

SOFTWARE

Open Access



An interpretable framework for clustering single-cell RNA-Seq datasets

Jesse M. Zhang¹, Jue Fan², H. Christina Fan², David Rosenfeld² and David N. Tse^{1*} 

Abstract

Background: With the recent proliferation of single-cell RNA-Seq experiments, several methods have been developed for unsupervised analysis of the resulting datasets. These methods often rely on unintuitive hyperparameters and do not explicitly address the subjectivity associated with clustering.

Results: In this work, we present DendroSplit, an interpretable framework for analyzing single-cell RNA-Seq datasets that addresses both the clustering interpretability and clustering subjectivity issues. DendroSplit offers a novel perspective on the single-cell RNA-Seq clustering problem motivated by the definition of “cell type”, allowing us to cluster using feature selection to uncover multiple levels of biologically meaningful populations in the data. We analyze several landmark single-cell datasets, demonstrating both the method’s efficacy and computational efficiency.

Conclusion: DendroSplit offers a clustering framework that is comparable to existing methods in terms of accuracy and speed but is novel in its emphasis on interpretability. We provide the full DendroSplit software package at <https://github.com/jessemzhang/dendrosplit>.

Keywords: Single-cell RNA-seq, Clustering, Feature selection, Interpretability

Background

In recent years, single-cell RNA-Seq has proven to be a powerful approach for studying biological samples in various settings [1]. Scientists have leveraged this technology to shed light on how cells differentiate [2–6], investigate known cell types [7–10], and discover new cell types and gene patterns [11–17]. These efforts have yielded a plethora of diverse datasets sharing characteristics such as missing entries (drop-out events) and high dimensionality. Additionally, technological breakthroughs such as droplet encapsulation, molecular barcoding, and cheap parallelization have produced datasets involving tens of thousands and even millions of cells [17–22]. After obtaining such datasets, scientists are often interested in clustering the high-dimensional points corresponding to individual cells, ideally recovering known cell populations while discovering new and perhaps rare cell types. While the definition of a cell type is not precise [23], biologists agree that gene expression levels are highly relevant.

With gene expression dictating protein expression (and hence cellular function), identifying the genes that distinguish a cell type is of paramount importance. Therefore from a computational perspective, there are two key problems in downstream analysis: 1) clustering and 2) feature selection, also known as differential expression.

General-purpose clustering algorithms such as *K*-means, DBSCAN [24], affinity propagation [25], and spectral clustering [26] have performed well for several single-cell datasets [27]. In order to achieve good performance, however, the datasets often need to be carefully preprocessed, and the algorithms require non-intuitive hyperparameter tuning. For example, both *K*-means and spectral clustering require choosing the desired number of clusters, DBSCAN requires choosing the max distance between two samples in the same neighborhood, and affinity propagation requires choosing both a preference parameter for determining which points are exemplars and a damping parameter for avoiding numerical oscillations. To address specific computational challenges of single-cell RNA-Seq datasets, researchers have developed a wide array of application-specific clustering algorithms [28–34] and packages for end-to-end analysis [21, 35–39].

*Correspondence: dntse@stanford.edu

¹Department of Electrical Engineering, Stanford, 94305 Stanford, California, USA

Full list of author information is available at the end of the article

Regardless of which set of these tools one uses, finding the right approach for clustering a specific dataset requires careful design of the computational workflow, but often finding a good combination of clustering algorithm and hyperparameters is time-consuming and difficult. Additionally, none of these approaches explicitly addresses the inherent subjectiveness behind clustering, which stems from the potential existence of subtypes and sub-subtypes.

With an emphasis on interpretability and ease of exploratory analysis, we introduce DendroSplit, a framework for clustering single-cell RNA-Seq data. In addition to speed, the framework has the following advantages:

- Gene-based justification for all decisions made when generating clusters
- Interpretable hyperparameters
- Ability to cheaply produce multiple clusterings for the same dataset
- Ease of incorporation into existing single-cell RNA-Seq workflows

At a high level, the approach leverages a feature selection algorithm to generate biologically meaningful clusters. The end-to-end DendroSplit workflow is illustrated in Fig. 1a. After preprocessing the $N \times M$ expression matrix X (where N and M represent the number of cells and genes, respectively), we generate the $N \times N$ distance matrix D . We use hierarchical clustering to iteratively group cells based on their pairwise distances, obtaining a dendrogram, a tree-like data structure illustrating how grouping was performed. The split step starts at the root of the tree. Each node in the dendrogram represents a potential partitioning of a larger cluster into two smaller ones. If this “split” results in two adequately separated clusters (according to a metric we call the **separation score**), the split is deemed valid and the algorithm continues on the new clusters. Otherwise, the algorithm terminates for the subtree below the node. After the split step, DendroSplit performs a pairwise comparison of the resulting clusters, repeatedly merging clusters until all clusters are sufficiently separated. The merge step counteracts the greedy nature of hierarchical clustering, allowing DendroSplit to compare clusters that may have incorrectly ended up far away from one another in the dendrogram. The overall approach involves two intuitive hyperparameters: the separation score threshold for accepting a split, and the separation score threshold for accepting a merge.

We use the term “framework” to underline how specific design choices for certain components in the workflow such as the separation score will result in different clustering “methods”. Our choice of separation score is motivated by a key assumption: **if two cell populations are**

of different types, then there should exist at least one gene that is differentially expressed between the two populations. Given a candidate split in the dendrogram, we perform a Welch’s t -test for each gene. The separation score is $-\log(p_{\min})$ where p_{\min} represents the smallest p -value achieved (Fig. 1b), and we will be using this definition of separation score for all experiments presented in this work. We demonstrate that the deterministic method outlined in Fig. 1 is applicable to a wide variety of single-cell datasets. We show how DendroSplit can help us investigate the most significant genes considered at each split or merge, providing insight for how clusters are generated. Finally, we show how DendroSplit can cheaply generate several clusterings for different hyperparameter values.

Some clustering approaches similar to DendroSplit exist in literature. For example, the most common method of generating clusters from a dendrogram involves simply cutting the dendrogram horizontally at some fixed height. This rigid approach often fails to generate meaningful clusters for more complex datasets. The Dynamic Tree Cut algorithm [40] adds significant flexibility and processes the dendrogram based on an adaptive cut. Though it does not explicitly use a dendrogram, the backSPIN algorithm [12] also uses cell-cell similarities to perform iterative splitting. Unlike DendroSplit, both of these algorithms require choosing unintuitive hyperparameter cutoffs based on nuanced criteria. The most similar clustering approach was used by Lake et al. [15] for analyzing their human brain single-cell dataset. Their approach fits into the DendroSplit framework, using a separation score based on random forests. This separation score, compared to the separation score mentioned above, has an element of randomness, is significantly more computationally expensive, and requires less intuitive hyperparameter choices.

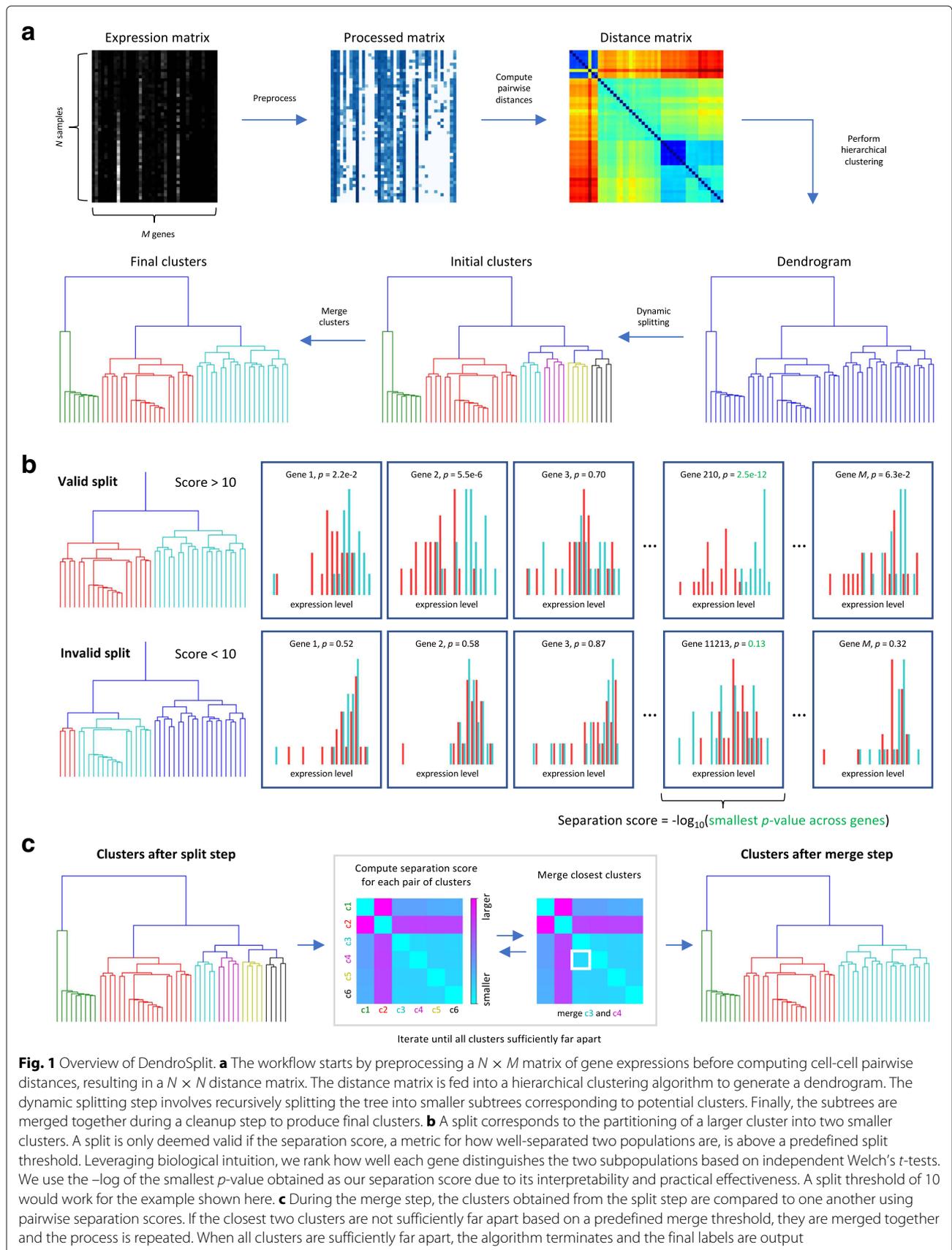
Implementation

Distance metric

For all single-cell datasets, we used the correlation distance. The correlation distance between \mathbf{x}_i and \mathbf{x}_j corresponding to cells i and j is

$$d(\mathbf{x}_i, \mathbf{x}_j) = 1 - r(\mathbf{x}_i, \mathbf{x}_j)$$

where r is the Pearson correlation coefficient. Therefore d is bounded between 0 and 2. This distance metric has the advantage of being agnostic to both shift and scale, making it robust to certain biases we would expect to vary across datasets. As a caveat, the distance metric has the disadvantage of depending on the number of zeros, and therefore the distance between two cells before and after gene filtering may be different due to removal of entries equal to 0. For all experiments in this work,



distance matrix computations were parallelized on 32 cores and computed using the `scikit-learn` Python package [41].

Hierarchical clustering

DendroSplit performs hierarchical clustering using the Scipy Python package [42]. One source of ambiguity for hierarchical clustering lies in the method for determining the distance between two clusters. We found that the “complete” method produces the best results, and this is the method used for all experiments reported below. For this method, the distance between two clusters is equal to the largest distance between a point from the first cluster and another point from the second cluster.

Separation score

The separation score effectively serves as a distance metric between two clusters, quantifying how different they are (see the supplementary material for further discussion). The cell-type assumption discussed in the “Background” section can also be phrased as: if two cell populations are of different types, then projection along one of the M axes should result in two distinguishable point clouds. For all experiments performed in this work, we defined the separation score between the $N_1 \times M$ population \mathbf{X} and the $N_2 \times M$ population \mathbf{Y} as

$$s(\mathbf{X}, \mathbf{Y}) = -\log_{10} \left(\min_i p(\mathbf{X}_i, \mathbf{Y}_i) \right)$$

where $p(\mathbf{X}_i, \mathbf{Y}_i)$ represents the p -value achieved using a Welch’s t -test for gene i . \mathbf{X}_i represents the i th column of \mathbf{X} corresponding to the expression of gene i in population \mathbf{X} . Welch’s t -test is similar to Student’s t -test but is more reliable when the two populations have unequal variance and size [43]. Compared to other differential expression approaches, Welch’s t -test is computationally cheap.

As an implementation note, if for a given split two or more genes have the exact same score, these genes are ranked by the magnitude of the t statistic. We note that because we are using Welch’s t -test rather than Student’s t -test, the degrees of freedom associated with each test is different, and hence outputting the largest t statistic is only approximately sound. Two genes may have the exact same score for larger datasets and for splits near the root of the dendrogram where p -values may be quite small, resulting in an underflow issue and a score of ∞ .

Handling singletons

In addition to the split and merge thresholds, the two major hyperparameters discussed in the “Background” section, the DendroSplit framework can also be customized using three minor hyperparameters. These three

hyperparameters are relevant for finding singletons (clusters containing one point), which are analogous to outliers.

Two of these hyperparameters are relevant for the split step. The first is the minimum cluster size. During a split, if one of the two candidate clusters contains less points than the minimum cluster size, that cluster is disbanded (each of its points are labeled as “Singleton”) and the algorithm continues on the other candidate. The second hyperparameter is the disband percentile. If a candidate split does not produce a subtree that meets the minimum cluster size requirement or if the candidate split does not achieve a high enough separation score, we look at the pairwise distances amongst samples in this final cluster. If all of them are greater than a certain percentile of distances in D , the original $N \times N$ distance matrix, then all points in this final cluster are marked as singletons. For all experiments performed in this work, the minimum cluster size was set to 2 (the smallest value) and the disband percentile was set to 50.

Before merging clusters, each singleton obtained during the split step is assigned to the same cluster as its nearest neighbor. If the distance between a singleton and its nearest neighbor is greater than a certain percentile of all pairwise distances in D , then the singleton remains unclassified. This percentile is the third minor hyperparameter and was set to 90 for all experiments performed in this work.

Hyperparameter sweeping

When choosing hyperparameters for DendroSplit, a relatively small split threshold such as 20 and a merge threshold set to half the split threshold often yields reasonable initial results. The DendroSplit approach can also rapidly generate several clusterings based on different split thresholds. Since DendroSplit saves the p -values and cell IDs considered at each split, we can obtain several split-step clustering results by exploiting the fact that the clusters generated with a smaller score threshold partition the clusters generated with a larger score threshold. The merge threshold can then be chosen by looking at pairwise separation scores between clusters.

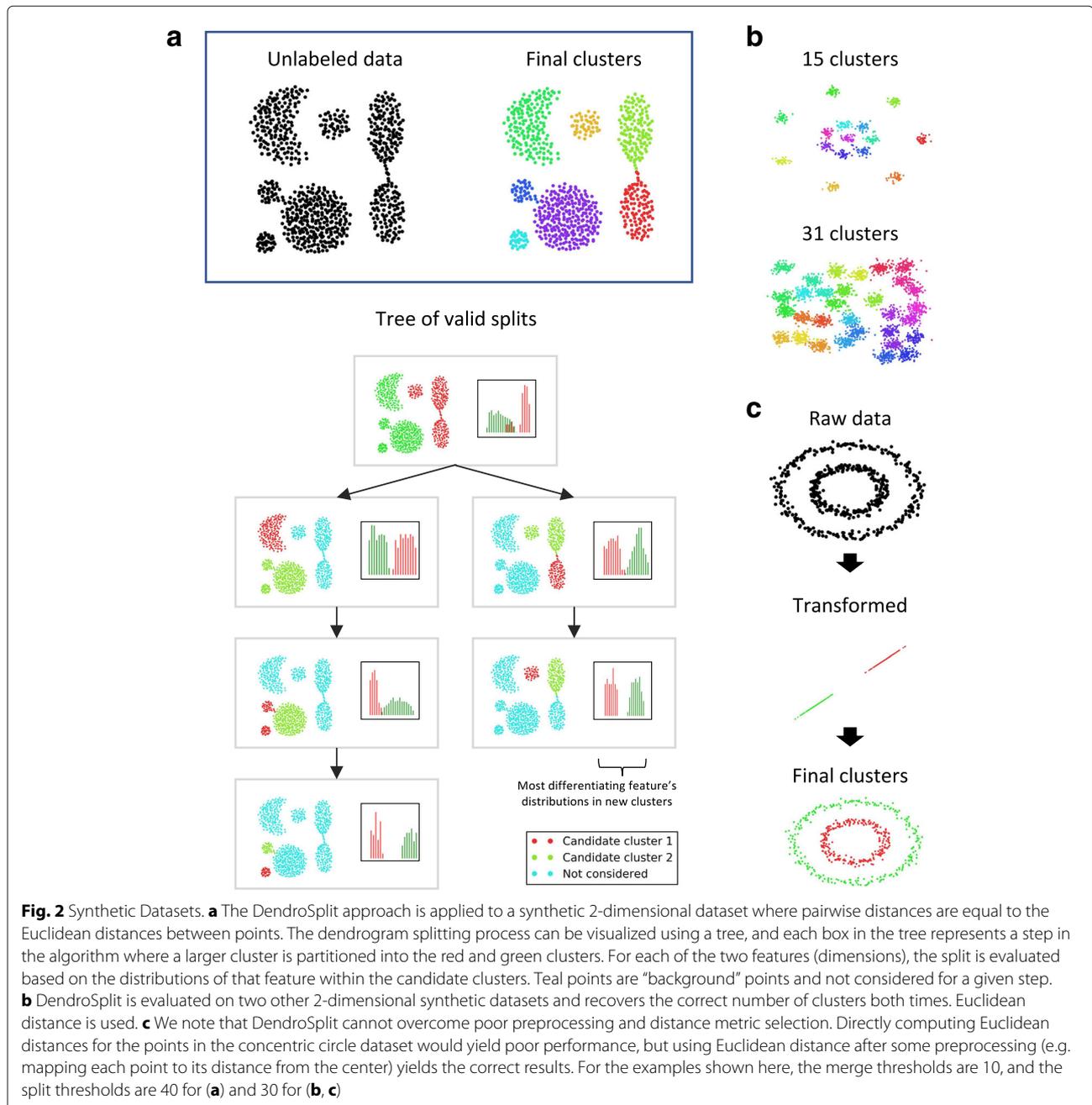
Results

Data preprocessing

For all single-cell datasets, we apply a logarithmic transformation $\log_{10}(X + 1)$ to the raw expression levels. We analyze 9 datasets in this paper. For each dataset, genes that have 0 expression across all cells were removed. Additionally, all datasets consisting of over 1000 cells undergo feature selection based on the method proposed by Macosko et al. [20] The M genes are sorted into

equal-sized bins depending on their mean expression values. Within a bin, genes are z -normalized based on their dispersions, where the dispersion for a gene is defined as the variance divided by the mean. Only genes corresponding to a z -score above a certain cutoff are retained. For the Zeisel et al. [12], Birey et al. [17], and Zheng et al. [21] datasets, we use DendroSplit's default setting of 5 bins with a z -cutoff of 1.5. For the Macosko et al. dataset, we first remove cells with less than 900 counts across all genes just like in Wang et al.'s [30]

approach, reducing the original 44808 cells to 11040. We then use Macosko et al.'s gene-filtering settings of 20 bins with a z -cutoff of 1.7. For the Zheng et al. dataset, reducing the number of genes results in several of the original 68579 cells having few and even 0 counts across all genes. We remove cells with less than 50 counts across all genes, resulting in 17426 cells, and again filtered out genes with 0 counts across all remaining cells. We also experimented with standardizing all log-transformed genes to have 0 mean and unit variance across all cells, but the



increased computational overhead did not yield better results.

Adjusted Rand index

The adjusted Rand index (ARI) is used to quantify how our clustering results match another given set of labels. The ARI ranges from 0 for poor matching to 1 for perfectly matched labels. For a set of n elements, we let X and Y represent two partitions of the elements. X_i represents the set of elements in partition i according to X . The adjusted Rand index is defined as:

$$ARI = \frac{\sum_{ij} \binom{n_{ij}}{2} - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}{\frac{1}{2} \left[\sum_i \binom{a_i}{2} + \sum_j \binom{b_j}{2} \right] - \left[\sum_i \binom{a_i}{2} \sum_j \binom{b_j}{2} \right] / \binom{n}{2}}$$

where n_{ij} is the number of elements in common between X_i and Y_j , $a_i = \sum_k n_{ik}$, and $b_j = \sum_k n_{kj}$.

Ground-truth datasets

To test the effectiveness of the DendroSplit framework, we first test the approach on datasets where the ground truth is known.

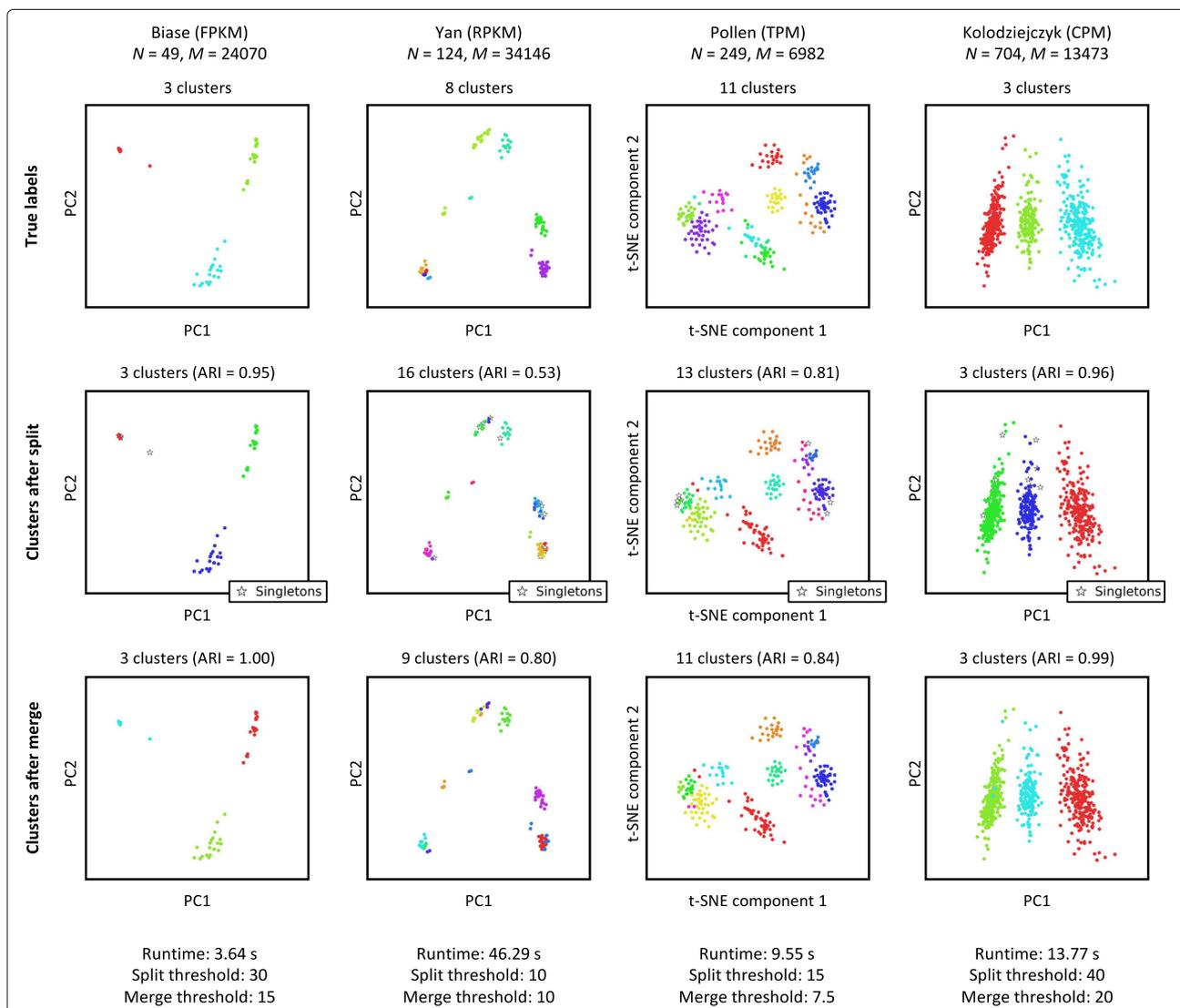


Fig. 3 Gold standard datasets. DendroSplit is evaluated on four single-cell RNA-Seq datasets where the labels are highly likely to be correct [2, 5, 8, 9, 34]. In addition to visual inspection, cluster quality is evaluated using the adjusted Rand index (ARI) based on the true labels. We observe here that the split step tends to generate more clusters than expected, shrinking the ARI. Additionally, due to how the dendrogram is constructed, a cell may end up in its own cluster and is consequently labeled as a “Singleton”. The merge step treats both these cases. The cells are visualized using either the first two principal components (PC) or the first two t-distributed stochastic neighbor embedding [63] (t-SNE) components. The reported runtimes include computation of the pairwise distance matrices

Synthetic datasets

Figure 2 shows the performance of DendroSplit on four synthetic datasets [44]. Since the 2-dimensional data points have clear, intuitive clustering structure, pairwise Euclidean distance is a natural choice. Figure 2a shows the exploratory power of DendroSplit on a toy dataset of oddly-oriented clusters. Because DendroSplit saves the information gathered at each valid split, we can easily investigate how the clustering was performed. At a given split, we can identify the points that went into each partition and look at the partition-specific distributions of the feature that validated the split. Thus the true advantage of DendroSplit is in its ability to justify its behavior with interpretable results. Figure 2b shows that DendroSplit has the power to uncover several clusters especially when the distance metric (Euclidean) suits the type of data (2-D Gaussian balls). Figure 2c emphasizes that, like other methods, DendroSplit cannot automatically overcome poor choices in preprocessing and distance metric selection.

Single-cell RNA-Seq datasets

Figure 3 shows the performance of DendroSplit on four single-cell RNA-Seq datasets featuring high-quality labels. Kiselev et al. [34] refers to these datasets as “gold standards”. We chose four datasets with varying amounts of cells, genes, and total clusters to understand how they affect the behavior of DendroSplit. We see that when N is on the order of 100s, the runtime is widely determined by M , the number of independent Welch’s t -tests that must be performed at every split. Figure 3 shows that for the Biase et al. [2], Pollen et al. [8], and Kolodziejczyk et al. [5] datasets, most of the final ARI is achieved after the split step. Therefore most of the information captured by the clusters lies in one of the dendrogram’s subtrees. Due to how the dendrogram is constructed, a cell may end up being split off into its own cluster and is temporarily labeled as a non-classified “Singleton”. The merge step cleans up singletons and small clusters, resulting in a higher ARI. For the Yan et al. [9] dataset, the ARI increases dramatically after the merge step. This is due to the fact

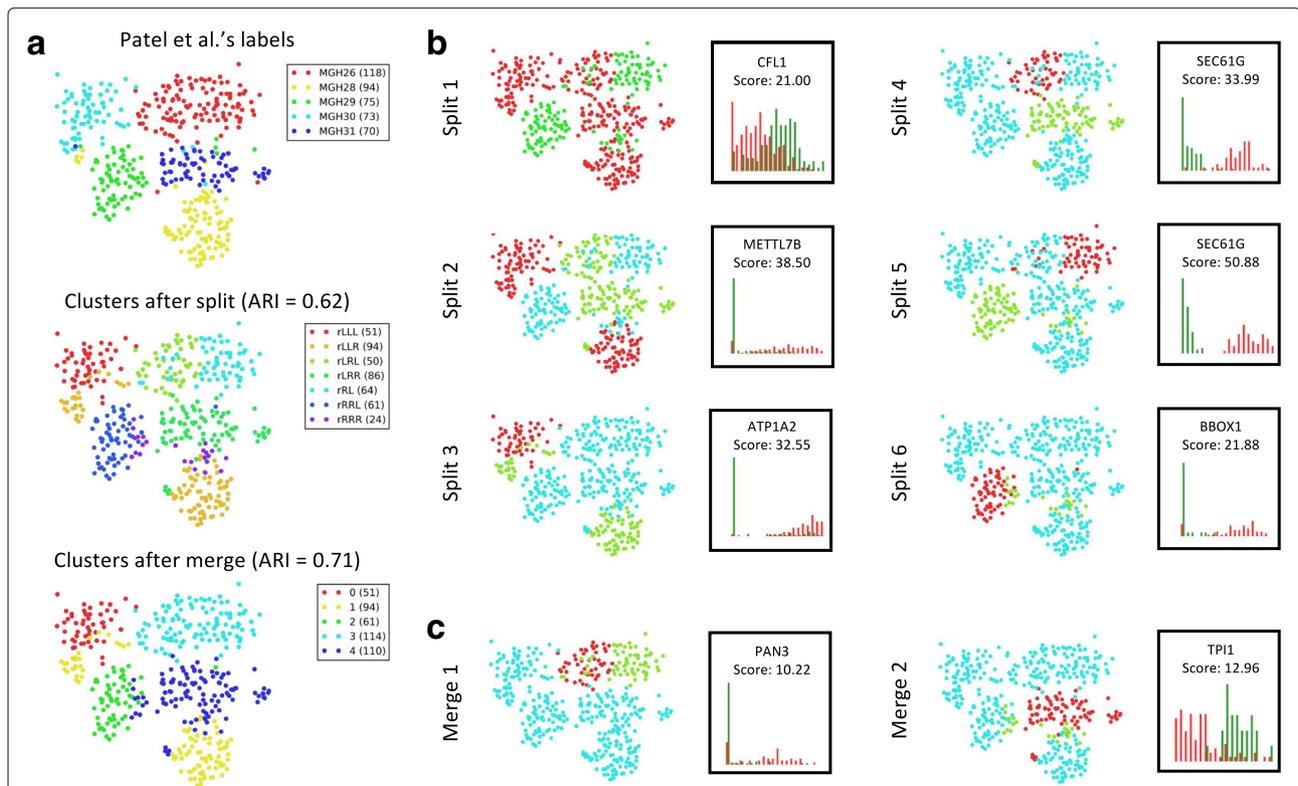


Fig. 4 Exploratory analysis on Patel et al. dataset. **a** DendroSplit is evaluated on Patel et al.’s dataset of 430 cells, 5948 features (genes) from five primary human glioblastomas [7]. Gene expression is quantified using TPM. The split and merge thresholds are 20 and 15, respectively, and the analysis takes 9.64 seconds to run. The numbers in the legends represent the number of points in the corresponding clusters. For the split step, the names of the clusters are generated based on the position of the subtrees in the dendrogram. “r” represents the root node, and “rRL” represents the subtree found at the left child of the right child of the root. **b** We can evaluate how cells were partitioned at each step of the split procedure, and DendroSplit can also show us the within-cluster distributions of the gene that validates the split. **c** We can also evaluate how clusters obtained after the split step were combined during the merge procedure, and DendroSplit can show us the distributions of the most distinguishing gene between two merged clusters

that after the split step, 1) 15 of the 124 cells ended up as singletons, and 2) splitting generated twice as many clusters as needed. In fact, for this dataset, dividing each true cluster into two equal-sized parts would result in an ARI of 0.74 when compared with the original labels. A more detailed visual analysis of the Yan et al. dataset is given in Additional file 1: Figure S1. Under certain conditions, some cells may remain in their own clusters even after the merge step (see “Implementation” section). These cells are analogous to outliers.

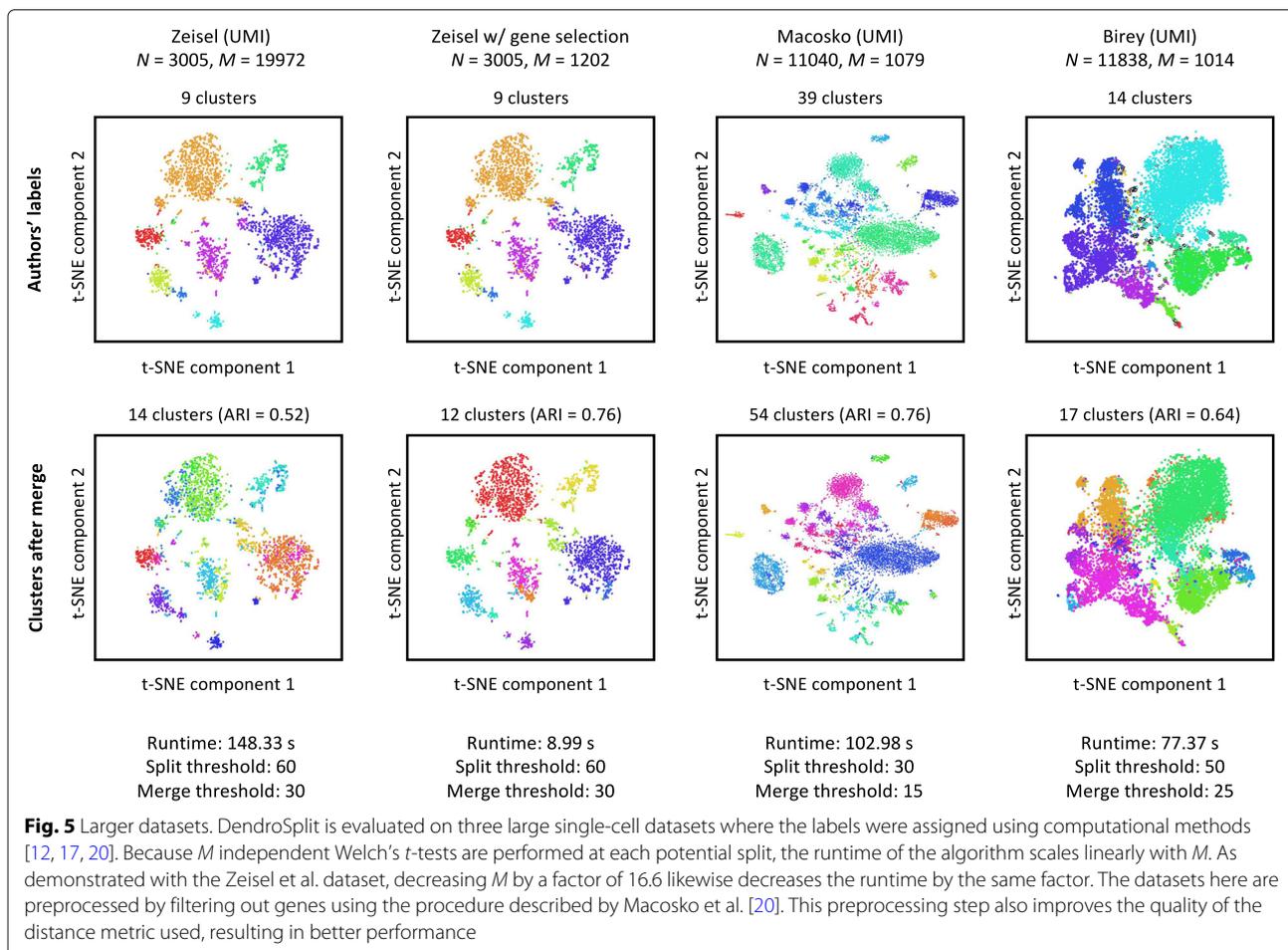
Exploratory analysis

We further demonstrate the exploratory power of DendroSplit on Patel et al.’s [7] dataset of 430 cells, 5948 features from five primary human glioblastomas. Without any further preprocessing, DendroSplit recovers five clusters corresponding to the five glioblastomas (Fig. 4a). Furthermore, DendroSplit can justify its findings by showing us the gene that plays the largest role in validating each split. Splits 4 and 5 in Fig. 4b show distinctively how *SEC61G*, for example, distinguishes MGH26 cells from MGH29 and MGH31 cells. The analysis also gives

insight on how the hierarchical clustering was performed. The cells from MGH26 were split in half during earlier stages of clustering, which is why they end up in separate superclusters at the root node. This is an artifact of the greedy nature of hierarchical clustering where clusters that should be close together may end up far apart. Merge 1 in Fig. 4c shows DendroSplit fixing this. At the same time, we see that *PAN3* may be a valid marker for distinguishing these two subtypes within MGH26 cells. Further analysis and perhaps side information would be needed to decide whether or not these two subtypes are truly different. DendroSplit handles the subjectiveness associated with clustering by showing the factors that contribute to its decisions.

Performance on larger single-cell datasets

We use DendroSplit to re-analyze three large single-cell RNA-Seq datasets that utilize unique molecular identifiers (UMIs) for quantifying genes [12, 17, 20]. Unlike for previous single-cell RNA-Seq datasets, the labels for these datasets were assigned using diverse computational methods. Figure 5 first shows that performing



a feature selection step prior to analysis with DendroSplit decreases the runtime dramatically. In fact, for the Zeisel et al. dataset, filtering out genes using the procedure described by Macosko et al. reduces both M and the runtime by a factor of 16.6. Additionally, the filtering out of noisy features improves the quality of the distance metric, and we see that the ARI improves dramatically. We also report that using a much smaller split threshold of 15 results in 43 non-singleton clusters. When compared with Zeisel et al.'s 47 subclasses, we achieve an ARI of 0.42. The gene filtering procedure is used for the all datasets presented in Figs. 5 and 6.

For the three datasets analyzed in Fig. 5, DendroSplit generates similar but not identical labels. Figure 6a shows that DendroSplit disagrees even more strongly on Zheng et al.'s dataset of 17426 cells, 908 features from fresh peripheral blood mononuclear cells (PBMCs). Noting that the merge step does not increase the ARI significantly, we focus on the split step labels. Although 15 valid splits were recorded, we investigate only the 5 shown in Fig. 6b. For the remaining splits, see Additional file 1: Figure S2. Split 1 was validated due to a lack of expression of several genes (*FCGR3A*, *LY86*, *FCN1*, and *IFI30*) in the red population, which we match to the authors' CD34+ cells.

Split 2 shows the separation of the red cells from the green cells based on high expression of *NKG7* and *GNLY*, markers for natural killer (NK) cells. The green cluster in split 5 likely corresponds to cytotoxic T cells based on increased expression of *GZMH*. The red cluster in split 9 shows greater expression of *CD79A* and may therefore represent B cells. Finally, the red cluster in split 14 does not have an obvious match with any of Zheng et al.'s set of labels. DendroSplit shows us that the existence of this cluster is justified based on increased expression of several genes including *FCGR3A*, *CFD*, and *LST1*. A one-versus-rest differential expression analysis based on independent Welch's t -tests (see Additional file 2: Table S1) further shows that *PSAP* and *SERPINA1* are also overexpressed, indicating that the red cluster (cluster 6 after the merge step in Fig. 6a) corresponds to some type of monocyte. We also repeat this analysis with the full 68579 cells, 20374 genes dataset, and the results are shown in Additional file 1: Figure S3.

Finally, Fig. 7 demonstrates the score threshold sweeping procedure for the Kolodziejczyk et al. and Zeisel et al. datasets. As observed in the experiments, larger datasets often require larger thresholds due to the t -statistic generally increasing with N .

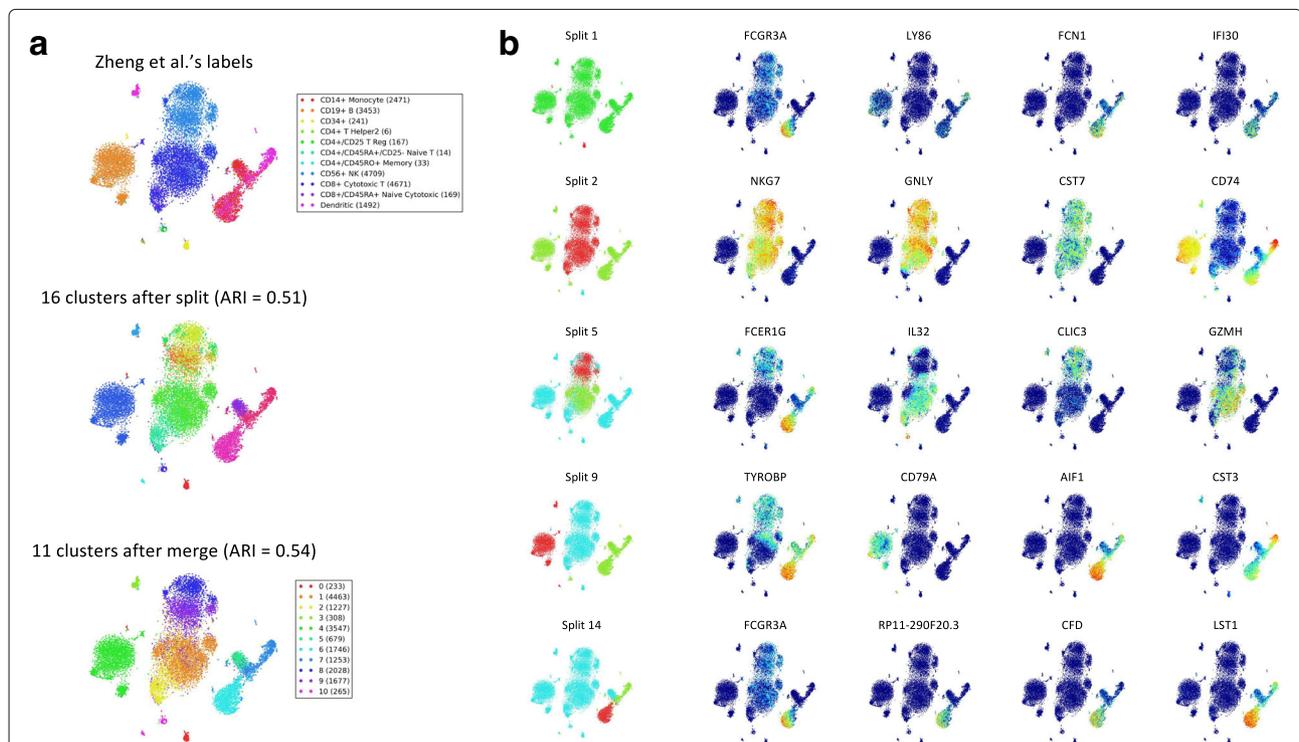
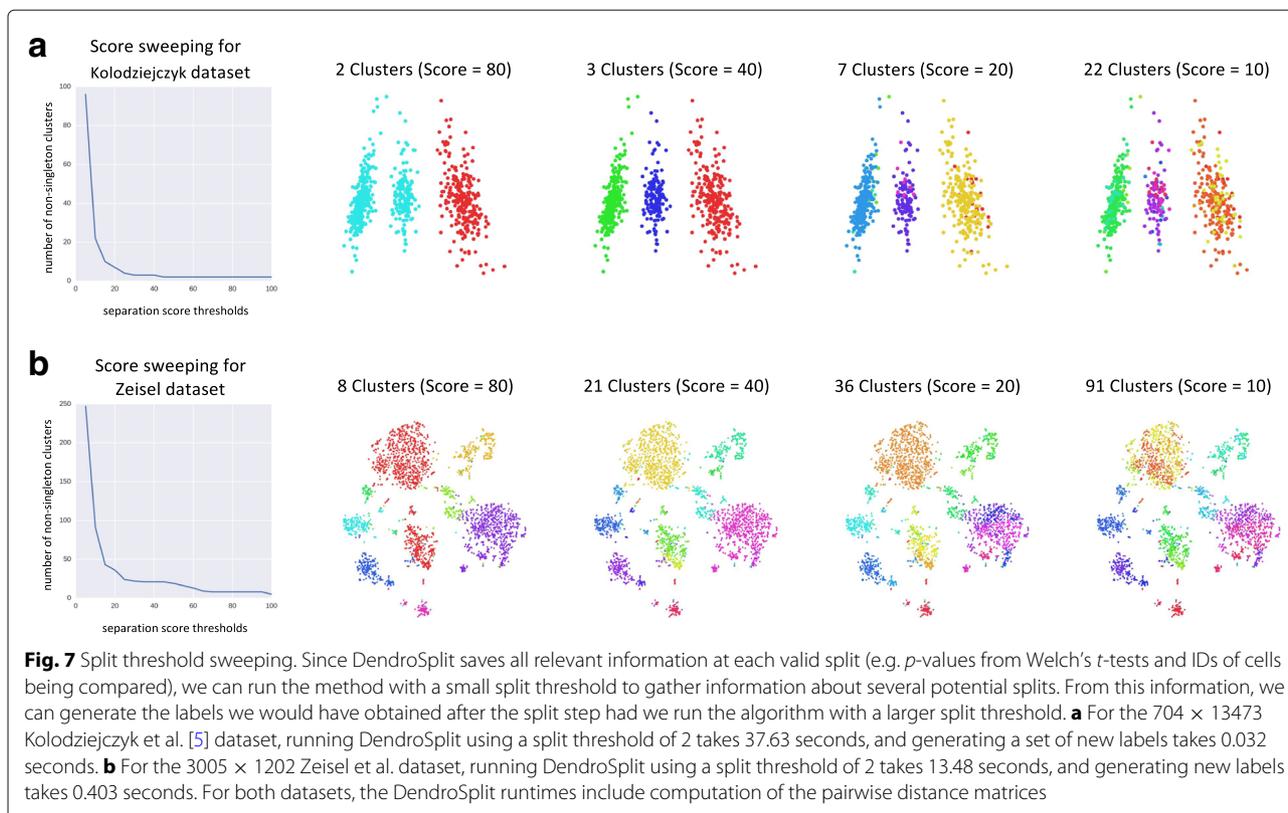


Fig. 6 Exploratory analysis on PBMC dataset. **a** After gene selection and removal of cells with less than 50 counts across all genes, DendroSplit generates clusters for Zheng et al.'s remaining dataset of 17426 cells, 908 features (genes) from fresh peripheral blood mononuclear cells (PBMCs) [21]. Gene expression is quantified using UMI counts. The split and merge thresholds are 200 and 100, respectively, and the analysis takes 119.97 seconds to run. **b** 5 of the 15 recorded valid splits are shown along with the expression levels of the top 4 genes used for validating each split. The reported runtimes include computation of the pairwise distance matrices



Conclusion

In this work, we presented a novel interpretable framework for tackling the single-cell RNA-Seq clustering problem. We demonstrated that a dendrosplit-splitting approach based on a separation score was key for uncovering the multiple layers of biological information within a dataset. In addition to recovering results from a diverse set of single-cell studies, we showed that the framework could cheaply produce several clusterings of the same dataset. Most importantly, the algorithm could justify each of its decisions in an interpretable way. Thus, DendroSplit is suitable as a backend algorithm for interactive analysis and interpretation.

With single-cell RNA-Seq technology improving, we can only expect increased cell throughput and larger datasets. While DendroSplit is able to generate clusters without expensive hyperparameter tuning, its optimal split and merge thresholds do depend on the size of the dataset since larger datasets tend to yield smaller p -values. To remove this size dependence, one could subsample a larger dataset to the same fixed size multiple times, run DendroSplit on each subsample, and ultimately report some consensus result. Another strategy for handling this dependence is in choosing a dataset-size-correcting statistical test rather than the naive Welch's t -test when computing the separation score.

For the analyses in this work, we used a separation score based on a computationally cheap method of performing differential expression and a simple definition of cell type. Separation scores based on more complex methods of evaluating differential expression such as those presented by [31, 45–51] may yield better results at the cost of greater computation. Additionally, just like for other clustering approaches, existing tools including those designed for outlier detection [13, 52], drop-out imputation [53], and correcting other sources of technical noise [54–62] can be easily incorporated into the DendroSplit framework by applying the desired correction procedures before the clustering step.

Availability and requirements

Project name: DendroSplit

Project home page: <https://github.com/jessem-zhang/dendrosplit>

Operating system(s): Platform independent

Programming language: Python 2.7

Other requirements: Python modules numpy 1.12.1, scipy 0.19.0, matplotlib 1.5.3, sklearn 0.18.1, networkx 1.11, community

License: Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International license

Any restrictions to use by non-academics: License needed

Additional files

Additional file 1: Supplementary figures and separation score analysis. A file containing a distance-metric interpretation of separability score along with the three supplementary figures described in the main text: 1) visual analysis of the Yan et al. dataset, 2) further analysis for the splits generated for the Zheng et al. dataset, 3) analysis on the full Zheng et al. dataset. (PDF 3707 kb)

Additional file 2: PBMC differential expression analysis. A one-versus-rest differential expression analysis based on independent Welch's *t*-tests for the Zheng et al. PBMC dataset of 17426 cells, 908 features described in the main text. (TSV 10 kb)

Abbreviations

ARI: Adjusted Rand index; backSPIN: Back sorting points into neighborhoods; DBSCAN: Density-based spatial clustering of applications with noise; NK: Natural killer cells; PBMC: Peripheral blood mononuclear cells; RNA-Seq: Ribonucleic acid sequencing

Acknowledgements

We thank Govinda Kamath and Vasilis Ntranos for feedback and useful discussions about analyzing single-cell RNA-Seq datasets.

Funding

JMZ and DNT are supported in part by the National Human Genome Research Institute of the National Institutes of Health under award number R01HG008164. JMZ performed part of this work as an employee of BD genomics.

Availability of data and materials

The [2, 5, 7–9, 12, 17, 20, 21] datasets were obtained from public repositories. The [5, 8] datasets were obtained from the GitHub repository for [30].

Authors' contributions

JMZ conceived the idea of splitting dendrograms and merging clusters using *p*-values, wrote the software package, performed analyses of data, interpreted results and wrote the manuscript. JF, HCF, DR, and DNT interpreted results, supervised the project and wrote the manuscript. All authors read and approved the final manuscript.

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

JMZ is a past employee of BD Genomics. JF, HCF, and DR are past and current employees of BD Genomics.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Department of Electrical Engineering, Stanford, 94305 Stanford, California, USA. ²BD Genomics, 94025 Menlo Park, California, USA.

Received: 22 November 2017 Accepted: 28 February 2018

Published online: 09 March 2018

References

- Yuan GC, et al. Challenges and emerging directions in single-cell analysis. *Genome Biol.* 2017;18:84. <https://doi.org/10.1186/s13059-017-1218-y>.
- Biase FH, Cao X, Zhong S. Cell fate inclination within 2-cell and 4-cell mouse embryos revealed by single-cell rna sequencing. *Genome Res.* 2014;24:1787–96.
- Trapnell C, et al. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol.* 2014;32:381–6.
- Goolam M, et al. Heterogeneity in oct4 and sox2 targets biases cell fate in 4-cell mouse embryos. *Cell.* 2016;165:61–74. <http://www.sciencedirect.com/science/article/pii/S0092867416300617>.
- Kolodziejczyk AA, et al. Single cell rna-sequencing of pluripotent states unlocks modular transcriptional variation. *Cell Stem Cell.* 2015;17:471–85. <http://www.sciencedirect.com/science/article/pii/S193459091500418X>.
- Treutlein B, et al. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell rna-seq. *Nature.* 2014;509:371–5. <https://doi.org/10.1038/nature13173>.
- Patel AP, et al. Single-cell rna-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science.* 2014;344:1396–401. <http://science.sciencemag.org/content/344/6190/1396>. <http://science.sciencemag.org/content/344/6190/1396.full.pdf>.
- Pollen AA, et al. Low-coverage single-cell mrna sequencing reveals cellular heterogeneity and activated signaling pathways in developing cerebral cortex. *Nat Biotech.* 2014;32:1053–8. <https://doi.org/10.1038/nbt.2967>.
- Yan L, et al. Single-cell rna-seq profiling of human preimplantation embryos and embryonic stem cells. *Nat Struct Mol Biol.* 2013;20:1131–9. <https://doi.org/10.1038/nsmb.2660>.
- Buettner F, et al. Computational analysis of cell-to-cell heterogeneity in single-cell rna-sequencing data reveals hidden subpopulations of cells. *Nat Biotech.* 2015;33:155–60. <https://doi.org/10.1038/nbt.3102>.
- Usoskin D, et al. Unbiased classification of sensory neuron types by large-scale single-cell rna sequencing. *Nat Neurosci.* 2015;18:145–53. <https://doi.org/10.1038/nn.3881>.
- Zeisel A, et al. Cell types in the mouse cortex and hippocampus revealed by single-cell rna-seq. *Science.* 2015;347:1138–42. <http://science.sciencemag.org/content/347/6226/1138>. <http://science.sciencemag.org/content/347/6226/1138.full.pdf>.
- Grun D, et al. Single-cell messenger rna sequencing reveals rare intestinal cell types. *Nature.* 2015;525:251–5. <http://doi.org/10.1038/nature14966>.
- Ting DT, et al. Single-cell (RNA) sequencing identifies extracellular matrix gene expression by pancreatic circulating tumor cells. *Cell Rep.* 2014;8:1905–18. <http://www.sciencedirect.com/science/article/pii/S2211124714007050>.
- Lake BB, et al. Neuronal subtypes and diversity revealed by single-nucleus rna sequencing of the human brain. *Science.* 2016;352:1586–90. <http://science.sciencemag.org/content/352/6293/1586>. <http://science.sciencemag.org/content/352/6293/1586.full.pdf>.
- Deng Q, Ramsköld D, Reinius B, Sandberg R. Single-cell rna-seq reveals dynamic, random monoallelic gene expression in mammalian cells. *Science.* 2014;343:193–6. <http://science.sciencemag.org/content/343/6167/193>. <http://science.sciencemag.org/content/343/6167/193.full.pdf>.
- Birey F, et al. Assembly of functionally integrated human forebrain spheroids. *Nature.* 2017;545:54–9. <http://doi.org/10.1038/nature22330>.
- Fan HC, Fu GK, Fodor SPA. Combinatorial labeling of single cells for gene expression cytometry. *Science.* 2015;347. <http://science.sciencemag.org/content/347/6222/1258367>. <http://science.sciencemag.org/content/347/6222/1258367.full.pdf>.
- Klein AM, et al. Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. *Cell.* 2015;161:1187–201. <http://www.sciencedirect.com/science/article/pii/S0092867415005000>.
- Macosko EZ, et al. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. *Cell.* 2015;161:1202–14. <http://doi.org/10.1016/j.cell.2015.05.002>.
- Zheng GXY, et al. Massively parallel digital transcriptional profiling of single cells. *Nat Commun.* 2017;8:14049. EP – <http://doi.org/10.1038/ncomms14049>.
- Cao J, et al. Comprehensive single-cell transcriptional profiling of a multicellular organism. *Science.* 2017;357:661–7. <http://science.sciencemag.org/content/357/6352/661>. <http://science.sciencemag.org/content/357/6352/661.full.pdf>.
- Trapnell C. Defining cell types and states with single-cell genomics. *Genome Res.* 2015;25:1491–8. <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC4579334/>.
- Ester M, Kriegel H-P, Sander J, Xu X, et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In: *Kdd. Portland: KDD-96; 1996.* p. 226–31.
- Dueck D, Frey BJ. Non-metric affinity propagation for unsupervised image categorization. In: *Computer Vision, 2007. ICCV 2007. IEEE 11th International Conference on. Rio de Janeiro: IEEE; 2007.* p. 1–8.

26. Ng AY, Jordan MI, Weiss Y. On spectral clustering: Analysis and an algorithm. In: *Advances in neural information processing systems*. Cambridge: MIT Press; 2002. p. 849–56.
27. Ntranos V, Kamath GM, Zhang JM, Pachter L, David NT. Fast and accurate single-cell rna-seq analysis by clustering of transcript-compatibility counts. *Genome Biol*. 2016;17:112.
28. Rostom R, Svensson V, Teichmann SA, Kar G. Computational approaches for interpreting scrna-seq data. *FEBS Lett*. 2017;17:112.
29. Pierson E, Yau C. Zifa: Dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biol*. 2015;16:241.
30. Wang B, Zhu J, Pierson E, Ramazzotti D, Batzoglou S. Visualization and analysis of single-cell rna-seq data by kernel-based similarity learning. *Nat Meth*. 2017;14:414–6. <http://doi.org/10.1038/nmeth.4207>.
31. Guo M, Wang H, Potter SS, Whitsett JA, Xu Y. Sincera: a pipeline for single-cell rna-seq profiling analysis. *PLoS Comput Biol*. 2015;11:e1004575.
32. Xu C, Su Z. Identification of cell types from single-cell transcriptomes using a novel clustering method. *Bioinformatics*. 2015;31:1974–80.
33. Žurauskienė J, Yau C. pcareduce: hierarchical clustering of single cell transcriptional profiles. *BMC Bioinformatics*. 2016;17:140. <https://doi.org/10.1186/s12859-016-0984-y>.
34. Kiselev VY, et al. Sc3: consensus clustering of single-cell rna-seq data. *Nat Meth*. 2017;14:483–6. <http://doi.org/10.1038/nmeth.4236>.
35. Satija R, Farrell JA, Gennert D, Schier AF, Regev A. Spatial reconstruction of single-cell gene expression data. *Nat Biotech*. 2015;33:495–502. <http://doi.org/10.1038/nbt.3192>.
36. Wolf FA, Angerer P, Theis FJ. Scanpy for analysis of large-scale single-cell gene expression data. *bioRxiv*. 2017;174029.
37. Lun AT, McCarthy DJ, Marioni JC. A step-by-step workflow for low-level analysis of single-cell rna-seq data with bioconductor. *F1000Research*. 2016;5:2122.
38. Qiu X, et al. Single-cell mrna quantification and differential analysis with census. *Nat Methods*. 2017;14:309–15.
39. McCarthy DJ, Campbell KR, Lun AT, Wills QF. Scater: pre-processing, quality control, normalization and visualization of single-cell rna-seq data in r. *Bioinformatics*. 2017;33:1179–86.
40. Langfelder P, Zhang B, Horvath S. Defining clusters from a hierarchical cluster tree: the dynamic tree cut package for r. *Bioinformatics*. 2007;24:719–20.
41. Pedregosa F, et al. Scikit-learn: Machine learning in Python. *J Mach Learn Res*. 2011;12:2825–30.
42. Jones E, Oliphant T, Peterson P, et al. SciPy: Open source scientific tools for Python. 2001. <http://www.scipy.org/>. Accessed July 2016.
43. Ruxton GD. The unequal variance t-test is an underused alternative to student's t-test and the mann–whitney u test. *Behav Ecol*. 2006;17:688–90.
44. Franti P. Clustering datasets. 2015. <http://cs.uef.fi/sipu/datasets/>. Accessed July 2017.
45. Love M, Anders S, Huber W. Differential analysis of count data—the *deseq2* package. *Genome Biol*. 2014;15:550.
46. Robinson MD, McCarthy DJ, Smyth GK. *edgeR*: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26:139–40.
47. Finak G, et al. *Mast*: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell rna sequencing data. *Genome Biol*. 2015;16:278.
48. Andrews TS, Hemberg M. Modelling dropouts allows for unbiased identification of marker genes in scRNAseq experiments. *bioRxiv*. 2016;065094.
49. Kharchenko PV, Silberstein L, Scadden DT. Bayesian approach to single-cell differential expression analysis. *Nat Methods*. 2014;11:740–2.
50. Fan J, et al. Characterizing transcriptional heterogeneity through pathway and gene set overdispersion analysis. *Nat Meth*. 2016;13:241–4. <http://doi.org/10.1038/nmeth.3734>.
51. Korthauer KD, et al. A statistical approach for identifying differential distributions in single-cell rna-seq experiments. *Genome Biol*. 2016;17:222.
52. Jiang L, Chen H, Pinello L, Yuan G-C. GiniClust: detecting rare cell types from single-cell gene expression data with gini index. *Genome Biol*. 2016;17:144.
53. Lin P, Troup M, Ho JW. Cidr: Ultrafast and accurate clustering through imputation for single-cell rna-seq data. *Genome Biol*. 2017;18:59.
54. Vallejos CA, Marioni JC, Richardson S. *Basics*: Bayesian analysis of single-cell sequencing data. *PLoS Comput Biol*. 2015;11:e1004333.
55. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical bayes methods. *Biostatistics*. 2007;8:118–27.
56. Benito M, et al. Adjustment of systematic microarray data biases. *Bioinformatics*. 2004;20:105–14.
57. Leek JT, Johnson WE, Parker HS, Jaffe AE, Storey JD. The *sva* package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*. 2012;28:882–3.
58. Leng N, et al. *Oefinder*: a user interface to identify and visualize ordering effects in single-cell rna-seq data. *Bioinformatics*. 2016;32:1408–10.
59. Brennecke P, et al. Accounting for technical noise in single-cell rna-seq experiments. *Nat Methods*. 2013;10:1093–5.
60. Illicic T, et al. Classification of low quality cells from single-cell rna-seq data. *Genome Biol*. 2016;17:29.
61. Risso D, Ngai J, Speed TP, Dudoit S. Normalization of rna-seq data using factor analysis of control genes or samples. *Nat Biotechnol*. 2014;32:896–902.
62. Grün D, Kester L, Van Oudenaarden A. Validation of noise models for single-cell transcriptomics. *Nat Methods*. 2014;11:637–40.
63. Maaten Lvd, Hinton G. Visualizing data using t-sne. *J Mach Learn Res*. 2008;9:2579–605.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at
www.biomedcentral.com/submit

