

# pipeComp, a general framework for the evaluation of computational pipelines, reveals performant single-cell RNA-seq preprocessing tools

Pierre-Luc Germain<sup>1,2,3,\*</sup>, Anthony Sonrel<sup>1,2</sup> & Mark D. Robinson<sup>1,2,\*</sup>

<sup>1</sup>*Department of Molecular Life Sciences, University of Zürich, Switzerland*

<sup>2</sup>*SIB Swiss Institute of Bioinformatics, Zürich, Switzerland*

<sup>3</sup>*D-HEST Institute for Neuroscience, Swiss Federal Institute of Technology (ETH), Zürich, Switzerland*

\*Correspondence to Pierre-Luc Germain ([pierre-luc.germain@hest.ethz.ch](mailto:pierre-luc.germain@hest.ethz.ch)) and Mark D. Robinson ([mark.robinson@imls.uzh.ch](mailto:mark.robinson@imls.uzh.ch))

## Abstract

1 The massive growth of single-cell RNA-sequencing (scRNAseq) and methods for its analysis still  
2 lacks sufficient and up-to-date benchmarks that would guide analytical choices. Moreover, current  
3 studies are often focused on isolated steps of the process. Here, we present a flexible R framework  
4 for pipeline comparison with multi-level evaluation metrics and apply it to the benchmark of  
5 scRNAseq analysis pipelines using datasets with known cell identities. We evaluate common steps  
6 of such analyses, including filtering, doublet detection (suggesting a new R package, *scDblFinder*),  
7 normalization, feature selection, denoising, dimensionality reduction and clustering. On the basis  
8 of these analyses, we make a number of concrete recommendations about analysis choices. The  
9 evaluation framework, *pipeComp*, has been implemented so as to easily integrate any other step  
10 or tool, allowing extensible benchmarks and easy application to other fields ([https://github.](https://github.com/plger/pipeComp)  
11 [com/plger/pipeComp](https://github.com/plger/pipeComp)).

## 12 Background

13 Single-cell RNA-sequencing (scRNAseq) and the set of attached analysis methods are evolving  
 14 fast, with more than 560 software tools available to the community [1], roughly half of which are  
 15 dedicated to tasks related to data processing such as clustering, ordering, dimension reduction  
 16 or normalization. This increase in the number of available tools follows the development of new  
 17 sequencing technologies and the growing number of reported cells, genes and cell populations [2].  
 18 As data processing is a critical step in any scRNAseq analysis, affecting downstream analysis and  
 19 interpretation, it is critical to evaluate the available tools.

20 A number of good comparison and benchmark studies have already been performed on various  
 21 steps related to scRNAseq analysis and can guide the choice of methodology [3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18].  
 22 However these recommendations need constant updating and often leave open many details of an  
 23 analysis. Another missing aspect of current benchmarking studies is their limitation to capture all  
 24 aspects of scRNAseq processing workflow. Although previous benchmarks already brought valuable  
 25 recommendations for data processing, some only focused on one aspect of data processing (e.g.,  
 26 [11]), did not evaluate how the tool selection affects downstream analysis (e.g., [14]) or did not  
 27 tackle all aspects of data processing, such as doublet identification or cell filtering (e.g., [15]). A  
 28 thorough evaluation of the tools covering all major processing steps is however urgently needed as  
 29 previous benchmarking studies highlighted that a combination of tools can have a drastic impact  
 30 on downstream analysis, such as differential expression analysis and cell-type deconvolution [15,3].  
 31 It is then critical to evaluate not only the single effect of a preprocessing method but also its  
 32 positive or negative interaction with all parts of a workflow.

33 Here, we develop a flexible R framework for pipeline comparison and evaluate the various  
 34 steps of analysis leading from an initial count matrix to a cluster assignment, which are critical in  
 35 a wide range of applications. We collected real datasets of known cell composition (Table 1) and  
 36 used a variety of evaluation metrics to investigate in a multilevel fashion the impact of various  
 37 parameters and variations around a core scRNAseq pipeline. Although we use some datasets based  
 38 on other protocols, our focus is especially on droplet-based datasets that do not include exogenous  
 39 control RNA (i.e. spike-ins); see Table 1 and Figure 1 for more details. In addition to previously-  
 40 used benchmark datasets with true cell labels [6,12], we simulated two datasets with a hierarchical  
 41 subpopulation structure based on real 10x human and mouse data using *muscat* [19]. Since graph-  
 42 based clustering [20] was previously shown to consistently perform well across several datasets

[6,12], we used the *Seurat* pipeline as the starting point to perform an integrated investigation of: 1) doublet identification, 2) cell filtering, 3) normalization, 4) feature selection, 5) dimension reduction, 6) clustering. We compared competing approaches and also explored more fine-grained parameters and variations on common methods. Importantly, the success of methods at a certain analytical step might be dependent on choices at other steps. Therefore, instead of evaluating each step in isolation, we developed a general framework for evaluating nested variations on a pipeline and suggest a multilevel panel of metrics. Finally, we evaluate several recent methods and provide concrete recommendations.

## Results

### A flexible framework for pipeline evaluation

The *pipeComp* package defines a pipeline as, minimally, a list of functions executed consecutively on the output of the previous one (Figure 2A). In addition, optional benchmark functions can be set for each step to provide standardized, multi-layered evaluation metrics. Given such a `PipelineDefinition` object, a set of alternative parameters (which might include different subroutines) and benchmark datasets, the `runPipeline` function then proceeds through all combinations of arguments, avoiding recomputing the same step twice and compiling evaluations (including running time) on the fly. Variations in a given parameter can be evaluated using all metrics from this point downward in the pipeline. This is especially important because end-point metrics, such as the adjusted Rand index (ARI) [21] for clustering, are not perfect. For example, although the meaning of an ARI score is independent of the number of true subpopulations [22], the number of clusters called is by far the most important determinant of the score: the farther it is from the actual number of subpopulations, the worse the ARI (Supplementary Figure 17). In this context, one strategy has been to cluster across various resolutions and only consider the results that have the right number of clusters [6]. While this has the virtue of making the results comparable, in practice the number of subpopulations is typically unknown and tools that operate well in this optimal context might not necessarily be best overall. Clustering results are also very sensitive and might not always capture improvements in earlier steps of the pipeline. We therefore favour monitoring several complementary metrics across multiple steps of the process.

A clustering output is relatively fragile to variations in the pipeline. We therefore wanted to capture whether the effect of a given parameter alteration is robust to changes in the rest of the

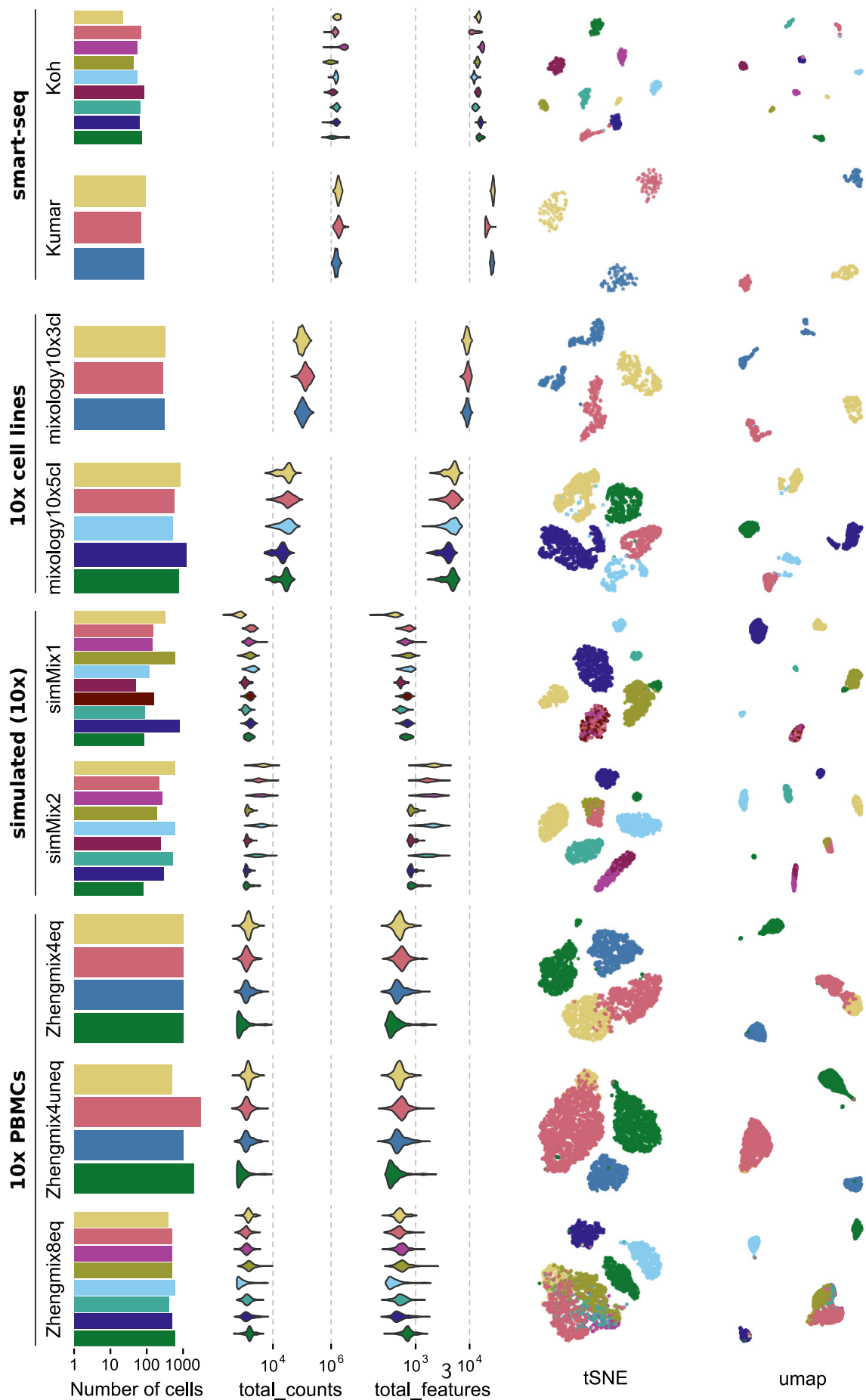
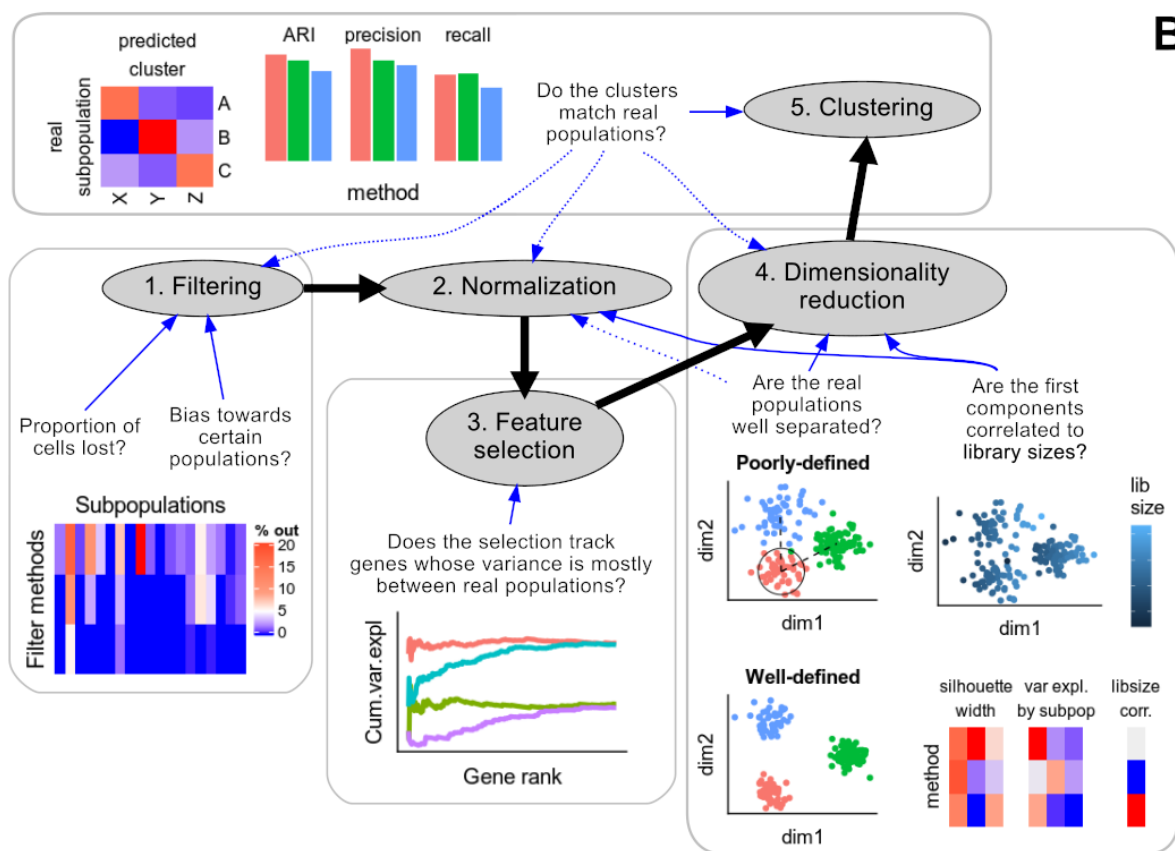
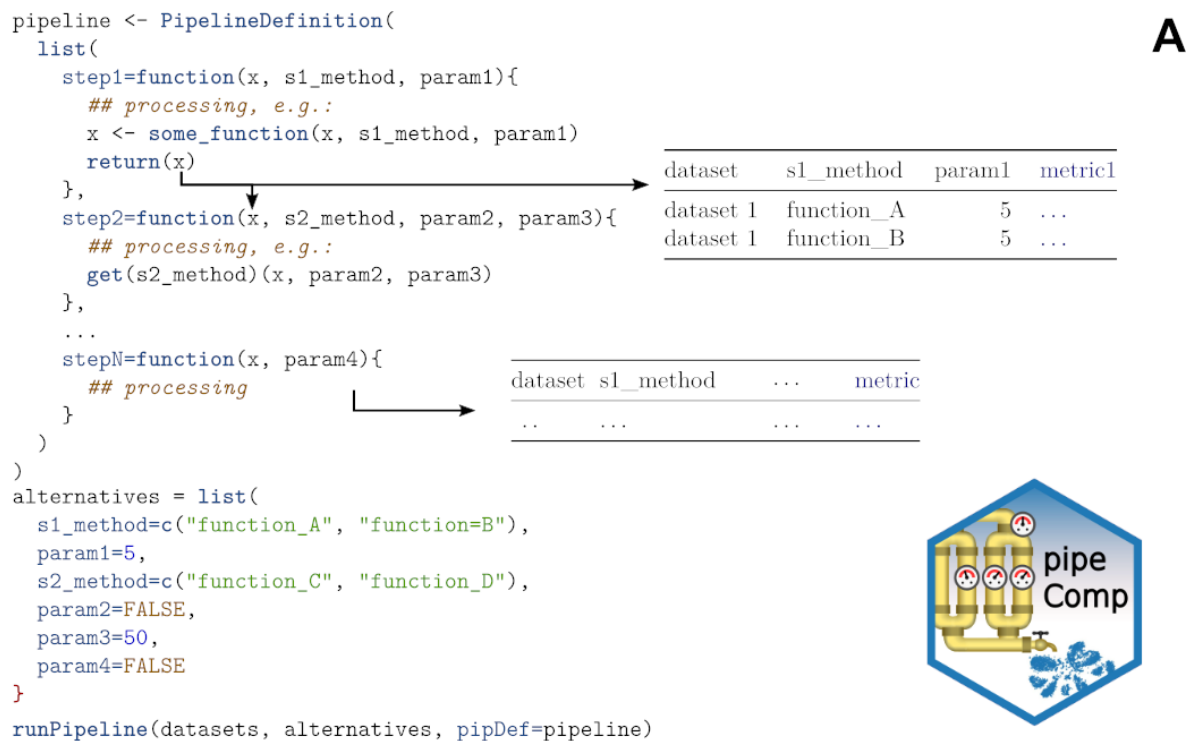


Figure 1: Overview of the benchmark datasets used.



**Figure 2: Overview of the pipeComp framework and its application to a scRNAseq clustering pipeline.** **A:** The package is built around a PipelineDefinition S4 class which defines a set of functions to be executed consecutively, as well as optional evaluation functions for each step. Each step can accept of number of parameters whose alternative values are provided as a list. All (or subsets of) combinations of parameters can then be simultaneously ran and evaluated using the runPipeline function. **B:** Scheme representing the application of *pipeComp* to evaluate a range of methods that are commonly used in scRNAseq studies and some of the metrics monitored at various steps.

73 pipeline, or rather specific to a set of other pipeline parameters. We proceeded in a step-wise  
74 fashion, first testing a large variety of parameters at the early steps of the pipeline along with only  
75 a set of mainstream options downstream, then selecting main alternatives and proceeding to a  
76 more detailed benchmark of the next step (Figure 2B).

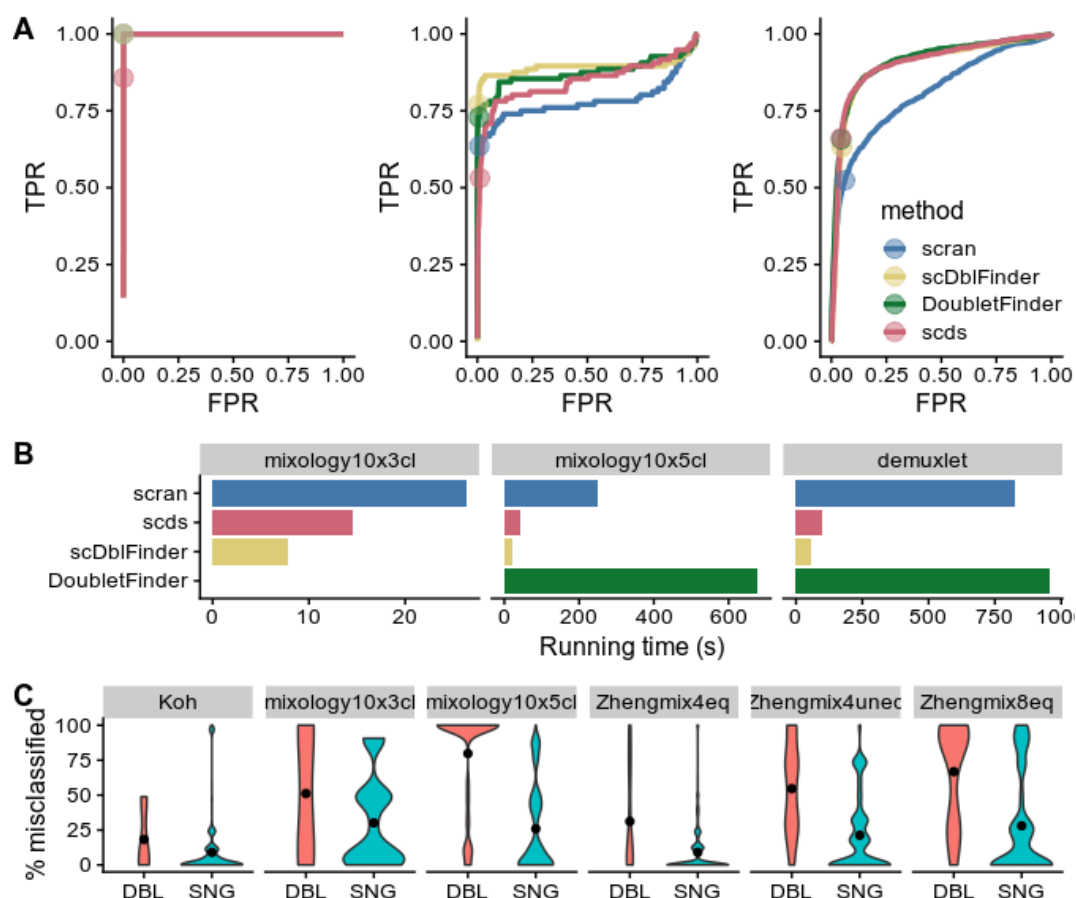
## 77 **Filtering**

### 78 **Doublet detection**

79 Doublets, defined as two cells sequenced under the same cellular barcode (e.g., being captured  
80 in the same droplet), are fairly frequent in scRNAseq datasets, with estimates ranging from 1 to  
81 10% depending on the platform and cell concentration used [23,24]. While doublets of the same  
82 cell type are relatively innocuous for most downstream applications due to their conservation of  
83 the relative expression between genes, doublets formed from different cell types or states are likely  
84 to be misclassified and could potentially distort downstream analysis. In some cases, doublets  
85 can be identified through their unusually high number of reads and detected features, but this is  
86 not always the case (Supplementary Figure 2). A number of methods were developed to identify  
87 doublets, in most cases by comparing each cell to artificially-created doublets [25,26,27]. We first  
88 evaluated the capacity of these methods to detect doublets using the two 10x datasets with cells  
89 of different genetic identity [12], using SNP genotypes as the ground truth. For the sole purpose  
90 of this section, we included an additional dataset with SNP information but lacking true cell labels  
91 [24]. Of note, SNP-based analyses also call doublets created by cells of the same cell type (but  
92 from different individuals) and which are generally described as homotypic (as opposed to neotypic  
93 or heterotypic doublets, i.e. doublets from different cell types). These homotypic doublets might  
94 not be identifiable from the mere gene counts, and their identification is not generally the primary  
95 aim of doublet callers since they are often considered innocuous and, when across individuals, can  
96 be identified through other means (e.g., SNPs). We therefore do not expect a perfect accuracy  
97 in datasets involving cells of the same type across individuals (as in the demuxlet dataset).

98 We tested *DoubletFinder* [25] and *scran*'s *doubletCells* [26], both of which use similarity  
99 to artificial doublets, and *scds* [27], which relies on a combination of co-expression and binary  
100 classification. *DoubletFinder* integrates a thresholding based on the proportion of expected dou-  
101 blets, while *scran* and *scds* return scores that must be manually thresholded. In these cases, we  
102 ensured that the right number of cells would be called doublets. In addition to these methods, we

reasoned that an approach such as *DoubletFinder* could be simplified by being applied directly on counts and by using a pre-clustering to create neotypic/heterotypic doublets more efficiently. We therefore developed a simple and fast Bioconductor package implementing this method for doublet detection, *scDbIFinder*, with the added advantage of accounting for uncertainty in the expected doublet rate and using meta-cells from the clusters to even include triplets (see Methods).



**Figure 3: Identification of doublet cells.** **A:** Receiver operating characteristic (ROC) curves of the tested doublet detection methods for three datasets with SNP-identified doublets. Dots indicate threshold determined by the true number of doublets. **B:** Running time of the different methods (DoubletFinder failed on one of the datasets). **C:** Rate of misclassification of the cells identified by *scDbIFinder* as doublets (DBL) or singlets (SNG), across a large range of clustering analysis. Even in datasets which should not have neotypic doublets, the cells identified as such tend to be misclassified.

While most methods accurately identified the doublets in the 3 cell lines dataset (mixology10x3cl), the other two datasets proved more difficult (Figure 3A). *scDbIFinder* achieved comparable or better accuracy than top alternatives while being the fastest method (Figure 3B). Across datasets, cells called as doublets tended to be classified in the wrong cluster more often than other

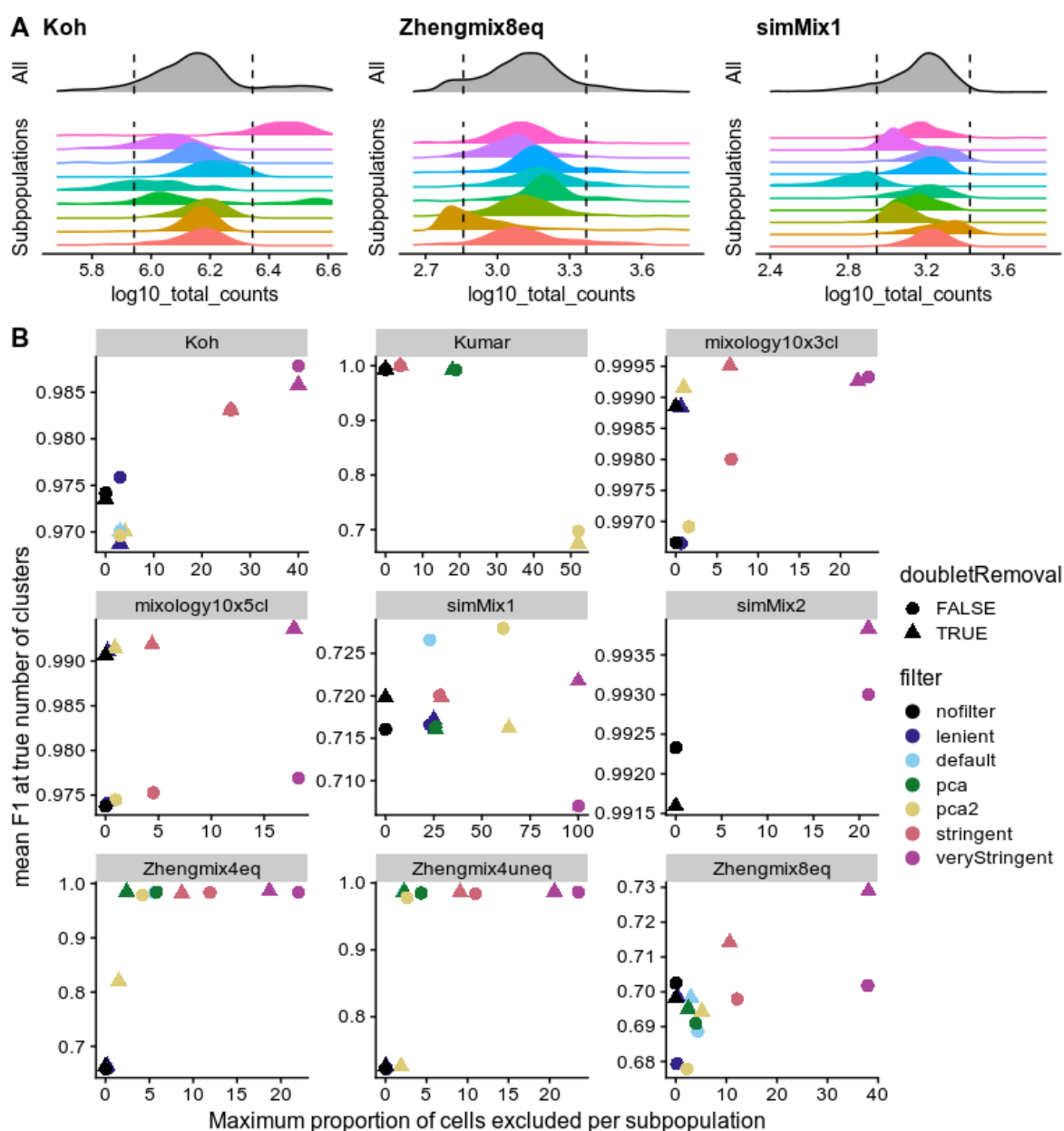
cells (Figure 3C). We also found that *scDbtFinder* improved the accuracy of the clustering across all benchmark datasets even when, by design, the data contained no heterotypic doublet (Figure 4).

## Excluding more cells is not necessarily better

Beyond doublets, a dataset might include low-quality cells whose elimination would reduce noise. This has for instance been demonstrated for droplets containing a high content of mitochondrial reads, often as a result of cell degradation and resulting loss of cytoplasmic mRNAs [28]. A common practice is to exclude cells that differ considerably from most other cells on the basis of such properties. This can for instance be performed through the *isOutlier* function from *scater* that measures, for a given control property, the number of median absolute deviations (MADs) of each cell from the median of all cells. Supplementary Figure 1 shows the distributions of some of the typical cell properties commonly used. Of note, these properties tend to be correlated, with some exceptions. For example, while a high proportion of mitochondrial reads is often correlated with a high proportion of the counts in the top features, there can also be other reasons for an over-representation of highly-expressed features (Supplementary Figure 3), such as an over-amplification in non-UMI datasets. In our experience, 10X datasets also exhibit a very strong correlation between the total counts and the total features even across very different cell types (Supplementary Figure 4). We therefore also measure the ratio between the two and treat cells strongly departing from this trend with suspicion.

Reasoning that the cells we wish to exclude are the cells that would be misclassified, we measured the rate of misclassification of each cell in each dataset across a variety of clustering pipelines, correcting for the median misclassification rate of the subpopulation and then evaluated what properties could be predictive of misclassification (Supplementary Figures 5-7). We could not identify any property or simple combination thereof that would be consistently predictive of misclassification; the only feature that consistently stood out across multiple datasets (the Zheng datasets) was that cells with very high read counts have a higher chance of being misclassified.

We next investigated the impact of filtering according to various criteria (see Methods). An important risk of excluding cells on the basis of their distance from the whole distribution of cells on some properties (e.g., library size) is that these properties tend to have different distributions across subpopulations. As a result, thresholds in terms of number of MADs from the whole distribution can lead to strong biases against certain subpopulations (Figure 4A). We therefore



**Figure 4: A:** Filtering on the basis of distance to the whole distribution can lead to strong bias against certain subpopulations. The dashed line indicates a threshold of 2.5 median absolute deviations (MADs) from the median of the overall population. **B:** Relationship between the maximum subpopulation exclusion rate and the average clustering accuracy per subpopulation across various filtering strategies. Of note, doublet removal appears to be desirable even when, due to the design, there are no heterotypic doublets in the data. The PCA methods refer to multivariate outlier detected as implemented in *scater* (see Methods for details).

143 examined the tradeoff between the increased accuracy of filtering and the maximum proportion of  
 144 cells excluded per subpopulation (Figure 4B and Supplementary Figure 8). Since filtering changes  
 145 the relative abundance of the different subpopulations, global clustering accuracy metrics such  
 146 as ARI are not appropriate here. We therefore calculated the per-subpopulation precision and  
 147 recall using the Hungarian algorithm [29] and monitored their mean F1 score. A first observation  
 148 was that, although the stringent filtering tended to be associated with an increase in accuracy, it  
 149 could also become deleterious and most of the benefits could be achieved without very stringent  
 150 filtering and minimizing subpopulation bias (Figure 4B). Applying the same filtering criteria on  
 151 individual clusters of cells (identified through *scran*'s `quickCluster` method) resulted in nearly  
 152 no cell being filtered out. This suggests that filtering on the global population tends to discard  
 153 cells of subpopulations with more extreme properties (e.g. high library size), rather than low-  
 154 quality cells. Finally, by changing the distributions, the doublet removal step in conjunction with  
 155 filtering sometimes resulted in a net decrease in the proportion of excluded cells while retaining  
 156 or improving accuracy. We therefore recommend the use of doublet removal followed by relatively  
 157 mild filtering, such as is implemented in our 'default' filtering (see Methods).

## 158 **Filtering features by type**

159 Mitochondrial reads have been associated with cell degradation and there is evidence that riboso-  
 160 mal genes can influence the clustering output, hiding other biological structure in the analysis [7].  
 161 We therefore investigated whether excluding one category of features or the other, or using only  
 162 protein-coding genes, had an impact on the ability to distinguish subpopulations (Supplementary  
 163 Figure 9). Removal of ribosomal genes robustly reduced the quality of the clustering, suggesting  
 164 that they represent real biological differences between subpopulations. Removing mitochondrial  
 165 genes and restricting to protein-coding genes had a very mild impact.

## 166 **Normalization and scaling**

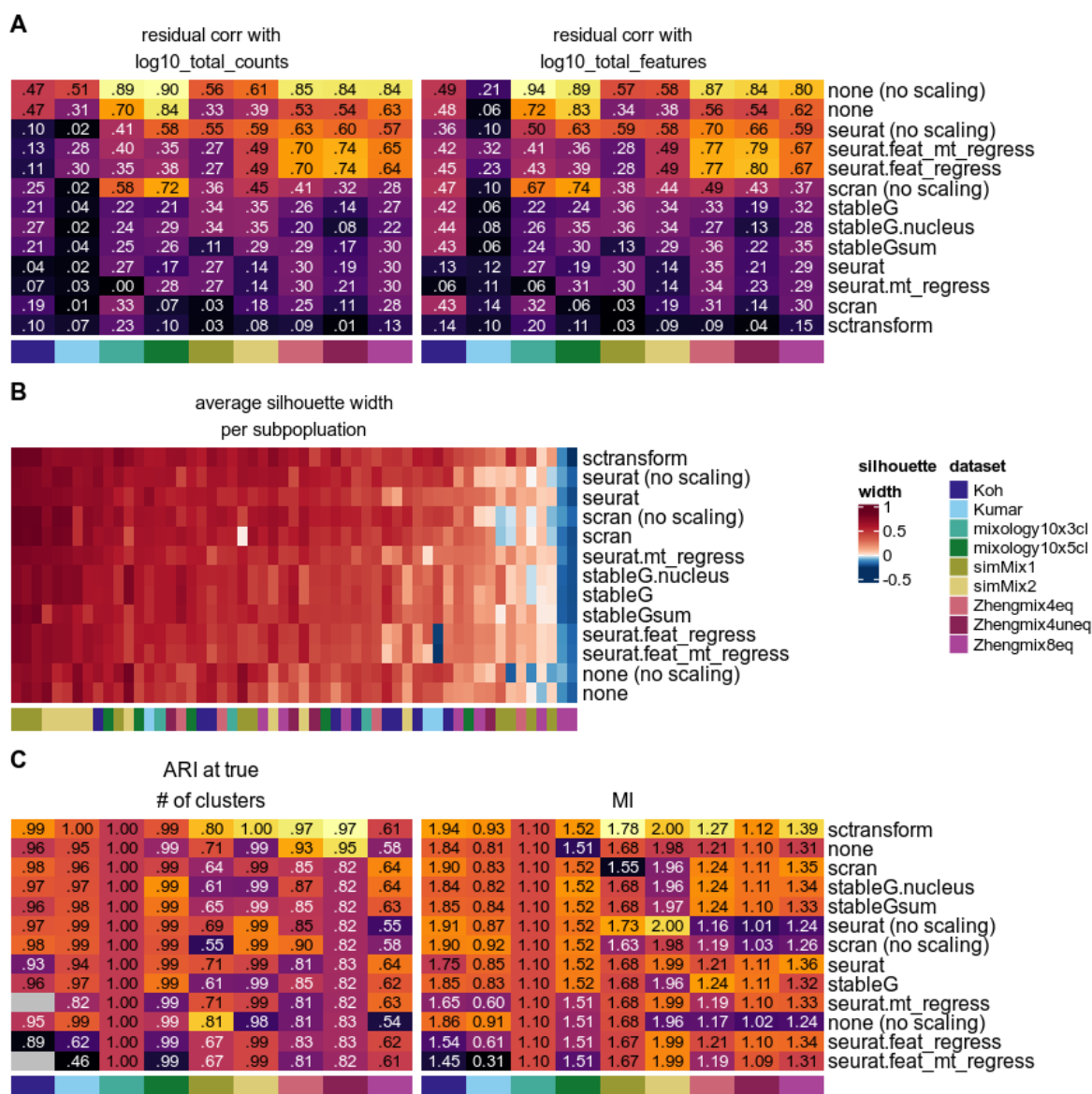
167 We next investigated the impact of different normalization strategies. Beside the standard log-  
 168 normalization included in *Seurat*, we tested *scran*'s pooling-based normalization [30], *sctransform*'s  
 169 variance-stabilizing transformation [31], and normalization based on stable genes [32,33]. In addition  
 170 to log-normalization, the standard *Seurat* clustering pipeline performs per-feature unit-variance  
 171 scaling so that the PCA is not too strongly dominated by highly-expressed features. We therefore  
 172 included versions of the different normalization procedures, with or without a subsequent scal-

ing (*sctransform*'s variance-stabilizing transformation involves an approach analogous to scaling).  
*Seurat*'s scaling function also includes the option to regress out the effect of certain covariates.  
We tested its performance by using the proportion of mitochondrial counts and the number of  
detected features as covariates. Finally, it has been proposed that the use of stable genes, in par-  
ticular cytosolic ribosomal genes, can be used to normalize scRNAseq<sup>[33]</sup>. We therefore evaluated  
a simple linear normalization based on the sum of these genes, as well as nuclear genes.

An important motivation for the development of *sctransform* was the observation that, even  
after normalization, the first principal components of various datasets tended to correlate with  
library size, suggesting an inadequate normalization<sup>[31]</sup>. However, as library size tends to vary  
across subpopulations, part of this effect can simply reflect biological differences. We therefore  
assessed to what extent the first principal component still retained a correlation with the library  
size and the number of detected features, removing the confounding covariation with the subpop-  
ulations (Figure 5A). The simple step of scaling tended to remove much of the correlation with  
these features and, in the absence of scaling, standard *Seurat* normalization resulted in fairly high  
correlation with technical covariates. *sctransform* led to the lowest correlation but most methods,  
including normalization based on stable genes, were able to remove most of the correlation when  
combined with scaling. The only exception is *Seurat* normalization with scaling regressing out  
covariates, which surprisingly increased the correlation in 8 of the 9 datasets.

We further evaluated normalization methods by investigating their impact on the separability  
of the subpopulations (Figure 5B-C). Since clustering accuracy metrics such as the ARI are very  
strongly influenced by the number of clusters, we complemented it with silhouette width<sup>[34]</sup> and  
mutual information. We found most methods (including no normalization at all) to perform fairly  
well in most of the subpopulations. Scaling tended to reduce the average silhouette width of some  
subpopulations and to increase that of some less distinguishable ones and was generally, but not  
always, beneficial on the accuracy of the final clustering. Regressing out covariates systematically  
gave poorer performance on all metrics. *sctransform* systematically outperformed other methods  
and, even though it was developed to be applied to data with unique molecular identifiers (UMI),  
it also performed fairly well with the Smart-seq protocol (Koh and Kumar datasets).

Finally, we monitored whether, under the same downstream clustering analysis, different  
normalization methods tended to lead to an over- or under-estimation of the number of clusters.  
Although some methods had a tendency to lead to a higher (e.g., *sctransform*) or lower (e.g., stable  
genes) number of clusters, the effect was very mild and not entirely systematic (Supplementary



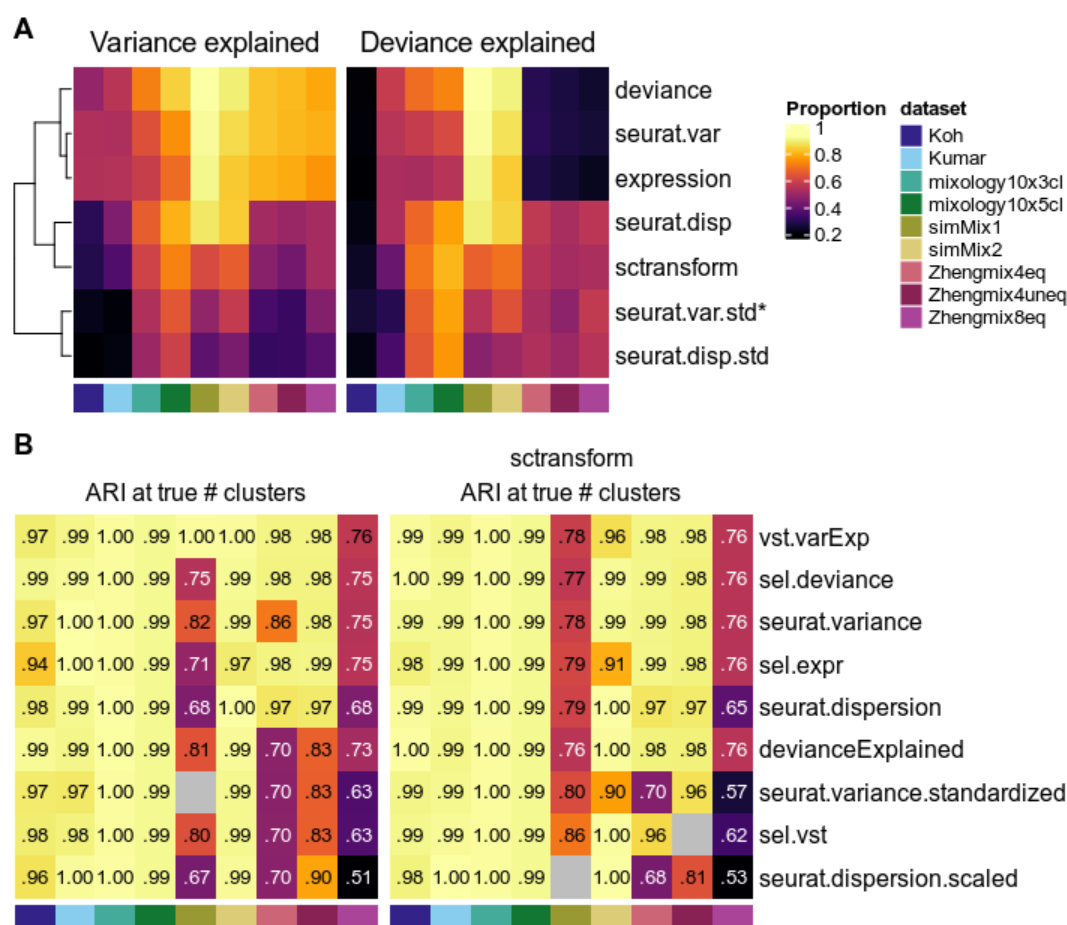
**Figure 5: Evaluation of normalization procedures.** ‘Seurat.feats\_mt\_regress’ regresses out the number of features and proportion of mitochondrial reads during scaling; ‘Seurat.feats\_regress’ regresses out the number of features only. Residual correlation of the first principal component with library size (left) and detection rate (right), after accounting for biological differences between subpopulations. **B:** Average silhouette width per true subpopulation, where higher silhouette width means a higher separability. **C:** Clustering accuracy (*Seurat* clustering), measured by the ARI at the true number of cluster (left) and by the average mutual information (MI) of the cluster assignment with true subpopulations.

205 Figure 10).

## 206 Feature selection

207 A standard clustering pipeline typically involves a step of highly-variable genes selection, which  
 208 is complicated by the digital nature and the mean-variance relationship of (sc)RNAseq. *Seurat*'s  
 209 earlier approaches involved the use of dispersion estimates standardized for the mean expression  
 210 levels, while more recent versions ( $\geq 3.0$ ) rely on a different measure of variance, again standardized.  
 211 While adjusting for the mean-variance relationship removes much of the bias towards highly-  
 212 expressed genes, it is plausible that this relationship may in fact sometimes reflects biological  
 213 relevance and would be helpful in classifying cell types. Another common practice in feature  
 214 selection is to use those with the highest mean expression. Recently, [35] instead suggested to use  
 215 deviance, while *sctransform* provides its own ordering of genes based on transformed variance.

216 Reasoning that a selection method should ideally select genes whose variability is higher be-  
 217 tween subpopulations than within, we first assessed to what extent each method selected genes  
 218 with a high proportion of variance or deviance explained by (real) subpopulation. As the pro-  
 219 portion of variability in a gene attributable to subpopulations can be measured in various ways,  
 220 we first compared three approaches: ANOVA on log-normalized count, ANOVA on *sctransform*  
 221 normalization, and deviance explained. The ANOVAs performed on a standard *Seurat* normaliza-  
 222 tion and on *sctransform* data were highly correlated (Supplementary Figure 11A). These estimates  
 223 were also in good agreement with the deviance explained, although lowly-expressed genes could  
 224 have a high deviance explained without having much of their variance explained by subpopulation  
 225 (Supplementary Figure 11B-D). We therefore compared the proportion of the cumulative vari-  
 226 ance/deviance explained by the top X genes that could be retrieved by each gene ranking method  
 227 (Supplementary Figures 12-13). We first focused on the first 1000 genes to highlight the differ-  
 228 ences between methods, although a higher number of selected genes decreased the differences  
 229 between methods (Supplementary Figures 12-14). The standardized measures of variability were  
 230 systematically worse than their non-standardized counterparts in selecting genes with a high pro-  
 231 portion of variance explained by subpopulation. Regarding the percentage of deviance explained  
 232 however, the standardized measures were often superior (Figure 6A and Supplementary Figures  
 233 12-13). Deviance proved the method of choice to prioritize genes with a high variance explained  
 234 by subpopulations (with mere expression level proving surprisingly good) but did not perform so  
 235 well to select genes with a high deviance explained.

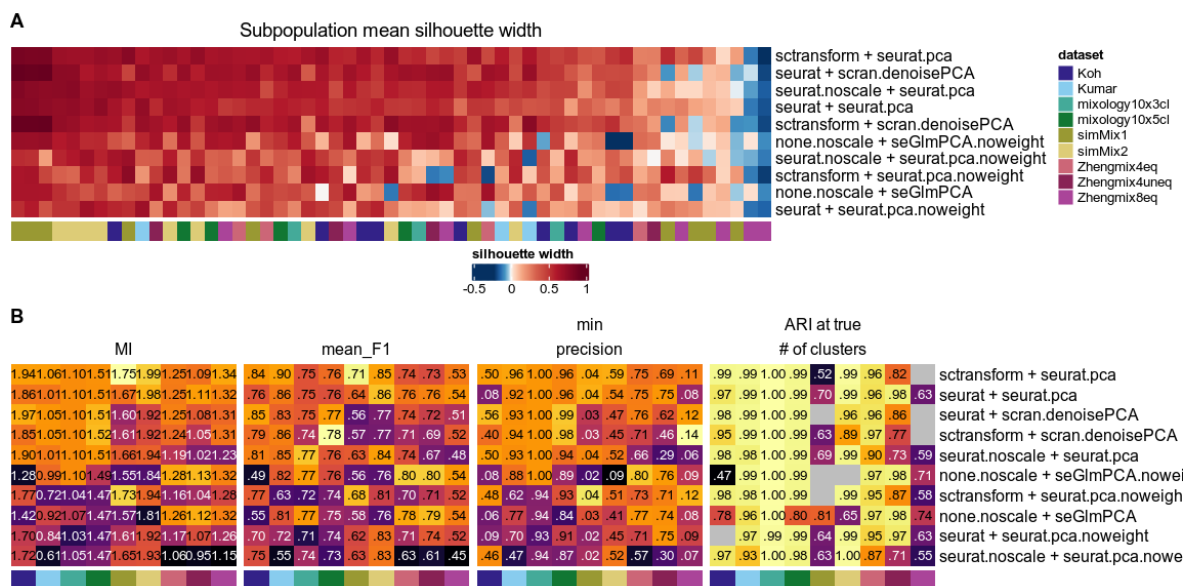


**Figure 6: Evaluation of feature selection methods.** **A:** Ability of different feature ranking methods to capture genes with a high proportion of variance (left) or deviance (right) explained by real subpopulations. The asterix denotes the default Seurat method. **B:** Accuracy of clusterings (at the true number of clusters) when selecting 1000 genes using the given methods. Based on standard *Seurat* normalization (left) or *sctransform* (right). The *vst.varExp* and *devianceExplained* methods correspond to the estimates used in **A** to evaluate the selection methods and were included here only for validation purpose.

We next evaluated how the use of different feature selection methods affected the clustering accuracy (Figure 6B). To validate the previous assay on the proportion of variance/deviance explained by real populations, we included genes that maximized these two latter measures. Interestingly, while these selections were on average the top-ranking methods, they were not systematically best for all datasets. The previous observations were reflected in the ARI of the resulting clustering (Figure 6B): non-standardized measures of variability, including mere expression level, tended to outperform more complex metrics. In general, we found deviance and unstandardized estimates of variance to provide the best results across datasets and normalization methods. Increasing the number of features selected also systematically led to an increase in the accuracy of the clustering, typically plateauing after 4000 features (Supplementary Figure 14).

## Dimensionality reduction

Since the various PCA approaches and implementations were recently benchmarked in a similar context [14], we focused on widely used approaches which had not yet been compared: *Seurat*'s PCA, *scran*'s denoisePCA, and GLM-PCA [35]. When relevant, we combined them with *sctransform* normalization. Given that *Seurat*'s default PCA weights the cell embeddings by the variance of each component, we also evaluated the impact of this weighting with each method.



**Figure 7: Evaluation of common dimensionality reduction methods. A:** Average silhouette width per subpopulation resulting from combinations of normalization and dimension reductions. **B:** Clustering accuracy, mean proportion of the variance in the first 5 components explained by real subpopulations (center), and median ARI of the resulting clustering (right).

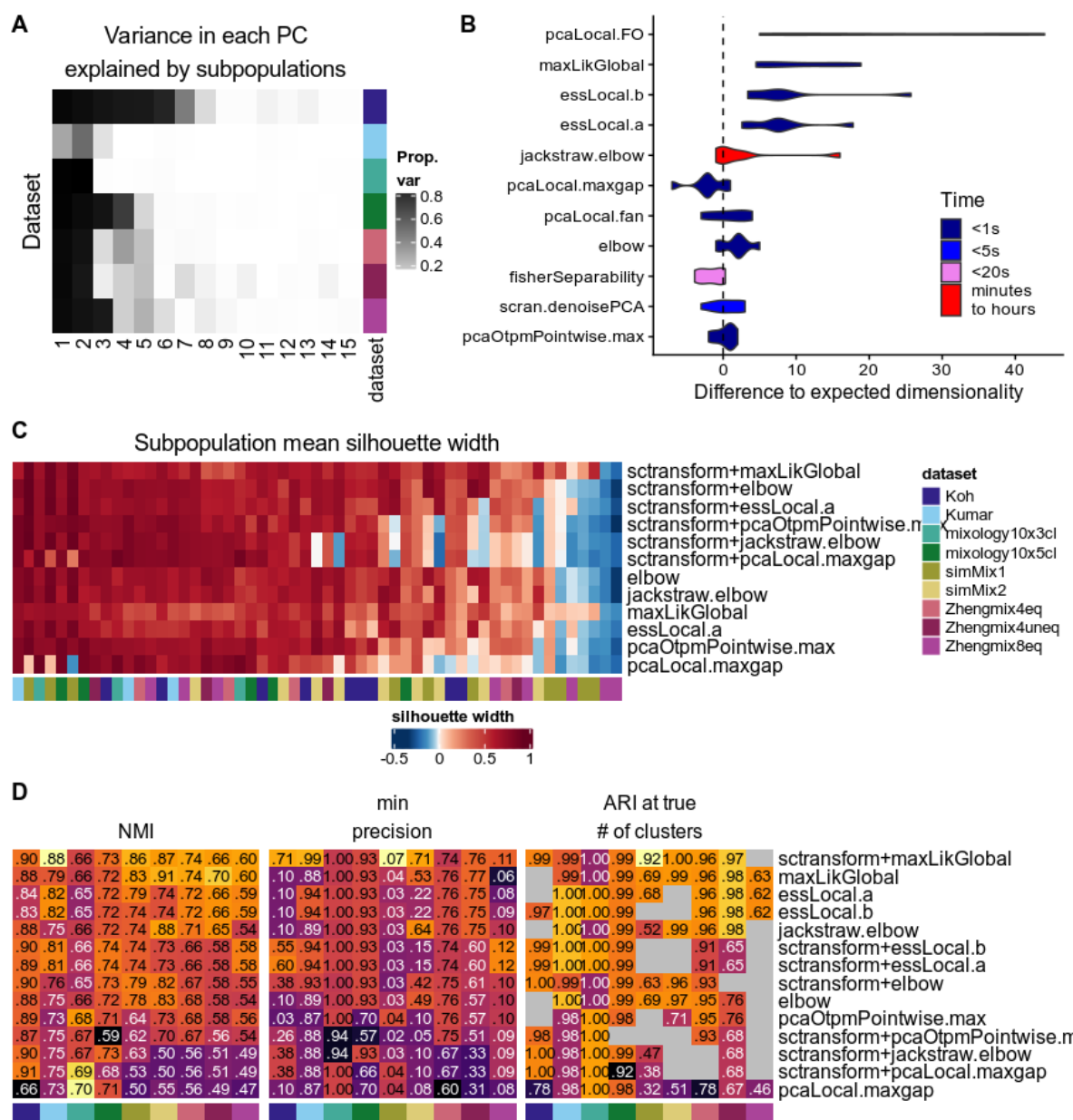
252 The impact of the choice of dimensionality reduction method was far greater than that of  
253 normalization or feature selection (e.g., Supplementary Figure 15). GLM-PCA tended to increase  
254 the average silhouette width of already well-defined subpopulations, but *Seurat*'s PCA procedure  
255 however proved superior on all metrics (Figure 7). Overall, weighting the principal components by  
256 their variance (as *Seurat* does) had a positive impact on silhouette widths and ARI scores.

## 257 Estimating the number of dimensions

258 A common step following dimension reduction is the selection of an appropriate number of di-  
259 mensions to use for downstream analysis. Since Euclidean distance decreases as the number of  
260 non-discriminating dimensions increases, there is usually a trade-off between selecting enough  
261 dimensions to keep most information and excluding smaller dimensions that may represent tech-  
262 nical noise or other unwanted sources of variation. Overall, increasing the number of dimensions  
263 robustly led to a decrease in the number of clusters (Supplementary Figures 10 and 15). This  
264 tended to affect the accuracy of the clustering (Supplementary Figure 16), although in both cases  
265 (number of clusters and ARI) *Seurat*'s resolution parameter had a much stronger impact.

266 Different approaches have been proposed to select the appropriate number of dimensions,  
267 from the visual identification of an 'Elbow' (inflection point) of the variance explained, to more  
268 complex algorithms. We evaluated the performance of dimensionality estimators implemented in  
269 the *intrinsicDimension* package [36], as well as common procedures such as the 'elbow' method  
270 (inflection point in the variance explained by each component), some scRNAseq-specific methods  
271 such as the JackStraw procedure [37] or *scan*'s denoisePCA [26], and the recent application of  
272 Fisher Separability analysis [38].

273 We compared the various estimates of dimensionality in their ability to retrieve the intrinsic  
274 number of dimensions in a dataset, based on *Seurat*'s weighted PCA space. As a first approx-  
275 imation of the true dimensionality, we computed the variance in each principal component that  
276 was explained by the subpopulations, which sharply decreased after the first few components in  
277 most datasets (Figure 8A). Figure 8B shows the difference between the dimension estimates of the  
278 above methods and that based on the subpopulations (i.e. from Figure 8A). Of note, the methods  
279 differ widely in view of their computing time (Figure 8B) and we saw no relationship between the  
280 accuracy of the estimates and the complexity of the method. Reasoning that over-estimating the  
281 number of dimensions is less problematic than under-estimating it, we kept the former methods for  
282 a full analysis of their impact on clustering (Figure 8C-D), when combined with *sctransform* or the



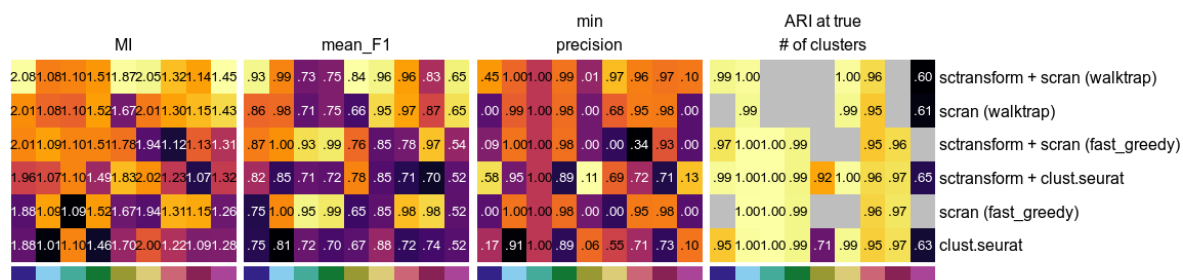
**Figure 8: Estimating dimensionality.** **A:** Estimated dimensionality using the proportion of variance in each component explained by the true subpopulations. **B:** Difference between ‘real’ dimensionality (from **A**) and dimensionality estimation methods, along with computing time. **C:** Average silhouette width per subpopulation across a selection of methods, combined with *sctransform* or standard *Seurat* normalization. **D:** Clustering accuracy across the normalization/dimensionality estimation methods.

standard Seurat normalization. Although most methods performed well on the various clustering measures, the global maximum likelihood based on translated Poisson Mixture Model (*maxLik-Global*) provided the dimensionality estimate that best separated the subpopulations (Figure 8C) and resulted in the best clustering accuracy (Figure 8D). This method systematically selected many more components than were associated with the subpopulations, suggesting that although these additional components appear individually uninformative, in combination they nevertheless contribute to classification.

## Clustering

The last step evaluated in our pipeline was clustering. Given previous works on the topic [6,7] and the success of graph-based clustering methods for scRNAseq, we restricted our evaluation to Seurat's method and two *scrn* SNN-based clustering approaches, based on random walks (*walktrap* method) or on the optimization of the modularity score (*fast\_greedy*). Again, the tested methods were combined with *Seurat*'s standard normalization and *sctransform*, otherwise using the parameters found optimal in the previous steps.

Since ARI is dominated by differences in the number of clusters (Supplementary Figure 17) and no single metric is perfect, we diversified them (Figure 9). Mutual information (MI) has the virtue of not decreasing when a true subpopulation is split into two clusters, which is arguably less problematic (and might well reflect unknown biological subgroups), but as a consequence it can be biased towards methods producing higher resolution clustering. Similarly, precision per true subpopulation is considerably more robust to differences in the number of clusters. We also tracked the mean F1 score and the ARI at the true number of clusters.



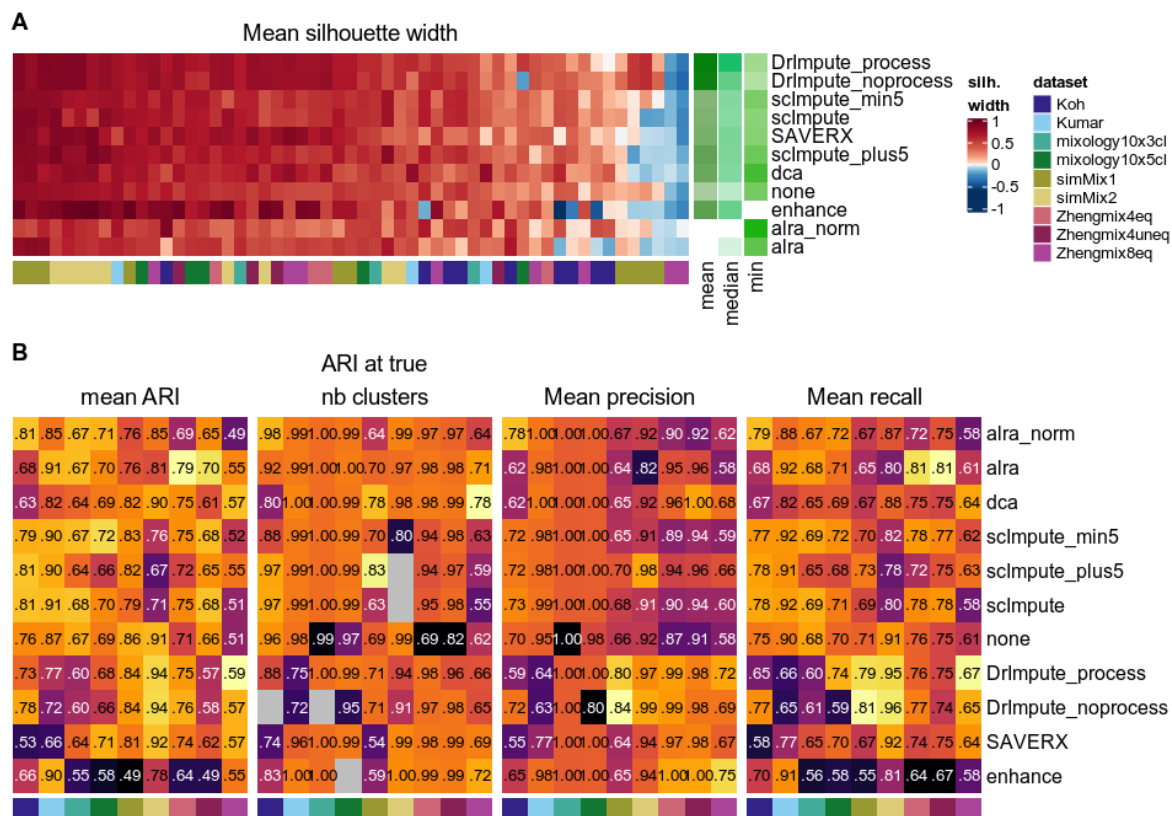
**Figure 9: Evaluation of clustering methods.** Clustering accuracy of common clustering tools in combination with *sctransform* and standard Seurat's normalization.

The MI score and minimum precision, which are largely independent of the estimated number of clusters, were overall higher for the *walktrap* method (Figure 9), while the mean F1 score

306 favored both *scran* methods (*walktrap* and *fast greedy*) over *Seurat*. Finally, the ARI score at the  
 307 true number of clusters, when available, showed similar performances. However, because *Seurat*'s  
 308 resolution parameter had a large impact on the number of clusters identified (Supplementary Fig-  
 309 ure 18), *Seurat* could always be coerced into producing the right number of clusters. Instead, the  
 310 number of clusters found by *scran* was considerably less influenced by the available parameters  
 311 (number of nearest neighbors or steps in the random walk - see Supplementary Figure 18), and  
 312 as a result *scran*-based clustering sometimes never produced a partitioning with the right number  
 313 of clusters. This observation, along with *scran*'s higher MI score, suggest that *scran* sometimes  
 314 simply divides a real subpopulation into two clusters (possibly tracking some unknown biolog-  
 315 ical differences) rather than committing misclassification errors. Overall, the *walktrap* method  
 316 appeared superior to the *fast greedy* algorithm and was generally less prone to misclassification  
 317 than *Seurat* clustering, although the latter offered more control over the resolution. Finally, some  
 318 poorly distinguishable subpopulations from both the Zhengmix8eq and simMix1 datasets remained  
 319 very inaccurately classified by all methods and in regards to all metrics.

## 320 **Further extensions to the pipeline: imputation/denoising**

321 The basic pipeline presented here can be extended with additional analysis steps while keeping the  
 322 same evaluation metrics. To demonstrate this, we evaluated various imputation or denoising tech-  
 323 niques based on their impact on classification. Since preliminary analysis showed that all methods  
 324 performed equally well or better on normalized data, we applied them after filtering and normal-  
 325 ization, but before scaling and reduction. Although some of the methods (e.g., *DRIImpute\_process*  
 326 and *alra\_norm*) did improve the separability of some more elusive subpopulations, no method had  
 327 a systematically positive impact on the average silhouette width across all subpopulations (Figure  
 328 10A). When restricting ourselves to clustering analyses that yielded the 'right' number of clusters,  
 329 all tested methods improved classification compared to a scenario with no imputation step ('none'  
 330 label, Figure 10B). However, the situation was not so straightforward with alternative metrics,  
 331 where some methods (e.g., *enhance*) consistently underperformed. 10X datasets, which are typ-  
 332 ically characterized by a lower per-cell coverage and feature detection rate, benefited more from  
 333 imputation and, in this context, *DRIImpute* and *dca* tended to show the best performance. On  
 334 the contrary, imputing on normalized counts originating from Smart-seq technology was instead  
 335 rather deleterious to the clustering accuracy.



**Figure 10: Imputation/denoising methods** Average silhouette width per subpopulation (**A**) and clustering accuracy (**B**) with or without (indicated as *none*) application of a denoising/imputation method.

## Discussion

### Concrete recommendations

On the basis of our findings, we can make a number of concrete analysis recommendations, also summarized in Figure 11:

#### 1. Filtering

- Doublet detection and removal is advised and can be performed at little computing cost with software such as *scDbIFinder* or *scds*.
- Distribution-based cell filtering fails to capture doublets and should use relatively lenient cutoffs (e.g., 5 MADs, or 3 MADs in at least 2 distributions) to exclude poor-quality cells while avoiding bias against some subpopulations.
- Features filtering based on feature type did not appear beneficial.

## 347 2. Normalization and scaling

- 348 • Most normalization methods tested yielded a fair performance, especially when com-  
349 bined with scaling, which tended to have a positive impact on clustering.
- 350 • *sctransform* offered the best overall performance in terms of the separability of the  
351 subpopulations, as well as removing the effect of library size and detection rate.
- 352 • The common practice of regressing out cell covariates, such as the detection rate or  
353 proportion of mitochondrial reads nearly always had a negative impact, leading to  
354 increased correlation with covariates and decreased clustering accuracy. We therefore  
355 advise against this practice.

## 356 3. Feature selection

- 357 • Deviance [35] offered the best ranking of genes for feature selection.
- 358 • Increasing the number of features included tended to lead to better classifications,  
359 plateauing from 4000 features in our datasets.

## 360 4. Denoising/imputation

- 361 • Denoising appeared beneficial to the identification of subpopulations in 10x datasets,  
362 but not in Smart-seq datasets.
- 363 • We found especially *alra* (with prior normalization), *DrlImpute* (with prior processing)  
364 and *dca* to offer the best performances.

## 365 5. PCA

- 366 • Similarly to previous reports [11], we recommend the irlba-based PCA using weighting  
367 of the components, as implemented in *Seurat*.
- 368 • Instead of the common elbow or jackstraw methods for deciding on the components  
369 to include, we recommend the *global maximum likelihood based on translated Poisson*  
370 *Mixture Model* method (e.g., implemented in the `maxLikGlobalDimEst` function of  
371 `intrinsicDimension`).

## 372 6. Clustering

- 373 • We found *scrn*-based *walktrap* clustering to show good performances.

- 374 • In cases where prior knowledge can guide the choice of a resolution, *Seurat* can be useful  
375 in affording manual control of it while, in the absence of such knowledge, *scrn*-based  
376 *walktrap* clustering provided reasonable estimates.

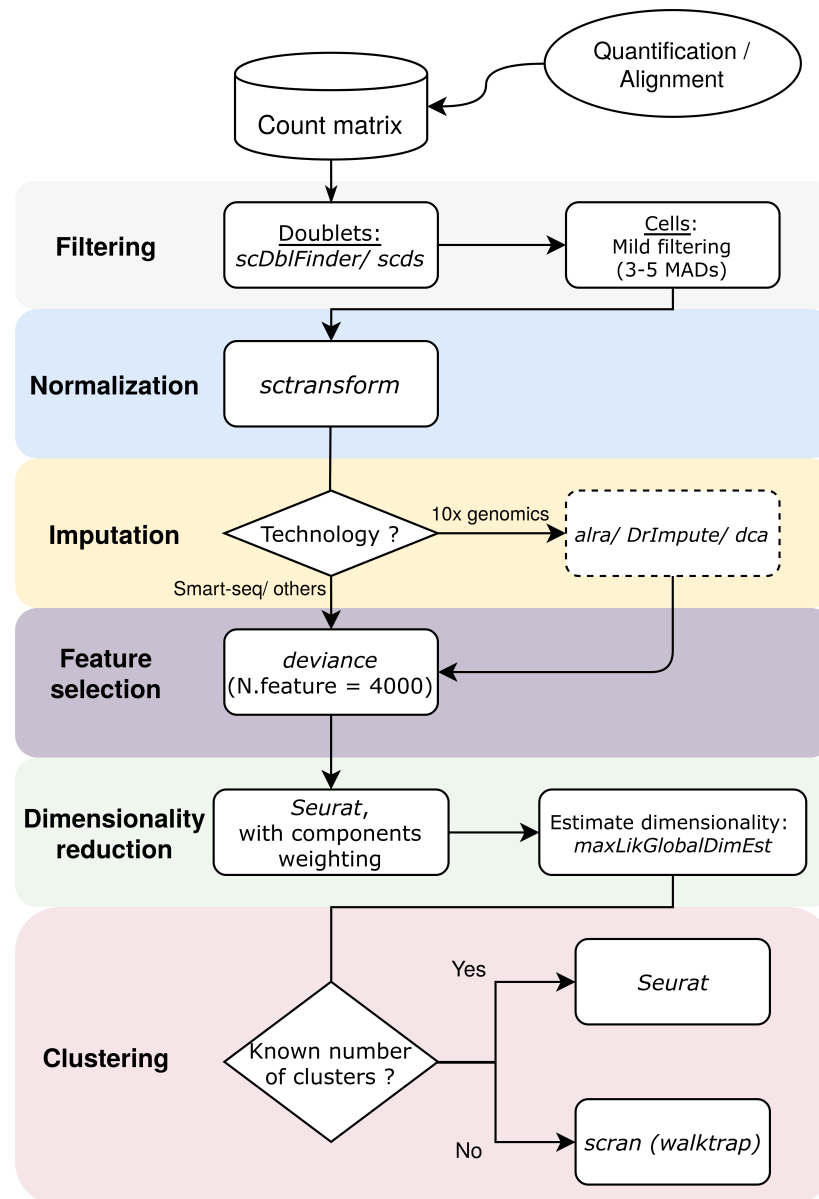


Figure 11: Recommendations of tools from filtering to clustering.

## 377 Limitations and open questions

378 In this study, we evaluated tools commonly used for the processing of scRNAseq with a focus on  
379 droplet-based datasets, namely from the 10x technology. Although this platform has been used

in almost half of the scRNAseq studies in 2019 [2], other popular technologies such as Drop-seq, InDrops or Smart-seq2/3 were not represented in the present benchmarking. Differences between such protocols, even among droplet-based technologies, can have a very large impact on cell capture efficiency, cell numbers and clustering [39,40,41]. Although most top-ranking methods in our comparison performed well on both Smart-seq and 10x datasets that we tested, future benchmarking efforts should strive to include less represented technologies. In addition, we did not compare any of the alignment and/or quantification methods used to obtain the count matrix, which was for instance discussed in [15]. Some steps, such as the implementation of the PCA, were also not explored in detail here as they have already been the object of recent and thorough study elsewhere [11,14]. We also considered only methods relying on Euclidean distance, while correlation was recently reported to be superior [42] and would require further investigation.

Here, we chose to concentrate on what could be learned from datasets with known cell labels (as opposed to labels inferred from the data, as in [39]). In contrast to Tian et al. [12], who used RNA mixtures of known proportions, we chose to rely chiefly on real cells and their representative form of variability. Given the limited availability of such well-described datasets however, several aspects of single-cell analysis could not be compared, such as batch effect correction or multi-dataset integration. For these aspects of scRNAseq processing that are critical in some experimental designs, we refer the reader to previous evaluations [13,43]. In addition, a focus on the identification of the subpopulations might fail to reveal methods that are instead best performing for tasks other than clustering, such as differential expression analysis or trajectory inference. Several informative benchmarks have already been performed on some of these topics [19,5,9,44,10,16]. Yet, such evaluations could benefit from considering methods not in isolation, but as parts of a connected workflow, as we have done here. We believe that the *pipeComp* framework is modular and flexible enough to integrate new steps in the pipeline, as shown by an example with imputation/denoising methods.

We developed *pipeComp* to address the need of a framework that can simultaneously evaluate the interaction of multiple tools. The work of Vieth and colleagues [15] already offered an important precedent in this respect, evaluating the interaction of various steps in the context of scRNAseq differential expression analysis, but did not offer a platform for doing so. Instead, the *CellBench* package was recently proposed to address a similar need [45], offering an elegant piping syntax to combine alternative methods at successive steps. A key additional feature of *pipeComp* is its ability to perform evaluation in parallel to the pipeline and offering on-the-fly evaluation at any

412 step. This avoid the need to store potentially large intermediate data for all possible permutations  
413 and thus allowing a combinatorial benchmark. In addition, the fact that benchmark functions  
414 are stored in the `PipelineDefinition` makes the benchmark more smoothly portable, making it  
415 easy to modify, extend with further methods, or apply to other datasets.

416 With respect to the methods themselves, we believe there is still space for improvement  
417 in some of the steps. Concerning cell filtering, we noticed that the current approach based on  
418 whole-population characteristic (e.g., MAD cut-off) can be biased against certain subpopulations,  
419 suggesting that more refined methods expecting multimodal distributions should be used rather  
420 than relying on whole-population characteristics. In addition, most common filtering approaches  
421 do not harness the relationship between cell QC properties. Finally, imputation had a varied impact  
422 on the clustering analysis and seemed to be linked to the technology that was used to generate  
423 the data. Also, the good performance of *DrlImpute* is in line with a previous study focused on the  
424 preservation of data structure in trajectory inference from scRNAseq [18] (*alra* was however not  
425 included in this previous benchmark). Our results also align with the hypothesis of this study on  
426 the respective strengths of linear and non-linear models and their use to different types of data,  
427 such as the highest performance of non-linear methods in developmental studies.

## 428 Conclusion

429 *pipeComp* is a flexible R framework for evaluating methods and parameter combinations in a  
430 complex pipeline, computing on-the-fly multilevel evaluation metrics. Applying this framework to  
431 scRNAseq clustering enabled us to make concrete recommendations on the steps of filtering, nor-  
432 malization, feature selection, denoising, dimensionality reduction and clustering. We demonstrate  
433 how a diversity of multilevel metrics can be more robust, more sensitive, and more nuanced than  
434 simply evaluating the final clustering performance. In addition, we provide a new and efficient  
435 Bioconductor package for doublet detection, *scDb1Finder*. We hope that the *pipeComp* frame-  
436 work can be applied to extend the current field of benchmarking, as well as to apply it to other  
437 range of methods.

## 438 **Methods**

### 439 **Code and data availability**

440 All analyses were performed through the pipeComp R package, which implements the pipeline  
441 framework described here. All code to reproduce the simulations and figures is available in the  
442 [https://github.com/markrobinsonuzh/scRNA\\_pipelines\\_paper](https://github.com/markrobinsonuzh/scRNA_pipelines_paper) repository, which also in-  
443 cludes the basic datasets with a standardized annotation.

444 The gene counts for the two mixology datasets were downloaded from the CellBench reposi-  
445 tory, commit 74fe79e. For the other datasets, we used the unfiltered counts from [6], available on  
446 the corresponding repository [https://github.com/markrobinsonuzh/scRNAseq\\_clustering\\_](https://github.com/markrobinsonuzh/scRNAseq_clustering_comparison)  
447 [comparison](https://github.com/markrobinsonuzh/scRNAseq_clustering_comparison). For simplicity, all starting *SingleCellExperiment* objects with standardized metadata  
448 are available on [https://github.com/markrobinsonuzh/scRNA\\_pipelines\\_paper](https://github.com/markrobinsonuzh/scRNA_pipelines_paper).

### 449 **Software and package versions**

450 Analyses were performed in R 3.6.0 (Bioconductor 3.9) and the following packages were installed  
451 from GitHub repositories: *Seurat* (version 3.0.0), *sctransform* (0.2.0), *DoubletFinder* (2.0.1),  
452 *scDbIFinder* (1.1.1), *scds* (1.0.0), *SAVERX* (1.0.0), *sclImpute* (0.0.9), *ALRA* (commit 7636de8),  
453 *DCA* (0.2.2), *DrlImpute* (1.2), *ENHANCE* (R version, commit 1571696), *SAVERX* (1.0.0). The  
454 code for the glmPCA was obtained from <https://github.com/willtownes/scrna2019> (com-  
455 mit 1ddcc30ebb95d083a685f12fe81d35dd1b0cb1b2).

### 456 **Simulated datasets**

457 The simMix1 dataset was based on the human PBMC CITE-seq data deposited under accession  
458 code GSE100866 of the Gene Expression Omnibus (GEO) website. Both RNA and ADT count  
459 data were downloaded from GEO and processed independently using *Seurat*. We then considered  
460 cells that were in the same cluster both in the RNA-based and ADT-based analyses to be real  
461 subpopulations and focused on the 4 most abundant ones. We then performed 3 sampling-based  
462 simulations with various degrees of separation using *muscat* [19] and merged the three simulations.  
463 The simMix2 dataset was generated from the mouse brain data published in [19] in a similar fashion.  
464 The specific code is available on [https://github.com/markrobinsonuzh/scRNA\\_pipelines\\_](https://github.com/markrobinsonuzh/scRNA_pipelines_paper)  
465 [paper](https://github.com/markrobinsonuzh/scRNA_pipelines_paper).

## 466 Default pipeline parameters

467 Where unspecified, the following default parameters or parameter sets were used:

- 468 • *scDbtFinder* was used for doublet identification
- 469 • the default filter sets (see below) were applied
- 470 • standard *Seurat* normalization was employed
- 471 • *Seurat* ( $\geq 3.0$ ) variable feature selection was employed, selecting 2000 genes
- 472 • standard *Seurat* scaling and PCA were employed
- 473 • To study the impact of upstream steps on clustering, various *Seurat* clustering analyses were
- 474 performed, selecting different numbers of dimensions (5, 10, 15, 20, 30, and 50) and using
- 475 various resolution parameters (0.005, 0.01, 0.02, 0.05, 0.1, 0.15, 0.2, 0.3, 0.4, 0.5, 0.8, 1,
- 476 1.2, 1.5, 2, 4). The range of resolution parameters was selected to ensure that the right
- 477 number of clusters could be obtained in all datasets.

## 478 Denoising/imputation

479 Most imputation/denoising methods were run with the default parameters with the following  
 480 exceptions; *DrlImpute* documentation advises to process the data prior to imputation by removing  
 481 lowly expressed genes and cells expressing less than 2 genes. As it is not clear if this step is a hard  
 482 requirement for the method to perform well, *DrlImpute* was run with and without prior processing  
 483 (*DrlImpute\_process* and *DrlImpute\_noprocess* labels, respectively). The *dca* method only accepts  
 484 integers counts, while two of the datasets had non-integer quantification of expected counts. For  
 485 these datasets, we rounded up the counts prior to imputation. *alra* is designed for normalized  
 486 data but as we are evaluating normalization downstream to imputation, we used the method on  
 487 both non-normalized (*alra* label) and normalized counts (*alra\_norm* label). *sclImpute* requires an  
 488 estimation of the expected number of clusters with the input data. As the estimation of the true  
 489 number of cluster may not be known by the user, we evaluated the tool using the true number of  
 490 clusters (*sclImpute* label) and using an over/under-estimation of this number (*sclImpute\_plus5* and  
 491 *sclImpute\_min5* labels, respectively). *ENHANCE* uses a k-nearest neighbor aggregation method  
 492 and automatically estimates the number of neighbors to merge prior to the imputation. With  
 493 the smallest datasets, this parameters was estimated to be 1, which lead to an early stop of the  
 494 function. In such cases, we manually set this parameter to 2 for the method to work.

## 495 Dimensionality estimates

496 For the methods that produce local dimensionality estimates, we used the maximum. For the  
 497 elbow method, we implemented an automatic procedure by taking the farthest point from a line  
 498 drawn between the variance explained by the first and last (i.e. 50th) components calculated. For  
 499 the JackStraw method, since the *Seurat* documentation advises not to use a p-value threshold  
 500 (which would typically yield a very large number of dimensions) but rather look for a drop in  
 501 significance, we applied the same farthest point algorithm on the  $\log_{10}(\text{p-values})$ , which in our  
 502 hands reproduced manual threshold selection.

## 503 Doublet detection method

504 Our doublet detection method is available at <https://github.com/plger/scDblFinder>. Briefly,  
 505 after reducing the data to the most expressed genes, we cluster the cells using the fast-greedy  
 506 algorithm, favouring overclustering. We then create artificial doublets by sampling the two cells  
 507 specifically from different clusters and sum the counts of each pair of cells. We also create meta-  
 508 cells from each cluster and use them to create additional doublets and triplets. We combine them  
 509 with the real cells, perform PCA and build a KNN graph using *BiocNeighbors*. We then calculate,  
 510 for each cell, the proportion of its neighbors that are artificial doublets, weighted by the distance.  
 511 This ratio serves as a doublet score, which is then thresholded by simultaneously minimizing the  
 512 error in classifying real vs artificial doublets and the deviation from the distribution of expected  
 513 doublet rate (accounting for homotypic doublets as done by *DoubletFinder*).

## 514 Filter sets

515 The *default* set of filters excludes cells that are outliers according to at least two of the following  
 516 thresholds:  $\log_{10}(\text{total\_counts}) > 2.5$  MADs or  $< 5$  MADs,  $\log_{10}(\text{total\_features}) > 2.5$  MADs or  $< 5$   
 517 MADs,  $\text{pct\_counts\_in\_top\_20\_features} > \text{or} < 5$  MADs,  $\text{featcount\_dist}$  (distance to expected ratio  
 518 of  $\log_{10}$  counts and features)  $> \text{or} < 5$  MADs,  $\text{pct\_counts\_Mt} > 2.5$  MADs and  $> 0.08$ .

519 The *stringent* set of filters uses the same thresholds, but excludes a cell if it is an outlier on  
 520 any single distribution.

521 The *lenient* set of filters excludes cells that are outliers on at least two distributions by at  
 522 least 5 MADs, except for  $\text{pct\_counts\_Mt}$  where the threshold is  $> 3$  MADs and  $> 0.08$ .

523 For cluster-wise filters, clusters were first identified with *scrn::quickCluster* and the filters

were then applied separately for each cluster.

The 'pca' and 'pca2' clusters refer to the multivariate outlier detection methods implemented in *scater*, running *runPCA* with `use_coldata=TRUE`, `detect_outliers=TRUE`. 'pca' uses all covariates, while 'pca2' uses only the  $\log_{10}(\text{counts})$ ,  $\log_{10}(\text{features})$ , proportion mitochondrial and proportion in the top 50 features.

## Variance and deviance explained

Unless specified otherwise, the variance in gene expression explained by subpopulations was calculated on the data normalized and transformed through *sctransform*. For each gene, we fitted a linear model using the subpopulation as only independent variable ( $\sim \text{subpopulation}$ ) and used the R-squared as a measure of the variance explained. The same method was used for principal components.

The deviance explained by subpopulations was calculated directly on counts using the `getDevianceExplained` function of *pipeComp*. The function uses *edgeR* to fit two models, a full ( $\sim \text{librarySize} + \text{population}$ ) and a reduced one ( $\sim \text{librarySize}$ ). For each gene, the deviance explained is then the difference between the deviance of each models, divided by the deviance of the reduced model. In the rare cases where this resulted in a negative deviance explained, it was set to zero.

To estimate the correlation between the principal components and covariates such as library size, we first fitted a linear model on the subpopulations and correlated the residuals of this model with the covariate of interest.

## Tables

**Table 1:** Overview of the benchmark datasets.

dataset	source	protocol	description
Koh	GSE85066	SMARTer	FACS purified H7 hESC in different differentiation stages
Kumar	GSE60749	SMARTer	Mouse ESC cultured in different conditions
Zhengmix4eq	[6]	10x	Mixtures of FACS purified PBMCs
Zhengmix4uneq	[6]	10x	Mixtures of FACS purified PBMCs
Zhengmix8eq	[6]	10x	Mixtures of FACS purified PBMCs
mixology10x3cl	[12]	10x	Mixture of 3 cancer cell lines from CellBench
mixology10x5cl	[12]	10x	Mixture of 5 cancer cell lines from CellBench
simMix1	-	10x-based	Simulation of 10 human cell subpopulations
simMix2	-	10x-based	Simulation of 9 mouse cell subpopulations

## 545 **Competing interests**

546 As developers of the *scDbIFinder* package, the authors have an interest in it performing well. The  
547 authors declare no further competing interests.

## 548 **Acknowledgements**

549 This work was supported by the Swiss National Science Foundation (grants 310030\_175841, CR-  
550 SII5\_177208) as well as the Chan Zuckerberg Initiative DAF (grant number 2018-182828), an  
551 advised fund of Silicon Valley Community Foundation. MDR acknowledges support from the  
552 University Research Priority Program Evolution in Action at the University of Zurich.

# References

1. Zappia, L., Phipson, B. & Oshlack, A. Exploring the single-cell RNA-seq analysis landscape with the scRNA-tools database. *PLoS Computational Biology* **14**, e1006245 (2018).
2. Svensson, V. & Beltrame, E. d. V. A curated database reveals trends in single cell transcriptomics. *bioRxiv*, 742304 (2019).
3. Cobos, F. A., Alquicira-hernandez, J. & Powell, J. Comprehensive benchmarking of computational deconvolution of transcriptomics data. *bioRxiv*. doi:10.1101/2020.01.10.897116.T (2020).
4. Cole, M. B. *et al.* Performance Assessment and Selection of Normalization Procedures for Single-Cell RNA-seq. en. *bioRxiv*. doi:10.1101/235382. (2019) (May 2018).
5. Dal Molin, A., Baruzzo, G. & Di Camillo, B. Single-cell RNA-sequencing: Assessment of differential expression analysis methods. *Frontiers in Genetics* **8**. doi:10.3389/fgene.2017.00062 (2017).
6. Duo, A., Robinson, M. D. & Soneson, C. A systematic performance evaluation of clustering methods for single-cell RNA-seq data. *F1000Research* **7**. doi:10.12688/f1000research.15666.2 (2018).
7. Freytag, S., Tian, L., Lönnstedt, I., Ng, M. & Bahlo, M. Comparison of clustering tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data. *F1000Research* **7**, 1–29 (2018).
8. Heiser, C. N. & Lau, K. S. A quantitative framework for evaluating single-cell data structure preservation by dimensionality reduction techniques. *bioRxiv*, 684340 (2019).
9. Jaakkola, M. K., Seyednasrollah, F., Mehmood, A. & Elo, L. L. Comparison of methods to detect differentially expressed genes between single-cell populations. *Briefings in bioinformatics* **18**, 735–743 (2017).
10. Soneson, C. & Robinson, M. D. Bias, robustness and scalability in single-cell differential expression analysis. *Nature Methods* **15**, 255–261 (2018).
11. Sun, S., Zhu, J., Ma, Y. & Zhou, X. Accuracy , robustness and scalability of dimensionality reduction methods for single-cell RNA-seq analysis. *Genome Biology* **20**, 1–21 (2019).
12. Tian, L. *et al.* scRNA-seq mixology: towards better benchmarking of single cell RNA-seq protocols and analysis methods. en. *bioRxiv*. doi:10.1101/433102. (2019) (Oct. 2018).
13. Tran, H. T. N. *et al.* A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome biology* **21**, 1–32 (2020).
14. Tsuyuzaki, K., Sato, H., Sato, K. & Nikaido, I. Benchmarking principal component analysis for large-scale single-cell RNA-sequencing. *Genome Biology* **21**, 1–17 (2020).
15. Vieth, B., Parekh, S., Ziegenhain, C., Enard, W. & Hellmann, I. A Systematic Evaluation of Single Cell RNA-Seq Analysis Pipelines. en. *bioRxiv*, 583013 (Mar. 2019).
16. Wang, T., Li, B., Nelson, C. E. & Nabavi, S. Comparative analysis of differential gene expression analysis tools for single-cell RNA sequencing data. *BMC Bioinformatics* **20.1**, 1–16 (2019).
17. Yip, S. H., Sham, P. C. & Wang, J. Evaluation of tools for highly variable gene discovery from single-cell RNA-seq data. *Briefings in Bioinformatics* **20**, 1583–1589 (2018).
18. Zhang, L. & Zhang, S. Comparison of computational methods for imputing single-cell RNA-sequencing data. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. doi:10.1109/TCBB.2018.2848633 (2018).
19. Crowell, H. L. *et al.* On the discovery of population-specific state transitions from multi-sample multi-condition single-cell RNA sequencing data. *bioRxiv*, 713412 (2019).
20. Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. en. *Nature Biotechnology* **33**, 495–502 (May 2015).
21. Hubert, L. & Arabie, P. Comparing partitions. *Journal of Classification* **2**, 193–218 (1985).
22. Steinley, D. Properties of the Hubert-Arabie Adjusted Rand Index. en. *Psychological Methods* **9**, 386–396 (2004).
23. Bloom, J. D. Estimating the frequency of multiplets in single-cell RNA sequencing from cell-mixing experiments. en. *PeerJ* **6**, e5578 (Sept. 2018).
24. Kang, H. M. *et al.* Multiplexed droplet single-cell RNA-sequencing using natural genetic variation. en. *Nature Biotechnology* **36**, 89–94 (Jan. 2018).

25. McGinnis, C. S., Murrow, L. M. & Gartner, Z. J. DoubletFinder: Doublet Detection in Single-Cell RNA Sequencing Data Using Artificial Nearest Neighbors. English. *Cell Systems* **0**. doi:10.1016/j.cels.2019.03.003. [https://www.cell.com/cell-systems/abstract/S2405-4712\(19\)30073-0](https://www.cell.com/cell-systems/abstract/S2405-4712(19)30073-0) (2019) (Apr. 2019).
26. Lun, A. T. L., McCarthy, D. J., Marioni, J. C. & McDavid, A. A step-by-step workflow for low-level analysis of single-cell RNA-seq data with Bioconductor. *F1000Research* **5**. doi:10.12688/f1000research.9501.2 (2016).
27. Bais, A. S. & Kostka, D. scds: computational annotation of doublets in single-cell RNA sequencing data. *Bioinformatics*, 1–9 (2019).
28. Ilicic, T. *et al.* Classification of low quality cells from single-cell RNA-seq data. *Genome Biology* **17.1**, 1–15 (2016).
29. Papadimitriou, C. H. & Steiglitz, K. *Combinatorial optimization: algorithms and complexity* 496. ISBN: 0486402584 (1998).
30. L. Lun, A. T., Bach, K. & Marioni, J. C. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. en. *Genome Biology* **17**. doi:10.1186/s13059-016-0947-7. <http://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0947-7> (2019) (Dec. 2016).
31. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. en. *bioRxiv*. doi:10.1101/576827. <http://biorxiv.org/lookup/doi/10.1101/576827> (2019) (Mar. 2019).
32. Lin, Y. *et al.* Evaluating stably expressed genes in single cells. en. *bioRxiv*, 229815 (Nov. 2018).
33. Deeke, J. M. & Gagnon-Bartsch, J. A. Stably expressed genes in single-cell RNA-sequencing: en. *bioRxiv*. doi:10.1101/475426. <http://biorxiv.org/lookup/doi/10.1101/475426> (2019) (Nov. 2018).
34. Rousseeuw, P. J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics* **20**, 53–65 (1987).
35. Townes, F. W., Hicks, S. C., Aryee, M. J. & Irizarry, R. A. Feature Selection and Dimension Reduction for Single Cell RNA-Seq based on a Multinomial Model. en. *bioRxiv*. doi:10.1101/574574. <http://biorxiv.org/lookup/doi/10.1101/574574> (2019) (Mar. 2019).
36. Johnsson, K., Soneson, C. & Fontes, M. Low Bias Local Intrinsic Dimension Estimation from Expected Simplex Skewness. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **37**, 196–202 (Jan. 2015).
37. Chung, N. C. & Storey, J. D. Statistical significance of variables driving systematic variation in high-dimensional data. *Bioinformatics* **31**, 545–554 (2015).
38. Albergante, L., Bac, J. & Zinovyev, A. *Estimating the effective dimension of large biological datasets using Fisher separability analysis in International Joint Conference on Neural Networks (IJCNN)* (IEEE, 2019), 1–8. doi:10.1109/IJCNN.2019.8852450. arXiv: arXiv:1901.06328v1.
39. Mereu, E. *et al.* Benchmarking Single-Cell RNA Sequencing Protocols for Cell Atlas Projects. *bioRxiv*, 630087 (2019).
40. Zhang, X. *et al.* Comparative Analysis of Droplet-Based Ultra-High-Throughput Single-Cell RNA-Seq Systems. *Molecular Cell* **73**, 130–142 (2019).
41. Salomon, R. *et al.* Droplet-based single cell RNAseq tools: a practical guide. *Lab on a Chip* **19**, 1706–1727 (2019).
42. Kim, T. *et al.* Impact of similarity metrics on single-cell RNA-seq data clustering. en. *Briefings in Bioinformatics* **20**, 2316–2326 (Nov. 2019).
43. Stuart, T. & Satija, R. Integrative single-cell analysis. *Nature Reviews Genetics* **20**, 257–272 (2019).
44. Saelens, W., Cannoodt, R., Todorov, H. & Saeys, Y. A comparison of single-cell trajectory inference methods. *Nature Biotechnology* **37**, 547–554 (2019).
45. Su, S. *et al.* CellBench: R/Bioconductor software for comparing single-cell RNA-seq analysis methods. en. *Bioinformatics*. doi:10.1093/bioinformatics/btz889. (2019).