legends:**

2 **Supplementary Fig. 1.** The performance of SCORE on the cell-line benchmarking
3 dataset was compared with the direct analysis on raw expression matrix within Seurat
4 pipeline (denoted as Raw) and other gene regulatory network (GRN) based methods
5 (i.e. CSN and SCENIC with continuous/binary features). The number of input
6 highlight variable genes is set as 14,000, and the colors in t-SNE plot denote cell lines
7 identity collected in the experiments.

8

9 **Supplementary Fig. 2.** The performance of SCORE on the cell-line benchmarking
10 dataset was compared with the direct analysis on raw expression matrix within Seurat
11 pipeline (denoted as Raw) and other gene regulatory network (GRN) based methods
12 (i.e. CSN and SCENIC with continuous/binary features). The number of input
13 highlight variable genes is set as 11,000, and the colors in t-SNE plot denote cell lines
14 identity collected in the experiments.

15

16 **Supplementary Fig. 3.** The performance of SCORE on the cell-line benchmarking
17 dataset was compared with the direct analysis on raw expression matrix within Seurat
18 pipeline (denoted as Raw) and other gene regulatory network (GRN) based methods
19 (i.e. CSN and SCENIC with continuous/binary features). The number of input
20 highlight variable genes is set as 8,000, and the colors in t-SNE plot denote cell lines
21 identity collected in the experiments.

22

23 **Supplementary Fig. 4.** The performance of SCORE on the cell-line benchmarking
24 dataset was compared with the direct analysis on raw expression matrix within Seurat
25 pipeline (denoted as Raw) and other gene regulatory network (GRN) based methods
26 (i.e. CSN and SCENIC with continuous/binary features). The number of input
27 highlight variable genes is set as 2,000, and the colors in t-SNE plot denote cell lines
28 identity collected in the experiments.

29

1 **Supplementary Fig. 5.** Comparison of marker gene analysis based on different data
2 integration results by SCORE and other methods, with the top 10 marker genes
3 displayed for each cluster. The columns in the matrices represent cells, and the rows
4 represent genes. The color key from blue to red denotes low to high expression levels,
5 respectively. The incorrectly clustered cells, which do not significantly express
6 marker genes of the assigned cluster, are marked by the black squares.

7

8 **Supplementary Fig. 6. Performance comparison between standard Seurat**
9 **pipeline (denoted as Raw) and SCORE on human fetal kidney dataset.**

10 **a,** UMAP visualization of cell clusters of kidney dataset based on standard Seurat
11 pipeline (denoted as Raw). Cells are colored by cell types (left) and developmental
12 weeks (right).

13 **b,** Feature plots displaying the expression levels of representative marker genes. The
14 color key from grey to red denotes low to high expression levels, respectively.

15 **c,** UMAP visualization of cell clusters of kidney dataset based on SCORE. Cells are
16 colored by cell types (left) and developmental weeks (right).

17 **d,** Feature plots displaying the expression levels of representative marker genes.

18 **e,** Schematic diagram of renal tubule.

19

20 **Supplementary Fig. 7. Performance comparison between standard Seurat**
21 **pipeline (denoted as Raw) and SCORE on human fetal heart dataset.**

22 **a,** UMAP visualization of cell clusters of heart dataset based on standard Seurat
23 pipeline (denoted as Raw). Cells are colored by cell types (left) and developmental
24 weeks (right).

25 **b,** Feature plots displaying the expression levels of representative marker genes. The
26 color key from grey to red denotes low to high expression levels, respectively.

27 **c,** UMAP visualization of cell clusters of heart dataset based on SCORE. Cells are
28 colored by cell types (left) and developmental weeks (right).

29 **d,** Feature plots displaying the expression levels of representative marker genes.

1 **e**, UMAP plots of subgroups of cardiomyocytes (CMs) based on SCORE (left). Cells
2 are colored by subgroups. Feature plots displaying the expression levels of
3 representative marker genes of subgroups (right).

4

5 **Supplementary Fig. 8. Performance comparison between standard Seurat**
6 **pipeline (denoted as Raw) and SCORE on human fetal gonad dataset.**

7 **a**, UMAP visualization of cell clusters of fetal gonad dataset based on standard Seurat
8 pipeline (denoted as Raw). Cells are colored by cell types (left) and developmental
9 weeks (right).

10 **b**, Feature plots displaying the expression levels of representative marker genes. The
11 color key from grey to red denotes low to high expression levels, respectively.

12 **c**, UMAP visualization of cell clusters of fetal gonad dataset based on SCORE. Cells
13 are colored by cell types (left) and developmental weeks (right).

14 **d**, Feature plots displaying the expression levels of representative marker genes.

15

16 **Supplementary Fig. 9. Performance comparison between standard Seurat**
17 **pipeline (denoted as Raw) and SCORE on human fetal cerebral cortex dataset.**

18 **a**, UMAP visualization of cell clusters of cerebral cortex dataset based on standard
19 Seurat pipeline (denoted as Raw). Cells are colored by cell types (left) and
20 developmental weeks (right).

21 **b**, Feature plots displaying the expression levels of representative marker genes. The
22 color key from grey to red denotes low to high expression levels, respectively.

23 **c**, UMAP visualization of cell clusters of cerebral cortex dataset based on SCORE.
24 Cells are colored by cell types (left) and developmental weeks (right).

25 **d**, Feature plots displaying the expression levels of representative marker genes.

26

27 **Supplementary Fig. 10. Performance comparison between standard Seurat**
28 **pipeline (denoted as Raw) and SCORE on human fetal prefrontal cortex dataset.**

29 **a**, UMAP visualization of cell clusters of prefrontal cortex dataset based on standard

1 Seurat pipeline (denoted as Raw). Cells are colored by cell types (left) and
2 developmental weeks (right).

3 **b**, Feature plots displaying the expression levels of representative marker genes. The
4 color key from grey to red denotes low to high expression levels, respectively.

5 **c**, UMAP visualization of cell clusters of prefrontal cortex dataset based on SCORE.
6 Cells are colored by cell types (left) and developmental weeks (right).

7 **d**, Feature plots displaying the expression levels of representative marker genes.
8

9 **Supplementary Fig. 11. Quality control of two human adult ileal samples.**

10 **a**, Histological sections of the two sampled human adult ileal tissues.

11 **b**, 10X dataset quality control information.
12

13 **Supplementary Fig. 12. Characterizations of the human adult ileal cell types.**

14 **a**, Boxplot displaying the differentiation potency inferred by the SCENT algorithm.

15 **b**, CMINs of representative clusters using the DEGs. The node sizes in CMIN
16 represent the PageRank score of each gene and the genes of top 20 PageRank scores
17 are displayed with their names in CMIN. The node colors denote the classification of
18 genes as DEGs (red) and the connecting Steiner genes (gray) in CMIN.
19

20 **Supplementary Fig. 13. Expression patterns of goblet cell and microfold cell
21 subgroups.**

22 **a**, Heatmap displaying the DEGs within each goblet cell subgroup. The top 10 marker
23 genes were listed on the right. The color key from purple to yellow denotes low to
24 high expression levels, respectively.

25 **b**, Heatmap displaying the DEGs within each microfold cell subgroup. The top 10
26 marker genes were listed on the right.
27

28 **Supplementary Table1. Cell information and DEGs of human adult ileal dataset.**
29

1 **Methods**

2 **Overview of SCORE**

3 For the convenience of users, SCORE is seamlessly compatible with the Seurat
4 pipeline. The input for SCORE workflow includes a single-cell gene expression
5 matrix and a PPI network. Both the cells and genes in the expression matrix could be
6 pre-filtered by users, while it is highly recommended that enough number of genes
7 (>5000) should be retained to achieve more robust analysis of SCORE, as shown via
8 the gold-standard cell line dataset (see main text and SI). The nodes of PPI network
9 should overlap considerably with the gene names in the expression matrix, and the
10 edges of the network shall represent the corresponding molecular interactions with
11 relatively high confidence. In the R implementation of SCORE, the procedure
12 supports the automatic download of PPI from public database such as Biogrid and
13 STRING, which is recommended in the standard workflow. The output of SCORE is a
14 cell-module matrix representing the activity of individual dynamic module within
15 each cell, which can be utilized for downstream visualization, clustering and cell
16 lineage analysis.

17

18 **Dissection of dynamic molecular networks**

19 To simultaneously reduce the false-positive interactions from correlation inference,
20 and prune data-irrelevant interactions of PPI network, SCORE constructs a weighted
21 molecular interaction network (WMIN) $\mathcal{G}(V, E, W)$ by combining the data-driven
22 and knowledge-based approaches. The vertex set V of the network only consists of
23 the genes in the input expression matrix, and the edge set E obtained from the input
24 PPI network. The weight W_{ij} on the edge E_{ij} is the Pearson correlation coefficients
25 between the corresponding nodes i and j calculated from the single-cell gene
26 expression matrix. To improve the interpretability of molecular networks and
27 highlight the co-expression features, by default we only keep the edges with positive
28 weight in the WMIN.

29

1 To achieve fast and robust identification of dynamic molecular networks, SCORE
2 utilizes the consensus detection of the weighted network community through random
3 walk approach. Given the weighted network $\mathcal{G}(V, E, W)$, a random walk on the
4 network is naturally induced, whose transition probability matrix (TPM) P is defined
5 by

$$P_{ij} = \frac{W_{ij}}{d_i}, \quad d_i = \sum_{j \in \mathcal{N}(i)} W_{ij},$$

6 Where $\mathcal{N}(i)$ denotes the neighbors of the node i . SCORE constructs an ensemble of
7 random walks with different step lengths on the WMIN, with the TPMs
8 $P^{l_1}, P^{l_2}, \dots, P^{l_R}$, where the power of matrix l_m denotes the length of time step. The
9 walktrap algorithm is applied to detect network community for each random walk,
10 respectively. The algorithm partitions the network in terms of the distance induced by
11 the random walk, based on the intuition that the random walker will be “trapped” in
12 the closely-connected sub-networks (termed as modules). All modules with molecule
13 number larger than 3 in each run will be kept as the final dynamic modules, resulting
14 in the module set $\{M_k\}_{k=1}^L$. It would be possible that certain modules occur
15 repeatedly in different runs of walktrap algorithm, indicating their stability to form a
16 closely-connected community. SCORE strengthens the weight of such modules
17 automatically in the subsequent analyses, by restoring all the modules without
18 deletion of repeated items.

19

20 **Quantification of the module activity**

21 For each detected module, SCORE utilizes AUCCell to quantify its activation level
22 within each individual cell. Given cell x , genes are ranked in descending order
23 according to their expression level in x . By default, the z-score is adopted in SCORE
24 to rank the genes in order to remove the effect of scaling. The recovery curve (ROC)
25 for module M_k is then derived by counting the top ranked genes enriched in M_k .
26 The activity measure $A_k(x)$ of module M_k in cell x (consists of the final output
27 matrix of SCORE) is defined as the area under the curve (AUC) for the top ranked

1 genes. Intuitively, modules with higher activity in the biochemical process tend to
2 possess the high-ranking gene expression level, therefore associate with higher
3 activity measure. The AUCell procedure is independent of gene expression unit or
4 normalization method, therefore achieving the effective removal of batch effect in the
5 single-cell experiments.

6

7 **Downstream Analysis**

8 The obtained module activity $A_k(x)$ matrix from SCORE (whose rows represent
9 modules and columns represent cells) can replace the raw gene expression matrix as
10 the input for downstream analysis, such as dimension reduction, clustering and
11 lineage inference. As shown in various datasets of the main text, the downstream
12 analysis based on SCORE module activity features outperforms the raw expression
13 matrix, in terms of clustering accuracy, development lineage trend, and removal of
14 experimental batch effects. Hence, the workflow of SCORE can be understood as the
15 extraction of biologically meaningful and robust features, guided by molecular
16 interactions in the single-cell transcriptome data.

17

18 For the convenience of downstream analysis, the R implementation of SCORE is
19 deeply fused with the workflow of Seurat v3.0 package. The input expression matrix
20 to SCORE can be a Seurat object with RNA assay, and the output module activity
21 features are returned as the Net assay in the same Seurat object. Users may
22 conveniently conduct dimension reduction and cell-clustering based on the Net assay,
23 and perform marker gene analysis based on the RNA assay, by switching the default
24 assay of Seurat object.

25

26 **Construction of Cell State-Specific Characteristic Molecular Interaction** 27 **Network (CMIN)**

28 To annotate a certain cell state, beyond the marker genes, **SCORE** constructs the
29 CMIN with the concept of Steiner Tree in graph theory. Given a graph $G = (V, E)$

1 and a subset of vertices $T \subset V$ (called terminal vertices), a Steiner tree $S \subset G$ is a
2 connected tree that spans through the given terminal vertices T . The Steiner tree S
3 may contain vertices not presented in T , known as the Steiner vertices, serving as the
4 interchange node to connect the vertices in T . In molecular interaction network, the
5 marker genes of a certain cell state are typically selected as the terminal vertices, and
6 the mediating Steiner vertices, although not necessarily differentially expressed, are
7 supposed to play important roles in formulating the specific cell state through
8 molecular interactions. Therefore, the Steiner tree provides a simplified representation
9 of original PPI network, and highlights the significance of non-marker interacting
10 genes surpassing traditional marker gene analysis.

11

12 In the downstream analysis of **SCORE**, given the specific cell cluster identified from
13 module activation features, we first construct the set of terminal genes T^* by
14 detecting the marker genes from two different levels,

15 1) Union of all genes in the **SCORE**-extracted modules that are differentially
16 activated in the cluster, denoted as differentially activated module genes (DAMGs);

17 2) Individual genes that are significantly up-regulated in the cluster, denoted as
18 differentially expressed genes (DEGs).

19 While the DEGs are commonly referred as the “markers” of cellular states, the
20 DAMGs also represent the key molecular interaction modules to mark the cell cluster
21 in the network resolution.

22

23 Next, to infer the CMINs that possibly formulate the cellular states rather than the
24 genes solely marking the cellular states, we propose to calculate a Steiner tree S^* that
25 spans the terminal gene set T^* with some optimal property, defined on the constructed
26 WMIN $\mathcal{G}(V, E, W)$ by **SCORE** in the first step of workflow. We require that the
27 separate DAMGs or DEGs in S^* are linked by the most relevant genes, as well as
28 through the most likely interaction path derived from the dataset. To this end, we
29 define the distance D_{ij} on the edge E_{ij} of WMIN by $D_{ij} = 1/(W_{ij} + \varepsilon)$, where

1 the small number ε is added to avoid zero in dominator. Highly correlated gene pairs
2 in PPI tend to possess much closer distance from the definition. Then S^* can be
3 optimized as the Steiner tree with the least sum of edge distances, which can be
4 tackled efficiently by the greedy algorithm.

5
6 For a better visualization of CMIN, in R implementation of **SCORE** we mark
7 DAMGs, DEGs, and Steiner connecting genes with different colors, and also use the
8 PageRank algorithm to measure the topological importance of the genes in optimal
9 Steiner tree S^* as shown by the size of the nodes. We can also provide any two
10 genesets to construct the CMIN.

11

12 **Settings in the benchmarking gold-standard dataset**

13 The gene expression matrix was processed as fragments per kilobase per million reads
14 (FPKM) as in the original literature. To test the robustness of different methods, we
15 first adopted the vst method in Seurat v3.0 package to select five groups of highly
16 variable genes (with the number of genes 2,000, 5,000, 8,000, 11,000, and 14,000,
17 respectively). We performed SCORE, as well as two other network or biological
18 information based methods, CSN and SCENIC to further compress and extract the
19 features, respectively, from the groups of highly variable genes (HVGs).

20

21 For SCORE, we downloaded the *Homo sapiens* PPI network (version 3.5.173) from
22 the BioGRID database. The top ranked genes included in the calculation of AUC
23 values varied with the sizes of input HVGs, with 250 and 200 for 2,000 and 5,000
24 variable genes, respectively, and 400 in other cases. The parameters in implementing
25 CSN and SCENIC were chosen with default values. The SCENIC yields both
26 continuous and binary features as the outputs.

27

28 The running time comparison was conducted on the 2.50GHz Xeon E5-2680 machine
29 with 128G RAM, 12 cores and Linux OS. For SCORE and SCENIC, the CPU core

1 number was set as 10. CSN was automatically paralleled with MATLAB 2019b. The
2 wall time of implementing each procedure were recorded as the running time in the
3 main text.

4

5 In the downstream clustering analysis, three methods, SNN, SIMLR and SC3 were
6 performed on the extracted features by different methods, and the adjusted rand index
7 (ARI) as well as the similarity matrix of SC3 were used to evaluate the accuracy and
8 the effect of batch removal. The direct analysis on raw expression matrix with
9 selected HVGs was also performed for the comparison. The true labels were the seven
10 collected cell lines identity (H1, GM12878, A549, HCT116, H1437, K562 and IMR90)
11 without batch information. For SNN, we tuned the resolution parameter to obtain the
12 optimal ARI value. As to SC3 and SIMLR, we set the number of clusters to 7. The
13 t-SNE plot was produced based on the top 10 principal components of the extracted
14 features.

15

16 **Human fetal datasets**

17 To evaluate the integration performance of SCORE, five human fetal datasets were
18 collected from our previously published studies, including fetal gonads (overies and
19 testis) (GEO number: GSE86146), heart (GEO number: GSE106118), kidney (GEO
20 number: GSE109488), prefrontal cortex (PFC) (GEO number: GSE104276), and
21 cerebral cortex (GEO number: GSE103723), spanning from 4 to 26 weeks of fetal
22 development. Importantly, we re-organized the five datasets using uniform pipeline
23 and format, which provided a rich and convenient resource for studying human fetal
24 development (<https://github.com/zorrodong/HECA>). In brief, barcode and UMI
25 information were extracted by UMI-tools form raw reads³³. After discarding the poly
26 A bases, TSO sequences and low-quality sequences, the clean reads were mapped to
27 GRCh38 reference using STAR aligner³⁴. We used featureCounts³⁵ to annotate the
28 mapped reads and quantified the UMI counts through UMI-tools. We provided the
29 pipeline for users (<https://github.com/zorrodong/HECA>).

1

2 To analyze the human fetal datasets, we first discard cells with gene number below
3 1,000 and UMI counts below 10,000. HVGs were chosen using Seurat (mean \geq 0.1,
4 dispersion \geq 0.1), and about 8,000 HVGs were selected for SCORE to perform the
5 evaluation. Two batch effect correction methods CCA and Harmony were used with
6 the recommended parameters. 30-50 reduced dimensions were used to perform
7 UMAP analysis and clustering using the graph-based method in Seurat. To accelerate
8 the speed, we used genesortR to conduct the differential expression analysis.

9

10 **Human adult ileal crypt dataset**

11 This study was approved by the Ethics Committee of Peking University Third
12 Hospital (License No. M2016170). All patients had signed written informed consent.
13 Ileal crypt samples were collected from 2 right-sided colon cancer patients
14 immediately after surgical resection. We dissociate the samples into single cells and
15 constructed the libraries using 10x Genomics (3' Library, Kit v3). Libraries were
16 sequenced using Illumina HiSeq 4000 platform with 150-bp paired-end reads.

17

18 We used Cellranger v3.1.0 (10X Genomics) to deal with the raw reads and quantify
19 the expression level. Next, the UMI count matrix were analyzed using Seurat pipeline.
20 We discarded cells with gene number below 1,000, UMI counts below 1,000, and
21 mitochondrial percentage above 30%. 8,000 HVGs were chosen for SCORE analysis
22 using FindVariableFeatures(nfeatures = 8000,selection.method = "vst"). The overall
23 dimensionality reduction and clustering were performed using all the obtained
24 modules.

25

26 **Data availability**

27 Human adult small intestine dataset are deposited in the GEO. SCORE is freely
28 available in <https://github.com/wycwycpku/RSCORE>. The five human fetal datasets
29 are available in <https://github.com/zorrodong/HECA>.

1

2

3 **Acknowledgements**

4 This project was supported by grants from the National Natural Science Foundation of
5 China (31625018 and 81521002 to F.T., 11825102 and 11421101 to T.L.).

6

7 **Author contributions**

8 F.T., T.L., W.F., J.D., P.Z. conceived the project; W.W., X.Z., Y.G., L.W., performed
9 the experiments; J.D., P.Z., Y.W., Y.C., H.X., J.L., J.Y., X.N.Z. conducted the
10 bioinformatics analyses; J.D., P.Z., T.Li, F.T., wrote the manuscript with the help of all
11 the authors.

12

13 **Competing interests**

14 The authors declare no competing financial interests.

1 **References:**

- 2 1 Arendt, D. *et al.* The origin and evolution of cell types. *Nature Reviews*
3 *Genetics* **17** (2016).
- 4 2 Achim, K. & Arendt, D. Structural evolution of cell types by step-wise
5 assembly of cellular modules. *Current Opinion in Genetics & Development* **27**,
6 102-108 (2014).
- 7 3 Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell.
8 *Nature Methods* **6**, 377-382 (2009).
- 9 4 Klein, A. M. *et al.* Droplet Barcoding for Single-Cell Transcriptomics Applied
10 to Embryonic Stem Cells. *Cell* **161**, 1187-1201 (2015).
- 11 5 Macosko, E. Z. *et al.* Highly Parallel Genome-wide Expression Profiling of
12 Individual Cells Using Nanoliter Droplets. *Cell* **161**, 1202-1214 (2015).
- 13 6 Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**,
14 1888-1902 e1821, doi:10.1016/j.cell.2019.05.031 (2019).
- 15 7 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene
16 expression data analysis. *Genome biology* **19**, 15,
17 doi:10.1186/s13059-017-1382-0 (2018).
- 18 8 Guo, M., Wang, H., Potter, S. S., Whitsett, J. A. & Xu, Y. SINCERA: A
19 Pipeline for Single-Cell RNA-Seq Profiling Analysis. *PLoS Comput Biol* **11**,
20 e1004575, doi:10.1371/journal.pcbi.1004575 (2015).
- 21 9 Aibar, S. *et al.* SCENIC: single-cell regulatory network inference and
22 clustering. *Nat Methods* **14**, 1083-1086, doi:10.1038/nmeth.4463 (2017).
- 23 10 Fan, J. *et al.* Characterizing transcriptional heterogeneity through pathway and
24 gene set overdispersion analysis. *Nature Methods* **13**, 241-+,
25 doi:10.1038/Nmeth.3734 (2016).
- 26 11 Dai, H., Li, L., Zeng, T. & Chen, L. Cell-specific network constructed by
27 single-cell RNA sequencing data. *Nucleic Acids Res* **47**, e62,
28 doi:10.1093/nar/gkz172 (2019).
- 29 12 Kiselev, V. Y., Andrews, T. S. & Hemberg, M. Challenges in unsupervised

- 1 clustering of single-cell RNA-seq data. *Nature Reviews Genetics* **20**, 273-282,
2 doi:10.1038/s41576-018-0088-9 (2019).
- 3 13 Oughtred, R. *et al.* The BioGRID interaction database: 2019 update. *Nucleic*
4 *Acids Research* **47**, D529-D541, doi:10.1093/nar/gky1079 (2019).
- 5 14 Szklarczyk, D. *et al.* STRING v11: protein-protein association networks with
6 increased coverage, supporting functional discovery in genome-wide
7 experimental datasets. *Nucleic Acids Research* **47**, D607-D613,
8 doi:10.1093/nar/gky1131 (2019).
- 9 15 Teschendorff, A. E. & Enver, T. Single-cell entropy for accurate estimation of
10 differentiation potency from a cell's transcriptome. *Nat Commun* **8**, 15599,
11 doi:10.1038/ncomms15599 (2017).
- 12 16 Ronen, J. & Akalin, A. netSmooth: Network-smoothing based imputation for
13 single cell RNA-seq. *F1000Research* **7**, 8,
14 doi:10.12688/f1000research.13511.3 (2018).
- 15 17 Li, H. *et al.* Reference component analysis of single-cell transcriptomes
16 elucidates cellular heterogeneity in human colorectal tumors. *Nature genetics*
17 **49**, 708-718, doi:10.1038/ng.3818 (2017).
- 18 18 Xu, C. & Su, Z. C. Identification of cell types from single-cell transcriptomes
19 using a novel clustering method. *Bioinformatics* **31**, 1974-1980,
20 doi:10.1093/bioinformatics/btv088 (2015).
- 21 19 Wang, B., Zhu, J., Pierson, E., Ramazzotti, D. & Batzoglou, S. Visualization
22 and analysis of single-cell RNA-seq data by kernel-based similarity learning.
23 *Nat Methods* **14**, 414-416, doi:10.1038/nmeth.4207 (2017).
- 24 20 Kiselev, V. Y. *et al.* SC3: consensus clustering of single-cell RNA-seq data.
25 *Nat Methods* **14**, 483-486, doi:10.1038/nmeth.4236 (2017).
- 26 21 Li, L. *et al.* Single-Cell RNA-Seq Analysis Maps Development of Human
27 Germline Cells and Gonadal Niche Interactions. *Cell stem cell* **20**, 891-892,
28 doi:10.1016/j.stem.2017.05.009 (2017).
- 29 22 Cui, Y. *et al.* Single-Cell Transcriptome Analysis Maps the Developmental

- 1 Track of the Human Heart. *Cell reports* **26**, 1934-1950 e1935,
2 doi:10.1016/j.celrep.2019.01.079 (2019).
- 3 23 Wang, P. *et al.* Dissecting the Global Dynamic Molecular Profiles of Human
4 Fetal Kidney Development by Single-Cell RNA Sequencing. *Cell reports* **24**,
5 3554-3567 e3553, doi:10.1016/j.celrep.2018.08.056 (2018).
- 6 24 Zhong, S. *et al.* A single-cell RNA-seq survey of the developmental landscape
7 of the human prefrontal cortex. *Nature* **555**, 524-528,
8 doi:10.1038/nature25980 (2018).
- 9 25 Fan, X. *et al.* Spatial transcriptomic survey of human embryonic cerebral
10 cortex by single-cell RNA-seq analysis. *Cell research* **28**, 730-745,
11 doi:10.1038/s41422-018-0053-3 (2018).
- 12 26 Korsunsky, I. *et al.* Fast, sensitive, and accurate integration of single cell data
13 with Harmony. *bioRxiv*, 461954, doi:10.1101/461954 (2018).
- 14 27 Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating
15 single-cell transcriptomic data across different conditions, technologies, and
16 species. *Nature biotechnology* **36**, 411-420, doi:10.1038/nbt.4096 (2018).
- 17 28 Bianco, C. *et al.* Role of Cripto-1 in Stem Cell Maintenance and Malignant
18 Progression. *American Journal of Pathology* **177**, 532-540 (2010).
- 19 29 Gehart, H. & Clevers, H. Tales from the crypt: new insights into intestinal
20 stem cells. *Nature Reviews Gastroenterology & Hepatology* **16**, 19-34,
21 doi:10.1038/s41575-018-0081-y (2019).
- 22 30 Parikh, K. *et al.* Colonic epithelial cell diversity in health and inflammatory
23 bowel disease. *Nature* **567**, 49-55, doi:10.1038/s41586-019-0992-y (2019).
- 24 31 Wang, N. & Teschendorff, A. E. Leveraging high-powered RNA-Seq datasets
25 to improve inference of regulatory activity in single-cell RNA-Seq data.
26 *bioRxiv*, 553040, doi:10.1101/553040 (2019).
- 27 32 Smillie, C. S. *et al.* Intra- and Inter-cellular Rewiring of the Human Colon
28 during Ulcerative Colitis. *Cell* **178**, 714-730.e722,
29 doi:https://doi.org/10.1016/j.cell.2019.06.029 (2019).

- 1 33 Smith, T., Heger, A. & Sudbery, I. UMI-tools: modeling sequencing errors in
2 Unique Molecular Identifiers to improve quantification accuracy. *Genome*
3 *Research* **27**, 491-499, doi:10.1101/gr.209601.116 (2017).
- 4 34 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**,
5 15-21, doi:10.1093/bioinformatics/bts635 (2013).
- 6 35 Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose
7 program for assigning sequence reads to genomic features. *Bioinformatics* **30**,
8 923-930, doi:10.1093/bioinformatics/btt656 (2014).
- 9

Fig. 1

bioRxiv preprint doi: <https://doi.org/10.1101/699959>; this version posted October 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

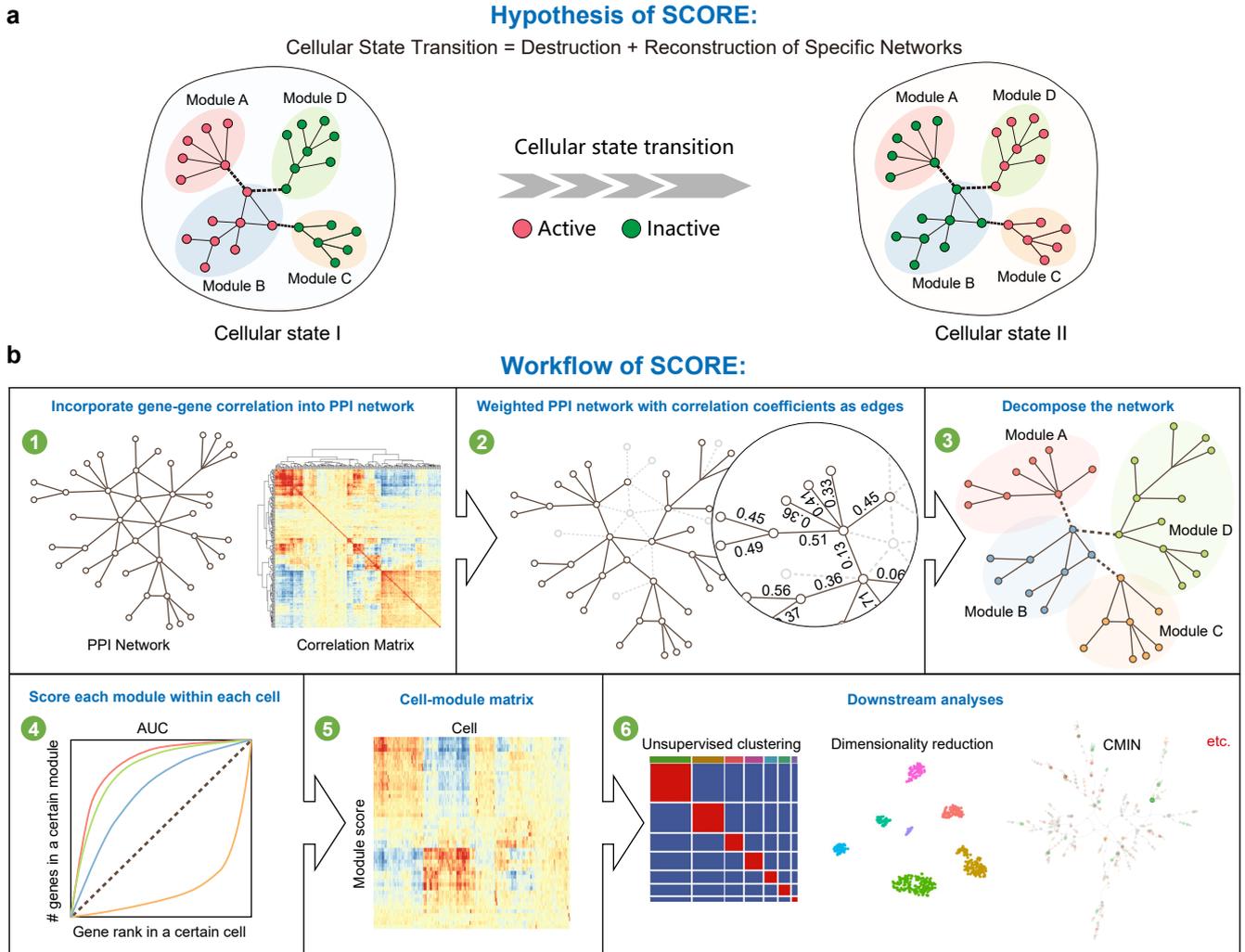


Fig. 2

bioRxiv preprint doi: <https://doi.org/10.1101/699959>; this version posted October 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

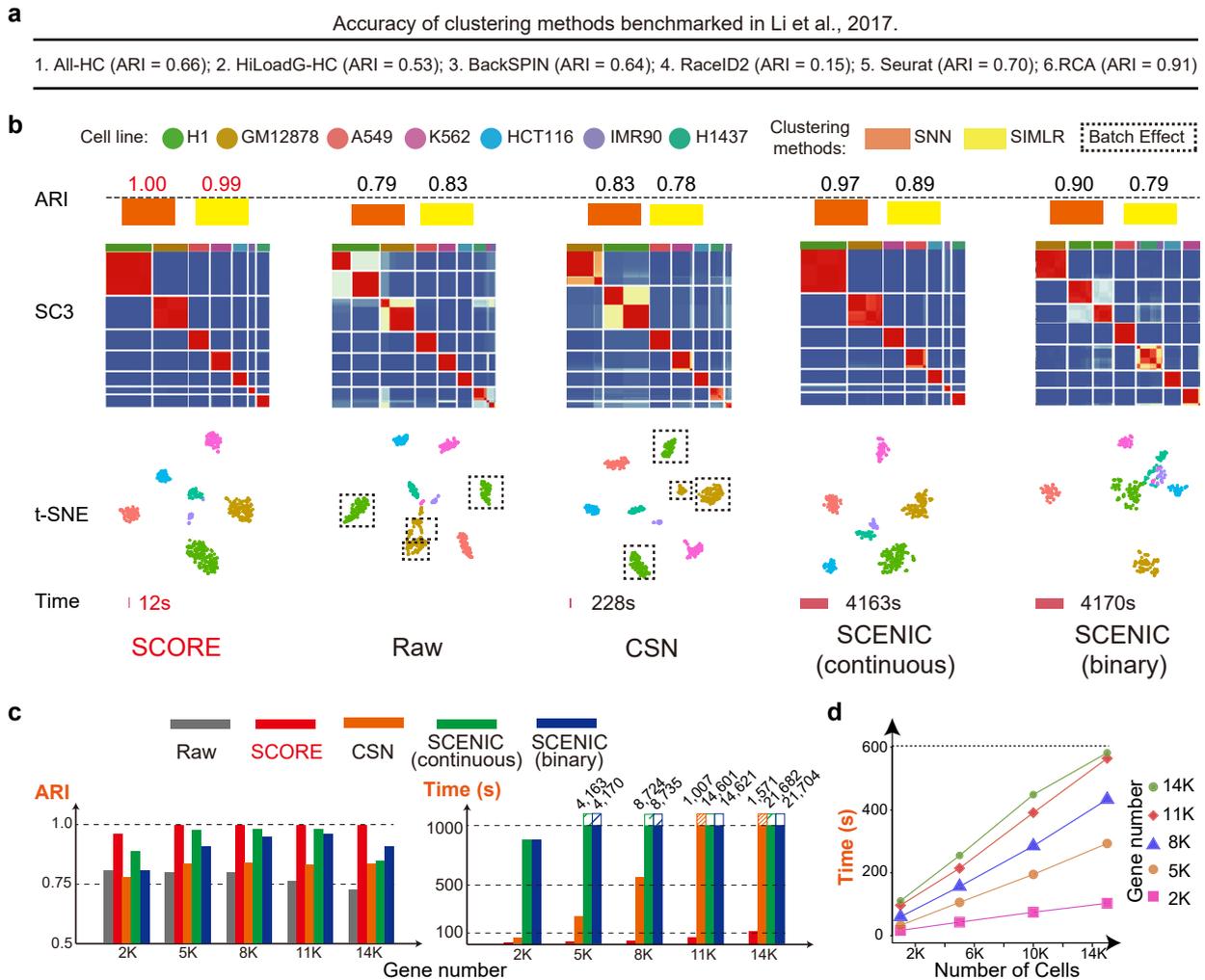


Fig. 3

bioRxiv preprint doi: <https://doi.org/10.1101/699959>; this version posted October 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

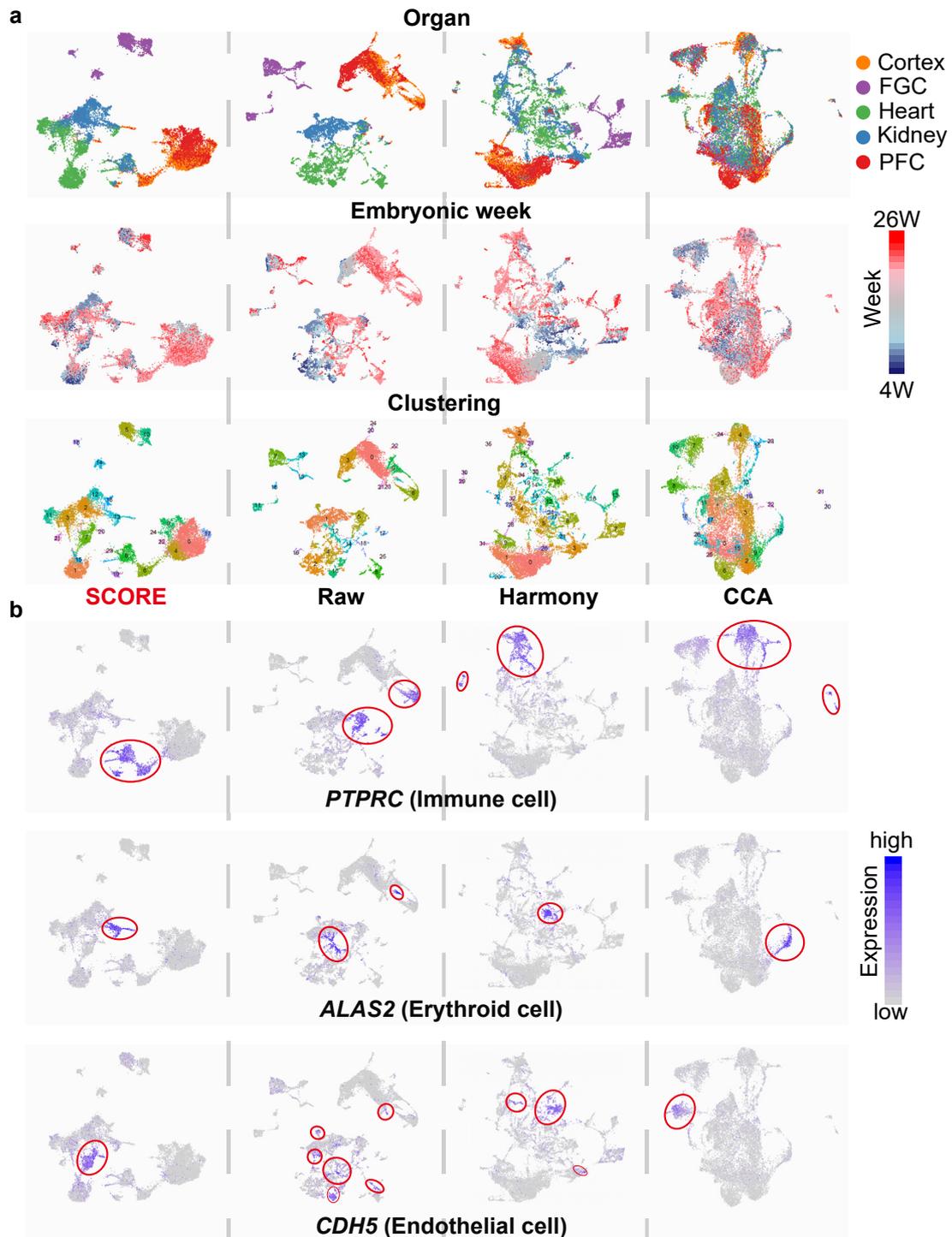


Fig. 4

bioRxiv preprint doi: <https://doi.org/10.1101/699959>; this version posted October 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

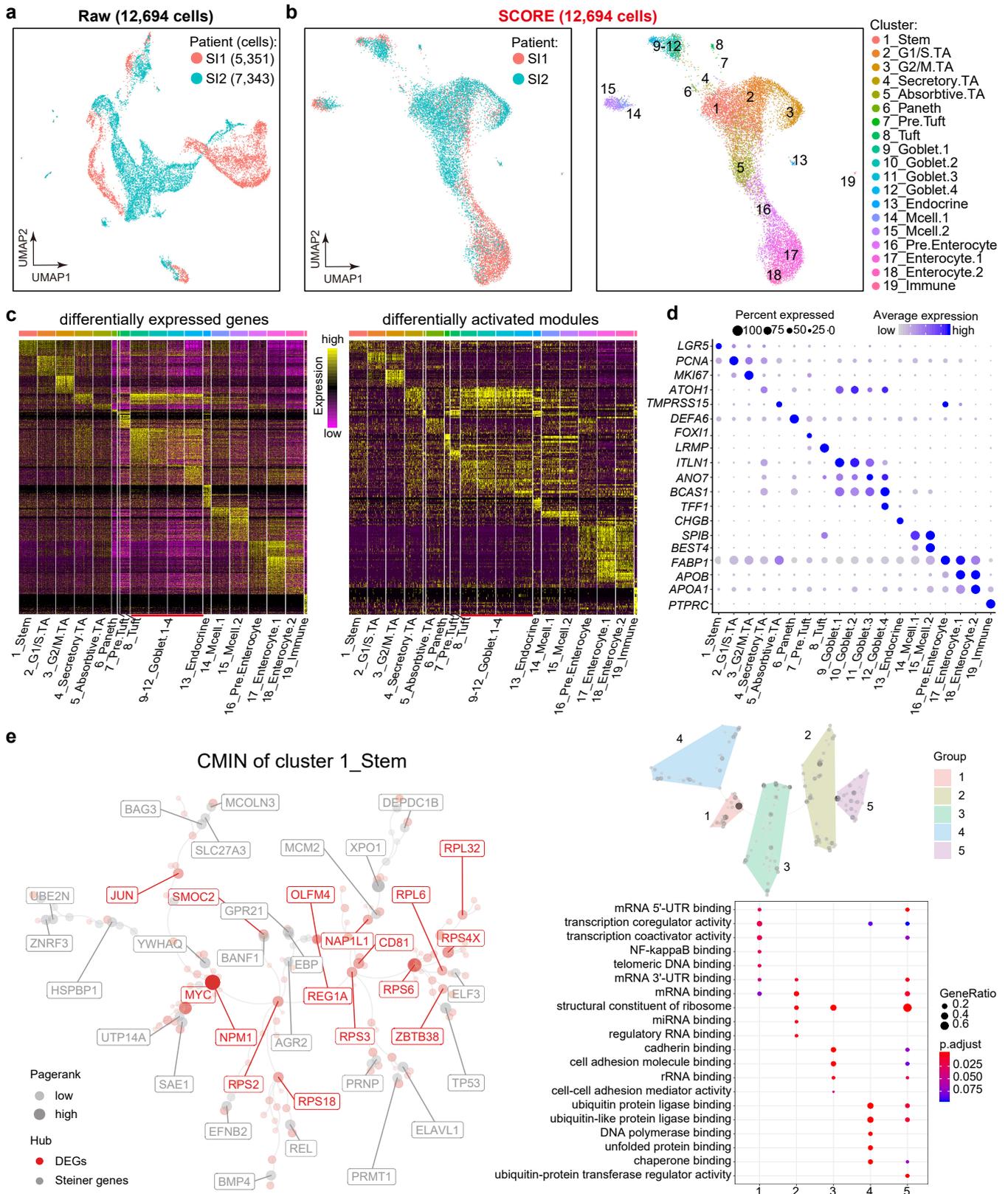


Fig. 5

bioRxiv preprint doi: <https://doi.org/10.1101/699959>; this version posted October 15, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

