## METHOD

# Benchmarking principal component analysis for large-scale single-cell RNA-sequencing

Koki Tsuyuzaki[1,2], Hiroyuki Sato[3], Kenta Sato[1,4] and Itoshi Nikaido[1,5*]

*Correspondence:
itoshi.nikaido@riken.jp
[1]Laboratory for Bioinformatics Research, RIKEN Center for Biosystems Dynamics Research, Japan
Full list of author information is available at the end of the article
†Equal contributor

**Abstract**

Principal component analysis (PCA) is an essential method for analyzing single-cell RNA-seq (scRNA-seq) datasets, but large-scale scRNA-seq datasets require long computational times and a large memory capacity.

In this work, we review 21 fast and memory-efficient PCA implementations (10 algorithms) and evaluate their application using 4 real and 18 synthetic datasets. Our benchmarking showed that some PCA algorithms are faster, more memory efficient, and more accurate than others. In consideration of the differences in the computational environments of users and developers, we have also developed guidelines to assist with selection of appropriate PCA implementations.

**Keywords:** Single-cell RNA-seq; Cellular heterogeneity; Dimension reduction; Principal component analysis; Online/Incremental algorithm; Randomized algorithm; Out-of-core; Sparse data format; R; Python; Julia

## Background

The emergence of single-cell RNA sequencing (scRNA-seq) technologies [1], has enabled the examination of many types of cellular heterogeneity. For example, cellular subpopulations consisting of various tissues [2–6], rare cells and stem cell niches [7], continuous gene expression changes related to cell cycle progression [8], spatial coordinates [9–11], and differences in differentiation maturity [12,13] have been captured by many scRNA-seq studies. As the measurement of cellular heterogeneity is highly dependent on the number of cells measured simultaneously, a wide variety of large-scale scRNA-seq technologies have been developed [14], including those using cell sorting devices [15–17], Fludigm C1 [18–21], droplet-based technologies (Drop-Seq [2–4], inDrop RNA-Seq [5,6], the 10X Genomics Chromium system [22]), and single-cell combinatorial-indexing RNA-sequencing (sci-RNA-seq [23]). Such technologies have encouraged the establishment of several large-scale genomics consortiums, such as the Human Cell Atlas [24–26], Mouse Cell Atlas [27], and Tabula Muris [28]. These projects are analyzing a tremendous number of cells by scRNA-seq and tackling basic life science problems such as the number of cell types comprising an individual, cell-type-specific marker gene expression and gene functions, and molecular mechanisms of diseases at a single-cell resolution.

Nevertheless, the analysis of scRNA-seq datasets poses a potentially difficult problem; the cell type corresponding to each data point is unknown *a priori* [1,29–35]. Accordingly, researchers perform unsupervised machine learning (UML) methods, such as dimensionality reduction and clustering, to reveal the cell type corresponding to each individual data point. In particular, principal component analysis (PCA [36–38]) is a commonly used UML algorithm applied across many situations.

Despite its wide use, there are several reasons why it is unclear how PCA should be conducted for large-scale scRNA-seq. First, because the widely used PCA algorithms and implementations load all elements of a data matrix into memory space, for large-scale datasets such as the 1.3 million cells measured by 10X Genomics Chromium [39] or the 2 million cells measured by sci-RNA-seq [23], the calculation is difficult unless the memory size of the user's machine is very large. Furthermore, the same data analysis workflow is performed repeatedly, with deletions or additions to the data or parameter changes for the workflow, and under such trial-and-error cycles, PCA can become a bottleneck for the workflow. Therefore, some fast and memory-efficient PCA algorithms are required.

Second, there are indeed some PCA algorithms that are fast and memory efficient, but their practicality for use with large-scale scRNA-seq datasets is not fully understood. Generally, there are trade-offs between the acceleration of algorithms by some approximation methods and the accuracy of biological data analysis. Fast PCA algorithms might overlook some important differential gene expression patterns. In the case of large-scale scRNA-seq studies aiming to find novel cell types, this property may cause a loss of clustering accuracy and not be acceptable.

Finally, actual computational time and memory efficiency are highly dependent on the specific implementation, including the programming language, the method for loading input files, and the data format. However, there is no benchmarking to evaluate these properties. Such information is directly related to the practicality of the software and is useful as a guideline for users and developers.

For the above reasons, in this research, we examine the practicality of fast and memory-efficient PCA algorithms for use with large-scale scRNA-seq datasets. This work provides four key contributions. First, we review the existing PCA algorithms and their implementations (Figure 1). Second, we present a benchmark test with selected PCA algorithms and implementations. To our knowledge, this is the first comprehensive benchmarking of PCA algorithms and implementations with large-scale scRNA-seq datasets. Third, we provide some original implementations of some PCA algorithms and utility functions for quality control (QC), filtering, and feature selection. All commands are implemented in a fast and memory-efficient Julia package. Finally, we propose guidelines for end-users and software developers.

## Results

### Review of PCA algorithms and implementations

PCA is widely used for data visualization [39–41], data QC [42], feature selection [13, 43–49], de-noising [50, 51], imputation [52–54], confirmation and removal of batch effects [55–57], confirmation and estimation of cell-cycle effects [58], rare cell type detection [59, 60], cell type and cell state similarity search [61], pseudotime inference [13, 62–66], and spatial reconstruction [9].

Additionally, principal component (PC) scores are also used as the input of other non-linear dimensionality reduction [67–73] and clustering methods [74–77] in order to preserve the global structure, avoid the "curse of dimensionality" [78–81], and save memory space. A wide variety of scRNA-seq data analysis tools actually include PCA as an internal function or utilize PC scores as input for down-stream analyses [22, 82–89].

We reviewed the existing PCA algorithms and implementations and classified the algorithms into six categories, namely, similarity transformation-based (SimT), downsampling-based (DS), singular value decomposition (SVD) update-based (SU), Krylov subspace-based (Krylov), gradient descent-based (GD), and random projection-based (Rand) (Additional file 1 [22, 42–44, 49–52, 55–61, 63, 65, 69, 74–77, 82, 85, 89–113]). We have listed 21 PCA implementations (comprising 10 algorithms) that are freely available and easy to download, install, and use for analyses. The correspondence of the reviewed PCA implementations and scRNA-seq studies are summarized in Table 1.

To extend the scope of the algorithms used in the benchmarking, we originally implemented some PCA algorithms in an out-of-core manner (Additional file 1). The pseudo-code and source code of all the algorithms benchmarked in this study are summarized in Additional file 2 and Additional file 3, respectively.

### Benchmarking of PCA algorithms and implementations

Next, we performed the benchmarking tests of the PCA algorithms and implementations. The results of the benchmarking are summarized in Figure 2 [69, 90, 92, 94–99, 107–109, 114, 115].

#### Real-world datasets

In consideration of the trade-offs among the large number of methods evaluated with our limited time, computational resources, and manpower, we carefully selected real-world datasets for the benchmarking. The latest scRNA-seq methods are divided into two categories, namely, full-length scRNA-seq methods and high-throughput scRNA-seq methods with specific cell dissociation and cellular/molecular barcoding technologies such as droplet-based and split-and-pool experiments [34, 35]. Because the number of cells measured by scRNA-seq has been increased by the latter technology, we selected the following four datasets generated by such technologies: human peripheral blood mononuclear cells (PBMCs), human pancreatic cells (Pancreas), mouse brain and spinal cord (BrainSpinalCord), and mouse cells from the cortex, hippocampus, and ventricular zone (Brain) (Table 2). These datasets have been used in many previous scRNA-seq studies [61, 76, 94, 116–122].

#### The accuracy of PCA algorithms

Here, we evaluate the accuracy of the various PCA algorithms by using the four real-world datasets. For the analyses of the PBMCs and Pancreas datasets, we set the result of `prcomp` as the gold standard, which is a wrapper function for performing SVD with LAPACK subroutines (Additional file 1). The other implementations are compared with this result (Figures 1b and 2). For the BrainSpinalCord and Brain datasets analyses, full-rank SVD by LAPACK is computationally difficult. According to the benchmarking guidelines developed by Mark D. Robinson's group [123], comparing the methods against each other is recommended when the ground truth cannot be defined. Therefore, we just compared the results of the methods against each other using several different criteria, such as the magnitude of the eigenvalues and the clustering accuracy.

First, we performed t-stochastic neighbor embedding (t-SNE [67, 68]) and uniform manifold approximation and projection (UMAP [71, 72]) for the results of each PCA

algorithm and compared the clarity of the cluster structures detected by the original studies (Figures 1b, 3, Additional file 4, and Additional file 5). For the BrainSpinalCord and Brain datasets, only downsampling, `IncrementalPCA` (*sklearn*), `orthiter`/`gd`/`sgd`/`halko`/`algorithm971` (*OnlinePCA.jl*), and `oocPCA_CSV` (*oocR-PCA*) could be performed, while the other implementations were terminated by out-of-memory errors on 96 and 128 GB RAM machines. For the PBMCS and Pancreas datasets, compared with the gold standard cluster structures, the structures detected by downsampling were unclear, and some distinct clusters determined by the original studies were incorrectly combined into single clusters. In the realistic situation when the cellular labels were unavailable *a priori*, the labels were exploratorily estimated by confirming differentially expressed genes, known markergenes, or related gene functions of clusters. In such a situation, downsampling may overlook subgroups hiding in a cluster.

We also performed four clustering algorithms on all the results of the PCA implementations and calculated the adjusted Rand index (ARI [124]) to evaluate clustering accuracy (Additional file 6). Here, we only show the result of Louvain clustering [125] (Figures 1b and 4). The ARI values show that the results of downsampling and `sgd` (*OnlinePCA.jl*) were worse compared with the gold standard or other implementations.

Next, we performed an all-to-all comparison between PCs from the gold standard and the other PCA implementations (Figures 1b, 5a, and Additional file 7). Because the PCs are unit vectors, when two PCs are directed in the same or opposite direction, their cross product becomes 1 or $-1$, respectively. Both the same and opposite direction vectors are mathematically identical in PCA optimization, and different PCA implementations may yield PCs with different signs. Accordingly, we calculated the absolute value of the cross product ranging from 0 to 1 for the all-to-all comparison and evaluated whether higher PCs, which correspond to lower eigenvalues, are accurately calculated. The Figure 5a and Additional file 7 show that the higher PCs based on downsampling, `orthiter`/`gd`/`sgd` (*OnlinePCA.jl*), and `PCA` (*dask-ml* [115]) become inaccurate as the dimensionality of a PC increases. The higher PCs of these implementations also appear noisy and unclear in pair plots of PCs between each implementation and seem uninformative (Additional file 8, Additional file 9, Additional file 10, and Additional file 11). In particular, the higher PCs calculated by downsampling and `sgd` (*OnlinePCA.jl*) are sometimes influenced by the existence of outlier cells (Additional file 8 and Additional file 9). When performing some clustering methods, such as $k$-means and Gaussian mixture model (GMM [126]) methods, such outlier cells are also detected as singleton clusters having only a single cell as their cluster member (Additional file 12). Contrary to these results, all the implementations of IRLBA and IRAM, as well as the randomized SVD approaches except for `PCA` (*dask-ml*), are surprisingly accurate regardless of the language in which they are written or their developers. Although `PCA` (*dask-ml*) is based on Halko's method and is nearly identical to the other implementations of Halko's method, this function uses the direct tall-and-skinny QR algorithm [127] (https://github.com/dask/dask/blob/a7bf545580c5cd4180373b5a2774276c2ccbb573/dask/array/linalg.py#L52), and this characteristic might be related to the inaccuracy of the implementations. Because there is no gold standard in the case of the BrainSpinalCord and Brain

datasets, we compared the eigenvectors of the PCA implementations in all possible combinations (Additional file 13) and found that the higher PCs of downsampling and `sgd` differed from those of the other PCA implementations.

Because gene-wise eigenvectors (i.e., loading vectors) are also retrieved from the data matrix and cell-wise eigenvectors (i.e., PCs), we also compared the loading vectors (Figure 5b and Additional file 14). We extracted the top 500 genes in terms of the largest absolute values of loading vectors and calculated the number of genes in common between the two loading vectors. As is the case with the eigenvectors, even for loading vectors, downsampling, `orthiter/gd/sgd` (*OnlinePCA.jl*), and `PCA` (*dask-ml* [115]) become inaccurate as the dimensionality of the PC increases. Because the genes with large absolute values for loading vectors are used as feature values in some studies [43–48], inaccurate PCA implementations may lower the accuracy of such an approach.

The distributions of the eigenvalues of downsampling, `IncrementalPCA` (*sklearn*), and `sgd` (*OnlinePCA.jl*) also differ from those of the other implementations (Figure 6).

*Calculation time, memory usage, and scalability*
We compared the computational time and memory usage of all the PCA implementations (Figure 7). For the BrainSpinalCord dataset, downsampling itself was faster than most of the PCA implementations, but other preprocessing steps, such as matrix transposition and multiplication of the transposed data matrix and loading vectors to calculate PCs, were slow and had high memory space requirements (Additional file 3). For the Brain dataset, downsampling became slower than most of the PCA implementations, and such a tendency is noticeable as the size of the data matrix increases, because downsampling is based on the full-rank SVD in LAPACK.

We also found that the calculation time of `PCA` (*dask-ml*) was not as fast in spite of its out-of-core implementation; for the BrainSpinalCord and Brain datasets, this implementation could not finish the calculation within three days in our computational environment. The other out-of-core PCA implementations, such as `IncrementalPCA` (*sklearn*), `orthiter/gd/sgd/halko/algorithm971` (*OnlinePCA.jl*), and `oocPCA_CSV` (*oocRPCA*), were able to finish those calculations.

We also systemically estimated the calculation time, memory usage, and scalability of all the PCA implementations using 18 synthetic datasets consisting of $\{10^2, 10^3, 10^4\}$ gene $\times$ $\{10^2, 10^3, 10^4, 10^5, 10^6, 10^7\}$ cell matrices (see Materials and methods). We evaluated whether the calculations could be finished or were interrupted by out-of-memory errors (Figure 1b). We also manually terminated a PCA process that was unable to generate output files within three days (i.e., *dask-ml*). All the terminated jobs are summarized in Additional file 15. To evaluate only the scalability and computability, we set the number of epochs (also known as passes) in `orthiter/gd/sgd` (*OnlinePCA.jl*) to one. However, in actual data analysis, a value several times larger should be used.

Figures 8 and 9 show the calculation time and the memory usage of all the PCA implementations, which can be scaled to a $10^4 \times 10^7$ matrix. `IncrementalPCA` (*sklearn*) and `oocPCA_CSV` (*oocRPCA*) were slightly slower than the other implementations (Figure 8), and this was probably because the inputs of these implementations were CSV files while the other implementations used compressed binary

files (Zstd). The memory usage of all the implementations were almost the same, except for `IncrementalPCA` (*sklearn*) and `oocPCA_CSV` (*oocRPCA*). `oocPCA_CSV` (*oocR-PCA*) has a parameter that controls the maximum memory usage (*mem*), and we set the value to 10 GB (Additional file 3). Indeed, the memory usage had converged to around 10 GB (Figure 9). This property is considered an advantage of this implementation; users can specify a different value to suit their computational environment.

### The relationship between file format and performance

We also counted the passes of the Brain matrix in the out-of-core implementations such as `oocPCA_CSV` (R, *oocRPCA*), `IncrementalPCA` (Python, *sklearn*), and `orthiter/gd/sgd/halko/algorithm971` (Julia, *OnlinePCA.jl*) (Figure 10a). In the `oocPCA_CSV` (R, *oocRPCA*), `IncrementalPCA` (Python, *sklearn*), the data matrix was passed to these function as the CSV format and in the other out-of-core implementations, the data matrix was firstly binarized and compressed in the Zstd file format. We found that the calculation time was correlated with the number of passes of the implementation. Furthermore, binarizing and data compression substantially accelerated the calculation time. This suggests that the data loading process is very critical for out-of-core implementation and that the overhead for this process has a great effect on the overall calculation time and memory usage.

Accordingly, using different data formats, such as CSV, Zstd, Loom [93], and hierarchical data format 5 (HDF5), provided by the 10X Genomics (10X-HDF5) for the Brain dataset, we evaluated the calculation time and the memory usage for the simple one-pass orthogonal iteration (qr(XW)), where qr is the QR decomposition, $X$ is the data matrix, and $W$ represents the 30 vectors to be estimated as the eigenvectors (Figure 10b). For this algorithm algorithm, incremental loading of large block matrices (e.g., 5000 rows) from a sparse matrix was faster than incremental loading of row vectors from a dense matrix, although the memory usage of the former was lower.

While it is not obvious that the usage of a sparse matrix accelerates the PCA with scRNA-seq datasets because scRNA-seq datasets are not particularly sparse compared with data from other fields (cf. recommender systems or social networks [128, 129]), we showed that it has the potential to speed up the calculation time for scRNA-seq datasets.

When all row vectors stored in 10X-HDF5 are loaded at once, the calculation is fastest, but the memory usage is also highest. Because the calculation time and the memory usage have a trade-off and the user's computational environment is not always high-spec, the block size should be optionally specified as a command argument. For the above reasons, we also developed `tenxpca`, which is a new implementation that performs Li's method for a sparse matrix stored in the 10X-HDF5 format. Using all the genes in the CSC matrix incrementally, `tenxpca` was able to finish the calculation in 1.3 hours with a maximum memory usage of 83.0 GB. This is the fastest analysis of the Brain dataset in this study.

In addition to `tenxpca`, some algorithms used in this benchmarking, such as orthogonal iteration, GD, SGD, Halko's method, and Li's method, are implemented as Julia functions and command line tools, which have been published as a Julia

package *OnlinePCA.jl* (Figure 11). When data are stored as a CSV file, they are binarized and compressed in the Zstd file format (Figure 11a) and then some out-of-core PCA implementations are performed. When data are in 10X-HDF5 format, Li's method is directly performed with the data by `tenxpca` (Figure 11b). We also implemented some functions and command line tools to extract row-wise/column-wise statistics such as mean and variance as well as highly variable genes (HVGs) [130] in an out-of-core manner. Because such statistics are saved as small vectors, they can be loaded by any programming language without out-of-core implementation and used for QC, and the users can select only informative genes and cells. After QC, the filtering command removes low-quality genes/cells and generates another Zstd file.

## Discussion

### Guidelines for users

Based on all the benchmarking results and our implementation in this work, we propose some user guidelines (Figure 12). Considering that bioinformatics studies combine multiple tools to construct a user's specific workflow, the programming language is an important factor in selecting the right PCA implementation. Therefore, we categorized the PCA implementations according to language (i.e., R [111], Python [112], and Julia [113]; Figure 12, column-wise). In addition to the data matrix size, we also categorized implementations according to the way they load data (in-memory or out-of-core) as well as their input matrix format (dense or sparse, Figure 12, row-wise). Here, we define the GC-value of a data matrix as the number of genes $\times$ the number of cells.

If the data matrix is not too large (e.g., GC $\leq 10^7$), the data matrix can be loaded as a dense matrix, and full-rank SVD in LAPACK is then accurate and optimal (in-memory & dense matrix). In such a situation, the wrapper functions for the full-rank SVD written in each language are suitable. However, if the data matrix is much larger (e.g., GC $\geq 10^8$), an alternative to the full-rank SVD is needed. Based on the benchmarking results, we recommend IRLBA, IRAM, Halko's method, and Li's method as alternatives to the full-rank SVD. For intermediate GC-values ($10^8 \leq$ GC $\leq 10^{10}$), if the data matrix can be loaded into memory as a sparse matrix, some implementations for these algorithms are available (in-memory & sparse matrix). In particular, such implementations are effective for large data matrices stored in 10X-HDF5 format using CSC format. Seurat2 [49] also introduces this approach by combining the matrix market format (R, *Matrix*) and `irlba` function (R, *irlba*). When the data matrix is dense and cannot be loaded into memory space (e.g., GC $\geq 10^{10}$), the out-of-core implementations, such as `oocPCA_CSV` (R, *oocRPCA*), `IncrementalPCA` (Python, *sklearn*), and `algorithm971` (Julia, *OnlinePCA.jl*), are useful (dense matrix & out-of-core). If the data matrix is extremely large and cannot be loaded into memory even if the data are formatted as a sparse matrix, out-of-core PCA implementations for sparse matrix are needed. Actually, R cannot load the Brain dataset, even if the data is formatted as a sparse matrix (https://github.com/satijalab/seurat/issues/1644). Hence, in such a situation, `tenxpca` can be used if the data is stored in the 10X-HDF5 format.

The PCA implementations examined in this work are affected by various parameters. For example, in `gd` and `sgd` (*OnlinePCA.jl*), the result is sensitive to the

value of learning parameters and the number of epochs. Therefore, a grid-search of such parameters is necessary (Additional file 17). When using `IncrementalPCA` (*sklearn*), the user specifies the chunk size of the input matrix, and a larger value slightly improves the accuracy of PCA (Additional file 16) and the calculation time (Figure 8), although there is a trade-off between these properties and memory usage (Figure 9). Both Halko's method and Li's method have a parameter for specifying the number of power iterations (*niter*), and this iteration step sharpens the distribution of eigenvalues and enforces a more rapid decay of singular values ( [114] and Additional file 3). In our experiments, the value of *niter* is critical for achieving accuracy, and we highly recommend a *niter* value of three or larger (Additional file 18). In some implementations, the default values of the parameters are specified as inappropriate values or cannot be accessed as a function parameter. Therefore, users should carefully set the parameter or select an appropriate implementation.

### Guidelines for developers

We have also established guidelines for developers. Many technologies such as data formats, algorithms, and computational frameworks and environments are available for developing fast, memory-efficient, and scalable PCA implementations (Additional file 19). Here, we focus on two topics.

The first topic is "loss of sparsity." As described above, the use of a sparse matrix can effectively reduce memory space and accelerate calculation, but developers must be careful not to destroy the sparsity of a sparse matrix. PCA with a sparse matrix is not equivalent to SVD with a sparse matrix; in PCA, all sparse matrix elements must be centered by the subtraction of gene-wise average values. Once the sparse matrix $X$ is centered ($X - X_{\mathrm{mean}}$), where $X_{\mathrm{mean}}$ has gene-wise average values as column vectors, it becomes a dense matrix and the memory usage is significantly increased. Obviously, the explicit calculation of the subtraction described above should be avoided. In such a situation, if multiplication of this centered matrix and a dense vector/matrix is required, the calculation should be divided into two parts, such as $(X - X_{\mathrm{mean}})W = XW - X_{\mathrm{mean}}W$, where $W$ represents the vectors to be estimated as eigenvectors, and these parts should be calculated separately. If one or both parts require more than the available memory space, such parts should be incrementally calculated in an out-of-core manner. There are actually some PCA implementations that can accept a sparse matrix, but they may require very long calculation times and large memory space because of a loss of sparsity (cf. `rpca` of *rsvd* https://github.com/cran/rsvd/blob/7a409fe77b220c26e88d29f393fe12a20a5f24fb/R/rpca.R#L158). To our knowledge, only `prcomp_irlba` in *irlba* (https://github.com/bwlewis/irlba/blob/8aa970a7d399b46f0d5ad90fb8a29d5991051bfe/R/irlba.R#L379), `irlb` in *Cell Ranger* (https://github.com/10XGenomics/cellranger/blob/e5396c6c444acec6af84caa7d3655dd33a162852/lib/python/cellranger/analysis/irlb.py#L118), `safe_sparse_dot` in *sklearn* (https://scikit-learn.org/stable/modules/generated/sklearn.utils.extmath.safe_sparse_dot.html), and `tenxpca` in *OnlinePCA.jl* (https://github.com/rikenbit/OnlinePCA.jl/blob/c95a2455acdd9ee14f8833dc5c53615d5e24b5f1/src/tenxpca.jl#L183) deal with this issue. Likewise, as an alternative to the centering calculation, `MaxAbsScaler` in *sklearn* (https://scikit-learn.org/

stable/modules/generated/sklearn.preprocessing.MaxAbsScaler.html) introduces a scaling method in which the maximum absolute value of each gene vector becomes one, thereby avoiding the loss of sparsity.

The second topic is "lazy loading." The out-of-core PCA implementations used in this benchmarking explicitly calculate centering, scaling, and all other relevant arithmetic operations from the extracted blocks of the data matrix. However, to reduce the complexity of the source code, it is desirable to calculate such processes as if the matrix was in memory and only when the data are actually required, so the processes are lazily evaluated on the fly. Some packages, such as DeferredMatrix in *BiocSingular* (R/Bioconductor, https://bioconductor.org/packages/devel/bioc/html/BiocSingular.html), *CenteredSparseMatrix* (Julia, https://github.com/jsams/CenteredSparseMatrix), *Dask* [115] (Python, https://dask.org), and *Vaex* (Python, https://vaex.io/), support lazy loading.

### Future perspective

In this benchmarking study, we found that PCA implementations based on full-rank SVD are accurate but cannot be scaled for use with high-throughput scRNA-seq datasets such as the BrainSpinalCord and Brain datasets, and alternative implementations are thus required. Some methods approximate this calculation by using truncated SVD forms that are sufficiently accurate as well as faster and more memory-efficient than full-rank SVD. The actual memory usage highly depends on whether an algorithm is implemented as out-of-core and whether sparse matrix can be specified as input. Some sophisticated implementations, including our *OnlinePCA.jl*, can handle such issues. Other PCA algorithms, such as downsampling and SGD, are actually not accurate, and their use risks overlooking cellular subgroups contained within scRNA-seq datasets. These methods commonly update eigenvectors with small fractions of the data matrix, and this process may overlook subgroups or subgroup-related gene expression, thereby causing the observed inaccuracy. Our literature review, benchmarking, special implementation for scRNA-seq datasets, and guidelines provide important resources for new users and developers tackling the UML of high-throughput scRNA-seq.

Although the down-stream analyses of PCA vary widely, and we could not examine all the topics of scRNA-seq analyses, such as rare cell-type detection [59, 60] and pseudotime analysis [13, 62–66], differences among PCA algorithms might also affect the accuracy of such analyses. Butler *et al.* showed batch effect removal can be formalized as canonical correlation analysis (CCA) [49], which is mathematically very similar to PCA. The optimization of CCA is also formalized in various ways, including randomized CCA [131] or SGD of CCA [132].

This work also sheds light on the effectiveness of randomized SVD. This algorithm is popular in population genetic studies [110]. In the present study, we also assessed its effectiveness with scRNA-seq datasets with high heterogeneity. This algorithm is relatively simple and some studies have implemented it from scratch (Table 1). Simplicity may be the most attractive feature of this algorithm.

There are also many focuses of recent PCA algorithms (Additional file 19). The randomized subspace iteration algorithm, which is a hybrid of Krylov and Rand methodologies, was developed based on randomized SVD [133, 134]. In pass-efficient

or one-pass randomized SVD, some tricks to reduce the number of passes have been considered [135, 136]. TeraPCA, which is a software tool for use in population genetics studies, utilizes the Mailman algorithm to accelerate the expectation–maximization algorithms for PCA [137, 138]. Townes *et al.* recently proposed the use of PCA for generalized linear models (GLM-PCA) and unified some PCA topics, such as log-transformation, size factor normalization, non-normal distribution, and feature selection, in their GLM framework [139, 140]. Although such topics are beyond the scope of the present work, the current discussion will be useful for the development and application of such methods above.

## Materials and methods

### Empirical datasets

The gene expression matrix and cell type labels for the PBMCs dataset and the Brain dataset [39] were downloaded from the 10X Genomics website (https://support.10xgenomics.com/single-cell-gene-expression/datasets/pbmc_1k_protein_v3 and https://support.10xgenomics.com/single-cell/datasets/1M_neurons, respectively). The gene expression matrix and cell type labels for the Pancreas dataset [40] and the BrainSpinalCord dataset [41] were retrieved from the GEO database (GSE84133 and GSE110823, respectively). For the Pancreas dataset, only the sample of GSM2230759 was used. The genes of all matrices with zero variance were removed because such genes are meaningless for PCA calculation. We also removed the ERCC RNA Spike-Ins, and the number of remaining genes and cells are summarized in Table 2. Additionally, we investigated the effect of feature selection on clustering accuracy (Additional file 20).

### Simulated datasets

All count datasets were generated by the R `rnbinom` (random number based on a negative binomial distribution) function with shape and rate parameters of 0.4 and 0.3, respectively. Matrices of $\{10^2, 10^3, 10^4\}$ genes $\times \{10^2, 10^3, 10^4, 10^5, 10^6, 10^7\}$ cells were generated.

### Benchmarking procedures

Assuming digital expression matrices of unique molecular identifier (UMI) counts, all the data files, including real and synthetic datasets, were in CSV format. When using the Brain dataset, the matrix stored in 10X-HDF5 format was converted to CSV using our in-house Python script (https://gist.github.com/kokitsuyuzaki/5b6cebcaf37100c8794bdb89c7135fd5).

After being loaded by each PCA implementation, the raw data matrix $X_{\mathrm{raw}}$ was converted to normalized values by count per median (CPMED [141–143]) normalization according to the formula $X_{\mathrm{cpmed}}(i,j) = \frac{X_{\mathrm{raw}}(i,j)}{\sum_{k=1}^{M} X_{\mathrm{raw}}(i,k)} * \mathrm{median}\,(\mathrm{Libsize})$, where $M$ is the number of columns and Libsize is the column-wise sum of counts of $X$. After normalization, $X_{\mathrm{cpmed}}$ was transformed to $X$ by the logarithm-transformation $X = \log_{10}(X_{\mathrm{cpmed}} + 1)$, where $\log_{10}$ is the element-wise logarithm. In all the randomized PCA implementation, random seed was fixed.

When $X_{\mathrm{raw}}$ was extremely large and could not be loaded into the memory space all at once, we prepared two approaches to perform PCA with $X$. When PCA

implementations are `orthiter`, `gd`, `sgd`, `halko`, or `algorithm971` (*OnlinePCA.jl*), each row-vector of $X_{\mathrm{raw}}$ is normalized using the pre-calculated Libsize by the `sumr` command, then log-transformed, and finally used for each of the PCA algorithms. When using other out-of-core PCA implementations such as `IncrementalPCA` (*sklearn*), `oocPCA_CSV` (*oocRPCA*), or `PCA` (*dask-ml*), there is no option to normalize and log-transform each row-vector of $X_{\mathrm{raw}}$, so we first calculated $X_{\mathrm{cpmed}}$ using our in-house Python script (https://gist.github.com/kokitsuyuzaki/5b6cebcaf37100c8794bdb89c7135fd5), which was then used for the input matrix of the PCA implementations.

We also investigated the effect of differences in normalization methods on the PCA results (Additional file 21). When performing each PCA implementation based on the truncated SVD, the number of PCs was specified in advance (Table 2).

Although it is unclear how many cells should be used in downsampling, one empirical analysis [94] suggests that 20,000 to 50,000 cells are sufficient for clustering and detecting subpopulations in the Brain dataset. Thus $50,000/1,300,000 \times 100 = 3.8\%$ of cells were sampled from each dataset and used for the downsampling method. When performing `IncrementalPCA` (*sklearn*), the row-vectors, which match the number of PCs, were extracted until the end of the lines of the files. When performing `irlb`(*Cell Ranger*), the loaded dataset was first converted to a scipy sparse matrix and passed to it because this function supports sparse matrix data stored in 10X-HDF5 format. When performing the benchmark, conversion time and memory usage were also recorded. When performing all the functions of *OnlinePCA.jl*, including `orthiter`/`gd`/`sgd`/`halko`/`algorithm971`, we converted the CSV data to Zstd format, and the calculation time and the memory usage were recorded in the benchmarking for fairness. For `orthiter`, `gd`, and `sgd` (*OnlinePCA.jl*), calculations were performed until they converged (Additional file 17). For all the randomized SVD implementations, the *niter* parameter value was set to 3 (Additional file 18). When performing `oocPCA_CSV`, the users can also use `oocPCA_BIN`, which performs PCA with binarized CSV files. The binarization is performed by the `csv2binary` function, which is also implemented in the *oocRPCA* package. Although data binarization accelerates the calculation time for PCA itself, we confirmed that `csv2binary` is based on in-memory calculation, and in our computing environment, `csv2binary` was terminated by an out-of-memory error. Accordingly, we only used `oocPCA_CSV`, and the CSV files were directly loaded by this function.

### Computational environment

All computations were performed on two-node machines with Intel Xeon E5-2697 v2 (2.70 GHz) processors and 128 GB of RAM, four-node machines with Intel Xeon E5-2670 v3 (2.30 GHz) processors and 96 GB of RAM, and four-node machines with Intel Xeon E5-2680 v3 (2.50 GHz) processors and 128 GB of RAM. Storage among machines was shared by NFS, connected using InfiniBand. All jobs were queued by the Open Grid Scheduler/Grid Engine (v2011.11) in parallel. The elapsed time and maximum memory usage were evaluated using the GNU `time` command (v1.7).

### Reproducibility

All the analyses were performed on the machines described above. We used R v3.5.0, Python v3.6.4, and Julia v1.0.1 in the benchmarking; for t-SNE and CSV conversion

of the Brain dataset, we used Python v2.7.9. The *Sklearn* (Python) package was used to perform *k*-means and GMM clustering methods. The igraph (R), nn2 (R), and Matrix (R) packages were used to perform Louvain clustering (Additional file 6). The hdbscan (Python) package was used to perform HDBScan clustering. The bht-sne (Python) package was used to perform t-SNE. Lastly, the umap (Python) package was used to perform UMAP. All the programs used to perform the PCA implementations in the benchmarking are summarized in Additional file 3. Orthogonal iteration, GD, SGD, Halko's method, and Li's method are implemented as `orthiter`, `gd`, `sgd`, `halko`, and `algorithm971`, respectively, which are the Julia functions or commands for *OnlinePCA.jl* (https://github.com/rikenbit/OnlinePCA.jl). We also published the script files used to perform the benchmarking (https://github.com/rikenbit/onlinePCA-experiments).

### Abbreviations

PCA: principal component analysis; scRNA-seq: single-cell RNA sequencing; sci-RNA-seq: single-cell combinatorial-indexing RNA-sequencing analysis; UML: unsupervised machine learning; QC: quality control; PC: principal component; EVD: eigenvalue decomposition; SVD: singular value decomposition; SimT: similarity transformation-based, DS: downsampling-based, SU: SVD update-based, Krylov: Krylov subspace-based, GD: gradient descent-based,Rand: Random projection-based, *Sklearn*: *scikit-learn*; SKL: sequential Karhunen-Loeve transform; IRLBA: augmented implicitly restarted Lanczos bidiagonalization; IRAM: implicitly restarted Arnoldi method; GD: gradient descent; SGD: stochastic gradient descent; t-SNE: t-stochastic neighbor embedding; UMAP: uniform manifold approximation and projection; FIt-SNE: Fourier transform-accelerated interpolation-based t-stochastic neighbor embedding; oocPCA: out-of-core PCA; GMM: Gaussian mixture model; ARI: adjusted Rand index; Zstd: Zstandard; UMI: unique molecular identifier; CSV: comma-separated values; HDF5: hierarchical data format 5; 10X-HDF5: HDF5 provided by 10X Genomics; CSC: compressed sparse column format; CSR: compressed sparse row format; CCA: canonical correlation analysis; GLM: generalized linear models; CPMED: Count per median; HVGs: highly variable genes

### Competing interests

The authors declare that they have no competing interests.

### Author's contributions

KT and HS surveyed the PCA algorithms and implementations. KT and IN designed the benchmarking test. KT and KS implemented the Julia program and performed all the analyses. KT retrieved and preprocessed the test dataset to evaluate the proposed method. All the authors have written, read, and approved the manuscript.

### Author details

[1]Laboratory for Bioinformatics Research, RIKEN Center for Biosystems Dynamics Research, Japan. [2]Japan Science and Technology Agency, PRESTO, Japan. [3]Department of Applied Mathematics and Physics, Graduate School of Informatics, Kyoto University, Japan. [4]Department of Biotechnology, Graduate School of Agricultural and Life Sciences, The University of Tokyo, Japan. [5]Bioinformatics Course, Master's/Doctoral Program in Life Science Innovation (T-LSI), School of Integrative and Global Majors (SIGMA), University of Tsukuba, Japan.

### References

1. Trapnell, C.: Defining cell types and states with single-cell genomics. Genome Research **25(10)**, 1491–1498 (2015)
2. Macosko, E.Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N., Martersteck, E.M., Trombetta, J.J., Weitz, D.A., Sanes, J.R., Shalek, A.K., Regev, A., McCarroll, S.A.: Highly parallel genome-wide expression profiling of individual cells using nanoliter dropltes. Cell **161**, 1202–1214 (2015)
3. Shekhar, K., Lapan, S.W., Whitney, I.E., Tran, N.M., Macosko, E.Z., Kowalczyk, M., Adiconis, Z., Levin, J.Z., Nemesh, J., Goldman, M., McCarroll, S.A., Cepko, C.L., Regev, A., Sanes, J.R.: Comprehensive classification of retinal bipolar neurons by single-cell transcriptomics. Cell **166**, 1308–1323 (2016)

4. Campbell, J.N., Macosko, E.Z., Fenselau, H., Pers, T.H., Lyubetskaya, A., Tenen, D., Goldman, M., Verstegen, A.M.J., Resch, J.M., McCarroll, S.A., Rosen, E.D., Lowell, B.B., Tsai, L.T.: A molecular census of arcuate hypothalamus and median eminence cell types. Nature Neuroscience **20(3)**, 484–496 (2017)

5. Klein, A.M., Mazutis, L., Akartuna, I., Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., Kirschner, M.W.: Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. Cell **161**, 1187–1201 (2015)

6. Baron, M., Veres, A., Wolock, S.L., Faust, A.L., Gaujoux, R., Vetere, A., Ryu, J.H., Wagner, B.K., Shen-Orr, S.S., Klein, A.M., Melton, D.A., Yanai, I.: A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. Cell Systems **3(4)**, 346–360 (2016)

7. Grun, D., Lyubimova, A., Kester, L., Wiebrands, K., Basak, O., sasaki, N., Clevers, H., Oudenaarden, A.: Single-cell messenger rna sequencing reveals rare intestinal cell types. Nature **525**, 251–255 (2015)

8. Buettner, F., Natarajan, K.N., Casale, F.P., Proserpio, V., Scialdone, A., Theis, F.J., Teichmann, S.A., Marioni, J.C., Stegle, O.: Computational analysis of cell-to-cell heterogeneity in single-cell rna-sequencing data reveals hidden subpopulations of cells. Nature Biotechnology **33(2)**, 155–160 (2015)

9. Durruthy-Durruthy, R., Gottlieb, A., Hartman, B.H., Waldhaus, J., Laske, R.D., Altman, R., Heller, S.: Reconstruction of the mouse otocyst and early neuroblast lineage at single-cell resolution. Cell **157**, 1–15 (2014)

10. Achim, K., Pettit, J.B., Saraiva, L.R., Gavriouchkina, D., Larsson, T., Arendt, D., Marioni, J.C.: High-throughput spatial mapping of single-cell rna-seq data to tissue of origin. Nature Computational Biology **33(5)**, 503–509 (2015)

11. Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., Regev, A.: Spatial reconstruction of single-cell gene expression data. Nature Biotechnology **33(5)**, 495–508 (2015)

12. Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokhare, P., Li, S., Morse, M., Lennon, N.J., Livak, K.J., Mikkelsen, T.S., Rinn, J.L.: The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nature Biotechnology **32**, 381–386 (2014)

13. Qiu, X., Mao, Q., Tang, Y., Wang, L., Chawla, R., Pliner, H.A., Trapnell, C.: Reversed graph embedding resolves complex single-cell trajectories. Nature Methods **14(10)**, 979–982 (2017)

14. Svensson, V., Tormo, R.V., Teichmann, S.A.: Exponential scaling of single-cell rna-seq in the past decade. Nature Protocols **13(4)**, 599–604 (2017)

15. Sasagawa, Y., Danno, H., Takada, H., Ebisawa, M., Tanaka, K., Hayashi, T., Kurisaki, A., Nikaido, I.: Quartz-seq2: a high-throughput single-cell rna-sequencing method that effectively uses limited sequence reads. BMC Genome Biology **19(29)** (2018)

16. Jaitin, D.A., Kenigsberg, E., Keren-Shaul, H., Elefant, N., Paul, F., Zaretsky, I., Mildner, A., Cohen, N., Jung, S., Tanay, A., Amit, I.: Massively parallel single cell rna-seq for marker-free decomposition of tissues into cell types. Science **343(6172)**, 776–779 (2014)

17. Hashimshony, T., Senderovich, N., Avital, G., Klochendler, A., de Leeuw, Y., Anavy, L., Gennert, D., Li, S., Livak, K.L., Rozenblatt-Rosen, O., Dor, Y., Regev, A., Yanai, I.: Cel-seq2: sensitive highly-multiplexed single-cell rna-seq. BMC Genome Biology **17(77)** (2016)

18. Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., Manno, G.L., Juréus, A., Marques, S., Munguba, H., He, L., Betsholtz, C., Rolny, C., Castelo-Branco, G., Hjerling-Leffler, J., Linnarsson, S.: Cell types in the mouse cortex and hippocampus revealed by single-cell rna-seq. Science **347(6226)**, 1138–1142 (2015)

19. Hashimshony, T., Senderovich, N., Avital, G., Klochendler, A., de Leeuw, Y., Anavy, L., Gennert, D., Li, S., Livak, K.J., Rozenblatt-Rosen, O., Dor, Y., Regev, A., Yanai, I.: Cel-seq2: sensitive highly-multiplexed single-cell rna-seq. Genome Biology **17(77)** (2016)

20. Shalek, A.K., Satija, R., Shuga, J., Trombetta, J.J., Gennert, D., Lu, D., Chen, P., Gertner, R.S., Gaublomme, J.T., Yosef, N., Schwartz, S., Fowler, B., Weaver, S., Wang, J., Ding, R., Raychowdhury, R., Friedman, N., Hacohen, N., Park, H., May, A.P., Regev, A.: Single cell rna seq reveals dynamic paracrine control of cellular variation. Nature **510(7505)** (2014)

21. Tasic, B., Menon, V., Nguyen, T.N., Kim, T.K., Jarsky, T., Yao, Z., Levi, B., Gray, L.T., Sorensen, S.A., Dolbeare, T., Bertagnolli, D., Goldy, J., Shapovalova, N., Pary, S., Parry, C., Lee, C., Smith, K., Bernard, A., Madisen, L., Sunkin, S.M., Hawrylycz, M., Koch, C., Zeng, H.: Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. Nature Neuroscience **19(2)**, 335–346 (2016)

22. Zheng, G.X.Y., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J., Gregory, M.T., Shuga, J., Montesclaros, L., Underwood, J.G., Masquelier, D.A., Nishimura, S.Y., Schnall-Levin, M., Wyatt, P.W., Hindson, C.M., Bharadwai, R., Wong, A., Ness, K.D., Beppu, L.W., Deeg, H.J., McFarland, C., Loeb, K.R., Valente, W.J., Ericson, N.G., Stevens, E.A., Radich, J.P., Mikkelsen, T.S., Hindson, B.J., Bielas, J.H.: Massively parallel digital transcriptional profiling of single cells. Nature Communications **8(14049)**, 1–12 (2017)

23. Cao, J., Spielmann, M., Qiu, X., Huang, X., Ibrahim, D.M., Hill, A.J., Zhang, F., Mundlos, S., Christiansen, L., Steemers, F.J., Trapnell, C., Shendure, J.: The single-cell transcriptional landscape of mammalian organogenesis. Nature (2019)

24. Consortium, T.H.: The human cell atlas white paper (2017)

25. Rozenblatt-Rosen, O., Stubbington, M.J.T., Regev, A., Teichmann, S.A.: The human cell atlas: from vision to reality. Nature **550**, 451–453 (2017)

26. Regev, A., Teichmann, S.A., Lander, E.S., Amit, I., Benoist, C., Birney, E., Bodenmiller, B., Campbell, P., Carninci, P., Clatworthy, M., Clevers, H., Deplancke, B., Dunham, I., Eberwine, J., Eils, R., Enard, W., Farmer, A., Fugger, L., Göttgens, B., Hacohen, N., Haniffa, M., Hemberg, M., Kim, S., Klenerman, P., Kriegstein, A., Lein, E., Linnarsson, S., Lundberg, E., Lundeberg, J., Majumder, P., Marioni, J.C., Merad, M., Mhlanga, M., Nawijn, M., Netea, M., Nolan, G., Pe'er, D., Phillipakis, A., Ponting, C.P., Quake, S., Reik, W., Rozenblatt-Rosen, O., Sanes, J., Satija, R., Schumacher, T.N., Shalek, A., Shapiro, E., Sharma, P., Shin, J.W., Stegle, O., Stratton, M., Stubbington, M.J.T., Theis, F.J., Uhlen, M., van Oudenaarden, A., Wagner,

A., Watt, F., Weissman, J., Wold, B., Xavier, R., Yosef, N., Participants, H.C.A.M.: Science forum: The human cell atlas. eLife, 37041 (2017)

27. Han, X., Wang, R., Zhou, Y., Fei, L., Sun, H., Lai, S., Saadatpour, A., Zhou, Z., Chen, H., Ye, F., Huang, D., Xu, Y., Huang, W., Jiang, M., Jiang, X., Mao, J., Chen, Y., Lu, C., Xie, J., Fang, Q., Wang, Y., Yue, R., Li, T., Huang, H., Orkin, S.H., Yuan, G.C., Chen, M., Guo, G.: Mapping the mouse cell atlas by microwell-seq. Cell **172(5)**, 1091–1107 (2018)

28. Consortium, T.T.M.: Single-cell transcriptomics of 20 mouse organs creates a tabula muris. Nature **562(7727)**, 367–372 (2018)

29. Wagner, A., Regev, A., Yosef, N.: Revealing the vectors of cellular identity with single-cell genomics. Nature Biotechnology **34(11)**, 1145–160 (2017)

30. Stegle, O., Teichmann, S.A., Marioni, J.C.: Computational and analytical challenges in single-cell transcriptomics. Nature Reviews Genetics **16(3)**, 133–145 (2015)

31. Bacher, R., Kendziorski, C.: Design and computational analysis of single-cell rna-sequencing experiments. BMC Genome Biology **17(63)** (2016)

32. Poulin, J.F., Tasic, B., Hjerling-Leffler, J., Trimarchi, J.M., Awatramani, R.: Disentangling neural cell diversity using single-cell transcriptomics. Nature Neuroscience **19(9)**, 1131–1141 (2016)

33. Kolodziejczyk, A.A., Kim, J.K., Svensson, V., Marioni, J.C., Teichmann, S.A.: The technology and biology of single-cell rna sequencing. Molecular Cell **58(4)**, 610–620 (2015)

34. Chen, G., Ning, B., Shi, T.: Single-cell rna-seq technologies and related computational data analysis. Frontiers in Genetics **10(317)** (2019)

35. Stuart, T., Satija, R.: Integrative single-cell analysis. Nature Reviews Genetics **20(5)**, 257–272 (2019)

36. Pearson, K.: On lines and planes of closest fit to systems of points in space. Philosophical Magazine **2(11)**, 559–572 (1901)

37. Hotelling, H.: Analysis of a complex of statistical variables into principal components. Journal of Educational Psychology **24**, 417–441 (1933)

38. Broa, R., K, S.A.: Principal component analysis. Royal Society of Chemistry **6(2812)**, 2812–2831 (2014)

39. Genomics, X.: 1.3 Million Brain Cells from E18 Mice. https://support.10xgenomics.com/single-cell/datasets/1M_neurons

40. Baron, M., Veres, A., Wolock, S.L., Faust, A.L., Gaujoux, R., Vetere, A., Ryu, J.H., Wagner, B.K., Shen-Orr, S.S., Klein, A.M., Melton, D.A., Yanai, I.: A single-cell transcriptomic map of the human and mouse pancreas reveals inter- and intra-cell population structure. Cell Systems **3(4)**, 346–360 (2016)

41. Rosenberg, A.B., Roco, C.M., Muscat, R.A., Kuchina, A., Sample, P., Yao, Z., Graybuck, L.T., Peeler, D.J., Mukherjee, S., Chen, W., Pun, S.H., Sellers, D.L., Tasic, B., Seelig, G.: Single-cell profiling of the developing mouse brain and spinal cord with split-pool barcoding. Science **360(6385)**, 176–182 (2018)

42. Cole, M.B., Risso, D., Wagner, A., DeTomaso, D., Ngai, J., Purdom, E., Dudoit, S., Yosef, N.: Performance assessment and selection of normalization procedures for single-cell rna-seq. Cell Systems **8(4)**, 315–328 (2019)

43. Taguchi, Y.-H.: Principal component analysis-based unsupervised feature extraction applied to single-cell gene expression analysis. In: 14th International Conference, ICIC 2018, pp. 816–826 (2018). China

44. Lin, Z., Yang, C., Zhu, Y., Duchi, J., Fu, Y., Wang, Y., Jiang, B., Zamanighomi, M., Xu, X., Li, M., Sestan, N., Zhao, H., Wong, W.H.: Simultaneous dimension reduction and adjustment for confounding variation. PNAS **113(51)**, 14662–14667 (2016)

45. Lasrado, R., Boesmans, W., Kleinjung, J., Pin, C., Bell, D., Bhaw, L., McCallum, S., Zong, H., Luo, L., Clevers, H., Vanden, B.P., Pachnis, V.: Lineage-dependent spatial and functional organization of the mammalian enteric nervous system. Science **356(6339)**, 722–726 (2017)

46. Wagner, F.: Go-pca: An unsupervised method to explore gene expression data using prior knowledge. PLOS ONE **10(11)**, 0143196 (2015)

47. Cerosaletti, K., Barahmand-Pour-Whitman, F., Yang, J., DeBerg, H.A., Dufort, M.J., Murray, S.A., Israelsson, E., Speake, C., Gersuk, V.H., Eddy, J.A., Reijonen, H., Greenbaum, C.J., Kwok, W.W., Wambre, E., Prlic, M., Gottardo, R., Nepom, G.T., Linsley, P.S.: Single-cell rna sequencing reveals expanded clones of islet antigen-reactive cd4+ t cells in peripheral blood of subjects with type 1 diabetes. Journal of Immunology **199(1)**, 323–325 (2017)

48. Li, J., Klughammer, J., Farlik, M., Penz, T., Spittler, A., Barbieux, C., Berishvili, E., Bock, C., Kubicek, S.: Single-cell transcriptomes reveal characteristic features of human pancreatic islet cell types. EMBO Reports **17(2)**, 178–187 (2016)

49. Butler, H.P, A., Smibert, P., Papalexi, E., Satija, R.: Integrated analysis of single cell transcriptomic data across conditions, technologies, and species. Nature Biotechnology **36**, 411–420 (2018)

50. Lun, A.T., McCarthy, D.J., Marioni, J.C.: A step-by-step workflow for low-level analysis of single-cell rna-seq data with bioconductor. F1000Research **Version2** (2016)

51. Ilicic, T., Kim, J.K., Kolodziejczyk, A.A., Bagger, F.O., McCarthy, D.J., Marioni, J.C., Teichmann, S.A.: Classification of low quality cells from single-cell rna-seq data. BMC Genome Biology **17(29)** (2016)

52. Dijk, D., Sharma, R., Nainys, J., Yim, K., Kathail, P., Carr, A.J., Burdziak, C., Moon, K.R., Chaffer, C.L., Pattabiraman, D., Bierie, B., Mazutis, L., Wolf, G., Krishnaswamy, S., Pe'er, D.: Recovering gene interactions from single-cell data using data diffusion. Cell **174(3)**, 716–729 (2018)

53. Li, W.V., Li, J.J.: An accurate and robust imputation method scimpute for single-cell rna-seq data. Nature Communication **9(997)** (2018)

54. Gong, W., Kwak, I.Y., Pota, P., Koyano-Nakagawa, N., Garry, D.J.: Drimpute: imputing dropout events in single cell rna sequencing data. BMC Bioinformatics **19(220)** (2018)

55. Büttner, M., Miao, Z., Wolf, F.A., Teichmann, S.A., Theis, F.J.: A test metric for assessing single-cell rna-seq batch correction. Nature methods **16(1)**, 43–49 (2019)

56. Shaham, U., Stanton, K.P., Zhao, J., Li, H., Raddassi, K., Montgomery, R., Kluger, Y.: Removal of batch effects using distribution-matching residual networks. Bioinformatics **33(16)**, 2539–2546 (2017)

57. Korsunsky, I., Fan, J., Slowikowski, K., Zhang, F., Wei, K., Baglaenko, Y., Brenner, M., Loh, P.-R., Raychaudhuri, S.: Fast, sensitive, and accurate integration of single cell data with harmony. bioRxiv (2018). doi:10.1101/461954

58. Scialdone, A., Natarajan, K.N., Saraiva, L.R., Proserpio, V., Teichmann, S.A., Stegle, O., Marioni, J.C., Buettner, F.: Computational assignment of cell-cycle stage from single-cell transcriptome data. Methods **85**, 54–61 (2015)

59. Tsoucas, D., Yuan, G.C.: Giniclust2: a cluster-aware, weighted ensemble clustering method for cell-type detection. BMC Genome Biology **19(1)** (2018)

60. Herman, J.S., Sagar, Grün, D.: Fateid infers cell fate bias in multipotent progenitors from single-cell rna-seq data. Nature methods **15**, 379–386 (2018)

61. Sato, K., Tsuyuzaki, K., Shimizu, K., Nikaido, I.: Cellfishing.jl: an ultrafast and scalable cell search method for single-cell rna sequencing. BMC Genome Biology **20(1)** (2019)

62. Diaz, A., Liu, S.J., Sandoval, C., Pollen, A., Nowakowski, T.J., Lim, D.A., Kriegstein, A.: Scell: integrated analysis of single-cell rna-seq data. Bioinformatics **32(14)**, 2219–2220 (2016)

63. Ji, Z., Ji, H.: Tscan: Pseudo-time reconstruction and evaluation in single-cell rna-seq analysis. Nucleic Acids Research **44(13)** (2016)

64. Shin, J., Berg, D.A., Zhu, Y., Shin, J.Y., Song, J., Bonaguidi, M.A., Enikolopov, G., Nauen, D.W., Christian, K.M., Ming, G.L., Song, H.: Single-cell rna-seq with waterfall reveals molecular cascades underlying adult neurogenesis. Cell Stem Cell **17(3)**, 360–372 (2015)

65. Street, K., Risso, D., Fletcher, R.B., Das, D., Ngai, J., Yosef, N., Purdom, E., Dudoit, S.: Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics. BMC Genomics **19(477)** (2018)

66. Campbell, K.R., Yau, C.: Probabilistic modeling of bifurcations in single-cell gene expression data using a bayesian mixture of factor analyzers. Wellcome Open Research **2(19)** (2017)

67. Maaten, L., Hinton, G.: Visualizing data using t-sne. Journal of Machine Learning Reseach, 2579–2605 (2008)

68. Maaten, L.: Accelerating t-sne using tree-based algorithms. Journal of Machine Learning Reseach, 3221–3245 (2014)

69. Linderman, G.C., Rachh, M., Hoskins, J.G., Steinerberger, S., Kluger, Y.: Fast interpolation-based t-sne for improved visualization of single-cell rna-seq data. Nature methods **16**, 243–245 (2019)

70. Lawrence, N.D.: Gaussian process latent variable models for visualisation of high dimensional data. In: In NIPS, p. 2004 (2003)

71. McInnes, L., Healy, J., Saul, N., Großberger, L.: Umap: Uniform manifold approximation and projection for dimension reduction. The Journal of Open Source Software **3(29)**, 861 (2018)

72. Becht, E., McInnes, L., Healy, J., Dutertre, C.A., Kwok, I.W.H., Ng, L.G., Ginhoux, F., Newell, E.W.: Dimensionality reduction for visualizing single-cell data using umap. Nature Biotechnology **37**, 38–44 (2019)

73. Weinreb, C., Wolock, S., Klein, A.M.: Spring: a kinetic interface for visualizing high dimensional single-cell expression data. Bioinformatics **34(7)**, 1246–1248 (2018)

74. Kiselev, V.Y., Kirschner, K., Schaub, M.T., Andrews, T., Yiu, A., Chandra, T., Natarajan, K.N., Reik, W., Barahona, M., Green, A.R., Hemberg, M.: Sc3: consensus clustering of single-cell rna-seq data. Nature methods **14(5)**, 483–486 (2017)

75. Wang, B., Zhu, J., Pierson, E., Ramazzotti, D., Batzoglou, S.: Visualization and analysis of single-cell rna-seq data by kernel-based similarity learning. Nature methods **14(4)**, 414–416 (2017)

76. Yang, Y., Huh, R., Culpepper, H.W., Lin, Y., Love, M.I., Li, Y.: Safe-clustering: Single-cell aggregated (from ensemble) clustering for single-cell rna-seq data. Bioinformatics (2018)

77. Zurauskiene, J., Yau, C.: pcareduce: hierarchical clustering of single cell transcriptional profiles. BMC Bioinformatics **17(140)** (2016)

78. Wagner, A., Regev, A., Yosef, N.: Revealing the vectors of cellular identity with single-cell genomics. Nature Biotechnology **34(11)**, 1145–1160 (2016)

79. Andrews, T.S., Hemberg, M.: Identifying cell populations with scrnaseq. Molecular Aspects of Medicine **59**, 114–122 (2018)

80. Kiselev, V.Y., Andrews, T.S., Hemberg, M.: Challenges in unsupervised clustering of single-cell rna-seq data. Nature Reviews **20(5)**, 273–282 (2019)

81. Oskolkov, N.: How to cluster in High Dimensions. https://towardsdatascience.com/how-to-cluster-in-high-dimensions-4ef693bacc6

82. McCarthy, D.J., Campbell, K.R., Lun, A.T., Wills, Q.F.: Scater: pre-processing, quality control, normalization and visualization of single-cell rna-seq data in r. Bioinformatics **33(8)**, 1179–1186 (2017)

83. Jenkins, D., Faits, T., Khan, M.M., Briars, E., Carrasco, P.S., Johnson, W.E.: singleCellTK: Interactive Analysis of Single Cell RNA-Seq Data. https://bioconductor.org/packages/release/bioc/html/singleCellTK.html (2018)

84. Tian, L., Su, S., Dong, X., Amann-Zalcenstein, D., Biben, C., Seidi, A., Hilton, D.J., Naik, S.H., Ritchie, M.E.: scpipe: A flexible r/bioconductor preprocessing pipeline for single-cell rna-sequencing data. PLOS Computational Biology **14(8)**, 1006361 (2018)

85. Yip, S.H., Wang, P., Kocher, J.A., Sham, P.C., Wang, J.: Linnorm: improved statistical analysis for single cell rna-seq expression data. Nucleic Acids Research **45(22)**, 179 (2017)

86. Finak, G., McDavid, A., Yajima, M., Deng, J., Gersuk, V., Shalek, A.K., Slichter, C.K., Miller, H.W., McElrath, M.J., Prlic, M., Linsley, P.S., Gottardo, R.: Mast: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell rna sequencing data. BMC Genome Biology **16(278)** (2015)

87. Demsar, J., Curk, T., Erjavec, A., Gorup, C., Hocevar, T., Milutinovic, M., Mozina, M., Polajnar, M., Toplak, M., Staric, A., Stajdohar, M., Umek, L., Zagar, L., Zbontar, J., Zitnik, M., Zupan, B.: Orange: Data mining toolbox in python. Journal of Machine Learning Research, 2349–2353 (2013)

88. Zhu, X., Wolfgruber, T.K., Tasato, A., Arisdakessian, C., Garmire, D.G., Garmire, L.X.: Granatum: a graphical single-cell rna-seq analysis pipeline for genomics scientists. BMC Genome Medicine **9(108)** (2017)

89. Azizi, E., Carr, A.J., Plitas, G., Cornish, A.E., Konopacki, C., Prabhakaran, S., Nainys, J., Wu, K., Kiseliovas, V., Setty, M., Choi, K., Fromme, R.M., Dao, P., McKenney, P.T., Wasti, R.C., Kadaveru, K., Mazutis, L., Rudensky, A.Y., Pe'er, D.: Single-cell map of diverse immune phenotypes in the breast tumor microenvironment. Cell **5(23)**, 1293–1308 (2018)

90. Golub, G.H., Loan, C.F.V.: Matrix Computations (Johns Hopkins Studies in the Mathematical Sciences), Fourth Edition. Johns Hopkins University Press, Baltimore (2012)

91. Senabouth, A., Lukowski, S., Alquicira, J., Andersen, S., Mei, X., Nguyen, Q., Powell, J.: ascend: R package for analysis of single cell rna-seq data. GigaScience **8(8)**, 087 (2019)

92. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondl, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in python. The Journal of Machine Learning Research **12**, 2825–2830 (2011)

93. Wolf, F.A., Angerer, P., Theis, F.J.: Scanpy: large-scale single-cell gene expression data analysis. BMC Genome Biology **19(15)** (2018)

94. Bhaduri, A., Nowakowski, T.J., Pollen, A.A., Kriegstein, A.R.: Identification of cell types in a mouse brain single-cell atlas using low sampling coverage. BMC Biology (2018)

95. Levy, A, M, K.: Sequential karhunen-loeve basis extraction and its application to images. IEEE Transactions on Image Processing **9(8)**, 1371–1374 (2000)

96. Bai, Z., Demmel, J., Dongarra, J., Ruhe, A., Vorst, H.V.D.: Templates for the Solution of Algebraic Eigenvalue Problems, A Practical Guide. Society for Industrial and Applied Mathematics, Philadelphia (1987)

97. Lehoucq, R., Maschhoff, K., Sorensen, D., Yang, C.: ARPACK SOFTWARE. https://www.caam.rice.edu/software/ARPACK/

98. Qiu, Y.: Spectra: C++ Library For Large Scale Eigenvalue Problems. https://spectralib.org

99. Larsen, R.M.: PROPACK homepage. http://sun.stanford.edu/~rmunk/PROPACK/

100. Baglama, J., Reichel, L.: Augmented implicitly restarted lanczos bidiagonalization methods. SIAM Journal on Scientific Computing **27(1)**, 19–42 (2005)

101. Lehoucq, R.B., Sorensen, D.C., Yang, C.: Arpack users' guide: Solution of large-scale eigenvalue problems with implicitly restarted arnoldi methods (1997)

102. Chen, J., Noack, A., Edelman, A.: Fast computation of the principal components of genotype matrices in julia. arXiv (2018). doi:arXiv:1808.03374v1

103. Balzano, L., Chi, Y., Lu, Y.M.: Streaming pca and subspace tracking: The missing data case. Proceedings of the IEEE **106(8)**, 1293–1310 (2018)

104. Oja, E.: A simplified neuron model as a principal component analyzer. Journal of Mathematical Biology **15**, 267–273 (1982)

105. Oja, E., Karhunen, J.: On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix author links open overlay panel. Journal of Mathematical Analysis and Applications **106(1)**, 69–84 (1985)

106. Oja, E.: Principal components, minor components, and linear neural networks. Neural Networks **5**, 927–935 (1992)

107. Halko, N., Martinsson, P.G., Tropp, J.A.: Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. SIAM Rev., Survey and Review **53(2)**, 217–288 (2011)

108. Halko, N., Martinsson, P.G., Shkolnisky, Y., M, T.: An algorithm for the principal component analysis of large data sets. SIAM Journal on Scientific Computing **33(5)**, 2580–2594 (2011)

109. Li, H., C, L.G., Szlam, A., Stanton, K.P., Kluger, Y., Tygert, M.: Algorithm 971: An implementation of a randomized algorithm for principal component analysis. ACM Transactions on Mathematical Software **43(3)** (2017)

110. Abraham, G., Inouye, M.: Fast principal component analysis of large-scale genome-wide data. PLOS ONE **9(4)**, 93766 (2014)

111. Ihaka, R., Gentleman, R.: R: A language for data analysis and graphics. Journal of Computational and Graphical Statistics **5(3)**, 299–314 (1996)

112. Rossum, G.: Python reference manual. Technical Report (1995)

113. Perkel, J.M.: Julia: come for the syntax, stay for the speed. Nature **572(7767)**, 141–142 (2019)

114. Erichson, N.B., Voronin, S., Brunton, S.L., Kutz, J.N.: Randomized matrix decompositions using r. Journal of Statistical Software **89(11)** (2019). doi:10.18637/jss.v089.i11

115. Rocklin, M.: Dask: Parallel computation with blocked algorithms and task scheduling. In: Huff, K., Bergstra, J. (eds.) Proceedings of the 14th Python in Science Conference, pp. 130–136 (2015)

116. Lacono, G., Mereu, E., Guillaumet-Adkins, A., Corominas, R., Cusco, I., Rodriguez-Esteban, G., Gut, M., Perez-Jurado, L.A., Gut, I., Heyn, H.: bigscale: an analytical framework for big-scale single-cell data. Genome Research **28(6)**, 878–890 (2018)

117. Aibar, S., Gonzalez-Blas, C.B., Moerman, T., Huynh-Thu, V.A., Imrichova, H., Hulselmans, G., Rambow, F., Marine, J.-C., Geurts, P., Aerts, J., Oord, J., Atak, Z.K., Wouters, J., Aerts, S.: Scenic: single-cell regulatory network inference and clustering. Nature methods **14**, 1083–1086 (2017)

118. Kisekev, V.Y., Yiu, A., Hemberg, M.: scmap: projection of single-cell rna-seq data across data sets. Nature methods **15**, 359–362 (2018)

119. Huang, M., Wang, J., Torre, E., Dueck, H., Shaffer, S., Bonasio, R., Murray, J.I., Raj, A., Li, M., Zhang, N.R.: Saver: gene expression recovery for single-cell rna sequencing. Nature methods **15**, 539–542 (2018)

120. Wang, D., Gu, J.: Vasc: Dimension reduction and visualization of single-cell rna-seq data by deep variational autoencoder. Genomics, Proteomics & Bioinformatics **16(5)**, 320–331 (2018)

121. Ding, J., Condon, A., Shah, S.P.: Interpretable dimensionality reduction of single cell transcriptome data with deep generative models. Nature Communications **2002** (2018)

122. Pliner, H.A., Shendure, J., Trapnell, C.: Supervised classification enables rapid annotation of cell atlases. Nature Methods (2019)

123. Weber, L.M., Saelens, W., Cannoodt, R., Soneson, C., Hapfelmeier, A., Gardner, P.P., Boulesteix, A.-L.,

Saeys, Y., Robinson, M.D.: Essential guidelines for computational method benchmarking. BMC Genome Biology **20(125)** (2019)

124. Hubert, L., Arabie, P.: Comparing partitions. Journal of Classification **2(1)**, 193–218 (1985)
125. Blondel, V.D., Guillaume, J.-L., Lambiotte, R., Lefebvre, E.: Fast unfolding of communities in large networks. arXiv (2008). doi:arXiv:0803.0476v2
126. Bishop, C.M.: Pattern Recognition and Machine Learning (Information Science and Statistics). Springer, New York City (2006)
127. Benson, A.R., Gleich, D.F., Demmel, J.: Direct qr factorizations for tall-and-skinny matrices in mapreduce architectures. Proceedings of the IEEE International Conference on Big Data (2013). doi:10.1109/BigData.2013.6691583
128. Koren, Y., Bell, R., Volinsky, C.: Matrix factorization techniques for recommender systems. IEEE Computer **42(8)**, 30–37 (2009)
129. Davis, T.: University of Florida Sparse Matrix Collection. https://sparse.tamu.edu
130. Yip, S.H., Sham, P.C., J, W.: Evaluation of tools for highly variable gene discovery from single-cell rna-seq data. Briefing in Bioinformatics, 011 (2018)
131. Mineiro, P., Karampatziakis, N.: A randomized algorithm for cca. arXiv (2014). doi:arXiv:1411.3409v1
132. Arora, R., Cotter, A., Livescu, K., Srebro, N.: Stochastic optimization for pca and pls. In: 2012 50th Annual Allerton Conference on Communication, Control, and Computing (Allerton), pp. 861–868 (2012)
133. Bose, A., Kalantzis, V., Kontopoulou, E., Elkady, M., Paschou, P., Drineas, P.: Terapca: a fast and scalable software package to study genetic variation in tera-scale genotypes. Bioinformtaics **btz157** (2019)
134. Musco, C., Musco, C.: Randomized block krylov methods for stronger and faster approximate singular value decomposition. arXiv (2015). doi:arXiv:1504.05477
135. Wang, S.: A practical guide to randomized matrix computations with matlab implementations. arXiv (2015). doi:arXiv:1505.07570v6
136. Yu, W., Gu, Y., Li, J., Liu, S., Li, Y.: Single-pass pca of large high-dimensional data. In: Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, pp. 3350–3356 (2017). Hong Kong
137. Agrawal, A., Chiu, A.M., Halperin, M.L.E., Sankararaman, S.: Scalable probabilistic pca for large-scale genetic variation data. bioRxiv (2019). doi:10.1101/729202
138. Liberty, E., Zucker, S.W.: The mailman algorithm: A note on matrix–vector multiplication. Information Processing Letters **109(3)**, 179–182 (2009)
139. 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. bioRxiv (2019). doi:10.1101/574574
140. Chen, M., Li, W., Zhang, W., Wang, X.: Dimensionality reduction with generalized linear models. In: Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence, pp. 1267–1272 (2013). Beijing
141. Shekhar, K., Lapan, S.W., Whitney, I.E., Tran, N.M., Macosko, E.Z., Kowalczyk, M., Adiconis, X., Levin, J.Z., Nemesh, J., Goldman, M., McCarroll, S.A., Cepko, C.L., Regev, A., Sanes, J.R.: Comprehensive classification of retinal bipolar neurons by single-cell transcriptomics. Cell **166(5)**, 1306–1323 (2016)
142. van Dijk, D., Sharma, R., Nainys, J., Yim, K., Kathail, P., Carr, A.J., Burdziak, C., Moon, K.R., Chaffer, C.L., Pattabiraman, D., Bierie, B., Mazutis, L., Wolf, G., Krishnaswamy, S., Peer, D.: Recovering gene interactions from single-cell data using data diffusion. Cell **174(3)**, 716–729 (2018)
143. Zheng, G.X., Terry, J.M., Belgrader, P., Ryvkin, P., Bent, Z.W., Wilson, R., Ziraldo, S.B., Wheeler, T.D., McDermott, G.P., Zhu, J., Gregory, M.T., Shuga, J., Montesclaros, L., Underwood, J.G., Masquelier, D.A., Nishimura, S.Y., Schnall-Levin, M., Wyatt, P.W., Hindson, C.M., Bharadwaj, R., Wong, A., Ness, K.D., Beppu, L.W., Deeg, H.J., McFarland, C., Loeb, W.J. K R andd Valente, Ericson, N.G., Stevens, E.A., Radich, J.P., Mikkelsen, T.S., Hindson, B.J., Bielas, J.H.: Massively parallel digital transcriptional profiling of single cells. Nature Communications **8(14049)** (2017)

**Figures**



**Figure 1 Overview of benchmarking in this work.** (a) Schematic overview of this work. (b) Evaluation metrics of the benchmarking with real-world datasets. (c) Evaluation metrics of the benchmarking with synthetic datasets.

| Implementation (*package*) | Algorithm | Category | Time complexity | Memory | Language | OS | Out-of-core | t-SNE / UMAP | Clustering | Higher PCs | Outlier robustness | Eigenvalues | Calculation time | Memory usage | Scalability (LGC) | Parameter tuning | Repository | Installation | Document type | Unit test | CI | Age (years) | Reference |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| prcomp | Golub-Kahan | SimT | O(NM min(N,M)) | O(NM) | R | LMW | – | + | + | + | + | + | – | – | 9 | + | CRAN | + | man/ | + | – | > 10 | [90] |
| PCA (*sklearn*, full) | Golub-Kahan | SimT | O(NM min(N,M)) | O(NM) | Python | LMW | – | + | + | + | + | + | – | – | 9 | + | GitHub | + | documentation.html | + | CircleCI | 8 | [90,92] |
| fit (*MultiVariateStats.jl*) | Golub-Kahan | SimT | O(NM min(N,M)) | O(NM) | Julia | LMW | – | + | + | + | + | + | – | – | 9 | + | GitHub | + | Documenter.jl | + | TravisCI | 5 | [90] |
| Downsampling | Golub-Kahan | DS | O(NM' min(N,M')) | O(NM') | ? | ? | – | – | – | – | – | – | + | + | – | ± | – | – | – | – | – | – | [94] |
| IncrementalPCA (*sklearn*) | SKL | SU | O(NM(K+B)²/B) | O(BM) | Python | LMW | + | + | + | + | + | + | + | + | > 11 | ± | PyPI | + | documentation.html | + | CircleCI | 5 | [95] |
| irlba (*irlba*) | IRLBA | Krylov | O(KM) | O(NM) | R | LMW | – | + | + | + | + | + | + | – | 9 | + | CRAN | + | vignettes | + | CRAN | 7 | [96] |
| svds (*RSpectra*) | IRAM | Krylov | O(K²M) | O(NM) | R | LMW | – | + | + | + | + | + | + | – | 9 | + | CRAN | + | vignettes | – | CRAN | 3 | [96,98] |
| propack.svd (*svd*) | IRLBA | Krylov | O(KM) | O(NM) | R | LMW | – | + | + | + | + | + | + | – | 9 | + | CRAN | + | man/ | – | CRAN | 8 | [96,99] |
| PCA (*sklearn*, arpack) | IRAM | Krylov | O(K²M) | O(NM) | Python | LMW | – | + | + | + | + | + | + | – | 10 | + | PyPI | + | documentation.html | + | CircleCI | 2 | [96,97] |
| irlb (*Cell Ranger*) | IRLBA | Krylov | O(KM) | O(NM) | Python | L | – | + | + | + | + | + | + | – | 9 | + | 10X(*) | ± | Support page | + | TravisCI | 0 | [96] |
| svds (*Arpack.jl*) | IRAM | Krylov | O(K²M) | O(NM) | Julia | LMW | – | + | + | + | + | + | + | – | 10 | + | GitHub | + | Sphinx-julia | + | TravisCI | 0 | [96,99] |
| orthiter (*OnlinePCA.jl*) | Orthogonal iteration | Krylov | O(KNM) | O(KM) | Julia | LMW | + | + | + | + | + | + | + | + | > 11 | ± | GitHub | ± | Documenter.jl | + | TravisCI | 0 | This paper |
| gd (*OnlinePCA.jl*) | GD | GD | O(KNM) | O(KM) | Julia | LMW | + | + | + | – | + | + | + | + | > 11 | – | GitHub | ± | Documenter.jl | + | TravisCI | 0 | This paper |
| sgd (*OnlinePCA.jl*) | SGD | GD | O(K²NM) | O(KM) | Julia | LMW | + | – | – | – | – | – | + | + | > 11 | – | GitHub | ± | Documenter.jl | + | TravisCI | 0 | This paper |
| rsvd (*rsvd*) | Halko's method | Rand | O(LNM) | O(NM) | R | LMW | – | + | + | + | + | + | + | – | 9 | + | CRAN | + | man/ | + | CRAN | 3 | [111] |
| oocPCA_CSV (*oocRPCA*) | Li's method | Rand | O(LNM) | O(BM) | R | LMW | + | + | + | + | + | + | + | + | > 11 | + | GitHub | ± | man/ | + | – | 2 | [69] |
| PCA (*sklearn*, randomized) | Halko's method | Rand | O(LNM) | O(NM) | Python | LMW | – | + | + | + | + | + | + | – | 10 | + | PyPI | + | documentation.html | + | CircleCI | 2 | [107,108] |
| randomized_svd (*sklearn*) | Li's method | Rand | O(LNM) | O(NM) | Python | LMW | – | + | + | + | + | + | + | – | 9 | + | PyPI | + | documentation.html | + | CircleCI | 7 | [109] |
| PCA (*dask-ml*) | Halko's method | Rand | O(LNM) | O(BM) | Python | LMW | + | + | + | – | + | + | + | + | 8 | + | PyPI | + | Sphinx | + | TravisCI | 1 | [112] |
| halko (*OnlinePCA.jl*) | Halko's method | Rand | O(LNM) | O(KM) | Julia | LMW | + | + | + | + | + | + | + | + | > 11 | + | GitHub | ± | Documenter.jl | + | TravisCI | 0 | This paper |
| algorithm971 (*OnlinePCA.jl*) | Li's method | Rand | O(LNM) | O(KM) | Julia | LMW | + | + | + | + | + | + | + | + | > 11 | + | GitHub | ± | Documenter.jl | + | TravisCI | 0 | This paper |

+: Good/Exists/Easy
±: Normal
–: Bad/Does not exists/Difficult
?: Unknown/Not evaluated

N : No. genes
M : No. cells
M' : No. of sampled cells
B : Block size (No. of sampled genes)
L : Random dimension

L : Linux
M : Mac OS
W : Windows

LGC : Minimum log₁₀(NM) of the jobs crashed by out-of-memory erros

* 10X Genomics website

**Figure 2 Summary of results.** (a) Theoretical properties summarized by our literature review. (b) Properties related to each implementation. (c) Performance evaluated by benchmarking with real-world and synthetic datasets. (d) User-friendliness evaluated by some metrics.

**Figure 3 The comparison of t-stochastic neighbor embedding (t-SNE) plots.** Comparison of multiple principal component analysis (PCA) implementations performed with empirical datasets: PBMCs ($10^2$ cells), Pancreas ($10^3$ cells), BrainSpinalCord ($10^5$ cells), and Brain datasets ($10^6$ cells). t-SNE was performed with the result of each PCA implementation.

**Figure 4 Clustering accuracy comparison.** Clustering accuracy was evaluated by the adjusted Rand index (ARI) of the result of Louvain clustering. Multiple principal component analysis (PCA) implementations were performed for PBMCs ($10^2$ cells), Pancreas ($10^3$ cells), BrainSpinalCord ($10^5$ cells), and Brain datasets ($10^6$ cells); Louvain clustering was performed for the PCA results. For each PCA result, Louvain clustering calculations were performed ten times. The cluster labels are the same as those of the respective original papers.

**Figure 5  Comparison of all combinations of eigenvectors.** Absolute values of the cross products of all combinations between the eigenvectors of the gold standard methods and those of the other principal component analysis (PCA) implementations were calculated. The closer the value is to 1, the closer the two corresponding eigenvectors are to each other. If two PCA results are equal without considering differences in sign, the matrix in this figure becomes an identity matrix.

**Figure 6 Comparison of eigenvalues.** Distribution of eigenvalues of all the principal component analysis (PCA) implementations for each real dataset.

**Figure 7 Comparison of the elapsed time and maximum memory usage for empirical datasets.** (a) Elapsed time and (b) memory usage of all principal component analysis (PCA) implementations calculated for each empirical dataset. We used our in-house Julia script to preprocess the Brain dataset.

**Figure 8   Comparison of the elapsed time for simulated datasets.** Synthetic datasets ($\{10^2, 10^3, 10^4\}$ gene $\times$ $\{10^2, 10^3, 10^4, 10^4, 10^5, 10^6, 10^7\}$ cell matrices) were randomly generated, and all the out-of-core principal component analysis (PCA) implementations were performed. In each panel, the logarithm of the number of cells is indicated along each $x$-axis, and the elapsed time (hours) is shown along each $y$-axis.

**Figure 9 Comparison of the maximum memory usage for simulated datasets.** Synthetic datasets ($\{10^2, 10^3, 10^4\}$ gene $\times$ $\{10^2, 10^3, 10^4, 10^4, 10^5, 10^6, 10^7\}$ cell matrices) were randomly generated, and all the out-of-core principal component analysis (PCA) implementations were performed. In each panel, the logarithm of the number of cells is indicated along the $x$-axis, and memory usage (GB) is shown along the $y$-axis.

**Figure 10   Relationships of the algorithms/implementations, the number of passes, and the file format with the elapsed time for performing principal component analysis (PCA) with the Brain dataset.** (a) Number of passes for the data matrix and the computation time for each algorithms/implementations were calculated. (b) Elapsed time and memory usage for one-pass orthogonal iteration were calculated.

**Figure 11  OnlinePCA.jl schematic.** Input CSV files were first saved as a binary file with the `csv2bin` command and analyzed with the `gd` and `sgd` commands, which perform incremental principal component analysis (PCA). When using the HDF5 file format defined by 10X Genomics, we converted the file to CSV format using an in-house Python script. Gene-wise or cell-wise summary statistics were calculated using the `sumr` command. Highly variable genes can also be calculated with the `hvg` command. Because the gene-wise and cell-wise summary statistics are expressed as small vectors, they can be used to perform precise data quality control (QC) with any programming language without out-of-core implementations. After QC, the `filtering` command removed low-quality genes and cells using a user-specified index. Combined with the small size vectors, some out-of-core PCA implementations, such as `orthiter/gd/sgd/halko/algorithm971` have commands to incrementally update eigenvectors from the row vector of the data matrix. The `tenxpca` command directly performed algorithm971 on 10X-HDF files.

**Figure 12 User guidelines.** Recommended PCA implementations categorized based on written language and matrix size.
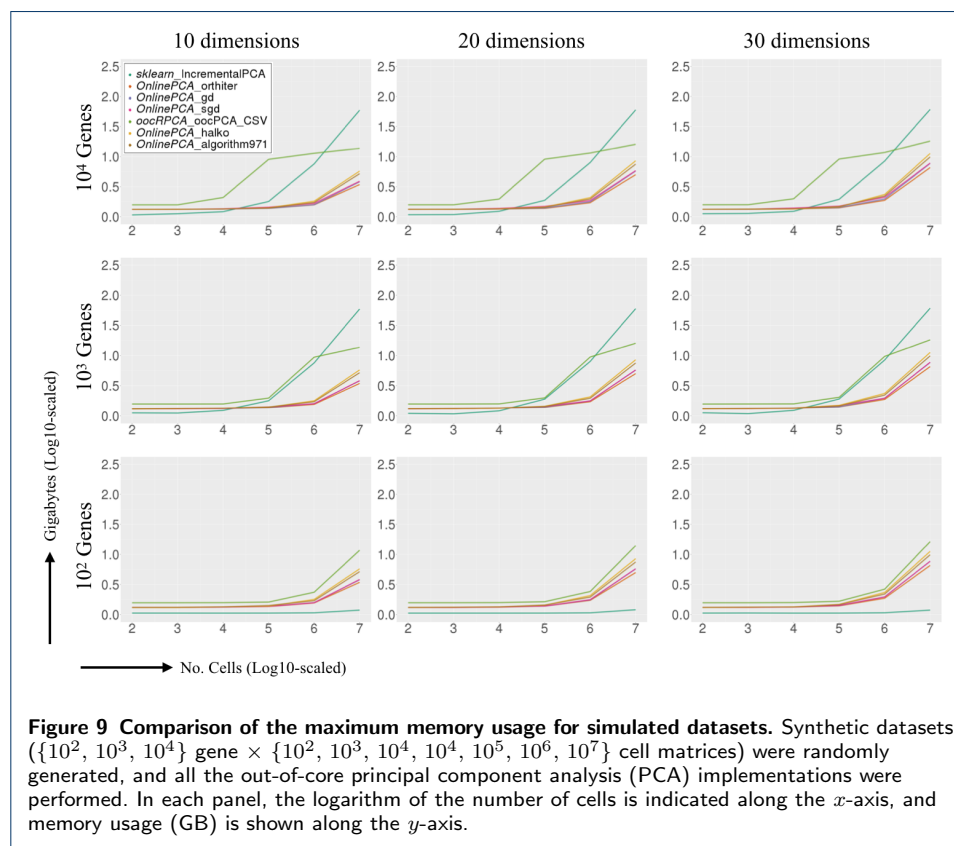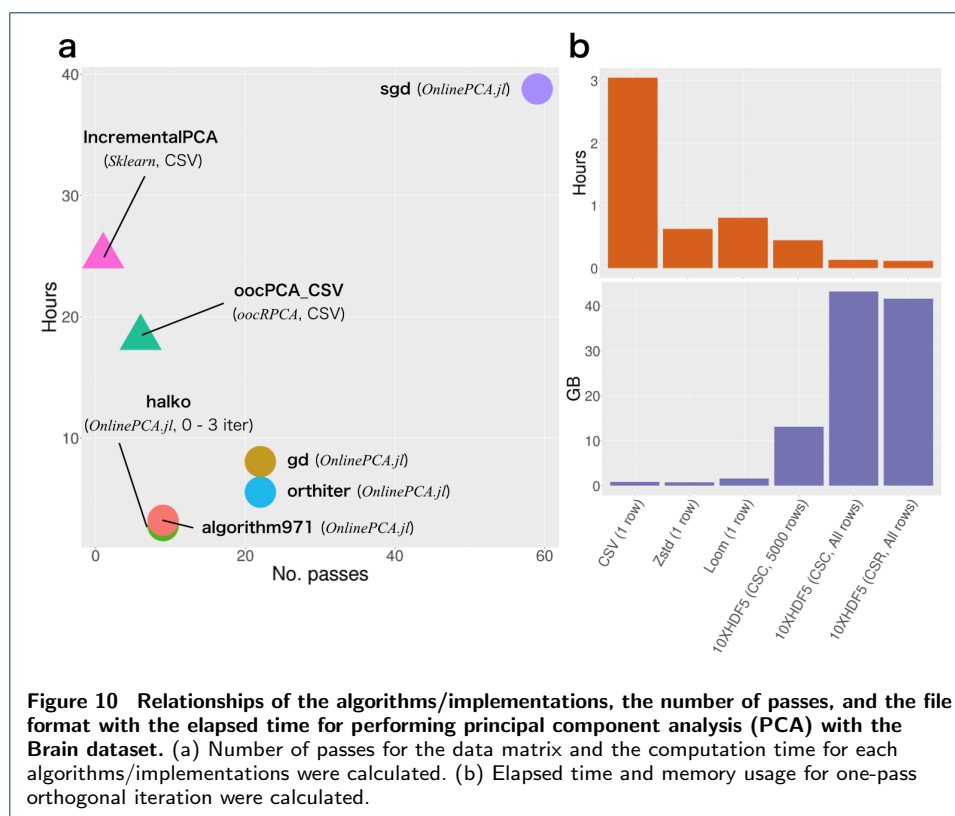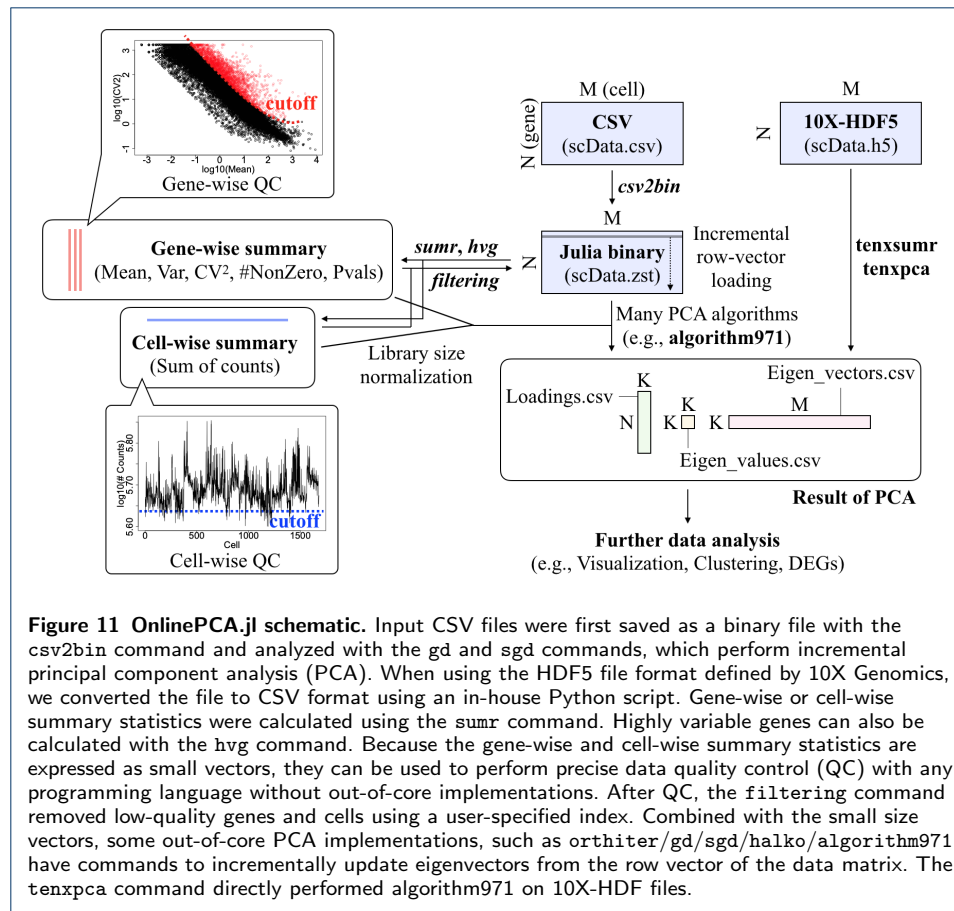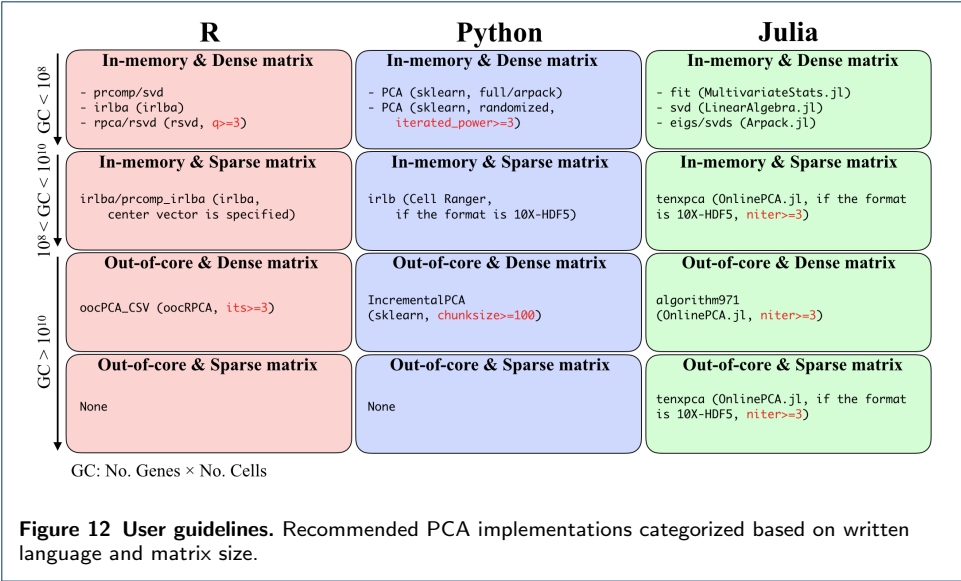
## Tables

**Table 1** Use cases of PCA implementations in scRNA-seq studies.

| scRNA-seq studies | PCA algorithms | Commands or functions used in the studies |
|---|---|---|
| In most cases | | `prcomp`/`svd` (R) |
| [13, 42, 43, 51, 52, 55, 56, 58, 60, 63, 65, 74, 77, 82, 85, 91, 93] | Golub-Kahan method | PCA (Python, *sklearn*) |
| Bhaduri *et al.*, [94] | Downsampling | Unknown |
| *Loompy* [93] | SKL | `IncrementalPCA` (Python, *sklearn*) |
| | IRLBA | PCA (Python, *sklearn*) |
| *Scanpy* [93] | SKL | `IncrementalPCA` (Python, *sklearn*) |
| | Halko's method | `TruncatedSVD` (Python, *sklearn*) |
| *Cell Ranger* [22] | IRLBA | `irlb` (Python, from scratch) |
| *Seurat2* [49] | IRLBA | `irlba` (R, *irlba*) |
| *Scran* [50] | Golub-Kahan method | `svd` (R) |
| | IRLBA | `irlba` (R, *irlba*) |
| *SAFE* [76] | IRLBA | `irlba` (R, *irlba*) |
| | Golub-Kahan method | `svds` (MATLAB) |
| *MAGIC* [52] | Halko's method | `randPCA` (MATLAB, from scratch) |
| | Halko's method | PCA (Python, *sklearn*) |
| *Harmony* [57] | IRLBA | `irlba` (R, *irlba*) |
| *Scater* [82] | Golub-Kahan method | `prcomp` (R) |
| | IRLBA | `irlba` (R, *irlba*) |
| *GiniClust2* [59] | IRLBA | `propack.svd` (R, *svd*) |
| *SIMLR* [75] | Halko's method | `fast.rsvd` (R, from scratch) |
| *SEQC* [89] | Golub-Kahan method | PCA (Python, *sklearn*) |
| | Halko's method | PCA (Python, *sklearn*) |
| *CellFishing.jl* [61] | Li's method | `rsvd` (Julia, from scratch) |

**Table 2** Real-world datasets for benchmarking

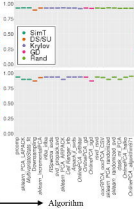| Dataset | No. Genes | No. Cells | No. Cell types | PCs used | File size (LogCPMED, CSV) | File size (Count, CSV) | File size (Count, Binary) |
|---|---|---|---|---|---|---|---|
| PBMCs | 17484 | 713 | 6 | PC1-3 | 45 MB | 24 MB | 2.1 MB |
| Pancreas | 17499 | 3605 | 14 | PC1-12 | 530 MB | 287 MB | 22 MB |
| Brain SpinalCord | 25893 | 156049 | 73 | PC1-16 | 9.3 MB | 7.5 GB | 197 MB |
| Brain | 18782 | 1306127 | 60 | PC1-20 | 290 GB | 58 GB | 3.2 GB |

**Additional Files**

Additional File 1 — Review of existing PCA algorithms and implementations. (PDF 271 KB)

Additional File 2 — Pseudo-code of all the PCA algorithms. (PDF 178 KB)

Additional File 3 — Source code of all the PCA implementations. (PDF 59 KB)

Additional File 4 — Results of t-SNE of all the PCA implementations. (PNG 623 KB)

Additional File 5 — Results of UMAP of all the PCA implementations. (PNG 368 KB)

Additional File 6 — Results of clustering methods of all the PCA implementations (PDF 3.6 MB)

Additional File 7 — Eigenvectors of all the PCA implementations (PBMCs and Pancreas). (PNG 308 MB)

Additional File 8 — Pair plots of all the PCA (PBMCs) implementations. (TAR.GZ 649 KB)

Additional File 9 — Pair plots of all the PCA (Pancreas) implementations. (TAR.GZ 4.9 MB)

Additional File 10 — Pair plots of all the PCA (BrainSpinalCord) implementations. (TAR.GZ 3.1 MB)

Additional File 11 — Pair plots of all the PCA (Brain) implementations. (TAR.GZ 5.8 MB)

Additional File 12 — Number of singleton clusters. (PNG 271 KB)

Additional File 13 — Eigenvectors of all the PCA implementations (BrainSpinalCord and Brain). (PNG 532 KB)

Additional File 14 — Loading vectors of all the PCA implementations (PBMCs and Pancreas). (PNG 349 KB)

Additional File 15 — Crashed jobs caused by out-of-memory errors. (TXT 882 B)

Additional File 16 — Parameter tuning of the `IncrementalPCA` implementations. (PDF 445 KB)

Additional File 17 — Parameter tuning of the orthogonal iteration, gradient descent, and stochastic gradient descent implementations. (PDF 1.3 MB)

Additional File 18 — Parameter tuning of the randomized SVD implementations. (PDF 734 KB)

Additional File 19 — Developer guidelines. (PNG 1.1 MB)

Additional File 20 — Effect of feature selection on clustering accuracy. (PDF 1 MB)

Additional File 21 — Comparison of normalizing size factors. (HTML 1.4 MB)

|  | SimT (LAPACK) *Skraam* | DS (Devonagiri) *p-comp* | SU (Frommenial) *Skraam* | Krylov *i/o* | GD (OrchitzPCA.fl (rgd)) | Rand (aquriv.777) |
|---|---|---|---|---|---|---|
| **a** PBMCs | | | | | | |
| Pancreas | | | | | | |
| **b** PBMCs | | | | | | |
| Pancreas | | | | | | |

**a** PBMCs     **b** Pancreas     **c** BrainSpinalCord

**d** Brain

SimT, DS/SU, Krylov, GD, Rand

Eigenvalues

Principal components

**a**

No. passes (x-axis), Hours (y-axis)

- sgd (*OnlinePCA.jl*)
- IncrementalPCA (*Sklearn*, CSV)
- oocPCA_CSV (*oocRPCA*, CSV)
- halko (*OnlinePCA.jl*, 0 - 3 iter)
- gd (*OnlinePCA.jl*)
- orthiter (*OnlinePCA.jl*)
- algorithm971 (*OnlinePCA.jl*)

**b**

Hours (top), GB (bottom)

Categories: CSV (1 row), Zstd (1 row), Loom (1 row), 10XHDF5 (CSC, 5000 cells), 10XHDF5 (CSC, All rows), 10XHDF5 (CSR, All rows)

**a**

Literature review of existing PCA algorithms and implementations

10 Algorithms
21 implementations

*Gold standard*
prcomp
(Only when the dataset is small)

4 Real Datasets          18 Synthetic Datasets

Accuracy          Time
Memory
Scalability

Fast, memory efficient, and scalable PCA,
QC, highly variable genes, and filtering

Recommendation for users and developers
Future extensions

**b**

kNN / LMU
Clustering accuracy

*Gold standard*      *+ Good*      *- Bad*

Higher PCs

*+ Good*      *- Bad*

Eigenvectors

Gene-wise QC

Gene-wise summary
(Mean, Var, CV², #NonZero, Pvals)

*sumr, hvg*

*filtering*

Cell-wise summary
(Sum of counts)

Library size
normalization



Cell-wise QC

M (cell)

N (gene)

**CSV**
(scData.csv)

*csv2bin*

**Julia binary**
(scData.zst)

N

Incremental
row-vector
loading

M

**10X-HDF5**
(scData.h5)

Z

*tenxsumr*
*tenxpca*

Many PCA algorithms
(e.g., **algorithm971**)

K

K

M

**Loadings.csv**

N

K

**Eigen_vectors.csv**

Eigen_values.csv

**Result of PCA**

**Further data analysis**
(e.g., Visualization, Clustering, DEGs)

Table columns are grouped under four headings: **a** Theoretical properties, **b** Implementation, **c** Performance, **d** User friendliness.

| Implementation (package) pcromp | Algorithm | Category | Time complexity | Memory | Language | OS | Out-of-core | t-SNE/UMAP | Clustering | Higher PC's | Point distances | Calculation time | Memory usage | Scalability (LCC) | Parameter tuning | Repository | Installation | Document type | Unit test | CI | Age (years) | Reference |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PCA (sklearn, full) | Gram–Kalman | SimT | $O(NM\min(N,M))$ | $O(NM)$ | Python | L,M,W | − | + | + | + | + | − | 3 | − | + | GitHub | many | documentation | many | TravisCI | >13 | [50,89] |
| fit (MultivariateStats.jl) | Gram–Kalman | SimT | $O(NM\min(N,M))$ | $O(NM)$ | Julia | L,M,W | − | + | + | + | + | − | 3 | − | + | GitHub | Documented | Documented | 5 | [90,91] |
| Downsampling | Gram–Kalman | I/S | $O(NM\min(N,M))$ | $O(NM)$ | ? | ? | ± | + | ± | − | ± | + | 0 | ± | ± | — | — | — | — | — | [84] |
| IncrementalPCA (sklearn) | SAL | SU | $O(NMt+Mt^2)$ | $O(MM'+Mt)$ | Python | L,M,W | + | + | + | + | + | ± | 0 | ± | + | GitHub | documentation | TravisCI | 9 | [85] |
| irlba (irlba) | IRLBA | Krylov | $O(NM)$ | $O(NM)$ | R | L,M,W | − | + | + | + | + | + | 3 | − | + | CRAN | vignettes | CRAN | 7 | [86] |
| prcomp (stats) | IRAM | Krylov | $O(NM)$ | $O(NM)$ | R | L,M,W | − | + | + | + | + | + | 3 | − | + | CRAN | Documented | CRAN | 3 | [92,98] |
| propack.svd (svd) | IRAM | Krylov | $O(NM)$ | $O(NM)$ | R | L,M,W | − | + | + | + | + | + | 3 | − | + | CRAN | Documented | CRAN | 8 | [88] |
| PCA (scran, pca) | IRLBA | Krylov | $O(NM)$ | $O(NM)$ | Python | L | − | + | + | + | + | − | 10 | − | − | GitHub | Support page | 4 | This paper |
| irlb (irlb.jl) | IRLBA | Krylov | $O(NM)$ | $O(NM)$ | Julia | L,M,W | − | + | + | + | + | + | 10 | − | + | GitHub | Sphinx | Julia | 4 | [96] |
| svds (svds) | IRAM | Krylov | $O(NM)$ | $O(NM)$ | Julia | L,M,W | − | + | + | + | + | + | 10 | − | + | GitHub | Documented | TravisCI | 5 | This paper |
| orthiter (OnlinePCA.jl) | Orthogonal iteration | Krylov | $O(NM)$ | $O(NM')$ | Julia | L,M,W | + | + | + | + | + | ± | 0 | + | − | GitHub | Documented | TravisCI | 3 | This paper |
| gd (OnlinePCA.jl) | GD | GD | $O(NM)$ | $O(NM')$ | Julia | L,M,W | + | + | + | + | + | ± | 0 | + | − | GitHub | Documented | TravisCI | 3 | This paper |
| sgd (OnlinePCA.jl) | SGD | GD | $O(NM)$ | $O(NM')$ | Julia | L,M,W | + | + | + | + | + | ± | 0 | + | − | GitHub | Documented | TravisCI | 3 | This paper |
| rsvd (rsvd) | Halko's method | Rand | $O(NM)$ | $O(NM)$ | R | L,M,W | − | + | + | + | + | + | 3 | − | ± | CRAN | many | CRAN | 3 | [117,119] |
| oocPCA_CSV (oocRPCA) | Halko's method | Rand | $O(NM)$ | $O(NM')$ | R | L,M,W | + | + | + | + | + | + | 0 | + | ± | GitHub | Documented | 3 | [110] |
| randomized_svd (sklearn) | Halko's method | Rand | $O(NM)$ | $O(NM)$ | Python | L,M,W | − | + | + | + | + | + | 3 | − | ± | GitHub | documentation | TravisCI | 9 | [112] |
| PCA (dask-ml) | Halko's method | Rand | $O(NM)$ | $O(NM')$ | Python | L,M,W | + | + | + | + | + | + | 0 | + | ± | GitHub | Sphinx | 3 | [113] |
| halko (OnlinePCA.jl) | Halko's method | Rand | $O(NM)$ | $O(NM')$ | Julia | L,M,W | + | + | + | + | + | + | 0 | + | − | GitHub | Documented | TravisCI | 3 | This paper |
| algorithm971 (OnlinePCA.jl) | Li's method | Rand | $O(NM)$ | $O(NM')$ | Julia | L,M,W | + | + | + | + | + | + | 0 | + | − | GitHub | Documented | TravisCI | 3 | This paper |

Legend:

- + Good/Fast/Easy
- ± Bad/Does not exist/Difficult
- − Bad/Slow
- ? Unknown/Not evaluated

N: No. genes
M: No. of samples/cells
M': No. of sampled cells (No. of sampled genes)
t: Random dimension

LCC: Minimum $\log_2(M)$ at the pos crashed by out-of-memory errors
* NIX Genomics website

**a** Elapsed Time  **b** Maximum Memory Usage

Legend: Data I/O + PCA (orange); Preprocess (e.g. binarization) (green)

PBMCs — Seconds / Gigabytes

Pancreas — Hours / Gigabytes

Brain SpinalCord — Hours / Gigabytes

Brain — Hours / Gigabytes

Out of Memory

x-axis labels:
picomo, sklearn_PCA_LAPACK, Multicore-TSNE, sklearn_incremental_PCA, Downsampling, Seurat, pbp, tSNE, mrtsne, Rtsne, sklearn_umap, sklearn_PCA_ARPACK, Asava-E_sva, OnlinePCA_jst, OnlinePCA_gd, OnlinePCA_sgd, oocPCA_exact, sklearn_PCA_randomized, sklearn_randomized_svd, irlba, irlba_PCA, OnlinePCA_halko, OnlinePCA_algorithm01

No. Cells (Log10-scaled)

|  | SimT | DS | SU | Krylov | GD | Rand |
|---|---|---|---|---|---|---|
|  | prcomp | prcomp (Downsampling) | *Sklearn* (Incremental) | *irlba* (irlba) | *OnlinePCA.jl* (sgd) | *OnlinePCA.jl* (algorithm971) |
| PBMCs | | | | | | |
| Pancreas | | | | | | |
| Brain SpinalCord | × | | | × | | |
| Brain | × | | | × | | |

× : Out-of-memory

| | R | Python | Julia |
|---|---|---|---|
| **GC < 10^9** | **In-memory & Dense matrix**<br>- prcomp/svd<br>- irlba (irlba)<br>- rpca/rsvd (rsvd, q>=3) | **In-memory & Dense matrix**<br>- PCA (sklearn, full/arpack)<br>- PCA (sklearn, randomized,<br>  iterated_power>=3) | **In-memory & Dense matrix**<br>- fit (MultivariateStats.jl)<br>- svd (LinearAlgebra.jl)<br>- eigs/svds (Arpack.jl) |
| **10^8 < GC < 10^10** | **In-memory & Sparse matrix**<br>irlba/prcomp_irlba (irlba,<br>  center vector is specified) | **In-memory & Sparse matrix**<br>irlb (Cell Ranger,<br>  if the format is 10X-HDF5) | **In-memory & Sparse matrix**<br>tenxpca (OnlinePCA.jl, if the format<br>is 10X-HDF5, niter>=3) |
| **GC > 10^10** | **Out-of-core & Dense matrix**<br>oocPCA_CSV (oocRPCA, its>=3) | **Out-of-core & Dense matrix**<br>IncrementalPCA<br>(sklearn, chunksize>=100) | **Out-of-core & Dense matrix**<br>algorithm971<br>(OnlinePCA.jl, niter>=3) |
| | **Out-of-core & Sparse matrix**<br>None | **Out-of-core & Sparse matrix**<br>None | **Out-of-core & Sparse matrix**<br>tenxpca (OnlinePCA.jl, if the format<br>is 10X-HDF5, niter>=3) |

GC: No. Genes × No. Cells