



RESEARCH ARTICLE

Comparison of clustering tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data [version 1; referees: 1 approved, 2 approved with reservations]

Saskia Freytag ^{1,2}, Luyi Tian^{2,3}, Ingrid Lönnstedt⁴, Milica Ng⁴, Melanie Bahlo ^{1,2}

¹Population Health and Immunity, Walter and Eliza Hall Institute of Medical Research, Parkville, Australia

²Department of Medical Biology, University of Melbourne, Parkville, Australia

³Molecular Medicine Division, Walter and Eliza Hall Institute of Medical Research, Parkville, Australia

⁴Bio21 Institute, CSL Limited, Parkville, Australia

v1 First published: 15 Aug 2018, 7:1297 (doi: [10.12688/f1000research.15809.1](https://doi.org/10.12688/f1000research.15809.1))
 Latest published: 15 Aug 2018, 7:1297 (doi: [10.12688/f1000research.15809.1](https://doi.org/10.12688/f1000research.15809.1))

Abstract

Background: The commercially available 10x Genomics protocol to generate droplet-based single-cell RNA-seq (scRNA-seq) data is enjoying growing popularity among researchers. Fundamental to the analysis of such scRNA-seq data is the ability to cluster similar or same cells into non-overlapping groups. Many competing methods have been proposed for this task, but there is currently little guidance with regards to which method to use.

Methods: Here we use one gold standard 10x Genomics dataset, generated from the mixture of three cell lines, as well as three silver standard 10x Genomics datasets generated from peripheral blood mononuclear cells to examine not only the accuracy but also robustness of a dozen methods.

Results: We found that some methods, including Seurat and Cell Ranger, outperform other methods, although performance seems to be dependent on the complexity of the studied system. Furthermore, we found that solutions produced by different methods have little in common with each other.

Conclusions: In light of this, we conclude that the choice of clustering tool crucially determines interpretation of scRNA-seq data generated by 10x Genomics. Hence practitioners and consumers should remain vigilant about the outcome of 10x Genomics scRNA-seq analysis.

Keywords

Clustering, Single-Cell RNA-seq, Benchmarking, 10x Genomics

Open Peer Review

Referee Status:

	Invited Referees		
	1	2	3
version 1 published 15 Aug 2018	 report	 report	 report

- 1 **Joshua W. K. Ho** , Victor Chang Cardiac Research Institute (VCCRI), Australia
- 2 **Shila Ghazanfar** , University of Sydney, Australia
- 3 **Stephanie Hicks** , Johns Hopkins Bloomberg School of Public Health (JHSPH), USA

Discuss this article

Comments (0)

Corresponding author: Saskia Freytag (freytag.s@wehi.edu.au)

Author roles: **Freytag S:** Conceptualization, Data Curation, Formal Analysis, Funding Acquisition, Investigation, Methodology, Software, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing; **Tian L:** Investigation, Writing – Review & Editing; **Lönnstedt I:** Conceptualization, Investigation, Methodology, Writing – Review & Editing; **Ng M:** Conceptualization, Funding Acquisition, Investigation, Methodology, Writing – Review & Editing; **Bahlo M:** Conceptualization, Investigation, Methodology, Project Administration, Supervision, Writing – Review & Editing

Competing interests: No competing interests were disclosed.

Grant information: We would like to thank the Australian Genome Research Facility and the Genomics Innovation Hub for their generous support of this project, including funding. This work was also supported by the Victorian Government's Operational Infrastructure Support Program and Australian Government NHMRC IRIIS. MB is funded by NHMRC Senior Research Fellowship 110297 and NHMRC Program Grant 1054618. *The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.*

Copyright: © 2018 Freytag S *et al.* This is an open access article distributed under the terms of the [Creative Commons Attribution Licence](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

How to cite this article: Freytag S, Tian L, Lönnstedt I *et al.* **Comparison of clustering tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data [version 1; referees: 1 approved, 2 approved with reservations]** *F1000Research* 2018, 7:1297 (doi: [10.12688/f1000research.15809.1](https://doi.org/10.12688/f1000research.15809.1))

First published: 15 Aug 2018, 7:1297 (doi: [10.12688/f1000research.15809.1](https://doi.org/10.12688/f1000research.15809.1))

Introduction

Single-cell RNA-sequencing (scRNA-seq) studies have opened the way for new data-driven definitions of cell identity and function. No longer is a cell's type determined by arbitrary hierarchies and their respective predefined markers. Instead, a cell's transcriptional and epigenomic profile can now be used¹ to accomplish this task. This is achieved using computational methods for scRNA-seq that characterize cells into novel and known cell types. Characterization consists of two steps: (i) unsupervised or semi-supervised clustering of same or similar cells into non-overlapping groups, and (ii) labeling clusters, i.e. determining the cell type, or related cell types, represented by the cluster. Here, we focus on the first step of this process.

Research into clustering has produced many algorithms for the task, including over 60 tools specifically designed for scRNA-seq². Due to the relative youth of the field, there are currently no rules guiding the application of these clustering algorithms. If tools' performances have been tested outside synthetic scenarios, testing seems to be confined to scenarios with limited biological variability. Furthermore, most tools were developed and consequently tested only on the Fluidigm C1 protocol, despite considerable differences in throughput capabilities and sensitivities³ in the different scRNA-seq platforms. Here we focus solely on

clustering performance on medium-sized scRNA-seq data generated by 10x Genomics as it is currently the most widely used platform. Commercially available scRNA-seq platforms, like 10x Genomics' Chromium, are being widely adopted due to their ease of use and relatively low cost per cell⁴. The 10x Genomics protocol uses a droplet-based system to isolate single cells. Each droplet contains all the necessary reagents for cell lysis, barcoding, reverse transcription and molecular tagging. This is followed by pooled PCR amplification and 3' library preparation, after which standard Illumina short-read sequencing can be applied⁵. Unlike other commercially available scRNA-seq protocols, like Fluidigm C1, 10x Genomics allows for sequencing of thousands of cells albeit at much shallower read depths per cell, and without allowing the use of fluorescence markers to establish cell identity. As such the 10x Genomics platform is particularly suited to detailed characterization of heterogeneous tissues.

Methods

In this study, we performed comprehensive evaluation of a dozen clustering methods (Table 1). We focused on analysis methods available in the R language, as this is one of the most commonly used programming languages for scRNA-seq data analysis. The exception to this is the 10x Genomics software

Table 1. Overview of the clustering tools included in this study, and several characteristics thereof.

Software	Year	Description	Properties	Ref
ascend version 0.4.0	2017	Adjusted RLE normalization followed by hierarchical clustering and merging	Gene filtering, limited documentation, medium user interaction	6
Cell Ranger version 2.0.0	2016	Graph-based clustering on the first 10 principal components	Gene filtering, detailed documentation, no user interaction	
CIDR version 0.1.5	2017	Imputation of potential dropout genes followed by hierarchical clustering on first 4 principal components	Gene filtering, imputation good documentation, high user interaction	7
countClust version 1.4.1	2014	Likelihood models to estimate a specified number of multinomial distributions	Requires number of clusters, limited documentation, medium user interaction	8
RaceID	2015	Two iterations of k-means clustering with merging of outlier cells and identification of rare cell types in last step	Gene filtering, limited documentation, little user interaction	9
RaceID2	2016	More advanced version of RaceID based on UMIs	Gene filtering, limited documentation, little interaction	10
RCA version 1.0	2017	Rudimentary filtering then projection onto reference datasets consisting of profiles of isolated cell types	Gene filtering, reference dataset required, limited documentation, medium user interaction	11
SC3 version 1.7.7	2016	Rudimentary filtering followed by ensemble clustering method	Gene filtering, good documentation, medium user interaction	12
scrn version 1.6.9	2016	Library scale normalization by cell pools followed by hierarchical clustering on rank correlation-based distances of marker genes	Gene filtering, requires marker genes, detailed documentation, medium user interaction	13
Seurat version 2.3.0	2015	Normalization using mitochondrial RNA followed by PCA of highly variable genes and then graph-based clustering	Gene filtering, detailed documentation, high user interaction	14
SIMLR version 1.4.1	2016	Multikernel learning finds best fit and forces blocks in similarity matrix to address dropouts then applies spectral clustering	Requires number of clusters, good documentation, little user interaction	15
TSCAN version 1.16.0	2016	<i>In-silico</i> pseudo time reconstruction with a cluster-based minimum spanning tree approach to order cells	Good documentation, little user interaction	16

Cell Ranger. Our evaluation comprised four core aspects: (i) accuracy of clustering solutions compared to a gold standard (near absolute truth, limited variability and complexity), (ii) performance of clustering methods using silver standard data (no absolute truth, realistic variability and complexity), (iii) stability of clustering solutions, and (iv) miscellaneous characteristics, such as time and practicality.

Data

Gold standard. Three human lung adenocarcinoma cell lines, HCC827, H1975 and H2228, were cultured separately. The cell lines were obtained from ATCC and cultured in Roswell Park Memorial Institute 1640 medium with 10% fetal bovine serum (FBS, catalog number: 11875-176; Thermo Fisher Gibco) and 1% penicillin-streptomycin. The cells were grown independently at 37°C with 5% carbon dioxide until near 100% confluence. Before mixing cell lines, cells were dissociated into single-cell suspensions in FACS buffer (phosphate-buffered saline (PBS); catalog number: 14190-144; Thermo Fisher Gibco) with 5% FBS, Corning, catalog number: 35-076-CV, stained with propidium iodide (catalog number: P21493; Thermo Fisher FluoroPure) and 120,000 live cells were sorted for each cell line by FACS (BD FACSAria III flow cytometer, BD FACSDiva software version 7.0; BD Biology) to acquire an accurate equal mixture of live cells from the three cell lines. The resulting mixture was then processed by the Chromium Controller (10x Genomics) using single Cell 3' Reagent Kit v2 (Chromium Single Cell 3' Library & Gel Bead Kit v2, catalog number: 120237; Chromium Single Cell A Chip Kit, 48 runs, catalog number: 120236; 10x Genomics) (see Table 2). Afterwards the library was sequenced using Illumina NextSeq500 and V4 chemistry (NextSeq 500/550 High Output Kit v2.5, 150 Cycles, catalog number: 20024907; Illumina) with 100bp paired end reads. RTA (version 1.18.66.3; Illumina) was used for base calling.

Silver standard. We consider three human peripheral blood mononuclear cells (PBMCs) scRNA-seq datasets to be the silver standard (Table 2). All datasets were generated using the 10X Genomics droplet system combined with Illumina sequencing. The Australian Genome Research Facility in partnership with CSL generated one dataset using the 10x Genomics Chromium system (Dataset 1). Two datasets were generated by 10x Genomics and are publicly available

(Datasets 2 and 3). Of these, one dataset was generated with an earlier version of the microfluidics instrument, the 10x Genomics GemCode Controller (Dataset 2). The second dataset was generated with the latest instrument, the 10x Genomics Chromium Controller (Dataset 3).

For the first dataset, PBMCs were isolated from whole blood obtained through the Australian Red Cross Blood Service in the following manner. First, 50ml of blood was diluted using 50ml of PBS (catalog number: D8537-500ml; Sigma-Aldrich). We then added 30ml of Ficoll-Paque medium (catalog number: Catalog: 17-1440-03; GE Healthcare). We then centrifuged at room temperature for 20 minutes at 400 g and carefully removed the interface layer containing PBMCs, located between the top plasma layer and middle layer (Heraeus Multifuge 3 S-R Centrifuge; Thermo Fisher Scientific). To remove the supernatant, we further centrifuged at 400 g for 10 minutes at room temperature. This process was repeated to remove the contaminating Ficoll medium or platelets. Finally, cells were resuspended in 20ml of cell culture media with 5% FBS (RPMI-1640 Medium, catalog number: R0884-500ml, Sigma-Aldrich) and counted (Nikon Eclipse TS100 Microscope; Nikon). The resulting mixture was then processed by the Chromium Controller (10x Genomics) using single Cell 3' Reagent Kit v2 (Chromium Single Cell 3' Library & Gel Bead Kit v2, catalog number: 120237; Chromium Single Cell A Chip Kit, 48 runs, catalog number: 120236; 10x Genomics). Afterwards the library was sequenced using HiSeq2500 (Illumina) and V4 chemistry (HiSeq PE Cluster Kit v4 cBot, catalog number: PE-401-4001; HiSeq SBS Kit V4 50 cycles, catalog number: FC-401-4002; Illumina) with 101bp paired end reads. RTA (version 1.18.66.3; Illumina) was used for base calling.

Preprocessing

We used the 10x Genomics software, Cell Ranger (version 2.0.0) to align, de-duplicate, filter barcodes and quantify genes for all datasets. Note that we aligned reads with Cell Ranger to the GRCh38 (version 90) genome annotation. Using the Bioconductor package *scater*¹⁷ (version 1.6.3), we then removed low quality data from cells with low library size or low number of expressed gene transcripts. We also removed cells with a high mitochondrial read proportion as this can indicate apoptosis, also known as programmed cell death. Stressed cells

Table 2. Properties of all benchmarking datasets used in the study.

Benchmark standard	Gold	Silver		
Dataset		Dataset 1	Dataset 2	Dataset 3
Tissue	Cell lines	PBMCs	PBMCs	PBMCs
Source	GSE111108	GSE115189	Website*	Website*
Instrument	Chromium	Chromium	GemCode	Chromium
Number of cells	1,039	3,372	2,690	4,337
Total genes detected	29,451	24,654	20,693	25,820
<i>After preprocessing</i>				
Number of cells	925	3,205	2,590	4,292
Mean counts per cell	114,426	3,818	2,605	4,528
Median genes detected per cell	8,499	1,158	877	1,318

*<https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k>

*<https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.2.0/pbmc4k>

undergoing apoptosis have an aberrant transcriptome profile in comparison to a living cell and have previously been acknowledged to adversely influence transcriptome studies¹³.

Establishing truth

Gold standard. By exploiting the genetic differences between the three different cell lines we were able to establish near absolute truth in the gold standard dataset. To this end we first called single nucleotide variants (SNVs) in publicly available bulk RNA-seq of the same cell lines (GSE86337)¹⁸. Drawing on these SNVs, we then apply *demuxlet*¹⁹ (version 0.0.1), which harnesses the natural genetic variation between the cell lines to determine the most likely identity of each cell. We observe almost complete concordance between the result from *demuxlet* and clustering of cells seen in dimension reduction visualizations of the data (compare [Supplementary Figure 1](#)).

Silver standard. For the silver standard data, we compared clustering solutions to a cell labeling approach by 10x Genomics⁵. This approach finds the cell type in a reference dataset which most closely resembles the expression in the cell. The reference dataset contains 11 isolated cell types sequenced using the 10x Genomics system. While this labeling does not constitute truth, it has been found to perform well in comparison with marker-based classification⁵. Furthermore, the proportions of cells assigned to the 11 cell types by the supervised labeling approach were consistent with the literature (see [Supplementary Table 1](#))^{20,21}.

Criteria for inclusion of clustering tool

We based our selection of method on the online list within www.scRNA-tools.org² in October 2017. We only considered methods with an R package that had sufficient documentation to enable easy installation and execution and had at least one preprint or publication associated with it. We also excluded any methods that required extensive prior information not provided in the package. We also excluded any methods that continually failed to run (e.g. *Linnorm*²² and *Monocle*²³). This resulted in the evaluation of 12 methods (see [Table 1](#)). Note that for some of the R packages the primary focus is not clustering, but the package authors explicitly describe how their packages can be applied to achieve clustering of the scRNA-seq data.

The aim of this study is to provide guidance for the use of clustering methods to non-experts. Hence, we use all clustering methods with their default parameters as this represents the most common use case. In the case of *countClust* and *SIMLR* parameters included the number of clusters, which we set to 3 and 8 for the gold standard and silver standard datasets, respectively. Marker genes were required for the analysis with *scran*, which we obtained by performing differential expression analyses on GSE86337 and an in-house dataset of isolated cell types in PBMCs²⁴ for the gold standard and silver standard datasets, respectively. Furthermore, we also followed upstream data handling, such as filtering of genes and normalization, as described in the documentation of the respective clustering method. We concede that it is possible that more care in the upstream data handling and selection of parameters could result in different

results. However, confronted with the extremely large number of parameter choices, we believe that this evaluation suffices to identify strengths and weaknesses of each method.

Methods for the comparison of clustering solutions

To evaluate the performance and similarity of different clustering solutions, we rely on three different metrics. We use the adjusted Rand index (ARI)²⁵ and the normalized mutual information (NMI)²⁶, two metrics routinely applied in the field of clustering, to assess the similarity of clustering solutions or their similarity to a known truth. Both metrics can take values from 0 to 1, with 0 signifying no overlap between two groupings and 1 signifying complete overlap. These metrics are also applicable in the absence of known cluster labels. Finally, we also use a homogeneity score²⁷. This score takes the value 1 when all of its clusters contain only data points that are members of a single known group. Values of this score closer to 0 indicate that clusters contain mixed known groups. Unlike NMI and ARI, this score requires knowledge of an underlying truth.

Let X be a finite set of size n . A clustering solution C is a set C_1, \dots, C_k of non-empty disjoint subsets of X such that their union equals X . Let $C' = C'_1, \dots, C'_l$ be a second clustering solution or the supervised labeling solution with the same properties. The contingency table $M = (m_{ij})$ of the pair of sets C, C' is a $k \times l$ matrix whose i, j -th entry equals the number of elements in the intersection of clusters C_i and C'_j :

$$m_{ij} = |C_i \cap C'_j|, 1 \leq i \leq k, 1 \leq j \leq l.$$

ARI

$$R_{adj}(C, C') = \frac{\sum_{i=1}^k \sum_{j=1}^l \binom{m_{ij}}{2} - t_3}{\left(\frac{1}{2}(t_1 + t_2) - t_3\right)},$$

Where $t_1 = \sum_{i=1}^k \binom{|C_i|}{2}$, $t_2 = \sum_{j=1}^l \binom{|C'_j|}{2}$ and $t_3 = \frac{2t_1 t_2}{n(n-1)}$. For ease of notation this is referred to as ARI in the text, dropping the reference to specific pairs of sets. Furthermore, we also distinguish between *ARI_truth* as a comparison of a clustering solution to an underlying known or suspected truth and *ARI_comp*, which refers to a comparison between two clustering solutions.

NMI

$$NMI_1 = \frac{I(C, C')}{\sqrt{(H(C)H(C'))}},$$

where $H(C) = I(C, C)$ is the entropy of C . Note that

$$I(C, C') = \sum_{i=1}^k \sum_{j=1}^l P(i, j) \log_2 \frac{P(i, j)}{P(i)P(j)},$$

where $P(i, j) = \frac{m_{ij}}{n}$ and $P(i) = \frac{|C_i|}{n}$, is the mutual information of C and C' .

Homogeneity. Now let us assume C' is the known and correct grouping of the cells. Then,

$$\text{Homogeneity} = \frac{I(C, C')}{H(C')}$$

Stability assessment

To test the robustness of different clustering methods we pursued a sampling strategy in terms of genes and cells. We used Dataset 3 for the cell robustness evaluation with regards to cells, since it had the most number of cells. Similarly, we used Dataset 1 for the robustness evaluation with regards to genes, since it had the most number of non-zero genes after filtering. The impact of different aligners and preprocessing was assessed using all appropriate combinations of programs.

Cells. We randomly sampled 3,000 cells in Dataset 3 (out of the total of 4,292 that were available after filtering), generating five (non-independent) datasets. For every combination of two datasets (10 combinations in total) we then investigated for each clustering method separately how often cells contained in all five sampled datasets were assigned to the same cluster using the ARI_comp.

Genes. We randomly filtered half of all genes in Dataset 1 (out of the total of 58,302 genes), generating 10 datasets. For every combination of two datasets (45 combinations in total) we then investigated for each clustering method separately how often cells were assigned to the same cluster using the ARI_comp.

Aligners and preprocessing pipelines. In order to assess the affect of using different preprocessing pipelines on the data, we applied the Bioconductor package `scPipe`²⁸ (version 1.0.6) to the raw data. Like Cell Ranger, `scPipe` can be used to align, deduplicate, filter barcodes and quantify genes. Since `scPipe` is modular, we tried it with both the `STAR`²⁹ (version 020201) and `Subread`³⁰ (version 1.5.2) aligners. In order to ensure comparability we aligned reads to the same GRCh38 genome annotation and repeated quality control with `scater`. We investigated the similarity of clustering solutions applied to the differently preprocessed and aligned versions of the same dataset by ARI_comp.

Run time

Each execution of a method on a dataset was performed in a separate R session. Each task was allocated as many CPU cores of a 24 core Intel(R) Xeon(R) CPU E5-2690 v3 @ 2.60GHz as specified by the default parameters. The `base::set.seed` was overridden in order to prevent stochasticity and thus give reduced unwanted variation in the results. Timings for each method include any preprocessing steps.

Influence assessment

We also investigated what properties of each cell's data were driving the clustering solutions produced by the different methods as well as the inferred cell labels. Properties of a cell's data refer to features such as the number of total reads that included

the cell's barcode, the total number of gene transcripts found for this cell, etc. To this end, we used linear mixed models where cell data properties were predicted using the indicators for cluster membership. We predicted cell data properties and not cluster membership for modeling ease. The adjusted R^2 of these models was used to assess which properties influenced the clustering solutions. Properties investigated included: (i) the total number of detected gene transcripts, (ii) the total read count, and (iii) the percentages of reads aligning respectively to ribosomal proteins, mitochondrial genes and ribosomal RNA.

Results

Evaluation of clustering tools

Gold standard data. For the gold standard data set consisting of three cell types, half of the tested clustering methods overestimated the true number of different cell types in the data. Methods with cluster number estimations close to the correct number of different cell types included methods with prior information, such as `SIMLR`, `countClust` and `scran`, as well as `ascend`, `Cell Ranger`, `RaceID` and `CIDR` (Figure 1). The clustering solutions produced by these methods, with the exception of `countClust`, largely reflect the cell types. This is indicated by `ARI_truth` > 0.8. The remaining methods overestimated the number of clusters by 2 to 85 clusters, with `SC3` and `RaceID2` representing the extremes, both estimating more than 20 clusters (see t-SNE plots in Supplementary Figure 1 for the impact). As a consequence of the greater number of estimated clusters, the `ARI_truth` of the other clustering methods is lower than 0.8. To see whether these methods split cell types into several clusters or instead assign cells types randomly to clusters, we also investigate the homogeneity of the clustering solutions with respect to the known labeling. Apart from `countClust` and `RCA`, all methods have extremely high homogeneity, indicating that they split cell types into more subtypes, rather than randomly creating more cell types, which is reassuring.

Silver standard data. We labeled the cells in each of the three silver standard datasets as one of 11 different PBMC cell populations. When using the `ARI_truth` to compare the likeness of the clustering solutions and the labels, no method produced solutions that were uniformly the most similar to the inferred labels (Figure 2). Two methods, `ascend` and `countClust`, tended to estimate smaller number of clusters and consequently did not agree with the labeling. Only `Seurat`, `SC3` and `Cell Ranger` achieved an `ARI_truth` above 0.4 for at least two out of the three silver standard datasets. All methods considerably improved their `ARI_truth` when we subset to more confidently labeled cells (see Supplementary Figure 2). `RCA` was particularly affected, showing much greater similarity for more confidently labeled cells. We also calculated the homogeneity of each method in each dataset with respect to the inferred labeling. Generally, most methods exhibited significantly lower performance on Dataset 2, which was generated with an older version of the 10x Genomics technology than that used to generate Datasets 1 and 3. Apart from `RCA` all methods had much lower accuracy than for the gold standard

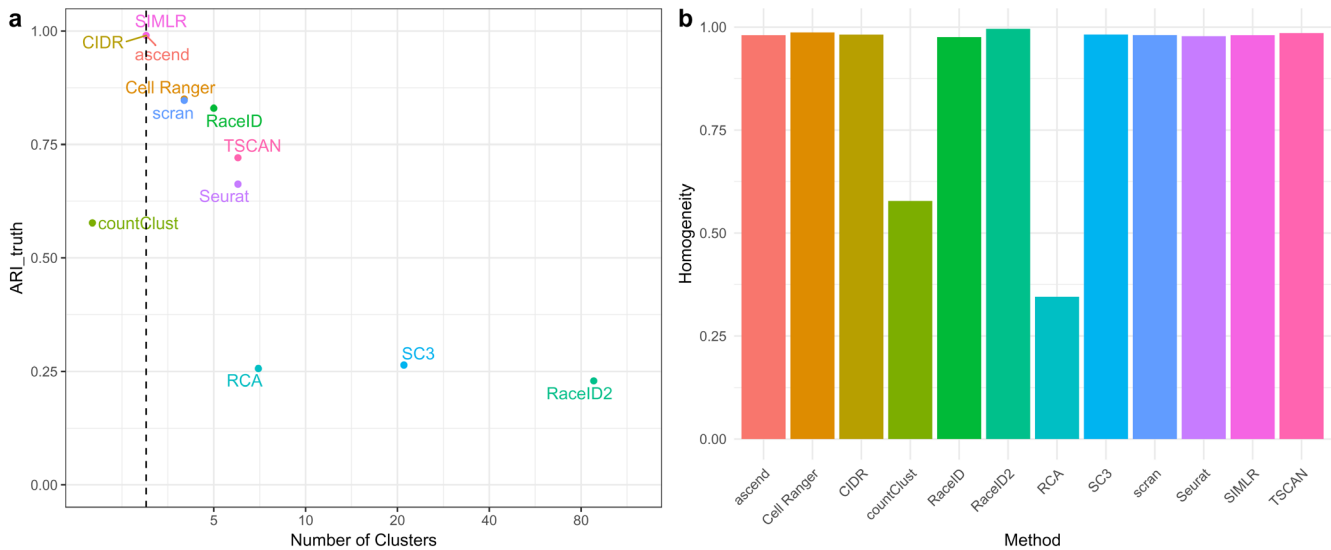


Figure 1. Performance on the gold standard dataset. (a) ARI_truth of each method with regards to the truth versus the number of clusters. The dashed line indicates the true number of clusters. (b) Homogeneity of clusters of each method, given the truth.

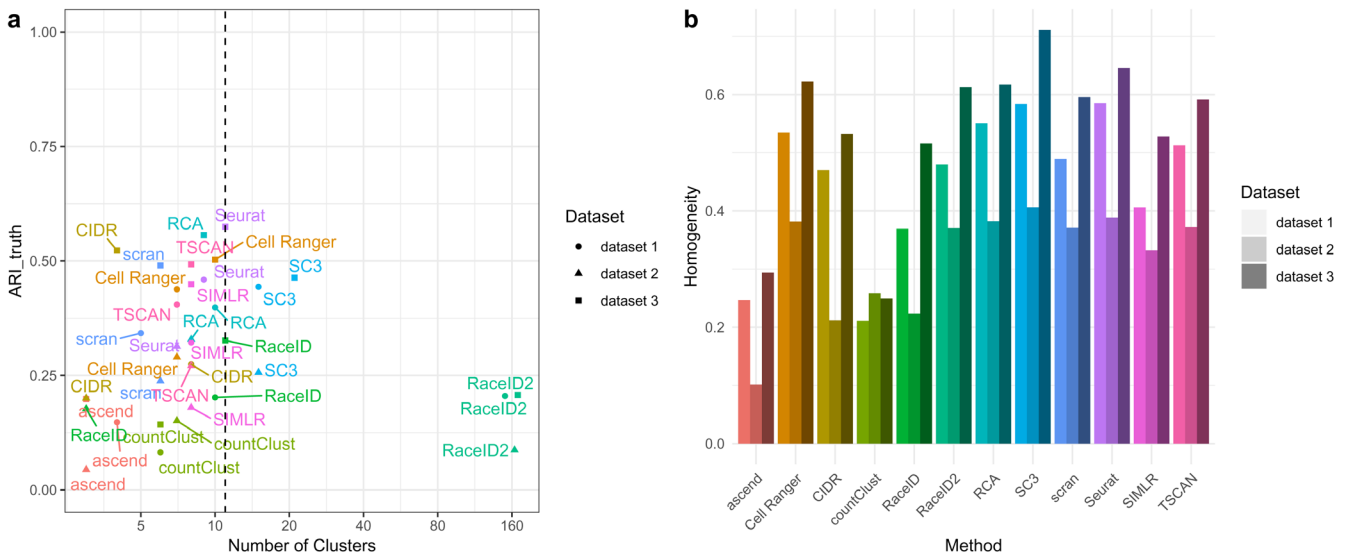


Figure 2. Performance on the three silver standard datasets. (a) ARI_truth of each method in each dataset, as indicated by different shapes, with regards to the supervised cell labeling versus the number of clusters. The dashed line indicates the number of cell populations estimated by the supervised cell labeling approach. (b) Homogeneity of clusters with regards to the inferred labeling for each method and each dataset. Different datasets are indicated by transparency.

data, indicating that most clusters represent mixtures of different inferred cell types. The exception is SC3’s clustering solution of Dataset 3, which achieved an homogeneity score above 0.7.

Interestingly, similar performance when compared to the labeling did not imply that cluster solutions were similar (compare Figure 3). In fact, only a few methods resulted in clustering solutions that were similar. For all three datasets, a group

of five methods (RCA, scran, Seurat, SIMLR and TSCAN) produced similar results, while the other seven methods appeared dissimilar between each other and to the set of five methods.

Stability. We evaluated the stability of the clustering methods by examining three different features: (i) filtering of cells, (ii) filtering of genes (Figure 4), and (iii) use of different aligners

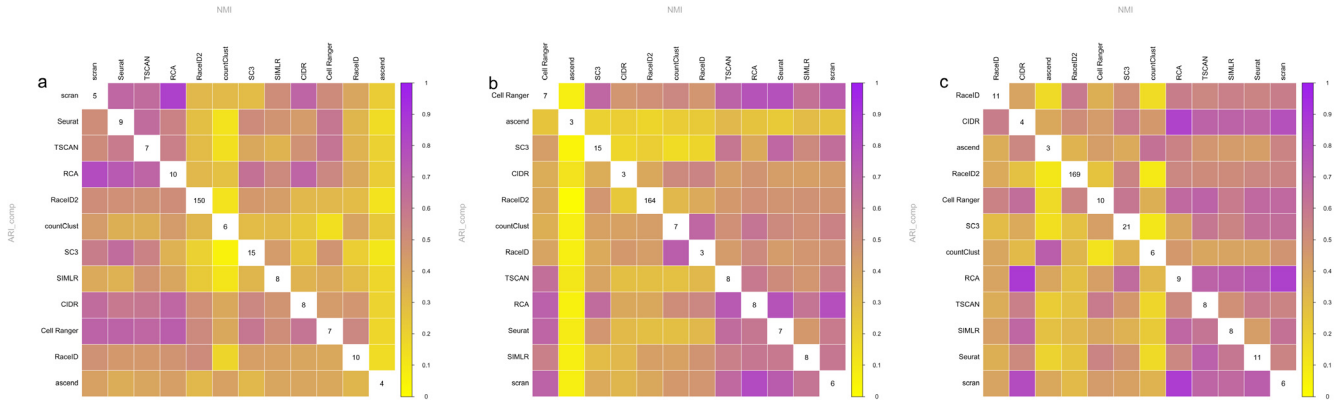


Figure 3. Similarity of all combinations of clustering methods as estimated by ARI_comp (lower triangle) and NMI (upper triangle) in (a) Dataset 1, (b) Dataset 2, and (c) Dataset 3. The similarity is indicated by the color; yellow indicating no similarity and purple indicating complete overlap. The diagonals give the number of clusters estimated by each respective method.

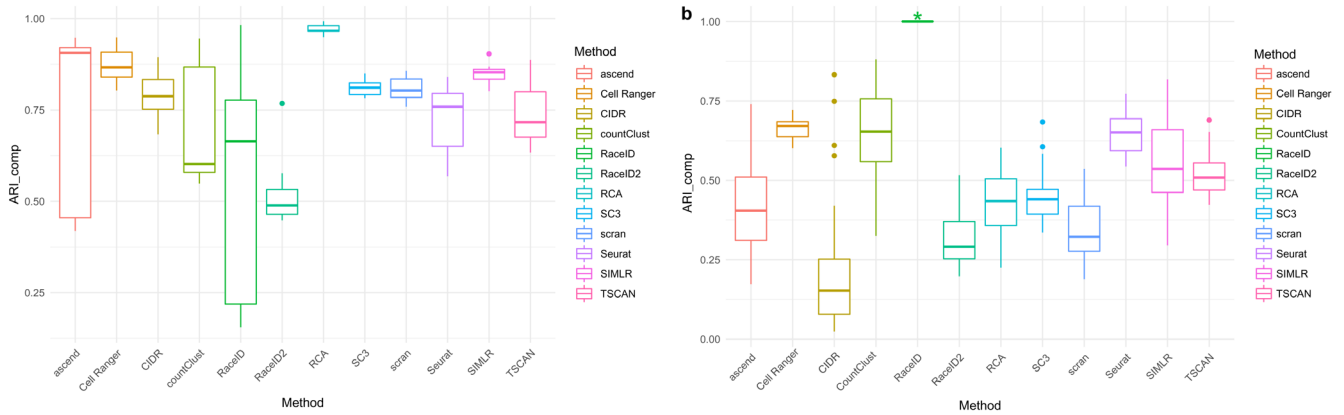


Figure 4. Tukey boxplots of ARI_comp results from the comparison of clustering solutions of the same method when (a) cell input was varied in Dataset 3 and (b) gene input was varied in Dataset 1. Note that RaceID only estimated 1 cluster when genes were varied.

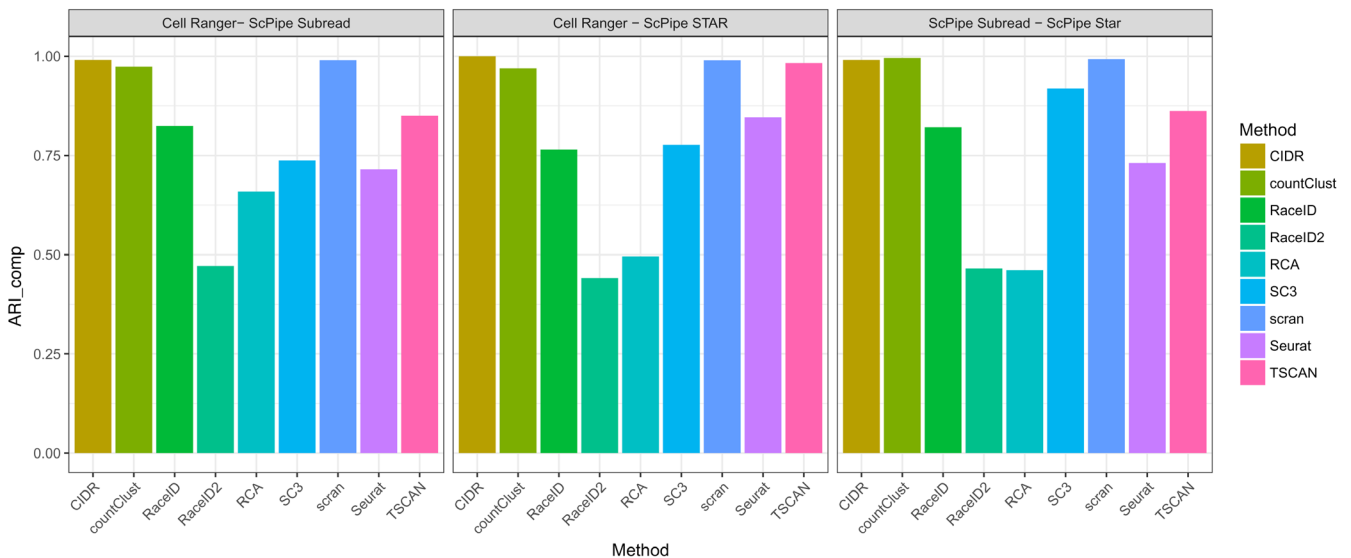


Figure 5. ARI_comp results from the comparison of clustering solutions of the same method used on datasets processed with different aligners. Only the subset of nine methods that worked in conjunction with all three aligners are shown.

(Figure 5). When assessing the stability with regards to input, countClust, RaceID and RaceID2 did not appear very robust. The robustness of ascend, countClust and RaceID was variable. Due to its reliance on reference profiles RCA is extremely robust, achieving ARI_comp above 0.9 consistently. Seurat, SIMLR and Cell Ranger demonstrated robustness with regards to input, but also exhibited robustness when changing gene filtering procedures (compare Figure 4b). CIDR appeared to be very sensitive to changes in gene filtering, which may be due to its imputation feature.

We also investigated how the stability of the clustering method was affected by the use of different aligners (Figure 5). In particular, we used Cell Ranger and ScPipe²⁸ with Subread³⁰, or STAR²⁹. We found that different aligners largely result in the same gene counts, but with some notable exceptions for processed pseudogenes (see Supplementary Figure 4, Supplementary Figure 5 and Supplementary Figure 6). Not all methods were able to be used in conjunction with scPipe. This included ascend and SIMLR, which failed to run, and Cell Ranger, which requires output from its own preprocessing pipeline. However we were able to evaluate eight methods. Apart from RaceID2 and RCA, all tested methods appeared robust.

Miscellaneous properties. Running time varies substantially between different methods. Methods like RaceID and RaceID2 take prohibitively long and thus do not lend themselves to interactive analysis when applied to 10x Genomics data (Figure 6). The fastest methods were RCA and TSCAN, with both taking less than 25 seconds on average for the entire dataset analysis. Their fast running time is due to both of these methods offering little

flexibility or intermediate results during their analysis (compare Table 1). By contrast Seurat's relatively long running time is partially due to extensive quality control during the analysis. Also note that methods differed in the quality of their documentation. For example, tools like Cell Ranger and Seurat offer detailed documentation, with many different use cases as well as tutorials. Tools, which are not found on Bioconductor, such as RaceID, RaceID2, ascend and RCA have more limited documentation.

Factors influencing clustering solutions

The variation in the percentage of reads aligning to ribosomal protein genes strongly predicted all clustering solutions as well as the inferred cell labels (see Figure 7). Expression of ribosomal protein genes has been successfully used to discriminate cell types belonging to different hematopoietic lineages³¹. Hence, it may be the case that overall mRNA amount of ribosomal protein genes can also serve as a discriminator. Furthermore, differences in abundance of ribosomal protein genes are likely to drive variation in PBMCs scRNA-seq datasets, as they typically account for a large proportion of reads (around 40% in all three datasets). In combination with ribosomal protein genes being less affected by dropout due to their relatively high expression, it is perhaps unsurprising that clustering solutions of all methods foremost reflect differences in the amount of ribosomal protein genes between cells.

Most methods' solutions were much more driven by the total number of features and total number of counts than the inferred solution. TSCAN was particularly affected ($R^2 = 0.52$), but for both ascend and RaceID2 similar effects were observed.

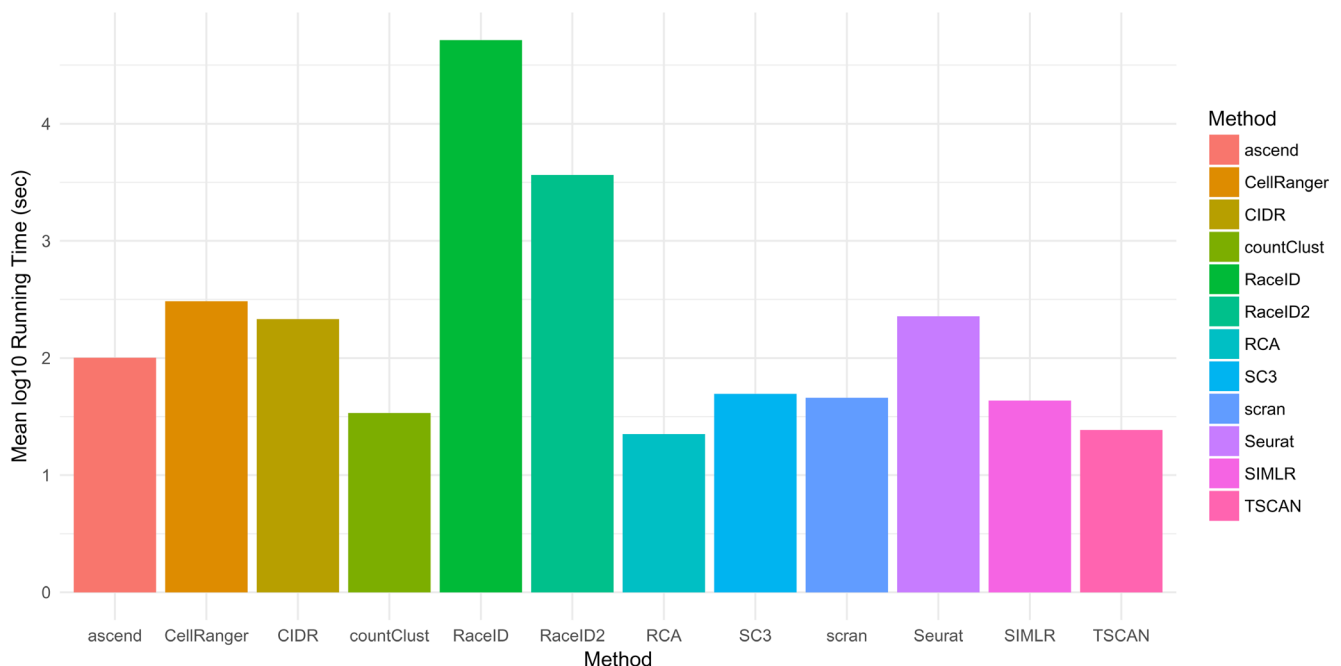


Figure 6. The bars indicate the average log₁₀ run time (in seconds) of all 12 methods on Dataset 1 with 29,151 genes over 10 iterations.



Figure 7. Radial plots describing the average effect of 5 cell features on the clustering solutions of different methods across the three silver standard datasets. For every method and every feature the adjusted R^2 of the linear model fitting the feature by the clustering solution is presented.

It can be speculated that this strong influence of total number of features and total number of count on their clustering solutions points to a failure to appropriately normalize the data.

Discussion

Most biological conclusions obtained from droplet-based scRNA-seq data crucially rely on accurate clustering of cells into homogeneous groups. Indeed, one can argue that it is the very act of clustering that unlocks the technology's potential for discovery. Therefore it is not surprising that according to several repositories, such as www.omicstools.org and www.scRNA-tools.org², many of the tools developed for scRNA-seq specifically focus on clustering. With so many choices, it is thus important to evaluate their performance for droplet based protocols, such as 10x Genomics, specifically.

In this study, we presented our evaluation of a dozen clustering method on scRNA-seq 10x Genomics data. The results of our

investigations will be useful for method users, as we provide clear and practical guidelines. Nonetheless, our evaluation has several limitations:

- Inclusion of methods limited to R packages and methods published before October 2017
- Parameter selection limited to defaults
- No assessment of robustness to noise and parameter changes
- No assessment of ability to discover rare cell populations
- Evaluation of more silver standard datasets from systems other than PBMCs
- No evaluation of quality of code and documentation
- No assessment of scalability of methods

Our evaluations suggest that *Seurat* and *Cell Ranger* provide the most stable and accurate clustering solutions for 10x Genomics scRNA-seq data. While *Seurat* performed slightly better, the choice between *Seurat* and *Cell Ranger*, in our opinion, should be informed by the user's familiarity with statistical concepts, which enable them to make the informed parameter choices required in *Seurat*. More generally we find that different clustering methods resulted in very different solutions. The good performance of *Seurat* as well as the vast difference between clustering methods have also been observed by Duò *et al.*³² in a benchmarking study including multiple scRNA-seq protocols. Our investigations suggest that biological differences between cells, such as cell type or state, and technical variation between cells (as well as combinations of biological differences and technical variation) all drive clustering. However, which aspects are captured by which clustering method remain to be confirmed. Our study merely pinpoints some of the drivers of performance, but not their origin, and thus cannot anoint an overall best method.

We recommend that practitioners and consumers of results generated from 10x Genomics scRNA-seq data alike remain vigilant about the outcome of their analysis, and acknowledge the variability and likelihood of undesired influences. The choice of clustering tool for scRNA-seq data generated by the 10x Genomics platform crucially determines interpretation. Hence, at least two clustering tools should be routinely applied to 10x Genomics scRNA-seq data in order to offer more than one subjective interpretation and hence increase robustness and confidence in any results.

Data availability

Repository: Gold Standard Dataset. Single cell profiling of 3 Human Lung Adenocarcinoma cell lines, GSE111108

Repository: Silver Standard Dataset 1. Single cell profiling of peripheral blood mononuclear cells from healthy human donor, GSE115189

Repository: Silver Standard Dataset 2. 3k PBMCs from a Healthy Donor, <https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k>

Repository: Silver Standard Dataset 3. 4k PBMCs from a Healthy Donor, <https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.2.0/pbmc4k>

Software availability

All code is available for download at: https://github.com/SaskiaFreytag/cluster_benchmarking_code.

Supplementary material

Supplementary Figure 1. T-SNE plot of the gold standard data after filtering and normalization with the package *scater*. Shapes indicate the cell identity as established by *demuxlet*. The colors indicate in (a) *RaceID2* clustering, (b) *SC3* clustering and in (c) *Seurat* clustering. It can be observed that these three programs present different degrees of complexity.

[Click here to access the data.](#)

Archived code at time of publication: [10.5281/zenodo.1324576](https://doi.org/10.5281/zenodo.1324576)

License: MIT License

Consent

Written informed consent for publication of the participant's transcriptomic information was obtained (Australian Red Cross Blood Service Supply Agreement 18-03VIC-07).

Author contributions

Freytag S: Conceptualization, Data Curation, Funding Acquisition, Formal Analysis, Investigation, Methodology, Software, Visualization, Writing – Original Draft Preparation, Writing – Review & Editing

Tian L: Investigation, Writing – Review & Editing

Lönnstedt I: Conceptualization, Methodology, Writing – Review & Editing

NG M: Conceptualization, Investigation, Funding Acquisition, Methodology, Writing – Review & Editing

Bahlo M: Supervision Conceptualization, Investigation, Funding Acquisition, Methodology, Writing – Review & Editing

Competing interests

No competing interests were disclosed.

Grant information

We would like to thank the Australian Genome Research Facility and the Genomics Innovation Hub for their generous support of this project, including funding. This work was also supported by the Victorian Government's Operational Infrastructure Support Program and Australian Government NHMRC IRIIS. MB is funded by NHMRC Senior Research Fellowship 110297 and NHMRC Program Grant 1054618.

The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Acknowledgements

We gratefully acknowledge the constructive comments and experimental work of Azadeh Seidi, Mark Biondo and Nicolas J. Wilson. Additionally, we want to acknowledge Mark Robinson for his great advice.

Supplementary Figure 2. ARI_truth of each method on all three silver standard datasets given the subset of the cells that reached respective minimum correlation using the cell labeling approach by Zheng *et al.*⁵.

[Click here to access the data.](#)

Supplementary Figure 3. Median of ARI_comp of each method when cell input is changed versus median of ARI_comp of each method when gene input is changed. See Figure 4 for variability associated with each methods' ARI_comp.

[Click here to access the data.](#)

Supplementary Figure 4. UpSeTR Venn diagram to compare of number genes detected in the gold standard dataset by different programs used for preprocessing. The vast majority of genes are detected by all programs.

[Click here to access the data.](#)

Supplementary Figure 5. Total counts for the same barcode as measured when processed with (a) Cell Ranger versus ScPipe Subread, (b) Cell Ranger versus ScPipe STAR, and (c) ScPipe STAR versus ScPipe STAR. Only barcodes are shown that appear in all three versions of the processed dataset.

[Click here to access the data.](#)

Supplementary Figure 6. Tukey boxplots showing the correlations between gene counts of particular gene category for all three comparisons between preprocessing programs used on gold standard dataset. Processed pseudogenes' counts seem to differ depending on program used.

[Click here to access the data.](#)

Supplementary Table 1. Proportion of cell types in different silver standard datasets as estimated by supervised cell labeling.

[Click here to access the data.](#)

References

- Tanay A, Regev A: **Scaling single-cell genomics from phenomenology to mechanism.** *Nature.* 2017; **541**(7637): 331–338.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Zappia L, Phipson B, Oshlack A: **Exploring the single-cell RNA-seq analysis landscape with the scRNA-tools database.** *PLoS Comput Biol.* 2018; **14**(6): e1006245.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Ziegenhain C, Vieth B, Parekh S, *et al.*: **Comparative Analysis of Single-Cell RNA Sequencing Methods.** *Mol Cell.* 2017; **65**(4): 631–643.e4.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Haque A, Engel J, Teichmann SA, *et al.*: **A practical guide to single-cell RNA-seq for biomedical research and clinical applications.** *Genome Med.* 2017; **9**(1): 75.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Zheng GX, Terry JM, Belgrader P, *et al.*: **Massively parallel digital transcriptional profiling of single cells.** *Nat Commun.* 2017; **8**: 14049.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Senabouth A, Lukowski S, Alquicira J, *et al.*: **ascend: R package for analysis of single cell RNA-seq data.** *bioRxiv.* 2017; 207704.
[Publisher Full Text](#)
- Lin P, Troup M, Ho JW: **CIDR: Ultrafast and accurate clustering through imputation for single-cell RNA-seq data.** *Genome Biol.* 2017; **18**(1): 59.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Dey KK, Hsiao CJ, Stephens M: **Visualizing the structure of RNA-seq expression data using grade of membership models.** *PLoS Genet.* 2017; **13**(3): e1006599.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Grün D, Lyubimova A, Kester L, *et al.*: **Single-cell messenger RNA sequencing reveals rare intestinal cell types.** *Nature.* 2015; **525**(7568): 251–5.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Grün D, Muraro MJ, Boisset JC, *et al.*: **De Novo Prediction of Stem Cell Identity using Single-Cell Transcriptome Data.** *Cell Stem Cell.* 2016; **19**(2): 266–277.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Li H, Courtois ET, Sengupta D, *et al.*: **Reference component analysis of single-cell transcriptomes elucidates cellular heterogeneity in human colorectal tumors.** *Nat Genet.* 2017; **49**(5): 708–718.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Kiselev VY, Kirschners K, Schaub MT, *et al.*: **SC3: consensus clustering of single-cell RNA-seq data.** *Nat Methods.* 2017; **14**(5): 483–486.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Lun AT, Bach K, Marioni JC: **Pooling across cells to normalize single-cell RNA sequencing data with many zero counts.** *Genome Biol.* 2016; **17**(1): 75.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Butler A, Hoffman P, Smibert P, *et al.*: **Integrating single-cell transcriptomic data across different conditions, technologies, and species.** *Nat Biotechnol.* 2018; **36**(5): 411–420.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Wang B, Ramazzotti D, De Sano L, *et al.*: **SIMLR: A Tool for Large-Scale Genomic Analyses by Multi-Kernel Learning.** *Proteomics.* 2018; **18**(2): 1700232.
[PubMed Abstract](#) | [Publisher Full Text](#)
- Ji Z, Ji H: **TSCAN: Pseudo-time reconstruction and evaluation in single-cell RNA-seq analysis.** *Nucleic Acids Res.* 2016; **44**(13): e117.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- McCarthy DJ, Campbell KR, Lun AT, *et al.*: **Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R.** *Bioinformatics.* 2017; **33**(8): 1179–1186.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Holik AZ, Law CW, Liu R, *et al.*: **RNA-seq mixology: designing realistic control experiments to compare protocols and analysis methods.** *Nucleic Acids Res.* 2017; **45**(5): e30.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Kang HM, Subramaniam M, Targ S, *et al.*: **Multiplexed droplet single-cell RNA-seq using natural genetic variation.** *Nat Biotechnol.* 2018; **36**(1): 89–94.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Sasaki Y, Darmochwal-Kolarz D, Suzuki D, *et al.*: **Proportion of peripheral blood and decidual CD4⁺CD25^{high} regulatory T cells in pre-eclampsia.** *Clin Exp Immunol.* 2007; **149**(1): 139–145.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
- Jing Y, Gravenstein S, Chaganty NR, *et al.*: **Aging is associated with a rapid**

- decline in frequency, alterations in subset composition, and enhanced Th2 response in CD1d-restricted NKT cells from human peripheral blood.** *Exp Gerontol.* 2007; **42**(8): 719–732.
[PubMed Abstract](#) | [Publisher Full Text](#)
22. Yip SH, Wang P, Kocher JA, *et al.*: **Linnorm: improved statistical analysis for single cell RNA-seq expression data.** *Nucleic Acids Res.* 2017; **45**(22): e179.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
23. Trapnell C, Cacchiarelli D, Grimsby J, *et al.*: **The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells.** *Nat Biotechnol.* 2014; **32**(4): 381–386.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
24. de Graaf CA, Choi J, Baldwin TM, *et al.*: **Haemopedia: An Expression Atlas of Murine Hematopoietic Cells.** *Stem cell reports.* 2016; **7**(3): 571–582.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
25. Hubert L, Arabie P: **Comparing partitions.** *J Classif.* 1985; **2**(1): 193–218.
[Publisher Full Text](#)
26. Studholme C, Hill DLG, Hawkes DJ: **An overlap invariant entropy measure of 3D medical image alignment.** *Pattern Recogn.* 1999; **32**(1): 71–86.
[Publisher Full Text](#)
27. Rosenberg A, Hirschberg J: **V-measure: A conditional entropy-based external cluster evaluation measure.** In *Proceedings of the 2007 joint conference on empirical methods in natural language processing and computational natural language learning (EMNLP-CoNLL)*. 2007.
[Reference Source](#)
28. Tian L, Su S, Amann-Zalcenstein D, *et al.*: **scPipe: a flexible R/Bioconductor preprocessing pipeline for single-cell RNA-sequencing data.** *bioRxiv.* 2017; 175927.
[Publisher Full Text](#)
29. Dobin A, Davis CA, Schlesinger F, *et al.*: **STAR: ultrafast universal RNA-seq aligner.** *Bioinformatics.* 2013; **29**(1): 15–21.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
30. Liao Y, Smyth GK, Shi W: **The Subread aligner: fast, accurate and scalable read mapping by seed-and-vote.** *Nucleic Acids Res.* 2013; **41**(10): e108.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
31. Guimaraes JC, Zavolan M: **Patterns of ribosomal protein expression specify normal and malignant human cells.** *Genome Biol.* 2016; **17**(1): 236.
[PubMed Abstract](#) | [Publisher Full Text](#) | [Free Full Text](#)
32. Duò A, Robinson MD, Soneson C: **A systematic performance evaluation of clustering methods for single-cell RNA-seq data [version 1; referees: 2 approved with reservations].** *F1000Res.* 2018; **7**: 1141.
[Publisher Full Text](#)

Open Peer Review

Current Referee Status:



Version 1

Referee Report 31 August 2018

doi:10.5256/f1000research.17256.r37228



Stephanie Hicks 

Johns Hopkins Bloomberg School of Public Health (JHSPH), Baltimore, MD, USA

Freytag et al. have produced a nice research article on assessing methods for clustering scRNA-seq data from the 10x Genomics platform. I was excited to read the article to learn about what they recommend using. I have made some suggestions below for improvements that are mostly related to providing more intuition and higher-level summaries. This is mostly because as a user of these methods, at the end of the paper, I still felt a little confused about which method the authors would recommend using. I hope the authors can update the article with some of the suggestions:

While the authors have provided detailed comparisons (running time, cluster stability, use of different aligners, different genes, etc), the biggest suggestion would be that the authors provide a higher-level summary of what the authors would suggest a user use to cluster his/her data. At the end of reading this paper, I felt a little overwhelmed at the amount of comparisons across various datasets. It's hard to look at Figs 1-7 and get an overall summary of which method to use. The authors do state in the abstract "We found that some methods, including Seurat and Cell Ranger, outperform other methods, although performance seems to be dependent on the complexity of the studied system", but it would be great if the authors could somehow provide a visual high-level summary of how they came to that conclusion, or elaborate in the discussion on that.

For the "gold standard" data, what was the percent of each human lung cell lines (HCC827, H1975, H2228) that were mixed together? Equal proportions? Was the reason you needed to use demuxlet was because the cell lines were mixed up for sequencing? It would be great if the authors could elaborate on the experimental design.

Is the "gold standard" data available with the SNVs called for each cell. It would be useful to have this count matrix and corresponding phenotypic information about each cell in a SingleCellExperiment object for others to have access to.

It would be great if the authors could include another example dataset with a batch effect in it or something with a slightly less clean design, given most datasets are not quite this "clean". Also, maybe different clustering methods would perform better / worse depending on they data contained rare vs common cell types or included more or less diversity.

There is a TENxPBMCsData package (<https://github.com/kasperdanielhansen/TENxPBMCData>) that has been submitted to Bioconductor (similar to the TENxBrainData). This includes all PBMC 10X datasets currently listed on their site and loads in a SingleCellExperiment object into R. For the Silver Standard

Datasets, you might incorporate this into your workflow.

How did you (or Cell Ranger) deal with empty droplets or swapped barcodes on the 10x platform? This seems relevant for discovering cell types using some form of clustering.

Supplemental Table 1 could use a caption and a label at the top saying "Supplemental Table 1". I had many tabs open with different supplemental figures and tables, and was getting confused about which was which one.

Why did Linnorm and Monocle "continually failed to run"? Did the authors contact the original authors of Linnorm and Monocle to determine if there was a problem with the actual software or if it was a problem with the implementation of the software? It would be great if the authors could elaborate.

I agree with this statement: " We concede that it is possible that more care in the upstream data handling and selection of parameters could result in different results." This is true for almost all benchmarking papers. Given the authors are working within the R/Bioconductor framework, it would be great if the authors could use something like SummarizedBenchmark (<http://bioconductor.org/packages/release/bioc/vignettes/SummarizedBenchmark/inst/doc/SummarizedBenchmark.html>) to keep track of these parameters.

Could the authors elaborate on how they decided which performance metrics to use?

What does this mean: "The impact of different aligners and preprocessing was assessed using all appropriate combinations of programs"? Could the authors be more specific?

I'm a little concerned about how much the solutions differ between methods and parameter choices. I understand the point of this paper is to make comparisons between already published methods, but as the authors are now very familiar with these methods, it would be great if they could provide some more practical guidance. What would the authors suggest using?

Fig 1 -- Could the authors hypothesize on why Seurat, TSCAN, RCA, SC3, RaceID, RaceID2 are estimating so many clusters? Also, why does countClust tend to underestimate the number of clusters? It would be great if the authors could provide some intuition.

Fig 3 -- If I'm understanding, ascend and countClust produce clusters that are very different than the rest?

Thank you to the authors for making their code publicly available!

Is the work clearly and accurately presented and does it cite the current literature?

Yes

Is the study design appropriate and is the work technically sound?

Yes

Are sufficient details of methods and analysis provided to allow replication by others?

Yes

If applicable, is the statistical analysis and its interpretation appropriate?

Yes

Are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions drawn adequately supported by the results?

Yes

Competing Interests: No competing interests were disclosed.

Referee Expertise: Statistics, genomics, analysis of single-cell data

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.

Referee Report 29 August 2018

doi:[10.5256/f1000research.17256.r37231](https://doi.org/10.5256/f1000research.17256.r37231)



Shila Ghazanfar 

School of Mathematics and Statistics, University of Sydney, Sydney, NSW, Australia

Freytag and colleagues provide a comprehensive comparison of clustering methods - specifically designed for scRNA-Seq data - on data collected using the popular droplet-based 10x Genomics platform. A total of four datasets, comprising a Gold standard mixture of cell lines as well as three Silver standard PBMC datasets, were compared in terms of accuracy, stability as well as other metrics like runtime and ease of use. Freytag et al also perform an analysis to try to determine the factors influencing the resulting clusterings for the Silver standard datasets.

It is a very challenging task to perform a comprehensive characterization and comparison of clustering methods on such types of high-dimensional data, due to the sheer number of choices that need to be made, the difficulty in establishing ideal performance, and the relative lack of ground truth. Freytag et al do a great job of addressing these challenges and working towards providing an overall recommendation of clustering methods for non-expert practitioners, while stressing the need for careful interpretation of such results.

With this in mind, I have some comments/suggestions, as well as a number of minor comments/suggestions, as follows:

****Comments to authors****

Linnorm and Monocle failed - expand on why? I understand that this is indeed a limitation especially for a non-expert practitioner, but it would be good to have an understanding towards what the issue might have been.

Could use a flowchart to summarise the study and various comparisons, as well which methods could no longer be compared (e.g. methods that could not work within the scPipe framework).

Different upstream data handling was performed for each clustering method. How much of a difference

was observed just due to this preprocessing, as opposed to the actual clustering step? I understand that each method provides their own preprocessing as *part* of the method, but at least some of these methods would have been developed with plate-based and/or non-UMI-based scRNA-Seq in mind, so may not be intended for the context of 10x Genomics data. Again I understand that you're comparing methods 'out of the box' but it would be insightful to see what differences there are. I suggest a figure like an upsetR plot for the genes/cells filtered and a correlation heatmap of the expression values themselves.

Could you summarise the distance metrics used in the clustering and if there is a general flavour to the clustering algorithm? e.g. hierarchical, k-means, density-based etc. How do these relate in terms of overall accuracy, stability and other metrics?

Stability assessment - mentions that half of the 58,302 genes were randomly selected, but Table 1 says 24,654 total genes detected. There's a big discrepancy between these two so please clarify; if half of the 58,302 genes were selected then a large proportion of genes would have identically zero rows. Also Table 1 shows Dataset 3 had the highest number of 'total genes detected', so how was Dataset 1 the one with "most number of non-zero genes after filtering"?

Run time section - What do you mean by 'overridden'? And for which aspects of the analysis steps was this done?

Figure 4 - These boxplots show ARI among multiple clustering solutions, so a method that gives a consistently bad result is still high (e.g. in this case the RCA method). Suggest an analogous set of boxplots but with ARI_truth, is there a similar variability observed, as seen in these boxplots?

Gene-wise stability analysis - I'm actually unsure how realistic this particular comparison is. It would be insightful to assess clusterings depending on different levels of gene filtering stringency (in the initial Cell Ranger read processing), or stringency on selection of features based on various criteria like highly variable genes.

Figure 7 - Please clarify how 'total number of features' is a cell-specific quantity. Do you mean total number of non-zero features? Was this analysis also performed on the Gold Dataset and what overall similarities could be observed?

Factors influencing clustering solutions - It would be interesting to consider the factors associated with 'correct' cluster assignment for cells. Optionally suggest to perform this for either the Gold Dataset or the Silver datasets and perform a logistic regression with the response being success/failure of a cell to belong to the cluster most associated with the 'true' cell type group. There is an added subtlety as far as matching clusters with cell type groups goes, but I think there are a few reasonable ways to perform this (e.g. assign candidate clusters to the 'true' groups by taking the higher proportion of cell overlap, and allow multiple candidate clusters to match to a single true group). Performing this kind of analysis could shed light on properties of cells that don't tend to cluster correctly, and if there is consistency in this across multiple disparate datasets.

****Minor comments****

Table 1 - countClust 'version' formatted with verbatim.

Table 1 - I would suggest the 'properties' column could be better presented in a checklist format, with ticks/crosses for fulfilling various criteria listed.

Section beginning "silver standard" - 10x is capitalised.

Supplementary Figure 1 - legend fallen off panel a), needs a higher resolution or larger points

NMI definition - trailing parenthesis in denominator

typo - assess the effect**

Figure 2a - I found this quite busy, hard to interpret. Suggest to add shading that covers the points for same method or to facet by dataset. I don't believe the ARI values are particularly comparable between datasets so I would prefer faceting by dataset.

Figure 3 - rows/columns are ordered differently between panels, what's driving this difference?

Supplementary Figure 3 was not mentioned in the main text

Supplementary Figure 4 is a two page pdf, with the first page blank

Figure 6 - Figure caption says Dataset 1 but reports 29,151 genes. Do you mean the Gold Dataset and 29,451 genes? If not, please clarify which data and how many genes.

Discussion - One instance of "Seurat" is missing verbatim format

Is the work clearly and accurately presented and does it cite the current literature?

Yes

Is the study design appropriate and is the work technically sound?

Yes

Are sufficient details of methods and analysis provided to allow replication by others?

Yes

If applicable, is the statistical analysis and its interpretation appropriate?

Yes

Are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions drawn adequately supported by the results?

Yes

Competing Interests: No competing interests were disclosed.

Referee Expertise: Statistics, statistical bioinformatics

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard, however I have significant reservations, as outlined above.

Referee Report 28 August 2018

doi:10.5256/f1000research.17256.r37232



Joshua W. K. Ho 

Victor Chang Cardiac Research Institute (VCCRI), Darlinghurst, NSW, Australia

This paper presents a well-designed and comprehensive evaluation of widely used clustering algorithms for medium-sized 10x Genomics scRNA-seq data. Clustering is a highly active area of research in scRNA-seq data analysis. With so many published clustering tools available, it is often difficult to choose the most appropriate tool. This paper attempts to address this problem by systematically comparing the performance of 12 commonly used clustering tools. The evaluation results should serve as an important guide to bioinformatics practitioners. This paper is a very useful contribution to the field.

Is the work clearly and accurately presented and does it cite the current literature?

Yes

Is the study design appropriate and is the work technically sound?

Yes

Are sufficient details of methods and analysis provided to allow replication by others?

Yes

If applicable, is the statistical analysis and its interpretation appropriate?

Yes

Are all the source data underlying the results available to ensure full reproducibility?

Yes

Are the conclusions drawn adequately supported by the results?

Yes

Competing Interests: No competing interests were disclosed.

Referee Expertise: Bioinformatics, single-cell transcriptomics

I have read this submission. I believe that I have an appropriate level of expertise to confirm that it is of an acceptable scientific standard.

The benefits of publishing with F1000Research:

- Your article is published within days, with no editorial bias
- You can publish traditional articles, null/negative results, case reports, data notes and more
- The peer review process is transparent and collaborative
- Your article is indexed in PubMed after passing peer review
- Dedicated customer support at every stage

For pre-submission enquiries, contact research@f1000.com

F1000Research