1   **CCSN: Single Cell RNA Sequencing Data Analysis by**

2   **Conditional Cell-specific Network**

3        **Lin Li[1,#,a], Hao Dai[1,2,#,b], Zhaoyuan Fang[1,*,c], Luonan Chen[1,2,3,4,*,d]**

4•   *[1]Key Laboratory of Systems Biology, CAS Center for Excellence in Molecular Cell*
5    *Science, Institute of Biochemistry and Cell Biology, Shanghai Institutes for Biological*
6    *Sciences, Chinese Academy of Sciences, University of Chinese Academy of Sciences,*
7    *Shanghai, 200031, China*
8•   *[2]Shanghai Research Center for Brain Science and Brain-Inspired Intelligence,*
9    *Shanghai, 201210, China*
10•  *[3]CAS Center for Excellence in Animal Evolution and Genetics, Chinese Academy of*
11   *Sciences, Kunming, 650223, China*
12•  *[4]School of Life Science and Technology, ShanghaiTech University, Shanghai, 201210,*
13   *China*
14•
15•  *To whom correspondence should be addressed.
16•  E-mail: lnchen@sibs.ac.cn (Luonan Chen), fangzhaoyuan@sibs.ac.cn (Zhaoyuan
17   Fang).
18•
19•  [#]These authors contributed equally to this paper as the first authors.

20

21

22

23

## Abstract

The rapid advancement of single cell technologies has shed new light on the complex mechanisms of cellular heterogeneity. However, compared with bulk RNA sequencing (RNA-seq), single-cell RNA-seq (scRNA-seq) suffers from higher noise and lower coverage, which brings new computational difficulties. Based on statistical independence, cell-specific network (CSN) is able to quantify the overall associations between genes for each cell, yet suffering from a problem of overestimation related to indirect effects. To overcome this problem, we propose the "conditional cell-specific network" (CCSN) method, which can measure the direct associations between genes by eliminating the indirect associations. CCSN can be used for cell clustering and dimension reduction on a network basis of single cells. Intuitively, each CCSN can be viewed as the transformation from less "reliable" gene expression to more "reliable" gene-gene associations in a cell. Based on CCSN, we further design network flow entropy (NFE) to estimate the differentiation potency of a single cell. A number of scRNA-seq datasets were used to demonstrate the advantages of our approach: (1) one direct association network for one cell; (2) most existing scRNA-seq methods designed for gene expression matrices are also applicable to CCSN-transformed degree matrices; (3) CCSN-based NFE helps resolving the direction of differentiation trajectories by quantifying the potency of each cell. CCSN is publicly available at http://sysbio.sibcb.ac.cn/cb/chenlab/soft/CCSN.zip.


**KEYWORDS:** Single cell analysis; Network flow entropy; Cell-specific network; Single cell network; Direct association; Conditional independence

## Introduction

49

50   With the development of high-throughput single-cell RNA sequencing (scRNA-seq),

51   novel cell populations in complex tissues [1-5] can be identified and the differentiation

52   trajectory of cell states [6-8] can be obtained, which opens a new way to understand the

53   heterogeneity and transition of cells [9-11]. However, compared to traditional bulk

54   RNA-seq data, the prevalence of high technical noise and dropout events is a major

55   problem in scRNA-seq [12-17], which raises substantial challenges for data analysis.

56   Many computational methods were proposed to improve the identification of new cell

57   types [18-21]. Meanwhile, imputation is an effective strategy to transform the dropouts

58   to the substituted values [22-26]. However, most of these methods mainly analyze

59   mRNA expression/concentrations, while the information of gene-gene interactions (or

60   their network) is ignored.

61        Recently, a network-based method, cell-specific network (CSN), was proposed to

62   perform network analysis for scRNA-seq data [27], which elegantly infers a network

63   for each cell and successfully transforms the noisy and "unreliable" gene expression

64   data to the more "reliable" gene association data. The network degree matrix (NDM)

65   derived from CSN can be further applied in downstream single cell analyses, which

66   performs better than traditional expression-based methods in terms of robustness and

67   accuracy. CSN is able to identify the dependency between two genes from single-cell

68   data based on statistical independence. However, CSN suffers from a problem of

69   overestimation on gene-gene associations, which include both direct and indirect

70   associations due to interactive effects from other genes in a network. In other words, a

71   gene pair without direct association can be falsely identified to have a link just because

72   they both have true associations with some other genes. Thus, the gene-gene network

73   of a cell constructed by CSN may be much denser than the real molecular network in

74   this cell, in particular when there are many complex associations among genes.

75       To overcome this shortcoming of CSN, we introduce a novel computational method

76   to construct a conditional cell-specific network (CCSN) from scRNA-seq data.

77   Specifically, CCSN identifies direct associations between genes by filtering out indirect

78   associations in the gene-gene network based on conditional independence. Thus, CCSN

79   can transform the original gene expression data of each cell to the direct and robust

80   gene-gene association data (or network data) of the same cell. In this paper, we first

81   demonstrate that the transformed gene-gene association data not only are fully

82   compatible with traditional analyses such as dimension reduction and clustering, but

83   also enable us to delineate the cell-specific network topology and its dynamics along

84   developmental trajectories. Then, by defining the network flow entropy (NFE) on the

85   gene-gene association data of each cell based on CCSN, we estimate the differentiation

86   potency of individual cells. We show that NFE can illustrate the lineage dynamics of

87   cell differentiation by quantifying the differentiation potency of cells, which is also one

88   of the most challenging tasks in developmental biology.

89

## 90   **Methods**

91   Assuming that $x$ and $y$ are two random variables, and $z$ is the third random variable. If

92   $x$ and $y$ are independent, then

93   $$p(x)p(y) = p(x,y) \tag{1}$$

94   where $p(x,y)$ is the joint probability distribution of $x$ and $y$; $p(x)$ and $p(y)$ are the

95   marginal probability distributions of $x$ and $y$.

96   If $x$ and $y$ with the condition $z$ are conditionally independent, then

97   $$p(x|z)p(y|z) = p(x,y|z) \tag{2}$$

98   where $p(x,y|z)$ is the joint probability distribution of $x$ and $y$ with the condition $z$,

4

99 $p(x|z)$ and $p(y|z)$ are conditionally marginal probability distributions. Note that

100 eqns. (1)-(2) are both necessary and sufficient conditions on mutual independence and

101 conditional independence, respectively.

102 Here we define

$$\rho_{xy} = p(x,y) - p(x)p(y). \quad (3)$$

$$\rho_{xy|z} = p(x,y|z) - p(x|z)p(y|z). \quad (4)$$

105 The original CSN method [27] uses $\rho_{xy}$ to distinguish the independency and

106 association between $x$ and $y$ (File S1 Note 1). However, if two independent variables $x$

107 and $y$ are both associated with a third random variable $z$, $\rho_{xy}$ cannot measure the direct

108 independency because there is an indirect association between $x$ and $y$. In other words,

109 the associations defined by CSN or eqn. (3) include both direct and indirect dependency,

110 thus resulting in the overestimation on gene-gene associations. To overcome this

111 problem of CSN, we develop a novel method, conditional cell-specific network

112 (CCSN), which measures the direct gene-gene associations based on the conditional

113 independency $\rho_{xy|z}$, i.e. eqn. (4), by filtering out the indirect associations in the

114 reconstructed network. The computational framework of CCSN is shown in **Figure 1**,

115 and is described in the next sections.

116 **Probability distribution estimation**

117 We numerically estimate the value of $\rho_{xy|z}$ by making a scatter diagram based on gene

118 expression data. Suppose there are m genes and n cells in the data. We depict the

119 expression values of gene $x$, gene $y$ and the conditional gene $z$ in a three-dimensional

120 space (Figure S1 A-G), where each dot represents one cell. First, we draw two parallel

121 planes which are orthogonal with $z$ axis near the dot $k$ to represent the upper and lower

122 bounds of the neighborhoods of $z_k$. And the number of dots in the space between the

123    two parallel planes (i.e. the neighborhood of $z_k$) is $n_z^{(k)}$ (Figure S1 D). Now we get a

124    subspace on condition of gene z. Then, we draw other four planes near the dot $k$, where

125    two planes are orthogonal with $x$ axis and the other two planes are orthogonal with $y$

126    axis. We can get the neighborhoods of $(x_k, z_k)$, $(y_k, z_k)$ and $(x_k, y_k, z_k)$ according

127    to the intersection space of six planes (Figure S1 E-G), where the numbers of dots

128    are $n_{xz}^{(k)}, n_{yz}^{(k)}$ and $n_{xyz}^{(k)}$, respectively. Then, we can get the estimation of probability

129    distributions:

130    $$p^{(k)}(x, y|z) \approx \frac{n_{xyz}^{(k)}}{n_z^{(k)}}, \qquad p^{(k)}(x|z) \approx \frac{n_{xz}^{(k)}}{n_z^{(k)}}, \qquad p^{(k)}(y|z) \approx \frac{n_{yz}^{(k)}}{n_z^{(k)}}$$

131    Based on eqn. (4), we construct a statistic

132    $$\rho_{xy|z}^{(k)} = \frac{n_{xyz}^{(k)}}{n_z^{(k)}} - \frac{n_{xz}^{(k)} n_{yz}^{(k)}}{n_z^{(k)2}} \tag{5}$$

133    to measure the conditional independence between gene $x$ and gene $y$ on the condition

134    of gene $z$ in cell $k$. And when gene $x$ and gene $y$ given gene z are conditionally

135    independent, the expectation $\mu_{xy|z}^{(k)}$ and standard deviation $\sigma_{xy|z}^{(k)}$ (File S1) of the

136    statistic $\rho_{xy|z}^{(k)}$ can be obtained:

137    $$\mu_{xy|z}^{(k)} = 0$$

138    $$\sigma_{xy|z}^{(k)} = \sqrt{\frac{n_{xz}^{(k)} n_{yz}^{(k)} \cdot \left(n_z^{(k)} - n_{xz}^{(k)}\right) \cdot (n_z^{(k)} - n_{yz}^{(k)})}{n_z^{(k)4}\left(n_z^{(k)} - 1\right)}}$$

139    Then, we normalize the statistic as

$$\hat{\rho}_{xy|z}^{(k)} = \frac{\rho_{xy|z}^{(k)} - \mu_{xy|z}^{(k)}}{\sigma_{xy|z}^{(k)}} \tag{6}$$

140    If gene $x$ and $y$ are conditionally independent on the condition of gene $z$, it can be proved

141    that the normalized statistic follows the standard normal distribution (File S1 Note 1

6

142    and Figure S2), and it is less than or equal to 0 when gene $x$ and $y$ are conditionally

143    independent (File S1 Note 2).

144

**Constructing conditional cell-specific network for each cell**

146    To estimate the conditional independency of gene $x$ and gene $y$ given the conditional

147    gene $z$ in cell $k$, we use the following hypothesis test:

148    $H_0(null\ hypothesis)$: gene $x$ and gene $y$ are conditionally independent given gene $z$ in

149    cell $k$.

150    $H_1(alternative\ hypothesis)$: gene $x$ and gene $y$ are conditionally dependent given

151    gene $z$ in cell $k$.

152    If $\hat{\rho}^{(k)}_{xy|z_g}$, the normalized statistic, is larger than $\mathcal{N}_\alpha$ (significance level α, $\mathcal{N}_\alpha$ is the

153    alpha quantile of the standard normal distribution), the null hypothesis will be rejected

154    and then $\omega^{(k)}_{xy|z} = 1$ ($\omega^{(k)}_{xy|z}$ is the edge weight of genes $x$ and $y$ on condition of gene $z$ ).

$$\omega^{(k)}_{xy|z} = \begin{cases} 1 & \text{gene x and y are directly dependent given gene z} \\ 0 & \text{gene x and y are conditionally independent given gene z} \end{cases} \qquad (7)$$

155    All gene pairs can be tested if they are conditionally independent given gene z in cell

156    k. And the conditional cell-specific network (CCSN) $C^{(k)}_z$ given conditional gene $z$ is

157    obtained for cell $k$.

158        Then, to estimate the direct association between a pair of genes in a cell,

159    theoretically we should use all the remaining m-2 genes as conditional genes, which is

160    computationally intensive. Suppose there are m genes in our analysis, then m*(m-1)/2

161    gene pairs should be tested. Fortunately, a molecular network is generally sparse, which

162    means that a pair of genes (i.e. genes x and y) are expected to have a very small number

163    of commonly interactive genes (as conditional genes z). In other words, numerically we

164    can use a small number of conditional genes to identify the direct association between

165    a pair of genes in a cell, which can significantly reduce the computational cost (File S1

166    Note 3, Table S1). For each gene pair in a cell, we choose G ($1 \leq G \leq m-2$) genes as the

167    conditional genes to test if the gene pair is conditionally independent or not. Generally,

168    the conditional genes may be the key regulatory genes in a biological process, such as

169    transcription factors and kinases. From a network viewpoint, these genes are usually

170    hub genes in the gene-gene network, and the network degrees of these genes would be

171    higher.

172        Practically, the conditional genes could be obtained from many available methods,

173    such as highly expressed genes, highly variable genes, key transcription factor genes,

174    or the hub genes in the CSN, and so on. For the CCSN method, the conditional gene

175    sets were defined by CSN. The following two steps were used to obtain the conditional

176    genes although other appropriate schemes can also be used:

177    1. For a given cell, we first construct a CSN without the consideration of conditional

178    genes, where the edge between gene x and gene y in cell k is determined by the

179    following hypothesis test:

180    $H_0$(*null hypothesis*):  gene x and gene y are independent in cell k.

181    $H_1$(*alternative hypothesis*):  gene x and gene y are dependent in cell k.

182    The statistic $\rho_{xy}$ can be used to measure the independency of genes x and y (File S1

183    Note1). If $\rho_{xy}$ is larger than a significant level, we will reject the null hypothesis and

184    edge$_{xy}$ (k) = 1, otherwise edge$_{xy}$ (k) = 0.

$$edge_{xy}^{(k)} = \begin{cases} 1 & gene\ x\ and\ y\ are\ dependent \\ 0 & gene\ x\ and\ y\ are\ independent \end{cases}$$

185

186    Then we use $D_z^{(k)}$ to measure the importance of conditional gene z in cell k:

$$D_z^{(k)} = \sum_{y=1, y \neq z}^{M} edge_{zy}^{(k)} \tag{8}$$

187    Eqn. (8) means that if a gene is connected to more other genes, this gene is more

188    important.

189    2. For a given cell k, we choose the top $G$ $(G \geq 1)$ largest 'importance' genes as the

190    conditional genes.

191        We assume that the conditional gene set is $\{z_g, g = 1,2,3,\cdots,G\}$, and the

192    conditional cell-specific network (CCSN) $C_{z_g}^{(k)}$ is obtained for cell $k$ given conditional

193    gene $z_g$. The CCSNs of the cell k on the condition of gene set $\{z_g, g = 1,2,3,\cdots,G\}$

194    are $\{C_{z_1}^{(k)}, C_{z_2}^{(k)}, \cdots, C_{z_G}^{(k)}\}$. Then, we use

195
$$\bar{C}_k = \frac{1}{G}\sum_{g=1}^{G} C_{z_g}^{(k)} = \left(c_{ij}^{(k)}\right) \tag{9}$$

196    to represent the degrees of gene-gene interaction network of cell $k$, where $c_{ij}^{(k)}$ for

197    $i, j = 1, \cdots, m$ is the (i,j) element of the matrix $\bar{C}_k$.

198        For scRNA-seq data with all $n$ cells, we can construct $n$ CCSNs, which can be used

199    for further dimension reduction and clustering. In other words, instead of the originally

200    measured gene expression data with $n$ cells, we use the $n$ transformed CCSNs for further

201    analysis. In addition, each CCSN is a network for a cell, which can be used for network

202    analysis (gene regulations and network biomarkers) on the basis of a single cell.

203    **Network degree matrix from CCSN**

204    CCSNs could be used for various biological studies by exploiting the gene-gene

205    conditional association network from a network viewpoint. We transform eqn. (9) to a

206    conditional network degree vector based on the following transformation

207
$$v_{ik} = \sum_{j=1}^{m} c_{ij}^{(k)} \tag{10}$$

208    Then, for $\{\bar{C}_1, \bar{C}_2, \cdots, \bar{C}_n\}$, an $m*n$ matrix CNDM is obtained.

209        CNDM $= (v_{ik})$   with $i = 1, \cdots, m; k = 1, \cdots, n$ $\qquad$ (11)

210        The matrix has the same dimension with the gene expression matrix (GEM), i.e.

211    GEM$=(x_{ik})$ with $i = 1, \cdots, m; k = 1, \cdots, n$, but CNDM can reflect the gene-gene

212  direct association in terms of interaction degrees. Moreover, this CNDM matrix after

213  normalization could be further analyzed by most traditional scRNA-seq methods for

214  dimension reduction and clustering analysis. The input/output settings as well as

215  application fields of our CCSN method are listed in File S1 Note 4.

216  **Network analysis of CCSN**

217  The relationship between gene pairs can be obtained by CCSN at a single cell level.

218  CCSN also provides a new way to build gene-gene interaction network for each cell.

219  And the CNDM derived from CCSN can be further used in dimension reduction,

220  clustering and network flow entropy analysis by many existing methods.

221  *Dimension reduction*

222  We used principal component analysis (PCA) [28] and t-distributed stochastic neighbor

223  embedding (t-SNE) [29] which respectively represent linear and nonlinear methods, to

224  perform dimension reduction on public scRNA-seq datasets with known cell types.

225  *Clustering*

226  To validate the good performance of CCSN in clustering analysis, several traditional

227  clustering methods such as K-means, Hierarchical cluster analysis, and K-medoids

228  were applied to clustering analysis. Furthermore, state-of-the-art scRNA-seq data

229  clustering methods such as SC3, SIMLR and Seurat [20, 30, 31] were also used for

230  comparison.

231  *Network flow entropy analysis.*

232  Quantifying the differentiation potency of a single cell is one of the important tasks in

233  scRNA-seq studies [15, 32, 33]. A recent study developed SCENT [34], which uses

234  protein-protein interaction (PPI) network and gene expression data as input to obtain

235  the potency of cells. However, SCENT depends on the PPI network, which may ignore

236  many important relationships between genes in specific cells. In this paper, we

237    developed network flow entropy (NFE) to estimate the differentiation potency of a cell

238    from its CSN or CCSN network, which is constructed for each cell. The normalized

239    gene expression profile and CSN/CCSN is used when we compute the network flow

240    entropy.

241       Estimating NFE requires a background network, which could be provided by CSN

242    or CCSN. Based on CSN or CCSN, we could know whether or not there is an edge

243    between gene i and gene j. We assume that the weight of an edge between gene $i$ and

244    gene $j$, $p_{ij}$ is proportional to the normalized expression levels of gene $i$ and gene $j$, that

245    is $p_{ij} \propto x_i x_j$ with $\sum_{j=1}^{m} p_{ij} = 1$. These weights are interpreted as interaction

246    probabilities. Then, we normalize the weighted network as a stochastic matrix, $P=(p_{ij})$

247    with

248    
$$p_{ij} = \frac{x_j}{\sum_{k \in E(i)} x_k} = \frac{x_j}{(Ax)_i} \qquad \text{for i, j=1, ..., m}$$

249    where $E(i)$ contains the neighbours of gene $i$, and $A$ is the CSN or CCSN ($A_{ij} = 1$ if

250    $i$ and $j$ are connected, otherwise $A_{ij} = 0$).

251    And then, we define the NFE as:

252    
$$\text{NFE} = -\sum_{i,j} x_i p_{ij} \log(x_i p_{ij}) \qquad (12)$$

253    where $x_i$ is the normalized gene expression of gene $i$. From the definition, NFE is

254    clearly different from network entropy.

255

256    **Data availability**

257    Twelve scRNA-seq datasets and one bulk RNA-seq dataset [15, 35-41] were used to

258    validate our CCSN method. The numbers of cells in these datasets range from 100 to

259    20,000. Table S2 gives a brief introduction of these datasets.

260

261 **Results**

262 **Visualization and clustering of scRNA-seq datasets with CNDM**

263 Characterizing the cell heterogeneity is one of the important tasks for scRNA-seq data

264 analysis. To test whether CCSN-transformed network data can help segregate cell types,

265 we performed dimension reduction and clustering on the CNDMs of gold-standard

266 scRNA-seq datasets, using algorithms widely employed in scRNA-seq studies. The

267 numbers of conditional genes used in CCSN construction were listed in Table S2.

268 For visualizing the structure of these datasets in a two-dimensional space, we used

269 the representative linear and nonlinear dimension reduction methods, principle

270 component analysis (PCA) [42] and t-distributed stochastic neighbor embedding (t-

271 SNE) [29], respectively. As shown in **Figure 2** and Figure S3, CNDMs can separate

272 different cell types clearly in the low-dimensional space by both PCA and t-SNE.

273 Notably, they generally perform even better than GEM (Figure 2, Figure S3). Hence,

274 the network data of CNDMs contain sufficient information for separating cell types in

275 scRNA-seq datasets.

276 To quantitatively evaluate the power of CNDMs in cell type identification, we

277 performed clustering on CNDMs and computed the adjusted random index (ARI) for

278 each dataset based on the background truth (File S1 Note 5). As shown in **Table 1** and

279 Figure S4, CNDMs perform obviously better than GEM on all datasets, either without

280 or with dimension reduction with t-SNE. These provide a strong support of the notion

281 that the CCSN-transformed network data are highly informative for characterizing

282 single cell populations. Interestingly, when further compared to NDM, CNDMs also

283 show a good performance (**Table 2** and Figure S5).

284 To further evaluate CCSN for larger datasets, the Tabula Muris droplet1 dataset [41]

12

285  comprising more than 20,000 cells from three tissues (bladder, trachea, and spleen)
286  were tested. The Seurat package was used to perform dimension reduction and
287  clustering analysis on the CNDM [31]. The cells are clearly segregated into three
288  dominant groups on the t-SNE map, which are largely defined by their cell origins (ARI
289  = 0.73 and Figure S6). This indicates that CCSN can be effectively extended to larger
290  datasets in addition to the relatively small gold standard datasets benchmarked above.

291  **CCSN reveals network structure and dynamics on a single cell basis**

292  In this paper, we apply CCSN to Wang dataset [39], which comes from a study of neural
293  progenitor cells (NPCs) that differentiate into mature neurons. The dataset contains six
294  time points over a 30-day period.

295  The CSN and CCSN are performed on a single cell (Day 0, RHB1742_d0) using
296  195 transcription factors which are differentially expressed across all the cell
297  subpopulations and all time points. In CCSN, two genes (*HMGB1* and *SOX11*) of high
298  coefficients of variation (CV) are chosen as the conditional genes. The results (**Figure
299  3**A) illustrate that the network of CCSN are much sparser than the network of CSN.
300  There are three modules in the CCSN, while there is only one dense network in the
301  CSN. Furthermore, three hub genes are obtained in three modules in CCSN. One of the
302  hub genes is *ASCL1* which plays an important role in neural development [13, 43]. Thus,
303  by removing indirect associations, CCSN can extract a more informative network
304  structure than CSN, which could improve the characterization of key regulatory factors
305  in individual cells.

306  CCSN also reveals the network dynamics over the differentiation trajectory. As
307  illustrated in Figure 3B, a core neural differentiation network composed of eight
308  regulatory genes is dynamically modulated through the temporal progression of NPC
309  differentiation. At day 0, the associations among these genes are the strongest,

310    consistent with the high potency of progenitor cells. As NPC differentiates, the network

311    becomes much sparser, suggesting more specified cell fates. In addition, when

312    constructing CCSN from all genes, the degrees of *MEIS2*, *PBX1* and *POU3F2* are also

313    larger in day 0 and quickly decreases afterwards (Figure 3C), indicating that these genes

314    are highly connected with other genes in NPCs, consistent with their known important

315    roles in early differentiation of neural progenitor cells [39].

316    Both theoretically and computationally, CCSN can also construct a gene-gene

317    network for a single bulk RNA-seq sample, in addition to a single cell. To validate this

318    biologically, we apply CCSN to the TCGA lung adenocarcinoma (LUAD) RNA-seq

319    dataset. The t-SNE plot based on CNDM reveals two obvious clusters, which

320    respectively corresponding to normal adjacent lung tissues and lung tumors (Figure

321    S8A), supporting the effective application of CCSN to bulk RNA-seq data as well.

322    Moreover, the EGFR pathway, a well-known oncogenic driver pathway for LUAD [44-

323    46], is densely connected in tumor samples but not in benign tissues, as illustrated in

324    the representative single-sample EGFR networks (Figure S8 B), and the CCSN degrees

325    of EGF and EGFR in each normal and tumor samples (Figure S8 C). These data

326    demonstrate that CCSN well extends to single sample bulk RNA-seq data analysis and

327    uncovers important biological connections related to disease states.

328    **CCSN-based network flow entropy analysis**

329    To quantify the differentiation state of cells, we further develop a new method "network

330    flow entropy" (NFE) to estimate the differentiation potency of cells by exploiting the

331    gene-gene network constructed by CCSN.

332    To assess the performance of NFE, we apply it to two datasets. In Wang dataset [39],

333    there are 484 cells with 6 stages (day 0, day 1, day 5, day 7, day 10, day 30) and the

334    CCSNs with one conditional gene are used to compute the network flow entropy. We

335    compared NPC (at Day 0 and Day 1) with mature neurons (at Day 30) (**Figure 4**A). In

336    Yang dataset [38], we compared the cells in day 10 with day 17 in differentiation of

337    mouse hepatoblasts (Figure 4B) and the CSN was used to compute the network flow

338    entropy. In both datasets, NFE assigns significantly higher scores to the progenitors

339    than the differentiated cells (one-sided Wilcox rank sum test, p-value = 3.756e-19 in

340    Yang dataset, p-value = 2.062e-12 on Wang dataset).

341        To further validate NFE, we generated a three-dimensional representation of the

342    cell-lineage trajectory for the Wang dataset. In the time-course differentiation

343    experiment of NPCs into neurons [39], NFE correctly predicted a gradual decrease in

344    differentiation potency (**Figure 5**). Therefore, NFE is effectively applicable to single

345    cell differentiation studies and highly predictive of developmental states and directions.

346

## Discussion

348    Estimating functional gene networks from noisy single cell data has been a challenging

349    task. Motivated by network-based data transformation, we have previously developed

350    CSN to uncover cell-specific networks and successfully applied it to extract

351    biologically important gene interactions. However, CSN does not distinguish direct and

352    indirect associations and thus suffers from the so-called overestimation problem. In this

353    study, we propose a more sophisticated approach termed CCSN, which constructs

354    direct gene-gene associations (network) of each cell by eliminating false connections

355    introduced by indirect effects.

356        CCSN can transform GEM to CNDM for downstream dimension reduction and

357    clustering analysis. These allow us to identify cell populations, generally better than

358    GEM in the datasets tested above. In addition, CCSN also shows good performance

359    when compared to CSN. Moreover, we can construct one direct gene-gene association

360     network by one cell based on CCSN. From the networks of the individual cells, we can

361     obtain the dynamically changed networks. In Figure 3C, the CCSNs of these cells

362     dynamically changed at different time points, and the network at day 0 shows the

363     strongest associations. Moreover, the hub genes of the networks constructed by CCSN

364     method may play an important role in biological processes. In Figure 3A, the hub genes

365     of three modules in the network constructed by CCSN play a vital role in neural

366     development. These clearly demonstrate the advantages of CCSN. Furthermore, we

367     develop a new NFE index which can accurately estimate the differentiation potency of

368     a single cell. And the results show that NFE performs well in distinguishing various

369     cells of differential potency.

370     Nonetheless, the computational cost of CCSN generally increases by G times

371     comparing with the original CSN due to G conditional genes. Thus, a parallel

372     computation scheme is desired to reduce the computation time. Also, CCSN is not

373     designed to construct the causal gene association networks, and the directions of the

374     gene associations cannot be obtained. These could be our future research topics.

## Author Contributions

376     L.L. and H.D. developed the methodology. L.L. executed the experiments. Z.Y.F.

377     helped the experiments and provided technical support. L.L., H.D., Z.Y.F. and L.N.C.

378     wrote and revised the manuscript. L.N.C. and Z.Y.F. supervised the work and critically

379     reviewed the paper. All authors have read and approved the final manuscript.

380

## Competing Interests

382     The authors have declared no competing interests.

383

## Acknowledgements

## References

[1] Yan L, Yang M, Guo H, Yang L, Wu J, Li R, et al. Single-cell RNA-Seq profiling of human preimplantation embryos and embryonic stem cells. Nat Struct Mol Biol 2013;20:1131-9.

[2] Treutlein B, Brownfield DG, Wu AR, Neff NF, Mantalas GL, Espinoza FH, et al. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. Nature 2014;509:371-5.

[3] Zeisel A, Munoz-Manchado AB, Codeluppi S, Lonnerberg P, La Manno G, Jureus A, et al. Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. Science 2015;347:1138-42.

[4] Fuzik J, Zeisel A, Mate Z, Calvigioni D, Yanagawa Y, Szabo G, et al. Integration of electrophysiological recordings with single-cell RNA-seq data identifies neuronal subtypes. Nat Biotechnol 2016;34:175-83.

[5] Scialdone A, Tanaka Y, Jawaid W, Moignard V, Wilson NK, Macaulay IC, et al. Resolving early mesoderm diversification through single-cell expression profiling. Nature 2016;535:289-93.

[6] Bendall SC, Davis KL, Amir el AD, Tadmor MD, Simonds EF, Chen TJ, et al. Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. Cell 2014;157:714-25.

[7] Nestorowa S, Hamey FK, Pijuan Sala B, Diamanti E, Shepherd M, Laurenti E, et al. A single-cell resolution map of mouse hematopoietic stem and progenitor cell differentiation. Blood 2016;128:e20-31.

[8] Woyke T, Doud DFR, Schulz F. The trajectory of microbial single-cell sequencing. Nat Methods 2017;14:1045-54.

[9] Jaitin DA, Kenigsberg E, Keren-Shaul H, Elefant N, Paul F, Zaretsky I, et al. Massively Parallel Single-Cell RNA-Seq for Marker-Free Decomposition of Tissues into Cell Types. Science 2014;343:776-9.

[10] Shapiro E, Biezuner T, Linnarsson S. Single-cell sequencing-based technologies will revolutionize

416    whole-organism science. Nature Reviews Genetics 2013;14:618-30.

417    [11] Grun D, Lyubimova A, Kester L, Wiebrands K, Basak O, Sasaki N, et al. Single-cell messenger

418    RNA sequencing reveals rare intestinal cell types. Nature 2015;525:251-5.

419    [12] Kuznetsov VA, Knott GD, Bonner RF. General statistics of stochastic process of gene expression in

420    eukaryotic cells. Genetics 2002;161:1321-32.

421    [13] Kim JK, Marioni JC. Inferring the kinetics of stochastic gene expression from single-cell RNA-

422    sequencing data. Genome Biol 2013;14:R7.

423    [14] Kharchenko PV, Silberstein L, Scadden DT. Bayesian approach to single-cell differential expression

424    analysis. Nature Methods 2014;11:740-U184.

425    [15] Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, et al. Computational

426    analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations

427    of cells. Nat Biotechnol 2015;33:155-60.

428    [16] Daigle BJ, Soltani M, Petzold LR, Singh A. Inferring single-cell gene expression mechanisms using

429    stochastic simulation. Bioinformatics 2015;31:1428-35.

430    [17] Vu TN, Wills QF, Kalari KR, Niu NF, Wang LW, Rantalainen M, et al. Beta-Poisson model for

431    single-cell RNA-seq data analyses. Bioinformatics 2016;32:2128-35.

432    [18] Tang H, Zeng T, Chen L. High-Order Correlation Integration for Single-Cell or Bulk RNA-seq Data

433    Analysis. Front Genet 2019;10:371.

434    [19] Jiang H, Sohn LL, Huang H, Chen L. Single cell clustering based on cell-pair differentiability

435    correlation and variance analysis. Bioinformatics 2018;34:3684-94.

436    [20] Kiselev VY, Kirschner K, Schaub MT, Andrews T, Yiu A, Chandra T, et al. SC3: consensus

437    clustering of single-cell RNA-seq data. Nat Methods 2017;14:483-6.

438    [21] Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S. Visualization and analysis of single-cell

439    RNA-seq data by kernel-based similarity learning. Nat Methods 2017;14:414-6.

440    [22] Huang M, Wang J, Torre E, Dueck H, Shaffer S, Bonasio R, et al. SAVER: gene expression recovery

441    for single-cell RNA sequencing. Nat Methods 2018;15:539-42.

442    [23] Li WV, Li JJ. An accurate and robust imputation method scImpute for single-cell RNA-seq data.

443    Nat Commun 2018;9:997.

444    [24] van Dijk D, Sharma R, Nainys J, Yim K, Kathail P, Carr AJ, et al. Recovering Gene Interactions

445    from Single-Cell Data Using Data Diffusion. Cell 2018;174:716-+.

446    [25] Gong WM, Kwak IY, Pota P, Koyano-Nakagawa N, Garry DJ. DrImpute: imputing dropout events

447    in single cell RNA sequencing data. Bmc Bioinformatics 2018;19.

448    [26] Talwar D, Mongia A, Sengupta D, Majumdar A. AutoImpute: Autoencoder based imputation of

449    single-cell RNA-seq data. Scientific Reports 2018;8.

450    [27] Dai H, Li L, Zeng T, Chen L. Cell-specific network constructed by single-cell RNA sequencing data.

451    Nucleic Acids Res 2019.

452    [28] Jolliffe IT, Cadima J. Principal component analysis: a review and recent developments. Philos Trans

453    A Math Phys Eng Sci 2016;374:20150202.

454    [29] van der Maaten L, Hinton G. Visualizing Data using t-SNE. Journal of Machine Learning Research

455    2008;9:2579-605.

456    [30] Wang B, Ramazzotti D, De Sano L, Zhu J, Pierson E, Batzoglou S. SIMLR: A Tool for Large-Scale

457    Genomic Analyses by Multi-Kernel Learning. Proteomics 2018;18.

458    [31] Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, 3rd, et al. Comprehensive

459    Integration of Single-Cell Data. Cell 2019;177:1888-902 e21.

460    [32] MacArthur BD, Lemischka IR. Statistical Mechanics of Pluripotency. Cell 2013;154:484-9.

461    [33] Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in single-cell

462    transcriptomics. Nat Rev Genet 2015;16:133-45.

463    [34] Teschendorff AE, Enver T. Single-cell entropy for accurate estimation of differentiation potency

464    from a cell's transcriptome. Nat Commun 2017;8:15599.

465    [35] Kolodziejczyk AA, Kim JK, Tsang JC, Ilicic T, Henriksson J, Natarajan KN, et al. Single Cell RNA-

466    Sequencing of Pluripotent States Unlocks Modular Transcriptional Variation. Cell Stem Cell

467    2015;17:471-85.

468    [36] Chu LF, Leng N, Zhang J, Hou Z, Mamott D, Vereide DT, et al. Single-cell RNA-seq reveals novel

469    regulators of human embryonic stem cell differentiation to definitive endoderm. Genome Biol

470    2016;17:173.

471    [37] Kim KT, Lee HW, Lee HO, Song HJ, Jeong da E, Shin S, et al. Application of single-cell RNA

472    sequencing in optimizing a combinatorial therapeutic strategy in metastatic renal cell carcinoma.

473    Genome Biol 2016;17:80.

474    [38] Yang L, Wang WH, Qiu WL, Guo Z, Bi E, Xu CR. A single-cell transcriptomic analysis reveals

475    precise pathways and regulatory mechanisms underlying hepatoblast differentiation. Hepatology

476    2017;66:1387-401.

477    [39] Wang J, Jenjaroenpun P, Bhinge A, Angarica VE, Del Sol A, Nookaew I, et al. Single-cell gene

478    expression analysis reveals regulators of distinct cell subpopulations among developing human neurons.

479    Genome Res 2017;27:1783-94.

480    [40] Gokce O, Stanley GM, Treutlein B, Neff NF, Camp JG, Malenka RC, et al. Cellular Taxonomy of

481    the Mouse Striatum as Revealed by Single-Cell RNA-Seq. Cell Reports 2016;16:1126-37.

482    [41] Tabula Muris C, Overall c, Logistical c, Organ c, processing, Library p, et al. Single-cell

483    transcriptomics of 20 mouse organs creates a Tabula Muris. Nature 2018;562:367-72.

484    [42] Baglama J, Reichel L (2005), 'Augmented implicitly restarted Lanczos bidiagonalization methods',

485    *Siam Journal on Scientific Computing*, pp. 19-42.

486    [43] Ming GL, Song H. Adult neurogenesis in the mammalian brain: significant answers and significant

487    questions. Neuron 2011;70:687-702.

488    [44] Ohsaki Y, Tanno S, Fujita Y, Toyoshima E, Fujiuchi S, Nishigaki Y, et al. Epidermal growth factor

489    receptor expression correlates with poor prognosis in non-small cell lung cancer patients with p53

490    overexpression. Oncol Rep 2000;7:603-7.

491    [45] Nicholson RI, Gee JM, Harper ME. EGFR and cancer prognosis. Eur J Cancer 2001;37 Suppl 4:S9-

492    15.

493    [46] Sharma SV, Bell DW, Settleman J, Haber DA. Epidermal growth factor receptor mutations in lung

494    cancer. Nat Rev Cancer 2007;7:169-81.

495

496

497

498

499

500

501

502

503

504

505

506

507    **Figure legends**

508    **Figure 1    Overview of CCSN**

509    The input data is gene expression matrix $E_{m*n}$ (The orange column represents the cell

510    k). (1) The normalized statistics $\hat{\rho}_{xy|z}^{(k)}$ of each gene pair gene x and gene y given a

511    conditional gene z for each cell k. $\hat{\rho}_{xy|z}^{(k)}$ can be used to measure the direct gene-gene

512    associations. (2) Hypothesis testing of the normalized statistic $\hat{\rho}_{xy|z}^{(k)}$. The significance

513    level of hypothesis testing is $\alpha$ and $\mathcal{N}_\alpha$ is the alpha quantile of the distribution. When

514    $\hat{\rho}_{xy|z}^{(k)} > \mathcal{N}_\alpha$, gene x and gene y are conditionally independent given the gene z in cell

515    k, $w_{xy|z}^{(k)} = 0$, else $w_{xy|z}^{(k)} = 1$. (3) Constructing conditional cell-specific network for

516    each gene pair for cell k and for the conditional gene set $\mathcal{Z} = \{z_1, z_2, \cdots, z_G\}$. (4)

517    Integrating the conditional cell-specific network of conditional gene set Z. For each cell,

518    we repeat the steps (1) – (4). Finally, we get a conditional degree matrix *CNDM* which

519    has the same dimension as gene expression matrix *E*. The *CNDM* can be used in

520    clustering, visualization and differentiation potency analysis.

521

522    **Figure 2    CNDM for visualization of scRNA-seq data**

523    The datasets are dimensionally reduced by t-SNE and cell types are encoded by

524    different colors.

525

526    **Figure 3    CCSN uncovers network topology and dynamics for single cells**

527    **A.** The cell specific network (CSN) and conditional cell specific network (CCSN) of

528    the same single cell from the Wang dataset. The same genes are used in network

529    construction. **B.** CCSNs of 8 core genes for representative single cells. **C.** CCSN

530    degrees of *MEIS2*, *PBX1* and *POU3F2* along six time points of the neuronal

531    differentiation.

532    **Figure 4    Network flow entropy** analyses for differentiated cells and progenitors

533    **A.** Network flow entropy between NPCs (at 0 and 1 day) and mature neurons (at 30

534    day). **B.** Network flow entropy between cells at day 10 and day 17 during differentiation

535    of mouse hepatoblasts. P-value is from one-sided Wilcoxon rank-sum test.

536    **Figure 5    The differentiation landscape of neural progenitor cells into mature**

537    **neurons**

538    The 3-dimensional plot shows the NFE of single cells gradually decrease along the

539    differentiation time-course of neural progenitor cells (day 0) into mature neurons (day

540    30). The z axis represents the NFE. The x axis and y axis are the two components of t-

541    SNE.

542

543    **Tables**

544 **Table 1   The comparison of CNDM and GEM in clustering of scRNA-seq data**

545 *Note*: The performance of clustering is evaluated by ARI. Hierarchical (t-SNE) and k-
546 means (t-SNE) indicates clustering after t-SNE.

547

548 **Table 2   The comparison of CNDM with NDM in clustering analysis**

549 *Note*: The performance of clustering is evaluated by ARI.
550

551 **Supplementary material**

552 **File S1   CCSN additional implementation details**

553 **Figure S1   Scatter diagram of the expression values of gene x, gene y and gene z**
554 **for cell k**

555 (A) the red plot k represents the cell k and x axis, y axis and z axis represent the

556 expression levels of gene x, gene y and gene z. gene z respectively. Gene z is set as

557 the conditional gene. n is the number of cells in the dataset. (B) The two parallel light

558 shadow planes $P_x^1, P_x^2$, where x-axis is orthogonal with two planes. The dots are

559 contained in the space between the two planes are the neighbors of $x_k$ and the

560 number of the dots is $n_x^{(k)}$. (C) The two parallel light shadow planes $P_y^1, P_y^2$, where y-

561 axis is orthogonal with two planes. The dots are contained in the space between the

562 two planes are the neighbors of $y_k$ and the number of the dots is $n_y^{(k)}$. (D) The two

563 parallel light shadow planes $\mathcal{P}_z^1, \mathcal{P}_z^2$, where z-axis is paralleled with the two planes.

564 The dots contained in the space between the two planes are the neighbors of $z_k$, and

565 the number of the dots is $n_z^{(k)}$. (E)The intersection of the four parallel light shadow

566 planes $P_x^1, P_x^2, \mathcal{P}_z^1, \mathcal{P}_z^2$ is the space which is surrounded by the green lines. The

567 number of dots which are contained in the space is $n_{xz}^{(k)}$. (F)The intersection of the

568    four parallel light shadow planes $P_y^1, P_y^2, \mathcal{P}_z^1, \mathcal{P}_z^2$ is the space which is surrounded

569    by the green lines. The number of dots which are contained in the space is $n_{yz}^{(k)}$. (G)

570    The intersection of the six parallel shadow planes $P_x^1, P_x^2, P_y^1, P_y^2, \mathcal{P}_z^1, \mathcal{P}_z^2$ is the

571    space which is surrounded by the green lines. The number of dots which are contained

572    in the space is $n_{xyz}^{(k)}$.

573

574    **Figure S2   The comparison of standard normal distribution and the distribution**

575    **of $\widehat{\boldsymbol{\rho}}_{xy|z}^{(k)}$**

576    The density function is calculated by kernel density estimation based on 20,000 plots,

577    and $n_x, n_y, n_z$ are equal to 0.2n. The gene x and gene y are conditional independent

578    given gene z.

579

580    **Figure S3   Performance comparison of GEM and CNDM**

581    PCA was applied for visualization and different colors represent different cell types.

582

583    **Figure S4   The clustering performance of CNDM and GEM**

584    K-means, hierarchical clustering algorithm (HCA) and K-medoids were used for

585    comparison. The data which was preprocessed by t-SNE was also performed to cluster.

586

587    **Figure S5   The clustering performance of CNDM and NDM**

588    K-means, hierarchical clustering algorithm (HCA) were used for comparison. The

589    data which was preprocessed by t-SNE was also performed to cluster.

590

591

592 **Figure S6    Visualization of 23,321 cells by t-SNE**

593 Different colors represent different tissues.

594

595 **Figure S7    The clustering performance of CNDM with different parameters**

596

597 **Figure S8    CCSN analysis of TCGA-LUAD dataset**

598 A. t-SNE plots are used for visualization based on CCSN. The normal samples and

599 tumor samples are represented by different colors. B. CCSNs of representative samples

600 for 18 genes involved in the EGFR pathway. C. Conditional network degrees of EGF

601 and EGGR in the normal samples and the tumor samples.
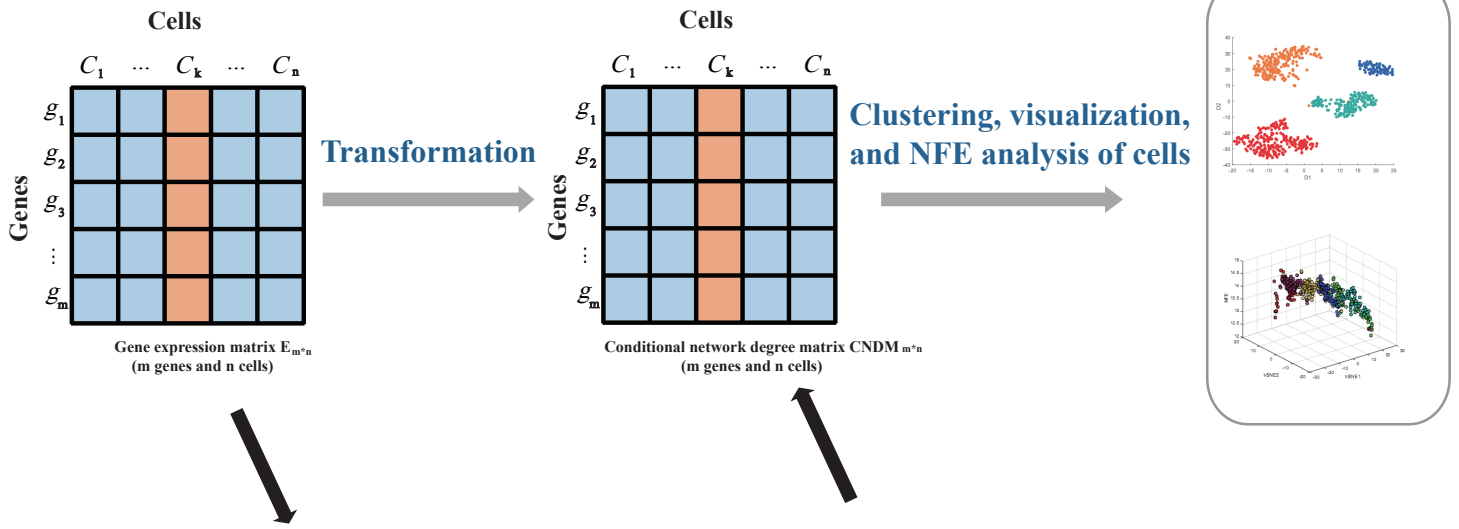
602

603 **Table S1    The running time of CCSN with different numbers of conditional**
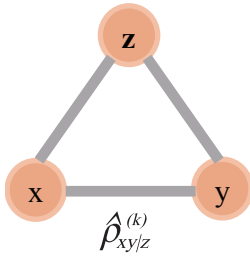
604 **genes**

605

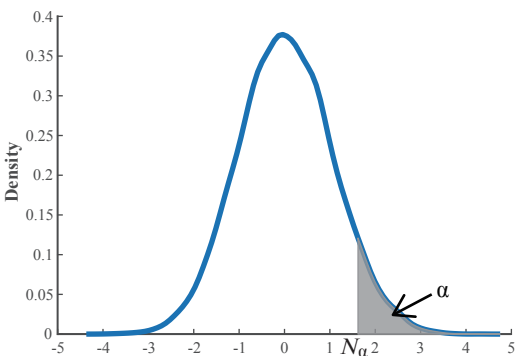606 **Table S2    Datasets used in this study**

607

**Transformation**

**Clustering, visualization, and NFE analysis of cells**

Gene expression matrix $E_{m*n}$
(m genes and n cells)

Conditional network degree matrix CNDM $_{m*n}$
(m genes and n cells)

(1) The normalized statistics of gene x and gene y given the conditional gene z for each cell k

(4) Integrate the conditional cell-specfic network of conditional gene set $z$

$\hat{\rho}_{xy|z}^{(k)}$

$\overline{C}_k = \frac{1}{G}\sum_{g=1}^{G} C_{z_g}^{(k)}$

(2) Hypothesis testing of the normalized statistic $\hat{\rho}_{xy|z}^{(k)}$

(3) Construct conditional cell-specific network for gene pairs for cell k and genes in the conditional gene set $\mathbf{Z} = \{z_1, z_2, \ldots, z_G\}$

**The distribution of the normalized statistic**

$C_{z_G}^{(k)}$

$C_{z_2}^{(k)}$

$C_{z_1}^{(k)}$

1. $\hat{\rho}_{xy|z}^{(k)} > N_\alpha : W_{xy|z}^{(k)} = 0$    2. $\hat{\rho}_{xy|z}^{(k)} <= N_\alpha : W_{xy|z}^{(k)} = 1$

A

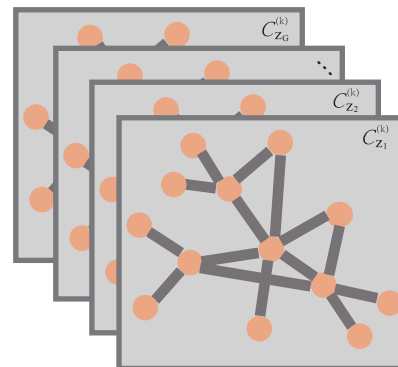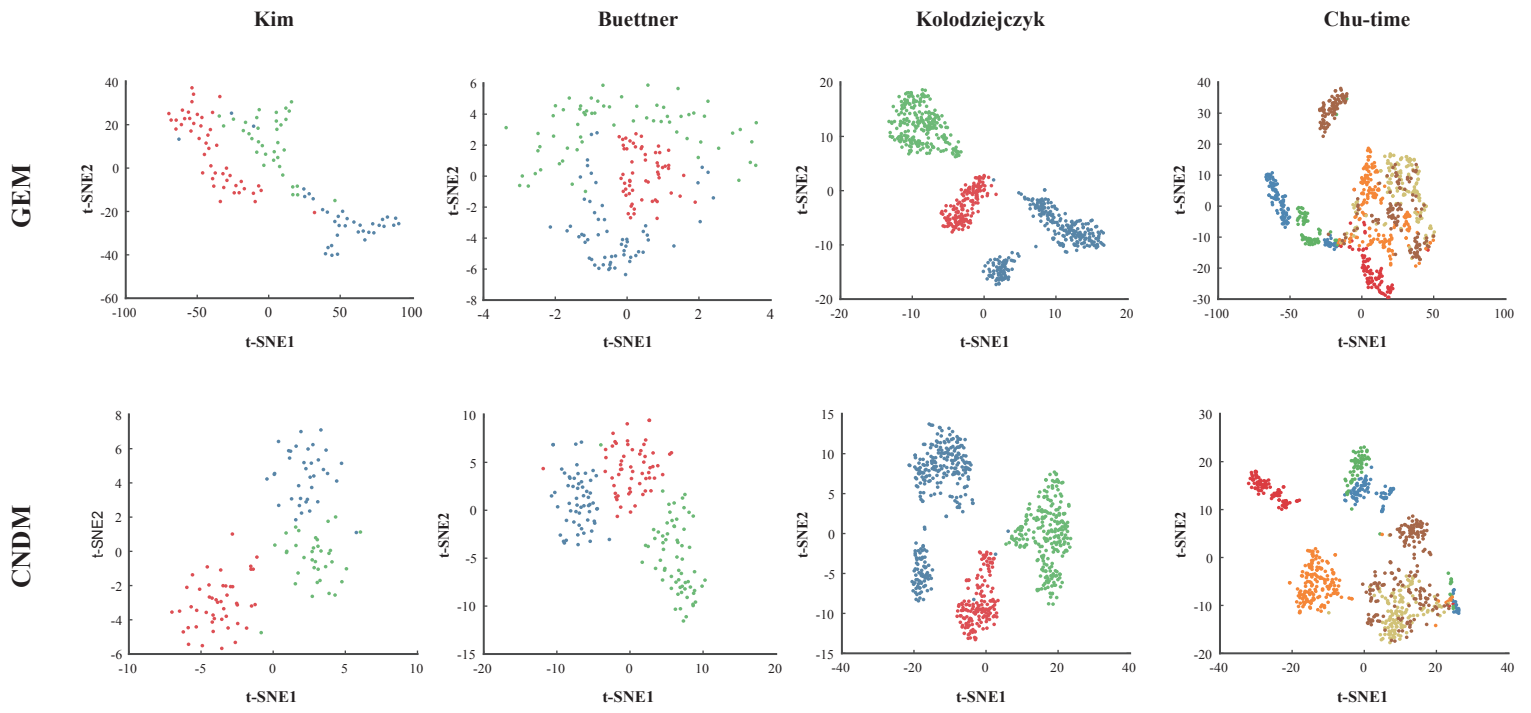CSN (RHB1742_d0)          CCSN (RHB1742_d0)



B



RHB1816_d0 (day 0)          RHB1221_d1 (day 1)          RHB1321_d5 (day 5)

RHB1825_d7 (day 7)          RHB1429_d10 (day 10)          RHB1502_d30 (day 30)

C

A



B

**Table 1    The comparison of CNDM and GEM in clustering of scRNA-seq data**

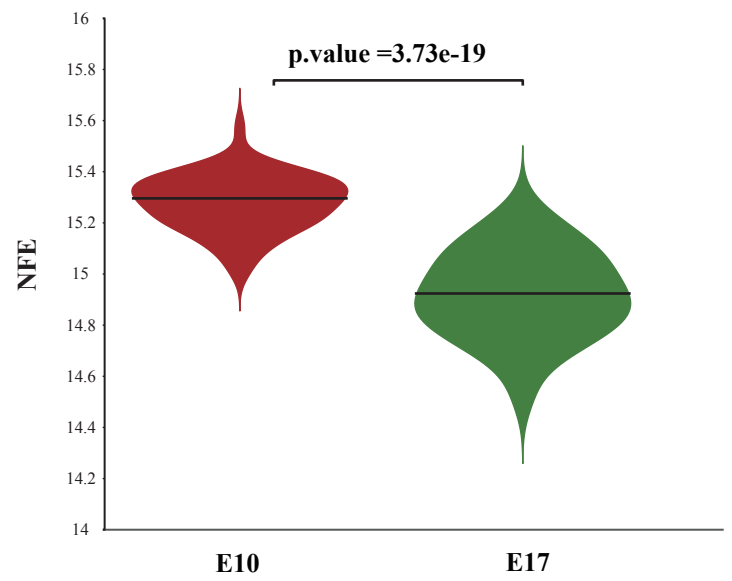|  |  | Buettner | Kolodziejczyk | Grokce | Chu-time | Chu-type | Kim |
|---|---|---|---|---|---|---|---|
| K-means | GEM | 0.29 | 0.54 | 0.42 | 0.17 | 0.22 | 0.20 |
|  | **CNDM** | **0.87** | **0.85** | **0.75** | **0.45** | **0.57** | **0.81** |
| Hierarchical | GEM | 0.32 | 0.49 | 0.47 | 0.22 | 0.22 | 0.12 |
|  | **CNDM** | **0.73** | **0.65** | **0.92** | **0.47** | **0.61** | **0.77** |
| K-means (t-SNE) | GEM | 0.41 | 0.87 | 0.43 | 0.33 | 0.55 | 0.53 |
|  | **CNDM** | **0.95** | **0.91** | 0.36 | 0.56 | **0.70** | **0.93** |
| Hierarchical (t-SNE) | GEM | 0.55 | 0.99 | 0.50 | 0.39 | 0.67 | 0.73 |
|  | **CNDM** | **0.95** | 0.99 | 0.39 | **0.61** | **0.80** | **0.95** |
| K-medoids | GEM | 0.23 | 0.29 | 0.40 | 0.33 | 0.33 | 0.79 |
|  | **CNDM** | **0.53** | **0.63** | **0.81** | 0.17 | **0.38** | **0.61** |
| SC3 | GEM | 0.89 | 1 | 0.56 | 0.66 | 0.78 | 0.89 |
|  | **CNDM** | **0.98** | 0.72 | **0.72** | 0.63 | **0.98** | **0.96** |
| SIMLR | GEM | 0.89 | 0.49 | 0.43 | 0.30 | 0.48 | 0.38 |
|  | **CNDM** | 0.63 | **0.52** | **0.85** | **0.58** | **0.54** | **0.95** |
| Seurat | GEM | 0.67 | 0.43 | 0.35 | 0.52 | 0.52 | 0.41 |
|  | **CNDM** | **0.90** | **0.56** | 0.32 | **0.56** | **0.69** | **0.84** |

*Note*: The performance of clustering is evaluated by adjusted random index (ARI). Hierarchical (t-SNE) and k-means (t-SNE) represent that the clustering analysis is performed after dimension-reduction by t-SNE

**Table 2 The comparison of CNDM with NDM in clustering analysis**

|  |  | Buettner | Kim | Wang | Grokce | Tabula Muris (Aorta) | Tabula Muris (Limb Muscle) |
|---|---|---|---|---|---|---|---|
| K-means | NDM | 0.50 | 0.50 | 0.30 | 0.79 | 0.21 | 0.58 |
|  | **CNDM** | **0.87** | **0.81** | **0.45** | **0.75** | **0.63** | **0.66** |
| Hierarchical | NDM | 0.69 | 0.59 | 0.38 | 0.95 | 0.12 | 0.65 |
|  | **CNDM** | **0.73** | **0.77** | **0.45** | **0.92** | **0.75** | **0.76** |
| K-means (t-SNE) | NDM | 0.83 | 0.84 | 0.61 | 0.38 | 0.46 | 0.62 |
|  | **CNDM** | **0.95** | **0.93** | **0.67** | 0.36 | **0.61** | **0.65** |
| Hierarchical (t-SNE) | NDM | 0.89 | **0.98** | 0.58 | 0.47 | 0.50 | 0.66 |
|  | **CNDM** | **0.95** | **0.95** | **0.72** | 0.39 | **0.50** | **0.66** |
| K-medoids | NDM | 0.26 | 0.49 | 0.31 | 0.60 | 0.35 | 0.14 |
|  | **CNDM** | **0.53** | **0.61** | 0.21 | **0.81** | **0.53** | **0.39** |
| SC3 | NDM | 0.67 | 1 | 0.70 | 0.45 | 0.29 | 0.66 |
|  | **CNDM** | **0.98** | 0.96 | **0.86** | **0.72** | **0.73** | **0.76** |
| SIMLR | NDM | 0.64 | 0.75 | 0.29 | 0.74 | 0.40 | 0.60 |
|  | **CNDM** | 0.63 | **0.95** | **0.60** | **0.85** | **0.70** | **0.71** |
| Seurat | NDM | 0.82 | 0.97 | 0.59 | 0.44 | 0.45 | 0.66 |
|  | **CNDM** | **0.90** | 0.84 | 0.59 | 0.32 | **0.76** | **0.75** |

*Note*: The performance of clustering is evaluated by adjusted random index (ARI). Hierarchical (t-SNE) and k-means (t-SNE) represent that the clustering analysis is performed after dimension-reduction by t-SNE