

Evaluating the Performance of Dropout Imputation and Clustering Methods for Single-cell RNA Sequencing Data

Junlin Xu^{1,#}, Lingyu Cui^{2,#}, Jujuan Zhuang³, Yajie Meng¹, Pingping Bing⁴, Bingsheng He⁴, Geng Tian^{5,6}, Taoyang Wu⁷, Bing Wang^{8,*}, Jialiang Yang^{5,6,*}

¹*College of Computer Science and Electronic Engineering, Hunan University, Changsha, Hunan 410082, China*

²*College of Life Science, Northeast Forestry University, Harbin, Heilongjiang 150000, China*

³*School of Science, Dalian Maritime University, Dalian, Liaoning 116026, China*

⁴*Academician Workstation, Changsha Medical University, Changsha 410219, China*

⁵*Geneis Beijing Co., Ltd., Beijing 100102, China*

⁶*Qingdao Geneis Institute of Big Data Mining and Precision Medicine, Qingdao, 266000, China*

⁷*School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, U.K.*

⁸*School of Electrical & Information Engineering, Anhui University of Technology, Anhui 243002, China*

The authors contributed equally to this study

* Corresponding author(s).

E-mail: yangjl@geneis.cn (Yang JL), wangbing@ustc.edu (Wang B)

Abstract

Recent advances in single-cell RNA sequencing (scRNA-seq) provide exciting opportunities for transcriptome analysis at the single-cell resolution. Clustering individual cells is a key step to reveal cell subtypes and infer cell lineage in scRNA-seq analysis. Although many dedicated algorithms have been proposed, clustering quality remains a computational challenge for scRNA-seq data, which becomes exacerbated due to excessive zero counts caused by various technical noise. To address this challenge, we assess the combinations of nine popular dropout

30 imputation methods and eight clustering methods using a collection of 10
31 well-annotated scRNA-seq datasets with different sample sizes. Our results show that
32 imputation algorithms do typically improve the performance of clustering methods, as
33 well as the quality of data visualization using t-Distributed Stochastic Neighbor
34 Embedding. However, the performance of a particular combination of imputation and
35 clustering methods may vary among datasets with different sizes. For example, the
36 combination of single-cell analysis via expression recovery and Sparse Subspace
37 Clustering (SSC) methods usually works well on small datasets, while the
38 combination of adaptively-thresholded low-rank approximation and single-cell
39 interpretation via multikernel learning (SIMLR) usually achieves the best
40 performance on large datasets.

41

42 **KEYWORDS:** Single-cell RNA sequencing; Dropout imputation; Cell clustering;
43 T-SNE; Adjusted Rand index

44

45 **Introduction**

46 Recent advances in single-cell sequencing provide a great opportunity for
47 understanding cell-specific gene expressions, cell lineage relationships, and various
48 important biological processes and functions at single-cell resolution [1-3]. Among
49 them, single-cell RNA sequencing (scRNA-seq) is widely used to quantify mRNA
50 expression in a single cell [4, 5]. However, effective analysis of scRNA-seq data
51 remains a challenging task, as they are typically much more complicated than
52 traditional sequencing data [6]. Indeed, because the amount of mRNA in a single cell
53 is very small, a million-fold amplification is usually required, which leads to greater
54 amplification noise [7]. Among these issues, the most common one includes dropout
55 events, referring to the value of certain genes in certain cells being zero or close to
56 zero.

57 A key step in scRNA-seq transcriptome profiling is to cluster individual cells to
58 reveal cell subtypes and/or subpopulations [8, 9]. To this end, a variety of

59 unsupervised clustering algorithms are proposed, ranging from simple k-means
60 clustering, hierarchical clustering [10] and its variants (e.g., RaceID [7], SC3 [11],
61 and CIDR [12]), to density-based spatial clustering [13], subspace clustering [14],
62 neural network [15, 16], ensemble clustering, and kernel-based methods (such as
63 SIMLR [17]). However, effective clustering remains a computational challenge due
64 to the high proportion of “dropouts” featured in scRNA-seq datasets. To address this,
65 several promising imputation methods specially designed for scRNA-seq data have
66 been developed [6, 18-21]. These methods are roughly divided into two categories:
67 similarity-based imputation methods, which use similarities between genes and
68 between cells to restore expression levels, and matrix-based imputation methods,
69 which are based on the postulate that the true expression matrix is a low-rank and
70 leverages various advanced techniques in matrix analysis [22].

71

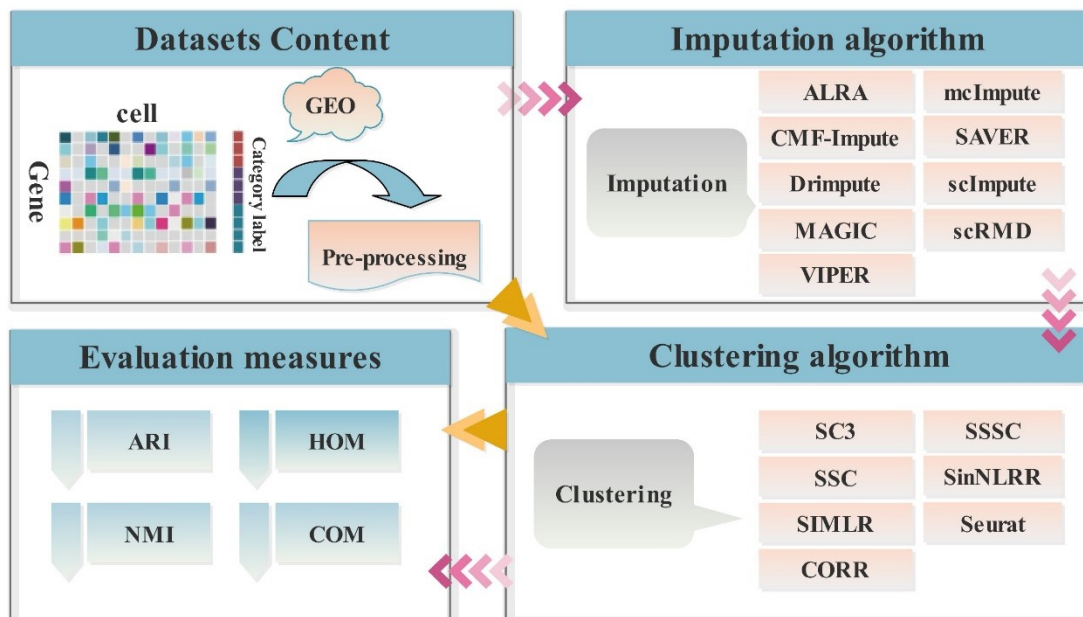
72 In this study, we present a critical review of nine promising imputation methods
73 and seven clustering techniques better suited for scRNA-seq data. Using a set of ten
74 well-annotated scRNA-seq datasets, we assessed the performance of various
75 combinations of these imputation and clustering approaches. Our results show that
76 imputation algorithms typically improve the performance of various clustering
77 methods, as well as the quality of data visualization using t-Distributed Stochastic
78 Neighbor Embedding (t-SNE). In addition, the performance of a particular
79 combination of imputation and clustering methods may vary among datasets with
80 different sizes. Therefore, it is critical to choose an appropriate combination of
81 imputation and clustering algorithms for obtaining high-quality clustering from
82 scRNA-seq data. Our results suggest that using single-cell analysis via expression
83 recovery (SAVER) + sparse subspace clustering (SSC) usually provides better
84 clustering results for small datasets with less than 100 cells. In contrast,
85 adaptively-thresholded low-rank approximation (ALRA) + single-cell interpretation
86 via multikernel learning (SIMLR) typically achieves the best performance for large
87 datasets with more than 1000 cells.

88

89 Results

90 The scRNA-seq data cluster evaluation framework

91 The cluster assessment framework outlined in **Figure 1** can be divided into the
 92 following four steps: (1) first, we preprocess the input gene expression matrix from a
 93 scRNA-seq dataset by removing rare genes and a logarithmic transformation (see
 94 Section 2); (2) next, we use the nine imputation algorithms reviewed in the last
 95 section to impute the processed expression matrix to obtain nine estimated expression
 96 matrices; (3) then, we use each of the seven clustering algorithms to cluster each of 10
 97 expression matrices (e.g., the original one and nine imputed ones); (4) finally, we
 98 compute the NMI, ARI, HOM, and COM scores to quantify the differences between
 99 the predefined annotations of cell types and the output cluster labels from each of the
 100 70 combinations of imputation and clustering algorithms.



101

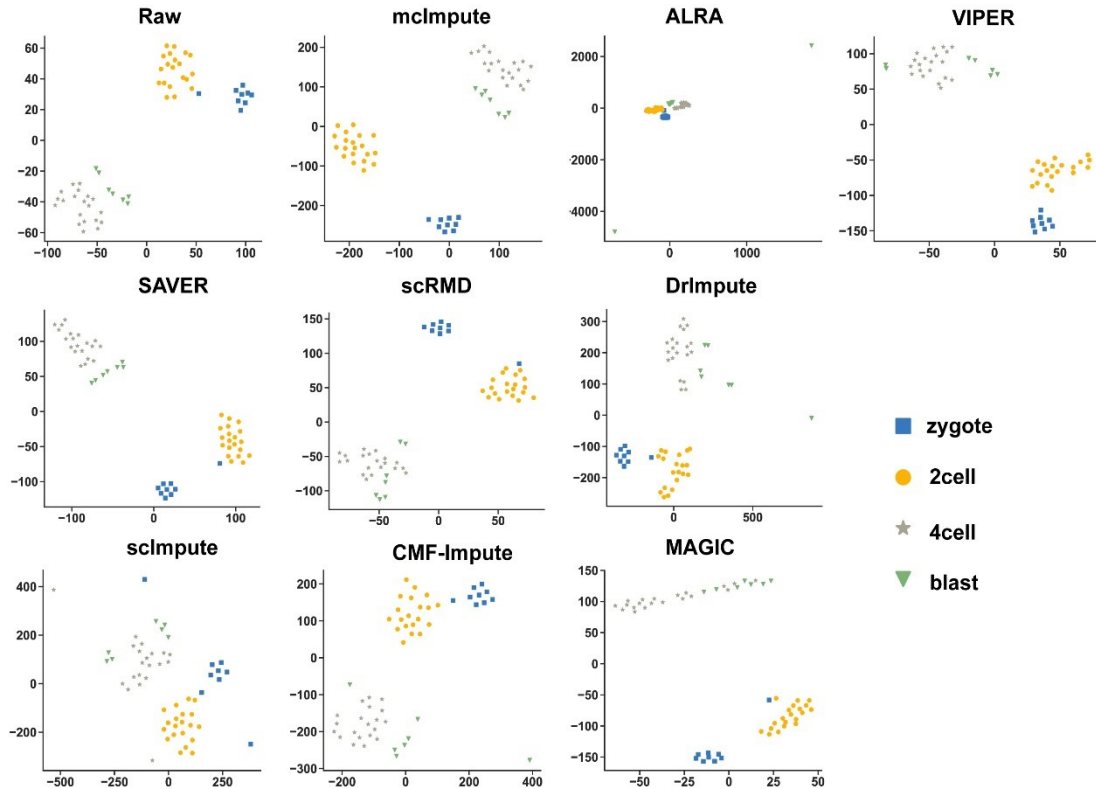
102 **Fig. 1.** Schematic workflow of the scRNA-seq dataset cluster evaluation framework. The framework is
 103 mainly divided into four parts: the collection of data sets, the direct analysis of the original data sets
 104 using different clustering algorithms, the clustering analysis of the data sets imputationed with various
 105 algorithms, and the evaluation of the clustering results.

106

107 Imputation on scRNA-seq data can often improve visualization

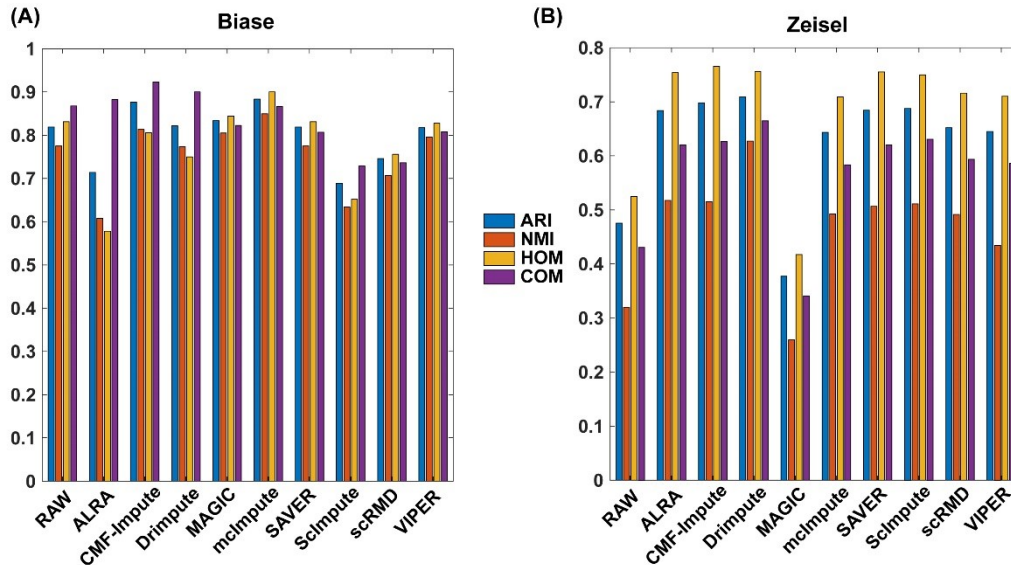
108 Finding an effective low-dimensional visualization of scRNA-seq data remains a key

109 computational challenge in single-cell data analysis. One popular dimensionality
 110 reduction visualization algorithm is t-SNE, which visualizes high-dimensional data by
 111 giving each data point a location in a two- or three-dimensional space. Since the
 112 t-SNE algorithm is not designed to handle the high rate of dropouts featured in
 113 scRNA-seq data, which may make this algorithm less suitable for some scRNA-seq
 114 datasets.



115
 116 **Fig. 2.** T-SNE visualization of cells from the Biase scRNA-seq dataset. Among them, the mcImpute
 117 algorithm improves the visualization of the Biase dataset. Note: Cells are color-coded by the cell type
 118 annotation of the original study.

119
 120 To assess the impact of imputation algorithms on visualization, we start with the
 121 smallest dataset in our collection, namely the Biase dataset [23]. As shown in **Figure**
 122 **2**, both MAGIC and scImpute did not improve the visualization of the Biase dataset.
 123 Indeed, cells of type “4-cell” and those of “blast” are often confused. We also perform
 124 k-means clustering on the resulting datasets transformed by t-SNE. The evaluation
 125 results summarized in **Figure 3(A)** indicate that CMF-Impute and McImpute
 126 algorithms improved the accuracy of clustering for the Biase dataset.

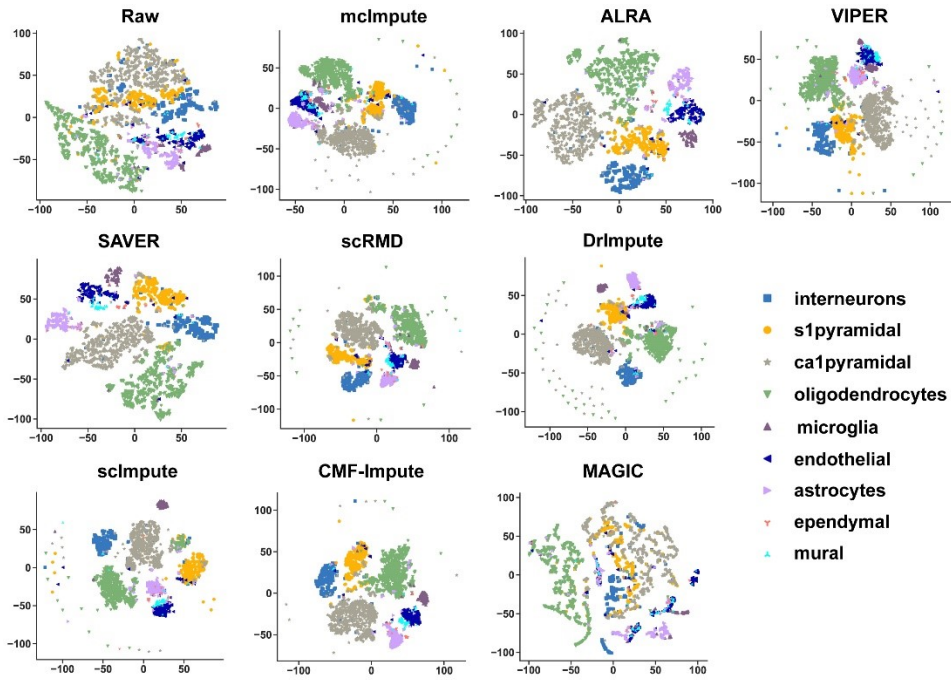


127
 128 **Fig. 3.** Benchmark of imputation algorithms on the t-SNE+k-means clustering of scRNA-seq dataset.
 129 (A) The use of different imputation algorithms to compare the four evaluation indicators obtained on
 130 the Biase scRNA-seq dataset through t-SNE+K-means clustering. CMF-Impute and McImpute
 131 algorithms improved the accuracy of clustering. (B) The use of different imputation algorithms to
 132 compare the four evaluation indicators obtained on the Zeisel scRNA-seq dataset through
 133 t-SNE+K-means clustering. Except for MAGIC, most of the imputation algorithms can improve the
 134 accuracy of clustering.

135

136 Next, we consider the largest dataset in our collection, namely the Zeisel dataset
 137 [24]. After applying various imputation algorithms to interpolate dropout events in
 138 this dataset, we use t-SNE to visualize the imputed dataset in **Figure 4**, where the
 139 cells are color-coded according to their types annotated in the original study. This
 140 demonstrates that t-SNE typically produces better decomposition or visualization
 141 results when imputation is applied. Furthermore, we perform K-means clustering on
 142 the datasets transformed by t-SNE. The evaluation results summarized in **Figure 3(B)**
 143 indicate that, except for MAGIC, all other imputation algorithms have a desirable
 144 effect on the visualization of the Zeisel dataset.

145 To systematically evaluate the impact of imputation algorithms on the
 146 performance of data visualization, we apply the above evaluation framework to each
 147 of the scRNA-seq datasets in **Table 1**. As shown in the evaluation results presented in
 148 Supplementary Table S1, various imputation algorithms contribute to the visualization
 149 performance of t-SNE. Among them, SAVER, DrImpute, and CMF-Impute



150

151 **Fig. 4.** T-SNE visualization of cells from the Zeisel scRNA-seq dataset. Except for MAGIC, all
 152 algorithms improve the visualization of the Zeisel dataset. Note: Cells are color-coded by the cell type
 153 annotation of the original study.

154

155

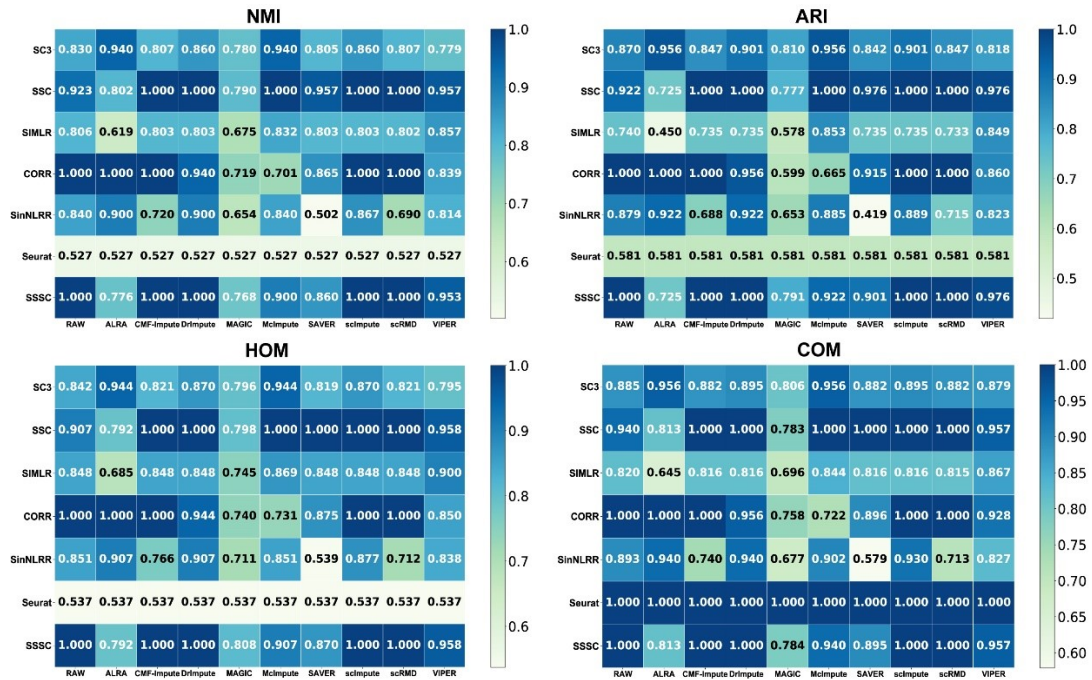
Table 1. scRNA-seq dataset for analysis and comparison.

No. of datasets	Names	No. of cells	No. of genes	No. of cell types
1	Biase	56	25,733	4
2	Deng	268	22,431	6
3	Goolam	124	41,427	5
4	Grun	251	23,459	4
5	Kolodziejczyk	704	38,615	9
6	Patel	430	5,948	5
7	Pollen	301	23,730	11
8	Usoskin	622	25,334	4
9	Yan	90	20,214	6
10	Zeisel	3,005	32,738	9

156

157 outperform other algorithms. Furthermore, the performance of some imputation
 158 algorithms may vary among the datasets, depending on the size of the datasets. For

159 example, ALRA performs significantly better compared to many other imputation
 160 algorithms in large datasets. However, ALRA does not work well on small datasets
 161 (such as the Deng and the Yan datasets), and all four evaluation indicators indicate
 162 poor performance. Therefore, it is important to choose an appropriate imputation
 163 algorithm for a given dataset, with the size of the dataset being a critical factor.



164
 165 **Fig. 5.** Comparison of the impact of different imputation algorithms on clustering performance on the
 166 Biase scRNA-seq dataset. Among the seven clustering algorithms, SSC, Corr, and SSSC have the best
 167 overall performance.

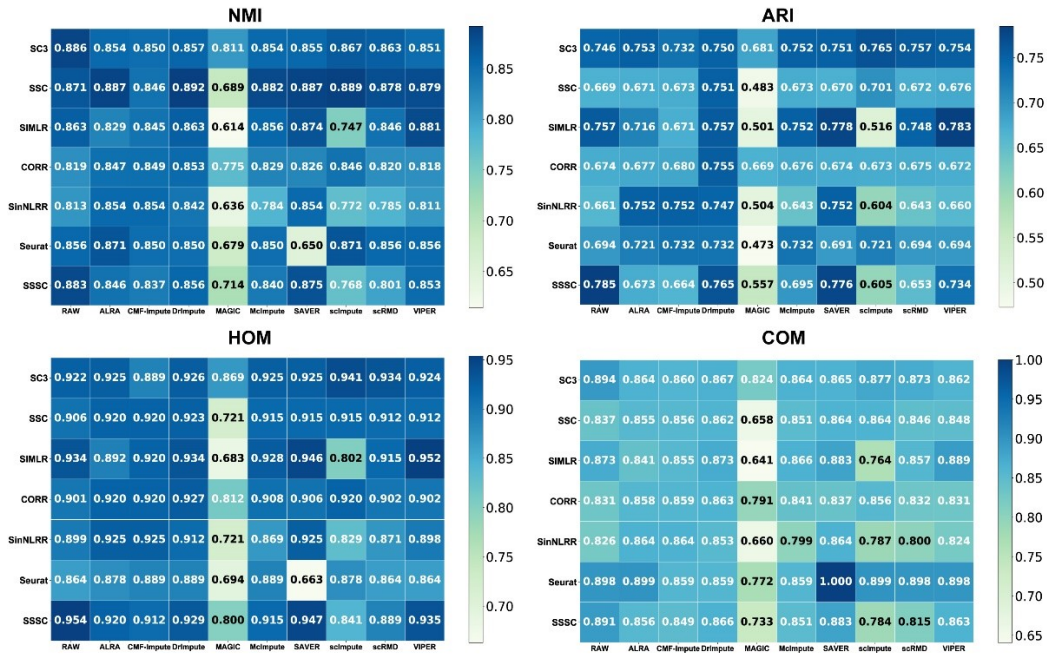
168

169 Imputation for most scRNA-seq data can improve clustering

170 To assess the impact of imputation algorithms on clustering methods, we apply the
 171 framework outlined in Figure 1 to three representative datasets in Table 1. The first
 172 one is the Biase dataset, which is a small dataset containing 56 cells of four types.

173 **Figure 5** shows the results for various combinations of imputation and clustering on
 174 the Biase dataset, which indicates that, in most cases, clustering after using various
 175 imputation algorithms can help improve performance. Among the seven clustering
 176 algorithms, SSC, Corr, and SSSC have the best overall performance.

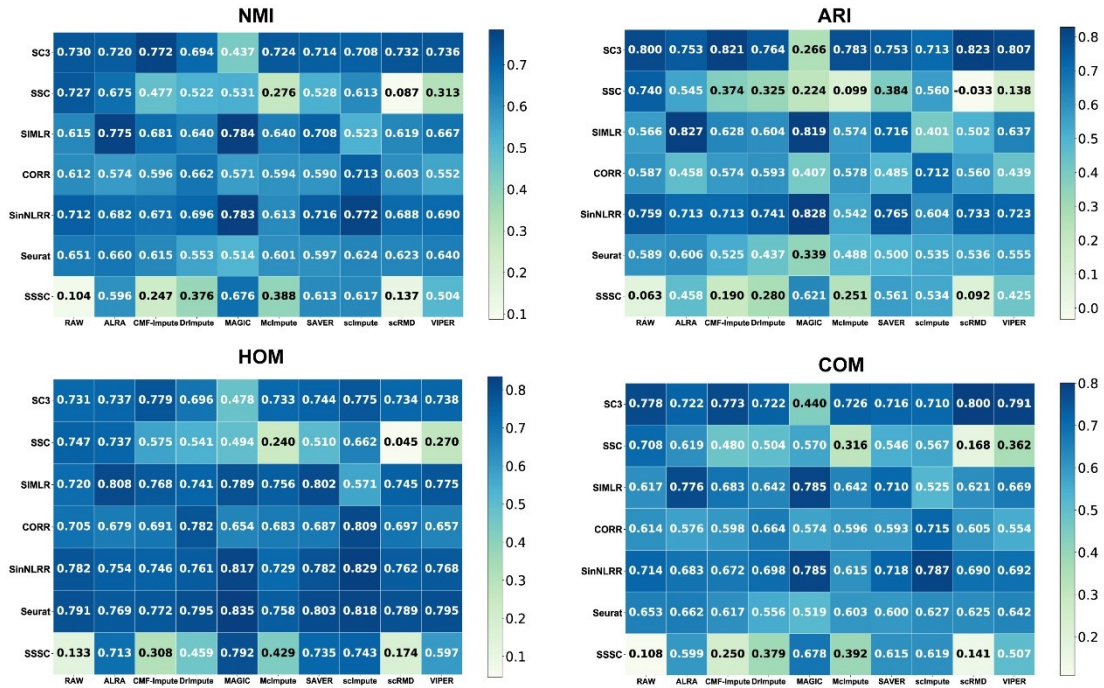
177



178
 179 **Fig. 6.** Comparison of the impact of imputation algorithms on clustering performance on the Pollen
 180 scRNA-seq dataset. Different imputation algorithms are selected for different sizes of data sets to
 181 improve the accuracy of clustering.

182 The dataset reported by Pollen [25] and colleagues is a medium dataset with 301
 183 cells from 11 different cell types. Compared with other datasets, the dropout rate of
 184 this dataset is relatively low. Consequently, clustering algorithms without imputation
 185 already achieve good performance. As shown in **Figure 6**, imputation algorithms do
 186 not necessarily lead to improvements in the performance of various clustering
 187 algorithms. On the contrary, some imputation algorithms (for example, MAGIC)
 188 indeed have an undesirable impact on the performance of most clustering algorithms.

189 The Zeisel dataset is a large dataset containing 3005 single cells from the mouse
 190 cortex and hippocampus, collected by the unique molecular identifier technology and
 191 divided into nine categories. **Figure 7** shows the performance of various
 192 combinations of imputation and clustering algorithms. For this dataset, most
 193 imputation algorithms can improve the performance of various clustering algorithms.
 194 For example, after using the ALRA algorithm to repair the missing values, the SIMLR
 195 algorithm has much better performance, achieving an ARI score of 0.827 and an NMI
 196 score of 0.774, which is much higher than that of other imputation algorithms.



197
198 **Fig. 7.** Comparison of the impact of imputation algorithms on clustering performance on the Zeisel
199 scRNA-seq dataset

200 Finally, we use the framework outlined in Figure 1 to assess the performance of
201 various combinations of imputation and clustering methods for all other datasets in
202 Table 1. As shown in Supplementary Table S2-S11, clustering methods generally have
203 a better performance when an imputation algorithm is applied. Since the performance
204 of each combination of imputation and clustering methods may vary among datasets
205 with different sizes. Based on these results, we observe that a combination of SAVER
206 and SSC usually performs well on small datasets (such as the Biase and the Yan
207 datasets [26]). In contrast, combining adaptively-thresholded low-rank approximation
208 ARLA and SIMLR typically achieves the best performance for large datasets (such as
209 the Zeisel dataset).

210
211 **Discussion**

212 In this study, we use 10 well-annotated scRNA-seq datasets and an objective
213 assessment framework to evaluate the performance of various combinations of
214 imputation and clustering algorithms that are suitable for scRNA-seq datasets. Our
215 empirical results show that imputation algorithms typically improve the performance

216 of various clustering methods, as well as the quality of data visualization using t-SNE.
217 However, the performance of a particular combination of imputation and clustering
218 methods may vary among datasets with different sizes. These results provide concrete
219 choices for how to choose an appropriate combination of imputation and clustering
220 algorithms for obtaining high-quality clustering from scRNA-seq data. Moreover, it
221 remains interesting to see how to utilize the insights obtained from here to design
222 better imputation and clustering algorithms.

223 In addition to this clustering problem, many other computational problems in
224 single-cell data analyses also face the same challenge derived from a high dropout
225 rate, such as standardization, differential expression analysis, and cell cycle
226 identification. Therefore, algorithms aiming to solve those problems could benefit
227 from various imputation techniques reviewed here, and the assessment framework
228 proposed here can be naturally extended to study their performance.

229

230 **Methods**

231 **Data preparation and preprocessing**

232 To determine the impact of different imputation algorithms on the clustering of
233 individual cells to their corresponding cell types, we collected 10 scRNA-seq datasets
234 with cell type annotations, which come from the National Center for Biotechnology
235 Information Gene Expression Omnibus (NCBI-GEO). Table 1 summarizes these 10
236 scRNA-seq datasets. At the same time, to reduce the technical noise in the scRNA-seq
237 datasets, genes expressed in less than or equal to two cells were filtered [27]. In order
238 to prevent the effect of highlighting genes with higher expression and weakening the
239 remaining genes, we used Equation 1 to perform a logarithmic transformation with
240 pseudo count 1 on the original expression data of single cells before analysis.

$$241 \quad X = \log_2(X + 1). \quad (1)$$

242

243 **Dropout imputation**

244 An important technical flaw in scRNA-seq data is the introduction of “dropout”
245 events [28, 29]. Deletion events usually refer to the incorrect quantification of genes
246 that are not expressed due to transcripts that are introduced during the reverse
247 transcription step or have low expression levels [30]. A large number of studies have
248 shown that simply deleting the less expressed genes and then normalizing them
249 cannot completely solve this issue in scRNA-seq data analysis. In order to better
250 perform downstream analysis of scRNA-seq data, a large number of missing value
251 repair algorithms have been proposed. Therefore, in this study, we analyzed in detail
252 the impact of nine better imputation algorithms on the clustering analysis of
253 scRNA-seq data.

254

255 *scImpute*

256 Li and Li proposed a three-step approach called scImpute to determine and impute
257 values that are affected by dropout events in scRNA-seq data [31]. Since this method
258 uses information of the same gene from similar cells to impute missing values, the
259 first step is to construct a candidate pool of neighboring cells for each cell, which is
260 achieved by principal component analysis (PCA) and spectral clustering. The second
261 step computes the dropout probability of each gene in each cell. To this end, the
262 expectation–maximization (EM) algorithm is utilized to estimate a gamma-normal
263 mixture model. In the final step, a separate regression model for each cell is
264 constructed to impute the expression of genes with high dropout probabilities, for
265 which information about the same genes in its neighboring cells identified in the first
266 step is used.

267 It is demonstrated that scImpute can automatically identify zero values of high
268 dropout probabilities and only perform imputation on these values without
269 introducing new deviations to the remaining data [31]. Furthermore, the method can
270 also detect outlier cells and exclude them from imputation. Evaluation based on
271 simulated and real human and mouse scRNA-seq datasets indicates that scImpute is
272 an effective tool for restoring transcriptome dynamics masked by dropouts. scImpute

273 can detect possible deletions, enhance the aggregation of cell subpopulations, improve
 274 the accuracy of differential expression analysis, and help the study of gene expression
 275 kinetics. Because scImpute requires the true number of cell subpopulations in the data
 276 a priori, this is not friendly to unknown structured data.

277

278 *DrImpute*

279 Gong et al. designed a simple and fast hot-platform imputation method called
 280 DrImpute to estimate missing events in scRNA-seq data [27]. Similar to scImpute,
 281 DrImpute performs cell clustering before imputation and also borrows information of
 282 the same gene from similar cells to impute missing values. DrImpute takes a
 283 consensus approach to obtain a more robust estimate. First, the clustering method
 284 used in DrImpute is single-cell consensus (SC3) clustering, which, as we review in
 285 the next subsection, is a consensus approach. Second, imputation is performed
 286 multiple times using different unit clustering results. Finally, the multiple estimates
 287 are averaged to get the final imputation. Specifically, let H be the number of cluster
 288 configurations (for example, the combination of distance metric and the number of
 289 clusters used in clustering), and C_1, C_2, \dots, C_H be the clustering results, one for each
 290 configuration. Assuming that the clustering of C_h is a true hidden cell classification
 291 result, the expected value of the dropout event can be obtained by averaging the
 292 entries in a given cell cluster:

$$293 \quad E(X_{ij}|C_h) = \text{mean}(X_{ij}|W), \quad (2)$$

294 where X is the input matrix, and W represents X_{ij} in the same cell group in cluster
 295 C_h .

296 Therefore, the final calculation of the estimated drop events X_{ij} and $E(X_{ij})$ can
 297 be calculated by a simple average:

$$298 \quad E(X_{ij}) = \text{mean}(E(X_{ij}|C)) = \frac{1}{H} \sum_{h=1}^H E(X_{ij}|C_h). \quad (3)$$

299 Experimental results show that DrImpute greatly improves several existing statistical
 300 tools including SC3, t-SNE, and Monocle, which cannot solve dropout events in the
 301 three most popular research areas in scRNA-seq analysis, namely cell clustering,

302 visualization, and lineage reestablishment. However, since the cluster number is
303 usually unknown in DrImpute, the results are not as accurate as expected.

304

305

306 *VIPER*

307 Both scImpute and DrImpute perform cell clustering before imputation, using cells
308 belonging to the same subgroup of cells for imputation. However, the cell
309 subpopulations used in this type of imputation algorithm are often not true cell
310 subpopulations, which cause serious deviations in the imputation results. To address
311 this concern, Chen and Zhou proposed a simple, accurate, unadjusted, and
312 computationally effective scRNA-seq imputation method called VIPER [32].

313 Compared to scImpute and DrImpute, VIPER mainly borrows information among
314 cells with similar expression patterns to estimate the expression measurement value in
315 a target cell. This is achieved by using a sparse non-generative regression model to
316 actively select the sparse local set of local neighborhoods that best predicts the target
317 cell. These sparse unit sets are selected in a progressive manner, and the attribution
318 weights associated with them are estimated in the final estimation step to ensure
319 robustness and computational scalability. In addition, VIPER uses cell type-specific
320 and gene-specific methods to model the deletion probability, which clearly illustrates
321 the uncertainty of the zero-value expression measurement in scRNA-seq. VIPER uses
322 efficient auxiliary programming algorithms to infer all modeling parameters from
323 existing data while maintaining low computational costs. The key feature of VIPER is
324 that the imputed data can retain the gene expression variability of the whole cell.
325 Compared to several existing imputation methods in several actual analysis
326 experiments based on scRNA-seq data, VIPER can obtain higher imputation accuracy.

327

328 *MAGIC*

329 One key challenge in many imputation algorithms for scRNA-seq data (e.g., scImpute
330 and DrImpute reviewed above) is to accurately find neighborhoods of similar cells. To

331 address this issue, van Dijk et al. developed a cell map imputation algorithm called
332 Markov affinity-based graph imputation of cells (MAGIC) [33], which mainly uses
333 data diffusion to share information among similar cells to denoise the cell count
334 matrix and fill in missing transcripts. Central to this approach is a Markov matrix M ,
335 which is derived using a Gaussian kernel and a normalization process from an
336 appropriate cell-by-cell distance matrix constructed from the input data. Once the
337 Markov matrix, whose (i,j) -entry represents the probability of transitioning from cell i
338 to cell j in a single diffusion step, is obtained, a data diffusion step is performed
339 through exponentiation of M to identify neighborhoods of similar cells. Then the
340 imputation step of MAGIC involves sharing information between cells in the resulting
341 neighborhoods through matrix multiplication.

342 Since MAGIC leverages the observation that cell phenotypes can often be roughly
343 embedded in a low-dimensional structure corresponding to the low-frequency trend of
344 the data containing the biological signal of interest, experiments have shown that it
345 can effectively alleviate the sparsity and noise caused by random mRNA captures and
346 reveal gene–gene relationships in scRNA-seq data. Moreover, unlike many other
347 imputation algorithms that only fill in “missing values,” MAGIC uses the value
348 diffusion between similar pixels along the affinity-based graph structure to correct the
349 entire data matrix and connect it to the basic manifold structure. However, imputation
350 on a low-dimensional space will likely eliminate gene expression variability across
351 cells and thus abolish a key feature of single-cell sequencing data.

352

353 *SAVER*

354 Huang et al. developed an expression restoration method for scRNA-seq data called
355 SAVER [34], which mainly uses information between genes and cells to estimate
356 missing values and improve the expression estimation of all genes. SAVER aims to
357 restore true gene expression patterns by eliminating technical differences and
358 retaining biological differences. It uses observed gene counts to form a prediction
359 model for each gene, and then uses the observed counts and the weighted average of

360 the prediction to estimate the true expression of the gene. Experimental results show
361 that SAVER can reliably restore cell-specific gene expression concentration,
362 cross-cell gene expression distribution, and gene-to-gene and cell-to-cell expression.
363 SAVER's powerful performance is attributed to its adaptive estimation of discrete
364 parameters at the gene level and its cross-validation-based model selection, which can
365 prevent unnecessary model complexity. But SAVER relies on a Markov chain Monte
366 Carlo algorithm to tune all parameters, which is computationally costly and might not
367 be scalable to large datasets

368

369 *ALRA*

370 Linderman et al. proposed a highly scalable method called ALRA to recover the true
371 expression level of scRNA-seq data [22]. It is a singular value decomposition (SVD)
372 followed by a thresholding scheme that takes advantage of the non-negativity of the
373 true expression matrix. A key assumption used in this approach is that the underlying
374 true matrix is non-negative and low-rank and contains many zeros, but none of these
375 zeros are associated with dropouts. The observed matrix from a scRNA-seq
376 experiment is then sparser as many values are incorrectly measured as zero due to the
377 dropout effect. Consequently, ALRA uses SVD to find the best k approximation of
378 this matrix, then transforms it into an imputed matrix where each element
379 corresponding to a dropout is not zero. Experimental results show that ALRA
380 improves the separation between cell types in the high-dimensional space of original
381 cells and restores the true expression of marker genes while retaining the biological
382 zero position. As ALRA has only one parameter, the approximate rank k of the matrix,
383 which is automatically selected based on the statistical information of the interval
384 between consecutive singular values, the method is widely applicable to various
385 scRNA-seq datasets.

386

387 *scRMD*

388 Chen et al. developed a single-cell RNA sequence imputation method based on robust

389 matrix decomposition (RMD) called scRMD [20]. A key postulate in this approach is
 390 that a gene expression matrix Y has the following decomposition:

$$391 \quad Y = L - S + E, \quad (4)$$

392 where L is a low-rank matrix, S is a sparse matrix, and E represents the combined
 393 effect of measurement errors and random fluctuations.

394 Moreover, this method also leverages the observation that the expression of gene
 395 i observed in cell j is less likely to be affected by the dropout if the value Y_{ij} is
 396 large enough. Formally, the index set of candidate dropouts can be represented by
 397 $\Omega = \{(i, j): y_{ij} \leq c\}$ for a threshold constant c . Then we have $\mathcal{P}_\Omega(S) \geq 0$, where
 398 the mask operator \mathcal{P}_Ω is defined as $\mathcal{P}_\Omega(s_{ij}) = s_{ij}$ if index (i, j) is contains in Ω
 399 and $\mathcal{P}_\Omega(s_{ij}) = 0$ otherwise. Then the scRMD model can be formulated as the
 400 following optimization problem:

$$401 \quad \min_{L, S} \frac{1}{2} \|Y - L + S\|_F^2 + \lambda_1 \|L\|_* + \lambda_2 \|S\|_1 \quad \text{s.t.} \quad \mathcal{P}_\Omega(S) \geq 0, \mathcal{P}_{\Omega^c}(S) = 0 \text{ and } L \geq$$

$$402 \quad 0 \quad (5).$$

403 Here, λ_1 and λ_2 are regularization parameters, $\Omega^c = \{(i, j): y_{ij} > c\}$. Moreover,
 404 $\|\cdot\|_*$, $\|\cdot\|_F$, and $\|\cdot\|_1$ represent the nuclear norm, Frobenious norm, and elementwise
 405 l_1 norm of a matrix, respectively. This optimization problem can be effectively
 406 solved by an alternating direction multiplier method (ADMM). Extensive data
 407 analysis shows that scRMD can accurately restore missing values and help improve
 408 downstream analyses, such as differential expression analysis and cluster analysis.

409

410 *mcImpute*

411 Mongia et al. presented an imputation algorithm based on matrix completion. The
 412 postulate used in this algorithm is similar to that in scRMD [18]; that is, the gene
 413 expression matrix Y is a low-rank matrix. Let \mathcal{P}_Y be the mask operator associated
 414 with Y , which is defined as $\mathcal{P}_Y(X_{ij}) = X_{ij}$ if $Y_{ij} > 0$ and $\mathcal{P}_Y(X_{ij}) = 0$ otherwise.

415 Then the optimization problem can be formulated as:

$$\min \|X\|_* \text{ s.t. } \|Y - \mathcal{P}_Y(X)\|_F^2 < \text{error}. \quad (6)$$

This problem is solved in mcImpute by iteratively limiting the singular values of the expression matrix. Compared with many other imputation algorithms, one distinguishing feature of mcImpute is that it does not assume any distribution of genes and maintains a complete biological silent expression value. Experiments using several real datasets show that mcImpute is competitive with other algorithms in improving the accuracy of cell clustering, identifying differentially expressed genes, enhancing the separability of cell types, and improving dimensionality reduction. But the complexity of mcImpute algorithm is very high and the running time is very long

CMF-Impute

Xu et al. proposed a novel method based on collaborative matrix decomposition to estimate missing items in a given scRNA-seq expression matrix [21]. A key step in the CMF-Impute algorithm is to find two characteristic matrices so that their product provides the best approximation to the original matrix. Specifically, for a gene expression matrix Y with g rows and n columns, the CMF-Impute algorithm seeks to find a k -dimensional cell feature matrix W and a k -dimensional gene feature matrix H such that $Y = WH^T$ and $k \ll \min(g, n)$. Noting that similar cells tend to have similar gene expression patterns, CMF-Impute explicitly incorporates a cell-to-cell similarity matrix S_c and a gene-to-gene similarity matrix S_g into its optimization formulation:

$$\min_{W,H} \|Y - WH^T\|_F^2 + \lambda_1 \|W\|_F^2 + \lambda_2 \|H\|_F^2 + \lambda_c \|S_c - WW^T\|_F^2 + \lambda_g \|S_g - HH^T\|_F^2, \quad (7)$$

where $\lambda_1, \lambda_2, \lambda_c, \lambda_g$ are regularization parameters.

Experiments on several simulated and real scRNA-seq datasets show that CMF-Impute improves the performance of existing cell clustering algorithms and methods for reconstructing cell-to-cell, gene-to-gene correlations, and inferring cell lineage trajectories.

445 **Clustering techniques**

446 Cell type identification based on single-cell sequencing data is one of the key
447 computational challenges in single-cell biology and has thus received widespread
448 attention [35, 36]. In this section, we review the application of eight clustering
449 methods to scRNA-seq data.

450

451 *t-SNE + K-means*

452 K-means clustering is one of the most frequently used cluster analysis methods. It
453 iteratively computes the mean of all data points of each class as the center point of the
454 class and assigns each data point into one of the k clusters whose mean is closest to
455 the given data point. As a consequence, two data points that are closer to each other
456 are more likely to be classified into the same category. The K-means algorithm and its
457 variants have been applied in a number of fields. Within the area of single-cell
458 technologies, many single-cell clustering algorithms use K-means, such as the SC3
459 [11] and pcaReduce [37]. However, one may argue that the most popular method is a
460 two-step combination of the t-SNE method and K-means clustering: First, t-SNE is
461 used to reduce the dimensionality of the data from a high-dimensional space to a 3d-
462 or 2d-plot. Next, the K-means method is used for clustering the processed dataset
463 whose dimension is reduced.

464

465 *SC3*

466 Kiselev et al. developed a method called SC3 for determining cell types based on
467 transcriptome profiles alone, achieving high accuracy and robustness by combining
468 multiple clustering solutions through a consensus approach [11]. The workflow of
469 SC3 can be grouped into the following five steps: (1) filter out rare genes and
470 common genes to reduce the dimensionality of the data; (2) construct three distance
471 matrices between cells using the Euclidean, Pearson, and Spearman metrics; (3)
472 transform all distance matrices using either PCA or by calculating the eigenvectors of
473 the associated Laplacian; the columns of the resulting matrices are then sorted in

474 ascending order by their corresponding eigenvalues; (4) perform K-means clustering
 475 on the first d eigenvectors of the transformed distance matrices; (5) obtain a
 476 hierarchical clustering from the consensus matrix, which is constructed using the
 477 cluster-based similarity partitioning algorithm (CSPA). Each clustering result is
 478 represented by a binary similarity matrix, and the consensus matrix is calculated by
 479 averaging all similarity matrices.

480 A major bottleneck of SC3 is its longer running time compared to other models.
 481 Furthermore, for dealing with even larger datasets, SC3 implements a hybrid
 482 approach that combines unsupervised and supervised methodologies, which have a
 483 possible limitation of rare cell types not being identified.

484

485 *SSC*

486 The SSC algorithm is a novel approach to the subspace clustering problem using
 487 sparse representation. It is particularly useful for clustering data drawn from multiple
 488 low-dimensional subspace embedded in a high-dimensional space, a feature common
 489 to many scRNA-seq datasets. The SSC algorithm solves the subspace clustering
 490 problem in two steps. The first one is to solve the following global sparse
 491 optimization problem:

$$492 \quad \min \|C\|_1 \text{ s.t. } X = XC \text{ and } \text{diag}(C) = 0, \quad (8)$$

493 where $X \in R^{D \times N}$ represents the input matrix.

494 The output matrix $C \in R^{N \times N}$ is a block diagonal matrix in which the nonzero
 495 block corresponding to data points in the same subspace. In the second step, this
 496 information about the membership of data points is utilized in the spectral clustering
 497 framework to obtain predicted labels.

498 Although SSC performs well in many applications, it ignores the constraint
 499 relationship between the coefficient matrix and the clustering result, which is a major
 500 shortcoming of the algorithm. To alleviate this problem, several improved SSC
 501 variants have been proposed, including the SSSC described below.

502

503

SSSC

504

As an improvement to SSC [38], the structural sparse subspace clustering (SSSC)

505

algorithm proposes a unified joint optimization framework, which not only obtains

506

the spectral clustering result through the sparse optimization step but also constrains

507

the coefficient matrix by the clustering result in turn. To this end, an incidence matrix

508

$Q = [q_1 \dots q_n] \in R^{N \times n}$ is constructed to associate each data point with the subspace

509

that contains it using the clustering result from a previous iteration. Since each data

510

point belongs to one subspace, we have $Q\mathbf{1}=\mathbf{1}$ and $rank(Q) = n$, where $\mathbf{1}$ is the

511

vector of all ones of appropriate dimension. Consequently, the set of all feasible

512

incidence matrices is:

513

$$\mathbb{Q} = \{Q \in \{0,1\}^{N \times n}: Q\mathbf{1} = \mathbf{1}, rank(Q) = n\}. \quad (9)$$

514

Using a subspace structure norm defined as:

515

$$\|C\|_Q = \sum_{i,j} |C_{ij}| \left(\frac{1}{2} \|q^i - q^j\|^2 \right), \quad (10)$$

516

where q^i and q^j are the i -th and j -th rows of the matrix Q .

517

The subspace clustering problem is reformulated in SSSC into solving the

518

following optimization problem:

519

$$\min_{C,Q} \|C\|_Q + \|C\|_1 \quad s.t. \quad X = XC, \text{diag}(C) = 0, \text{ and } Q \in \mathbb{Q}. \quad (11)$$

520

This optimization problem can be solved efficiently via a combination of an

521

alternating direction method of multipliers (ADMM) with spectral clustering.

522

523

Seurat

524

Butler et al. developed a comprehensive R package, which is an indispensable tool in

525

the field of single cell RNA-seq analysis. [39]. This toolkit provides a number of

526

functions including t-SNE dimensionality reduction analysis, cluster analysis,

527

differential expression, construction of developmental trajectories, mark gene

528

recognition and so on. For this work, we are mainly interested in the cluster analysis

529

module, which is used to identify cell subtypes. Instead of a direct cluster analysis

530

applied to all cells, Seurat first performs a PCA to select the principal components

531

with the largest contribution, and then uses the selected principal components to

532 perform cluster analysis. The clustering algorithm includes original Louvain
 533 algorithm (The default), Louvain algorithm with multilevel refinement, SLM, and
 534 Leiden. The t-SNE dimensionality reduction technique is also employed to display
 535 expression distributions of cells in a 2d-plot, where cells in the same cluster are coded
 536 by the same color.

537

538 *SIMLR*

539 Wang et al. developed a novel similarity-learning framework called SIMLR [17],
 540 which learns an appropriate distance metric from combining multiple types of
 541 distances between cells. SIMLR assumes that cells in the same subpopulation are
 542 more similar, and the similarity matrix should have an approximate block diagonal
 543 structure in which the number of blocks is determined by the number of separable
 544 subpopulations of the input cells. In the default implementation of SIMLR, Gaussian
 545 kernels, which generate the best empirical performance among a number of candidate
 546 kernels, take the form:

$$547 \quad K(c_i, c_j) = \frac{1}{\epsilon_{ij}\sqrt{2\pi}} \exp\left(-\frac{\|c_i - c_j\|_2^2}{2\epsilon_{ij}^2}\right). \quad (12)$$

548 Here $\|c_i - c_j\|_2$ is the Euclidean distance between cell i and cell j , and the
 549 variance ϵ_{ij} can be calculated with different scales:

$$550 \quad \mu_i = \frac{\sum_{l \in KNN(c_i)} \|c_i - c_l\|_2}{k}, \quad \epsilon_{ij} = \frac{\sigma(\mu_i + \mu_j)}{2}, \quad (13)$$

551 where $KNN(c_i)$ represents cells that are top k neighbors of cell i .

552 SIMLR uses learned similarities to visualize cells, reduce the dimensionality of
 553 the input data, and cluster cells into subgroups, giving priority to genes with the
 554 highest variability that can explain differences in the entire population. Since the
 555 implementation of SIMLR needs the number of clusters as an input, it is not suitable
 556 for analyzing data with unknown structure.

557

558 *SinNLRR*

559 Based on similarity learning, Zheng et al. proposed a scRNA-seq cell-type detection
 560 method called SinNLRR [40], where non-negative and low-rank structures on the

561 similarity matrix are imposed. This leads to an optimization problem with the form:

$$562 \quad \min \frac{1}{2} \|X - XC\|_F + \lambda \|C\|_* \quad s. t. \quad C \geq 0, \quad (14)$$

563 where X is the input matrix, and C is a coefficient matrix in which the entry $C_{i,j}$
564 denotes the confidence of cells i and j in the same subpopulation.

565 SinNLRR applies the alternating direction method of the multiplier (ADMM) to
566 solve the optimization problem and proposes an adaptive penalty selection method to
567 avoid sensitivity to parameters. The learned similarity matrix can be visualized, and
568 Laplace scores can be used to prioritize gene markers. SinNLRR is benchmarked with
569 ten human and mouse scRNA-seq datasets, whose sizes range from dozens to
570 thousands of cells. SinNLRR obtained stronger robustness and more accurate results
571 using different datasets. At present, the main goal of SinNLRR is to reduce the
572 running time on large-scale scRNA-seq data.

573

574 *Corr*

575 By introducing a new similarity measure named differentiability correlation, Jiang et
576 al. proposed a hierarchical clustering-based algorithm called Corr to predict cell types
577 [32]. Differentiability correlation evaluates the similarity between any two cells by
578 using the correlation between the gene expression profiles of two cells and
579 incorporating information from all other cells. Since the relationship of cell-specific
580 gene expression patterns over the whole cell population is considered, this novel
581 measure turns out to be more robust against cell heterogeneity and data noise. Using
582 the framework of hierarchical clustering, Corr incorporates factorial ANOVA in
583 optimal cluster number determination, which allows the number of clusters to be
584 automatically determined. Corr is benchmarked with several real scRNA-seq datasets,
585 with outstanding performance and a correct cluster number obtained for each dataset.

586

587 **Evaluations of clustering**

588 To benchmark the performance of various clusters through the imputation algorithm
589 used, four clustering evaluation indicators are chosen to quantify the clustering

590 performance on each scRNA-seq dataset. Formally, let $P = \{p_1, p_2, \dots, p_m\}$ and
 591 $T = \{t_1, t_2, \dots, t_n\}$ represent the real cell type from the dataset and the cell type
 592 generated by clustering algorithm, respectively. The dissimilarity between P and T
 593 can then be measured by one of the following four indicators: normalized mutual
 594 information (NMI) [41], adjusted Rand index (ARI) [42], homogeneity (HOM) [43],
 595 and completeness (COM) [43].

596 597 *NMI*

598 The NMI score between P and T is defined as:

$$599 \quad NMI(P, T) = \frac{MI(P, T)}{\sqrt{H(P)H(T)}}, \quad (15)$$

600 where $H(P)$ and $H(T)$ denote the entropy of P and T , respectively, and $MI(P, T)$
 601 represents the mutual information between them. It is well known that NMI has an
 602 upper bound of 1 and lower bound of 0.

603 604 *ARI*

605 To define the ARI score, cell pairs in the dataset are classified into one of the
 606 following four types: the number of cell pairs that are in the same cluster in both P
 607 and T denoted by N_{11} ; the number of cell pairs that are in different clusters in both
 608 P and T by N_{00} ; the number of cell pairs that are in the same cluster in P but in
 609 different clusters in T by N_{10} ; the number of cell pairs that are in different clusters
 610 in P but in the same cluster in T by N_{01} . The ARI score between P and T is then
 611 defined as:

$$612 \quad ARI(P, T) = \frac{2(N_{00}N_{11} - N_{01}N_{10})}{(N_{00} + N_{01})(N_{01} + N_{11}) + (N_{00} + N_{10})(N_{10} + N_{11})}. \quad (16)$$

613 The ARI score is bounded above by 1 and equals 0 when the Rand index is the
 614 same as its expected value (under the generalized hypergeometric distribution for
 615 randomness).

616 617 *HOM*

618 One desired property for the output clustering T is that it satisfies the homogeneity

619 criteria; that is, each cluster in T contains only cells from a single cluster in P . To
620 measure how close the clustering T to this ideal situation, the homogeneity score
621 (HOM) is defined as:

$$622 \quad \text{HOM}_P(T) = 1 - \frac{H(P|T)}{H(P)}. \quad (17)$$

623 Here $H(P|T)$ is the entropy of P conditioned on T . Note that for a perfectly
624 homogeneous clustering T , we have $H(P|T) = 0$, and hence its homogeneity score is
625 1. We also use the convention that when P contains only one cluster, that is, $H(P) =$
626 0, the homogeneity score is always 1. Then the HOM score is between 0 and 1, with
627 close to 1 being desirable. However, being homogeneous alone is not sufficient for
628 good clustering. For instance, the trivial clustering T in which each cluster contains
629 only one cell always has a homogeneity score of 1. To deal with such cases, we will
630 consider one additional score based on completeness.

631

632 *COM*

633 To some extent, completeness is a property that is symmetrical to homogeneity. That
634 is, if each cluster in P contains only units from a single cluster in T , then cluster T
635 is complete. To measure how close the clustering T to this ideal situation, the
636 completeness score (COM) is defined as:

$$637 \quad \text{COM}_P(T) = 1 - \frac{H(T|P)}{H(T)}. \quad (18)$$

638 Similar to the homogeneity score, here we also use the convention that the
639 completeness score is 1 when T contains only one cluster. Furthermore, the
640 completeness score is between 0 and 1, with 1 being desirable. Note that the
641 homogeneity score and the completeness score run roughly in opposition: a high
642 homogeneity score often means a low completeness score. Hence, a clustering that is
643 high on both the homogeneity and the completeness scores is truly desirable because
644 it indicates that the clustering is indeed rather consistent with the golden standard.

645

646 **Authors' contributions**

647 Junlin Xu: Writing - original draft, Conceptualization, Data curation,
648 Methodology, Visualization. Lingyu Cui: Writing - review & editing,
649 Investigation, Formal analysis. Jujuan Zhuang: Investigation, Data curation. Yajie
650 Meng: Visualization, Writing - review & editing. Pingping Bing: Investigation,
651 Data curation. Bingsheng He: Investigation, Data curation. Geng Tian: Data
652 curation, Validation. Taoyang Wu: Investigation, Writing - review & editing.
653 Jialiang Yang: Writing - review & editing, Conceptualization, Supervision. All
654 authors read and approved the final manuscript.

655

656 **Competing interests**

657 JY and GT are currently employed by Geneis (Beijing) Co. Ltd.; All other authors
658 have declared no competing interests.

659

660 **Acknowledgements**

661 This work was supported by the Hunan Provincial Innovation Foundation for
662 Postgraduate of China (Grant No. CX20200434) awarded to Junlin Xu, and the
663 National Natural Science Foundation of China (Grant No. 62172004) to Bing Wang.
664 We thank LetPub (www.letpub.com) for its linguistic assistance during the
665 preparation of this manuscript.

666

667 **References**

- 668 1. Jaitin DA, Kenigsberg E, Keren-Shaul H et al. Massively parallel single-cell RNA-seq for
669 marker-free decomposition of tissues into cell types, *Science* 2014;343:776-779.
- 670 2. Kolodziejczyk AA, Kim JK, Svensson V et al. The technology and biology of single-cell RNA
671 sequencing, *Molecular cell* 2015;58:610-620.
- 672 3. Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in single-cell
673 transcriptomics, *Nature Reviews Genetics* 2015;16:133-145.
- 674 4. Zappia L, Phipson B, Oshlack A. Exploring the single-cell RNA-seq analysis landscape with
675 the scRNA-tools database, *PLoS computational biology* 2018;14:e1006245.
- 676 5. Yang J, Liao B, Zhang T et al. Editorial: Bioinformatics Analysis of Single Cell Sequencing Data
677 and Applications in Precision Medicine, *Front Genet* 2019;10:1358.

- 678 6. Zhang L, Zhang S. Comparison of computational methods for imputing single-cell
679 RNA-sequencing data, *IEEE/ACM transactions on computational biology and bioinformatics*
680 2018.
- 681 7. Grün D, Muraro MJ, Boisset J-C et al. De novo prediction of stem cell identity using
682 single-cell transcriptome data, *Cell stem cell* 2016;19:266-277.
- 683 8. Buettner F, Natarajan KN, Casale FP et al. Computational analysis of cell-to-cell
684 heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells, *Nature*
685 *biotechnology* 2015;33:155-160.
- 686 9. Peng L, Tian X, Tian G et al. Single-cell RNA-seq clustering: datasets, models, and algorithms,
687 *RNA Biol* 2020;17:765-783.
- 688 10. Tasic B, Menon V, Nguyen TN et al. Adult mouse cortical cell taxonomy revealed by single
689 cell transcriptomics, *Nature neuroscience* 2016;19:335-346.
- 690 11. Kiselev VY, Kirschner K, Schaub MT et al. SC3: consensus clustering of single-cell RNA-seq
691 data, *Nature methods* 2017;14:483-486.
- 692 12. Lin P, Troup M, Ho JW. CIDR: Ultrafast and accurate clustering through imputation for
693 single-cell RNA-seq data, *Genome biology* 2017;18:59.
- 694 13. Dey KK, Hsiao CJ, Stephens M. Visualizing the structure of RNA-seq expression data using
695 grade of membership models, *PLoS genetics* 2017;13:e1006599.
- 696 14. Zhuang J, Cui L, Qu T et al. A streamlined scRNA-Seq data analysis framework based on
697 improved sparse subspace clustering, *IEEE Access* 2021;PP:1-1.
- 698 15. Kim DH, Marinov GK, Pepke S et al. Single-cell transcriptome analysis reveals dynamic
699 changes in lncRNA expression during reprogramming, *Cell stem cell* 2015;16:88-101.
- 700 16. Macosko EZ, Basu A, Satija R et al. Highly parallel genome-wide expression profiling of
701 individual cells using nanoliter droplets, *Cell* 2015;161:1202-1214.
- 702 17. Wang B, Zhu J, Pierson E et al. Visualization and analysis of single-cell RNA-seq data by
703 kernel-based similarity learning, *Nature methods* 2017;14:414-416.
- 704 18. Mongia A, Sengupta D, Majumdar A. McImpute: Matrix completion based imputation for
705 single cell RNA-seq data, *Frontiers in genetics* 2019;10:9.
- 706 19. Badsha MB, Li R, Liu B et al. Imputation of single-cell gene expression with an autoencoder
707 neural network, *Quantitative Biology* 2020:1-17.
- 708 20. Chen C, Wu C, Wu L et al. scRMD: Imputation for single cell RNA-seq data via robust matrix
709 decomposition, *Bioinformatics* 2020;36:3156-3161.
- 710 21. Xu J, Cai L, Liao B et al. CMF-Impute: an accurate imputation tool for single-cell RNA-seq
711 data, *Bioinformatics* 2020;36:3139-3147.
- 712 22. Linderman GC, Zhao J, Kluger Y. Zero-preserving imputation of scRNA-seq data using
713 low-rank approximation, *bioRxiv* 2018:397588.
- 714 23. Biase FH, Cao X, Zhong S. Cell fate inclination within 2-cell and 4-cell mouse embryos
715 revealed by single-cell RNA sequencing, *Genome research* 2014;24:1787-1796.
- 716 24. Zeisel A, Muñoz-Manchado AB, Codeluppi S et al. Cell types in the mouse cortex and
717 hippocampus revealed by single-cell RNA-seq, *Science* 2015;347:1138-1142.
- 718 25. Pollen AA, Nowakowski TJ, Shuga J et al. Low-coverage single-cell mRNA sequencing
719 reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex,
720 *Nature biotechnology* 2014;32:1053.
- 721 26. Yan L, Yang M, Guo H et al. Single-cell RNA-Seq profiling of human preimplantation

722 embryos and embryonic stem cells, *Nature structural & molecular biology* 2013;20:1131.

723 27. Gong W, Kwak I-Y, Pota P et al. DrImpute: imputing dropout events in single cell RNA
724 sequencing data, *BMC bioinformatics* 2018;19:1-10.

725 28. Zhang L, Zhang S. PBLR: an accurate single cell RNA-seq data imputation tool considering
726 cell heterogeneity and prior expression level of dropouts, *bioRxiv* 2018:379883.

727 29. Qiu P. Embracing the dropouts in single-cell RNA-seq analysis, *Nature communications*
728 2020;11:1-9.

729 30. Ye P, Ye W, Ye C et al. scHinter: imputing dropout events for single-cell RNA-seq data with
730 limited sample size, *Bioinformatics* 2020;36:789-797.

731 31. Li WV, Li JJ. An accurate and robust imputation method scImpute for single-cell RNA-seq
732 data, *Nature communications* 2018;9:1-9.

733 32. Jiang H, Sohn LL, Huang H et al. Single cell clustering based on cell-pair differentiability
734 correlation and variance analysis, *Bioinformatics* 2018;34:3684-3694.

735 33. van Dijk D, Nainys J, Sharma R et al. MAGIC: A diffusion-based imputation method reveals
736 gene-gene interactions in single-cell RNA-sequencing data, *bioRxiv* 2017:111591.

737 34. Huang M, Wang J, Torre E et al. SAVER: gene expression recovery for single-cell RNA
738 sequencing, *Nature methods* 2018;15:539-542.

739 35. Kim T, Chen IR, Lin Y et al. Impact of similarity metrics on single-cell RNA-seq data clustering,
740 *Briefings in bioinformatics* 2019;20:2316-2326.

741 36. Petegrosso R, Li Z, Kuang R. Machine learning and statistical methods for clustering
742 single-cell RNA-sequencing data, *Briefings in bioinformatics* 2020;21:1209-1223.

743 37. Yau C. pcaReduce: hierarchical clustering of single cell transcriptional profiles, *BMC*
744 *bioinformatics* 2016;17:140.

745 38. Elhamifar E, Vidal R. Sparse subspace clustering: Algorithm, theory, and applications, *IEEE*
746 *transactions on pattern analysis and machine intelligence* 2013;35:2765-2781.

747 39. Satija R, Butler A, Hoffman P. Seurat: Tools for Single Cell Genomics, R package version
748 2018;2.

749 40. Zheng R, Li M, Liang Z et al. SinNLRR: a robust subspace clustering method for cell type
750 detection by non-negative and low-rank representation, *Bioinformatics* 2019;35:3642-3650.

751 41. Strehl A, Ghosh J. Cluster ensembles---a knowledge reuse framework for combining
752 multiple partitions, *Journal of machine learning research* 2002;3:583-617.

753 42. Hubert L, Arabie P. Comparing partitions, *Journal of classification* 1985;2:193-218.

754 43. Rosenberg A, Hirschberg J. V-measure: A conditional entropy-based external cluster
755 evaluation measure. In: *Proceedings of the 2007 joint conference on empirical methods in natural*
756 *language processing and computational natural language learning (EMNLP-CoNLL)*. 2007, p.
757 410-420.

758

759 **Figure legends**

760 **Figure 1. Schematic workflow of the scRNA-seq dataset cluster evaluation**
761 **framework**

762

763 **Figure 2. T-SNE visualization of cells from the Biase scRNA-seq dataset**

764 *Note:* Cells are color-coded by the cell type annotation of the original study.

765

766 **Figure 3. Benchmark of imputation algorithms on the t-SNE+k-means clustering**
767 **of scRNA-seq dataset**

768 **(A)** The use of different imputation algorithms to compare the four evaluation
769 indicators obtained on the Biase scRNA-seq dataset through t-SNE+K-means
770 clustering. **(B)** The use of different imputation algorithms to compare the four
771 evaluation indicators obtained on the Zeisel scRNA-seq dataset through
772 t-SNE+K-means clustering.

773

774 **Figure 4. T-SNE visualization of cells from the Zeisel scRNA-seq dataset**

775 *Note:* Cells are color-coded by the cell type annotation of the original study.

776

777 **Figure 5. Comparison of the impact of different imputation algorithms on**
778 **clustering performance on the Biase scRNA-seq dataset**

779

780 **Figure 6. Comparison of the impact of imputation algorithms on clustering**
781 **performance on the Pollen scRNA-seq dataset**

782

783 **Figure 7. Comparison of the impact of imputation algorithms on clustering**
784 **performance on the Zeisel scRNA-seq dataset**

785

786 **Supplementary material**

787 **Supplementary Table S1. Evaluate T-SNE visualization performance of different**
788 **impute algorithms on 10 scRNA-seq datasets**

789

790 **Supplementary Table S2. Compare the performance of clustering algorithms on**
791 **10 scRNA-seq datasets**

792

793 **Supplementary Table S3. Compare the effects of ARLA on various clustering**
794 **performances on 10 scRNA-seq datasets**

795

796 **Supplementary Table S4. Compare the effects of CMF-Impute on various**
797 **clustering performances on 10 scRNA-seq datasets**

798

799 **Supplementary Table S5. Compare the effects of DrImpute on various clustering**
800 **performances on 10 scRNA-seq datasets**

801

802 **Supplementary Table S6. Compare the effects of MAGIC on various clustering**
803 **performances on 10 scRNA-seq datasets**

804

805 **Supplementary Table S7. Compare the effects of mcImpute on various clustering**
806 **performances on 10 scRNA-seq datasets**

807

808 **Supplementary Table S8. Compare the effects of SAVER on various clustering**
809 **performances on 10 scRNA-seq datasets**

810

811 **Supplementary Table S9. Compare the effects of scImpute on various clustering**
812 **performances on 10 scRNA-seq datasets**

813

814 **Supplementary Table S10. Compare the effects of scRMD on various clustering**
815 **performances on 10 scRNA-seq datasets**

816

817 **Supplementary Table S11. Compare the effects of VIPER on various clustering**
818 **performances on 10 scRNA-seq datasets**

819