

Ntranos *et al.*

METHOD

Fast and accurate single-cell RNA-Seq analysis by clustering of transcript-compatibility counts

Vasilis Ntranos^{1†}, Govinda M. Kamath^{2†}, Jesse Zhang^{2†}, Lior Pachter^{*3} and David N. Tse^{*2,1}

* Correspondence:
lpachter@math.berkeley.edu
dntse@stanford.edu

³Departments of Mathematics and Molecular and Cell Biology, University of California, Berkeley,
²Department of Electrical Engineering, Stanford University,
Full list of author information is available at the end of the article
†Equal contributor

Abstract

Current approaches to single-cell transcriptomic analysis are computationally expensive and require assay-specific modeling which limit their scope and generality. We propose a novel method that departs from standard analysis pipelines, comparing and clustering cells based not on their transcript or gene quantifications but on their transcript-compatibility read counts. In re-analysis of two landmark yet disparate single-cell RNA-Seq datasets, we show that our method is up to two orders of magnitude faster than previous approaches, provides accurate and in some cases improved results, and is universal, being directly applicable to data from a wide variety of assays.

Keywords: Single-cell RNA seq, clustering, pseudoalignment, transcript compatibility counts

Introduction

Single-cell RNA-Seq (scRNA-Seq) has proved to be a powerful tool for probing cell states [1–5], defining cell types [6–9], and describing cell lineages [10–13]. These applications of scRNA-Seq all rely on two computational steps: quantification of gene or transcript abundances in each cell and clustering of the data in the resulting abundance \times cell expression matrix [14, 15]. There are a number of challenges in both of these steps that are specific to scRNA-Seq analysis. While methods for transcript/gene abundance estimation from bulk RNA-Seq have been extensively tested and benchmarked [16], the wide variety of assay types in scRNA-Seq [17–25] have required a plethora of customized solutions [2, 6, 7, 9, 11–13, 24, 26–37] that are difficult to compare to each other. Furthermore, the quantification methods used all rely on read alignment to transcriptomes or genomes, a time consuming step that will not scale well with the increasing numbers of reads predicted for scRNA-Seq [15, 38]. Clustering based on scRNA-Seq expression matrices can also require domain specific information, e.g. temporal information [33] or functional constraints [37] so that in some cases hand curation of clusters is performed after unsupervised clustering [7].

In [39], a method of collapsing bulk read alignments into “equivalence classes” of reads was introduced for the purpose of estimating alternative splicing isoform frequencies from bulk RNA-Seq data. Each equivalence class consists of all the reads that are compatible with the same set of transcripts. (See Figure 1 for an example.) The collapsing of reads into equivalence classes was initially introduced to allow for significant speedup of the E-step in the expectation-maximization (EM) algorithm used in some RNA-Seq quantification programs [40, 41], as the read counts

in the equivalence classes, or *transcript-compatibility counts* (TCC), correspond to the sufficient statistics for a standard RNA-Seq model [42]. In other words, the use of transcript-compatibility counts was an intermediate computation step towards quantifying transcript abundances. In this paper we instead consider the direct use of such counts for the comparison and clustering of scRNA-Seq cells. Figure 2 shows an outline of a method we have developed for clustering and analyzing scRNA-Seq data; the key idea is to base clustering not on the quantification of transcripts or genes but on the transcript-compatibility counts for each cell.

To better understand the relevance of transcript-compatibility counts, consider their relationship to “gene-level” counts used in many RNA-Seq analyses. In the same way that “genes” represent groupings of transcripts [43], equivalence classes as introduced by [39] are also groups of transcripts. However while the former is a biologically motivated construction, the latter is technical, consisting of groupings that capture the extent of ambiguous multiple mappings among reads. The lack of direct biological interpretation of equivalence classes makes transcript-compatibility counts less intuitive; however, as we will show, there are two significant advantages to working with them: 1) unlike transcript or gene-level quantifications, transcript-compatibility counts can be computed without a read-generation model, and hence a single clustering pipeline based on transcript-compatibility counts can be used across a wide range of scRNA-Seq assays; 2) transcript-compatibility counts can be computed by pseudoalignment [41], a process that does not require read alignment and can be done extremely efficiently^[1]. Hence the transcript-compatibility based approach is fast, accurate and universal.

To demonstrate both the general applicability of our method as well as its accuracy, we re-analyzed data from two recently published scRNA-Seq papers: the pseudotemporal ordering of primary human myoblasts by [12] and the cell classification in the mouse cortex and hippocampus by [7]. We show that not only are we able to recapitulate the analyses of the papers two orders of magnitude faster than previously possible, but we also provide a refinement of the published results, suggesting that our approach is both fast and accurate. We show that the speedup of our approach will be important as single-cell RNA sequencing scales to more cells and illustrate the advantages of the universality of our approach via its application to a multitude of assays without the requirement for numerous parameters. Existing pipelines lack either efficiency or universality, or both.

Results

Transcript-compatibility counts from pseudoalignments

To demonstrate the effectiveness of transcript-compatibility counts for scRNA-Seq analysis, we first examined how efficiently they can be computed. While transcript-compatibility counts can be extracted from read alignments (e.g. in SAM/BAM

^[1]In our implementation of the method, we use kallisto to compute transcript-compatibility counts via pseudoalignment (we avoid the quantification step that is usually performed when running kallisto altogether). A naive pipeline based on using kallisto to quantify transcript abundances and then clustering cells based on the quantifications would not be accurate if there is a mismatch between its read-generation model and the single-cell data; see Results section for an example.

format), they do not require the full information contained in alignments. Instead, we examined the speedup possible with pseudoalignment [41], which obtains for each read the set of transcripts it is compatible with and therefore can be directly used to obtain transcript-compatibility counts.

Figure 3 shows the speed of obtaining transcript-compatibility counts via pseudoalignment in comparison to the time required to quantify RNA-Seq data with other approaches. The key result relevant for single-cell analysis is the scalability of pseudoalignment for obtaining transcript-compatibility counts (Figure 3 and Supplementary Figure 2). The fixed extra cost for aligning (rather than pseudoaligning) reads for each cell is small, but when extrapolated to hundreds of thousands of cells becomes a significant (computational) cost.

Pseudotime for differentiating human myoblasts

The recently published Monocle software [12] that builds on the Cufflinks program [44] is rapidly becoming a standard tool for scRNA-Seq analysis. We therefore sought to compare our approach to Monocle, and in order to do so began with a re-analysis of the data in [12]. Figure 4 shows the temporal ordering of differentiating primary human myoblasts using transcript-compatibility counts clustering based on the Jensen-Shannon metric and the affinity propagation algorithm (see Methods). We note that unlike Cufflinks, which consists of an explicit model of RNA-Seq suitable for the data in [12] but not necessarily for other assays, our transcript-compatibility counts make no assumption about the nature of the data. Furthermore, while the re-analysis appears to match that of [12], affinity propagation with different parameters provided a more refined clustering, possibly capturing seven stages of myoblast differentiation (see also Supplementary Figure 3).

A central idea in pseudo-temporal ordering of cells relies upon the construction of a minimum spanning tree (MST) over the pairwise distances of their corresponding gene expression vectors [45]. This attempts to capture the trajectory of a hypothetical cell that gradually “moves” through different cellular states or differentiation stages in a high-dimensional gene expression space. Our results show that the same concept can be applied to transcript-compatibility counts. A key step in Monocle is to *first* reduce the dimensionality of the data by independent component analysis (ICA) and *then* compute the MST based on Euclidean distances on the plane. Here we take a different approach and compute the MST on “cluster centers” in high dimensions (See Methods). Both approaches aim to battle the biological and technical noise that is inevitably introduced in scRNA-Seq experiments. Even though we could have used Monocle directly on transcript-compatibility counts, the design and comparison of specialized tools is beyond the scope of this paper.

Figure 4d validates the three primary clusters and the pseudo-temporal ordering obtained by our method based on three key myoblast differentiation markers, *MYOG*, *CDK1* and *PDGFRA* (see Supplementary Figure 4 for an additional set of genes taken from [12]). Interestingly, the expression of these genes gradually evolves over the pseudo-temporally ordered clusters, capturing both the underlying differentiation trajectory of proliferating cells to myoblasts, and the corresponding branching towards mesenchymal cells, as was observed in [12].

Finally, we should point out that although the three primary clusters of [12] are evident in our results, they are not identical. This naturally raises the question

of whether clustering on (high-dimensional) transcript-compatibility counts could possibly lead to cell mis-classification. Our results show that this is not the case. In Figure 5 we investigated one cell that seemed to have been severely mis-classified by our method as a differentiating myoblast while it was identified as a proliferating cell by Monocle. However, an analysis of the expression levels of 12 marker genes obtained from [12] shows that this cell displays more similarity to differentiating myoblasts than proliferating cells. Overall our results seem to suggest that transcript-compatibility counts, being directly obtained from sequenced reads, might constitute a less noisy representation of the “transcriptomic state” of a cell compared to the one obtained by quantifying its gene expression.

Cell classification in the mouse cortex and hippocampus

The re-analysis of [12] shows that clustering of transcript-compatibility counts can be useful on a single dataset, but we believe that the true power of our approach lies in its universality. In contrast to the standard quantification pipeline, obtaining transcript-compatibility counts does not require a read-generation model; our method can be directly applied to a wide range of scRNA-Seq datasets and transcript-compatibility counts can be used to analyze sequenced reads without any assay-specific information. To make this point, we re-analyzed a recent large scRNA-Seq experiment published earlier this year [7] that uses an assay based on unique molecular identifiers (UMI). In contrast to [12] where paired-end reads were sampled from fragments covering the entire length of the transcripts, [7] used single-end reads that were only obtained from the 3'-end of the transcripts.

Zeisel et al. [7] examined a very diverse population of 3005 cells obtained from the cortical and hippocampal regions of the mouse brain. In order to analyze this complex dataset, the authors developed a state-of-the-art hierarchical bi-clustering method called BackSPIN (based on SPIN [30]) and were able to identify 47 distinct sub-populations of cells within nine major brain cell types. This fine-grained analysis also revealed a previously unknown post-mitotic oligodendrocyte sub-class, referred to as Oligo1 in [7].

Figure 6 shows the clusters obtained by applying our method to the above dataset and compares our method’s clustering accuracy to various quantification-based methods. In order to systematically assess the clustering accuracy, we iteratively sub-sampled cells from two different cell types at random and evaluated the ability of each method to distinguish between these types. Since the development of specialized clustering algorithms is orthogonal to our paper, we compared based on the same clustering algorithm throughout (see Methods). Our results indicate that transcript-compatibility counts can be more accurate than standard model-based RNA-Seq quantification tools (such as eXpress) that try to estimate the underlying read-generation model from the data. Our transcript-compatibility counts based method is in fact able to achieve similar accuracy with the assay-specific quantification approach used in [7] (that explicitly takes into account the significant 3'-end bias in this dataset). Clustering transcript abundance quantifications output by kallisto results in lower accuracy due to the mismatch between kallisto’s read-generation model and this dataset, further emphasizing the importance of using transcript-compatibility counts which are computed without using any such model.

Quite remarkably, our method (via affinity propagation on all cells) was further able to recover the Oligo1 cluster of cells, showing that transcript-compatibility counts can indeed capture distinct cell signatures without actually quantifying their gene expression (Figure 6, Methods). Overall, in our experiments we observed that unsupervised clustering of transcript-compatibility counts typically yielded more than 47 clusters, which was also the case in [7]. Some of our clusters were very small, probably capturing outlier cells, while others seemed to be further splitting the 47 cell subtypes identified in [7].

To further investigate this, we focused on another oligodendrocyte sub-population, referred to as Oligo3 in [7]. As reported in [7], Oligo3 cells were almost exclusively observed in the somatosensory cortex and were identified by the authors as being in an intermediate stage of maturation – in between premyelinating and myelinating oligodendrocytes. Even though the Oligo3 cells appear to be well-clustered together, as visualized by t-SNE (Figure 7a), affinity propagation on transcript-compatibility counts with various parameters consistently separated them into two sub-clusters. Our results in Figure 7b seem to suggest that a sub-population of Oligo3 cells (captured by one of our sub-clusters) expresses an unusual signature of endothelial/vascular genes on top of the expected myelin related genes. Interestingly, similar findings have been reported recently in [37], suggesting a possible (experimental) contamination of several oligodendrocyte cells in the dataset at hand.

Conclusions

The extraordinary developments in single-cell RNA-Seq technology over the past few years have demonstrated that “single-cell resolution” is not just a gimmick but an unprecedented tool for probing transcriptomes that can reveal the inner-workings of developmental programs and their resulting tissues. However the computational challenges of scRNA-Seq analysis, already very high due to the large number of cells to analyze, have been further exacerbated by the smorgasbord of assays that each introduce unique technical challenges.

The new method we have proposed and evaluated in this paper, namely analysis of scRNA-Seq based on transcript-compatibility counts, offers a universal, efficient and accurate solution for extracting information from scRNA-Seq experiments. In the same way that single-cell analysis can be viewed as the ultimate resolution for transcriptomics, transcript-compatibility counts are the most direct way to “count” reads. While we have focused on clustering of cells in this paper, we believe that transcript-compatibility counts may have applications in many other sequencing-based assays, and that further development of methods based on such counts offers a fruitful avenue of exploration.

The ability to obtain transcript-compatibility counts by pseudoalignment is a benefit that has its own implications and applications. For example, the speed of pseudoalignment facilitated quick experimentation with our method, and in assessing our accuracy on different datasets one discovery was that much less sampling than is currently performed is necessary to cluster cells. In the re-analysis of [7], we found that the main results, namely the clustering of cells and identification of cell types, were achievable with only 1% of the data (see Figure 8a and Supplementary Figure 1). This observation has significant implications for scRNA-Seq as

it suggests that for clustering of cells, low-coverage sequencing may be sufficient thus allowing for larger experiments with more cells. Moreover, this low-coverage clustering performance can be achieved using our method, which is not tailored to the specific scRNA-Seq assay.

Methods

The code used to generate the results presented in this paper is available online [46].

Computation of transcript-compatibility counts

We utilized the "pseudo" option of the kallisto RNA-Seq program which computes equivalence classes of reads after pseudoalignment. We used kallisto version 0.42.3 with k set to kallisto's default value of 31.

Transcript-compatibility counts based on UMI information

The dataset of [7] has reads with unique molecular identifiers (UMIs). UMIs are typically used in scRNA-Seq to correct for PCR bias; biological copies of a transcript (distinct molecules) can be identified based on their UMIs. This information can be utilized in generating the transcript-compatibility counts from equivalence classes. Instead of counting all the reads in each equivalence class, we only count the reads with distinct UMIs. Transcript-compatibility counts with UMIs are shown in Figure 6b and Figure 8a (represented as "TCC with UMI" in the figures).

Clustering Methodology

On obtaining the transcript-compatibility counts for each cell, we normalize by the total number of mapped reads to obtain a probability distribution called the transcript-compatibility count distribution or TCC distribution. We then compute the square-root of the Jensen-Shannon divergence [47] between the TCC distributions for each pair of cells. As a distance metric which satisfies the triangle [48] inequality, the square-root of Jensen-Shannon divergence is a natural choice for computing pairwise distances between two probability distributions. However, the results obtained here are not contingent on using the square root of Jensen-Shannon divergences as the measure of distances, and quite similar results are obtained when we use other distances between probability distributions such as the ℓ_1 distance to compute the pairwise distance matrix, (see Supplementary Figure 1). ℓ_1 distance (which is just twice the total-variation distance) in fact seems to perform better than Jensen-Shannon distance for low coverage (Supplementary Figure 1b). In contrast, Euclidean distance (ℓ_2 distance) seems to perform much worse (see Supplementary Figure 1). The fact that Euclidean distance is not a good distance metric to measure distances between probability distributions is widely documented (see for instance [49]).

All clustering carried out in this paper were done using off-the-shelf clustering methods.

We used spectral clustering using the pairwise distance matrices when we know the number of clusters in the data. This includes Figures 6b, 8a, and Supplementary Figure 1b with 2 clusters for the pairwise distance matrix (from TCC distributions) obtained for the data from [7].

The clustering method used when the number of clusters is not known is affinity propagation [50]. This is an unsupervised clustering algorithm based on message passing, which needs a pairwise similarity matrix as input. The pairwise similarity matrix is computed as the negative of the pairwise distance matrix that was computed.

To evaluate the clustering accuracy of our method in Figure 6*b*, we performed binary classification tests using the labels reported in [7] as the ground truth. In particular, we randomly sub-sampled two different types of cells and evaluated the ability of each pipeline to separate them into two clusters via spectral clustering. We performed these binary classification tests between 1) the sub-classes Oligo1 (45 cells) and Oligo4 (106 cells), 2) the cell types Astrocytes (198 cells) and Interneurons (290 cells), and 3) the more general cell types neurons (1628 cells) and non-neurons (1377 cells). The error rates for each test were obtained by randomly sampling 22, 99 and 200 cells from each of the two labels respectively, averaged over 10 monte-carlo iterations.

For clustering the dataset of [7], we used affinity propagation with preference value set to the median of the similarity scores and the damping parameter set to 0.5. On doing this, we obtained 89 clusters. Of the 89 clusters obtained, cluster number 22 had the largest match with the set of cells the authors labeled as Oligo1 (which was the new type of cells discovered in [7]). 24 out of the 28 cells in the cluster were labeled Oligo1 by [7]. There were a total of 45 cells labeled Oligo1 in [7] out of the total of 3005 considered. This is investigated in Figure 6*c*.

Also, affinity propagation with different parameters seems to split the class labelled Oligo3 in [7] into 2 classes. This is investigated in Figure 7, where the two classes considered were classes obtained with parameters set as before.

For clustering the dataset of [12], we used affinity propagation with preference parameter set to 1.3 and damping parameter set to 0.95 to obtain three clusters in Figure 4. To obtain 8 clusters on the dataset of [12], we used affinity propagation with preference parameter set to 0.6 and damping parameter set to 0.95 to obtain 8 clusters. Then after collapsing any cluster with less than 5 cells into the cluster closest to it, we obtain the seven clusters investigated in figure 5.

Partial order on clusters

On the [12] data set, for the seven clusters obtained, we first find the centroid TCC distribution of each cluster as the mean TCC distribution of all cells in the cluster. Then, we compute the pairwise Jensen-Shannon distances between the centroid TCC distributions (cluster centers). We then run a minimum weight spanning tree on the complete graph between the cluster centers with weights given by the computed pairwise distances. This gives us a partial-order on the clusters, which is investigated in Figure 4 and Supplementary Figure 3.

Visualization of cells and clusters

We used t-SNE [51] to visualize the cells and clusters in Figures 6*a*, 8*b*, and Supplementary Figure 1*a*.

The left panel of Figure 4*d*, 5*a*, 5*b* and Supplementary Figure 4*a* was created using an implementation [52] of the diffusion map algorithm of [53].

Competing interests

The authors declare that they have no competing interests.

Author's contributions

VN,GMK and JZ conceived the idea of clustering without quantification, performed analyses of data, analyzed and interpreted results and wrote the manuscript. DNT and LP interpreted results, supervised the project and wrote the manuscript.

Acknowledgements

We thank Pál Melsted for implementing the pseudo command in kallisto. This is the command that allows for direct output of transcript-compatibility counts via pseudoalignment. Thanks to Bo Li for useful discussions about single-cell RNA-Seq assays and their biases. GMK and JZ are supported by the Center for Science of Information, an NSF Science and Technology Center, under grant agreement CCF-0939370. VN is supported in part by the Center for Science of Information and in part by a gift from Qualcomm Inc. LP is supported in part by the National Human Genome Research Institute of the National Institutes of Health under award number R01HG006129. DNT is supported in part by the Center of Science of Information and in part by the National Human Genome Research Institute of the National Institutes of Health under award number R01HG008164.

Author details

¹Department of Electrical Engineering and Computer Sciences, University of California, Berkeley, ²Department of Electrical Engineering, Stanford University, ³Departments of Mathematics and Molecular and Cell Biology , University of California, Berkeley,

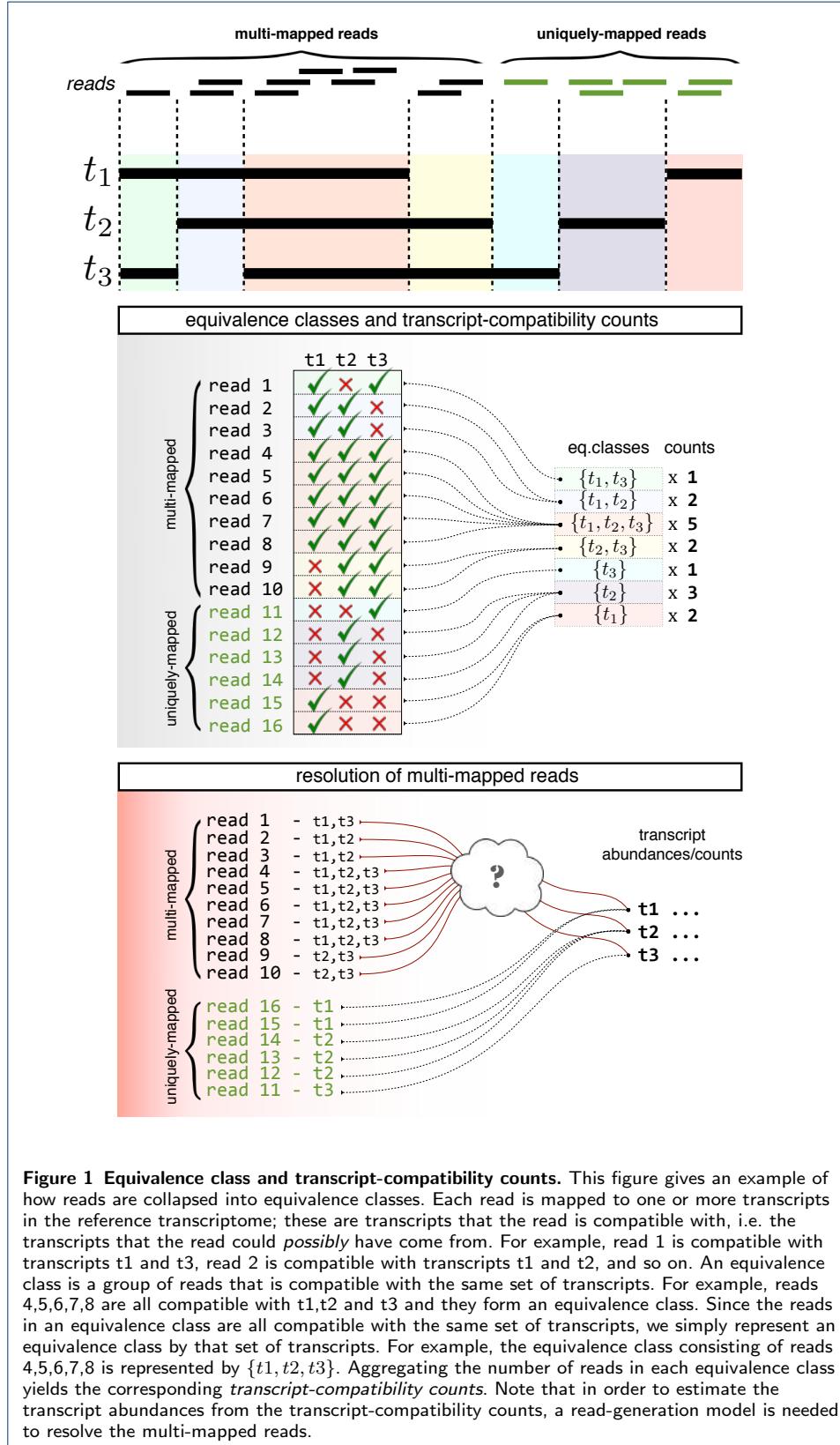
References

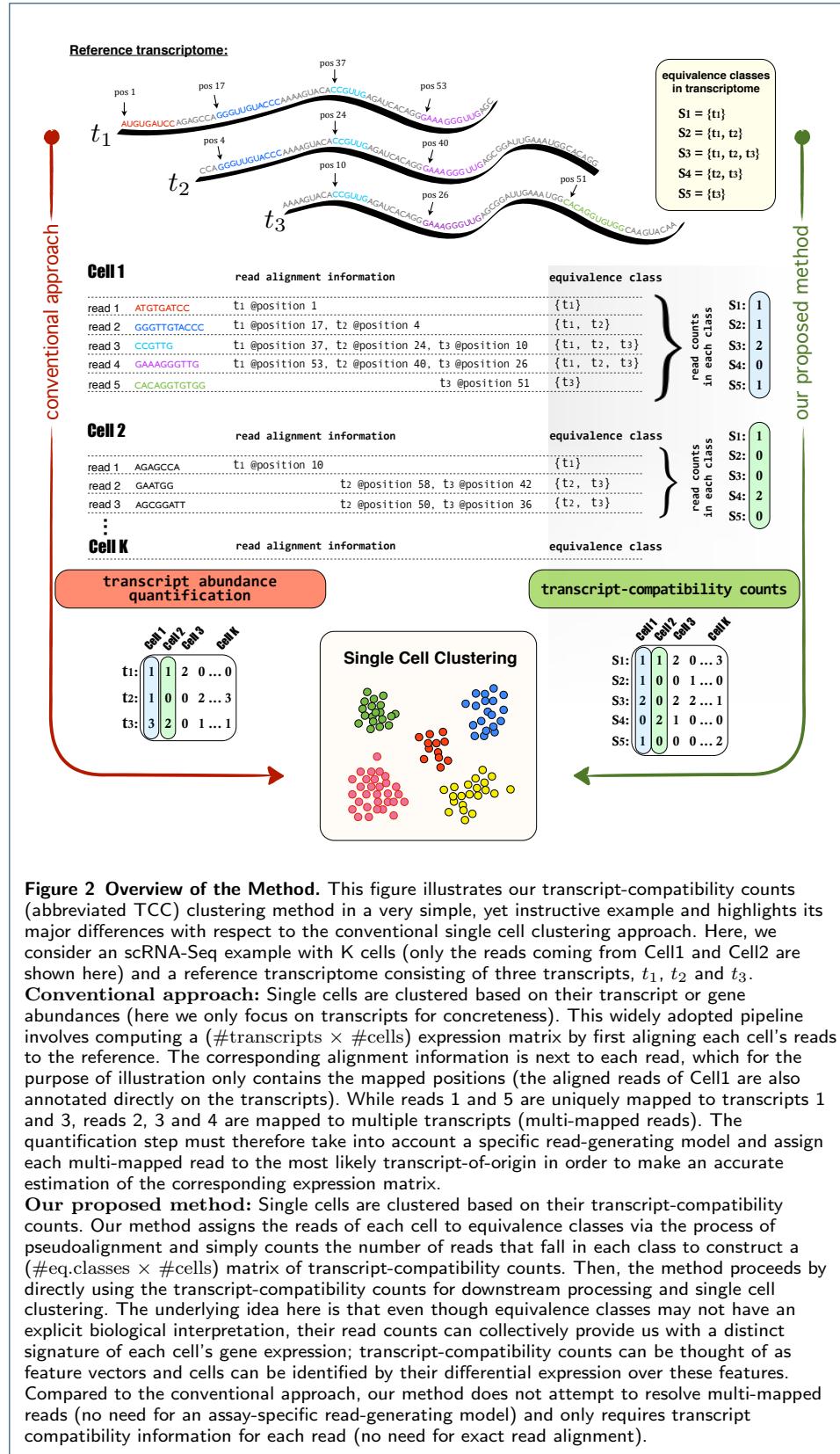
1. Patel, A.P., Tirosh, I., Trombetta, J.J., Shalek, A.K., Gillespie, S.M., Wakimoto, H., Cahill, D.P., Nahed, B.V., Curry, W.T., Martuza, R.L., Louis, D.N., Rozenblatt-Rosen, O., Suvà, M.L., Regev, A., Bernstein, B.E.: Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. *Science* **344**(6190), 1396–1401 (2014). doi:10.1126/science.1254257
2. Pollen, A.A., Nowakowski, T.J., Chen, J., Retallack, H., Sandoval-Espinosa, C., Nicholas, C.R., Shuga, J., Liu, S.J., Oldham, M.C., Diaz, A., Lim, D.A., Leyrat, A.A., West, J.A., Kriegstein, A.R.: Molecular Identity of Human Outer Radial Glia during Cortical Development. *Cell* **163**(1), 55–67. doi:10.1016/j.cell.2015.09.004
3. Gaublomme, J.T., Yosef, N., Lee, Y., Gertner, R.S., Yang, L.V., Wu, C., Pandolfi, P.P., Mak, T., Satija, R., Shalek, A.K., Kuchroo, V.K., Park, H., Regev, A.: Single-Cell Genomics Unveils Critical Regulators of Th17 Cell Pathogenicity. *Cell* **163**(6), 1400–1412. doi:10.1016/j.cell.2015.11.009
4. Kowalczyk, M.S., Tirosh, I., Heckl, D., Rao, T.N., Dixit, A., Haas, B.J., Schneider, R.K., Wagers, A.J., Ebert, B.L., Regev, A.: Single-cell RNA-seq reveals changes in cell cycle and differentiation programs upon aging of hematopoietic stem cells. *Genome Research* **25**(12), 1860–1872 (2015). doi:10.1101/gr.192237.115. <http://genome.cshlp.org/content/25/12/1860.full.pdf+html>
5. Lande-Diner, L., Stewart-Ornstein, J., Weitz, C.J., Lahav, G.: Single-cell analysis of circadian dynamics in tissue explants. *Molecular Biology of the Cell* **26**(22), 3940–3945 (2015). doi:10.1091/mbc.E15-06-0403
6. Usoskin, D., Furlan, A., Islam, S., Abdo, H., Lönnerberg, P., Lou, D., Hjerling-Leffler, J., Haeggstrom, J., Kharchenko, O., Kharchenko, P.V., Linnarsson, S., Ernfors, P.: Unbiased classification of sensory neuron types by large-scale single-cell RNA sequencing. *Nat Neurosci* **18**(1), 145–153 (2015). doi:10.1038/nn.3881
7. Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., La Manno, G., Jurèus, 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). doi:10.1126/science.aaa1934. <http://www.sciencemag.org/content/347/6226/1138.full.pdf>
8. Burns, J.C., Kelly, M.C., Hoa, M., Morell, R.J., Kelley, M.W.: Single-cell RNA-Seq resolves cellular complexity in sensory organs from the neonatal inner ear. *Nat Commun* **6** (2015). doi:10.1038/ncomms9557
9. Grün, D., Lyubimova, A., Kester, L., Wiebrands, K., Basak, O., Sasaki, N., Clevers, H., van Oudenaarden, A.: Single-cell messenger RNA sequencing reveals rare intestinal cell types. *Nature* **525**(7568), 251–255 (2015)
10. Kafri, R., Levy, J., Ginzberg, M.B., Oh, S., Lahav, G., Kirschner, M.W.: Dynamics extracted from fixed cells reveal feedback linking cell growth to cell cycle. *Nature* **494**(7438), 480–483 (2013). doi:10.1038/nature11897
11. Bendall, S.C., Davis, K.L., Amir, E.-a.D., Tadmor, M.D., Simonds, E.F., Chen, T.J., Shenfeld, D.K., Nolan, G.P., Pe'er, D.: Single-Cell Trajectory Detection Uncovers Progression and Regulatory Coordination in Human B Cell Development. *Cell* **157**(3), 714–725 (2014). doi:10.1016/j.cell.2014.04.005
12. Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, 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. *Nat Biotech* **32**(4), 381–386 (2014). doi:10.1038/nbt.2859
13. 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. *Nat Biotech* **33**(2), 155–160 (2015). doi:10.1038/nbt.3102
14. Shapiro, E., Biezuner, T., Linnarsson, S.: Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nat Rev Genet* **14**(9), 618–630 (2013). doi:10.1038/nrg3542
15. Stegle, O., Teichmann, S.A., Marioni, J.C.: Computational and analytical challenges in single-cell transcriptomics. *Nat Rev Genet* **16**(3), 133–145 (2015). doi:10.1038/nrg3833
16. Oshlack, A., Robinson, M., Young, M.: From RNA-seq reads to differential expression results. *Genome Biology* **11**(12), 220 (2010). doi:10.1186/gb-2010-11-12-220
17. Tang, F., Barbacioru, C., Wang, Y., Nordman, E., Lee, C., Xu, N., Wang, X., Bodeau, J., Tuch, B.B., Siddiqui, A., Lao, K., Surani, M.A.: mRNA-Seq whole-transcriptome analysis of a single cell. *Nat Meth* **6**(5), 377–382 (2009). doi:10.1038/nmeth.1315

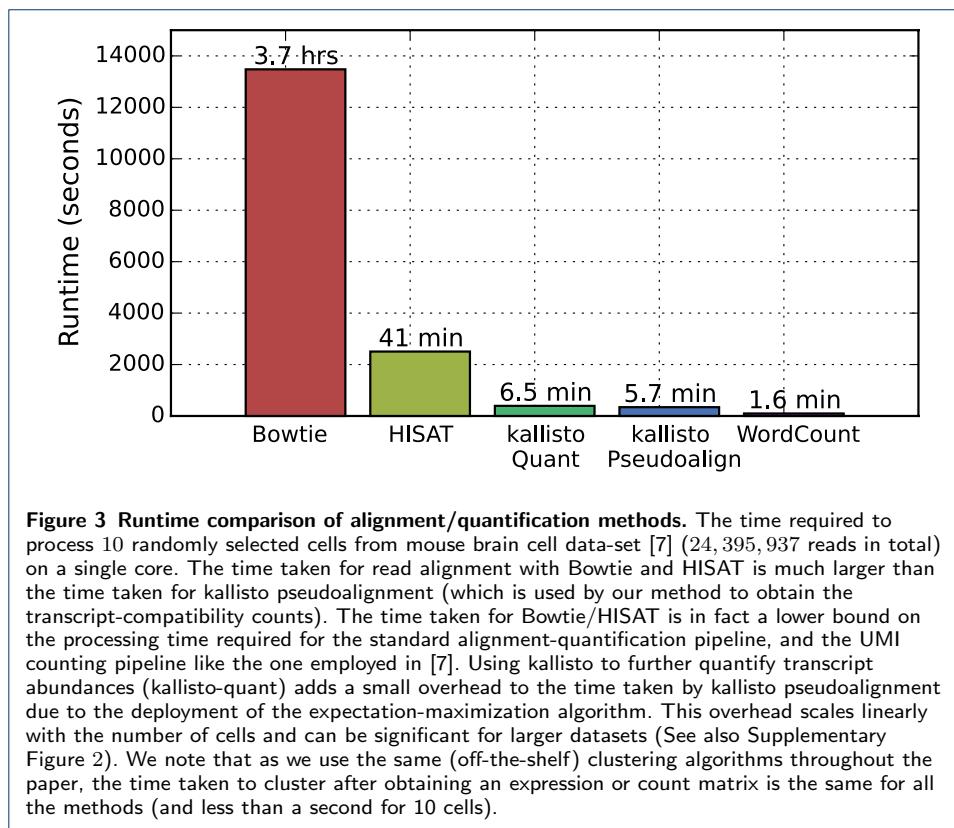
18. Islam, S., Kjällquist, U., Moliner, A., Zajac, P., Fan, J.-B., Lönnerberg, P., Linnarsson, S.: Characterization of the single-cell transcriptional landscape by highly multiplex RNA-seq. *Genome Research* **21**(7), 1160–1167 (2011). doi:10.1101/gr.110882.110
19. Ramsköld, D., Luo, S., Wang, Y.-C., Li, R., Deng, Q., Faridani, O.R., Daniels, G.A., Khrebtukova, I., Loring, J.F., Laurent, L.C., Schroth, G.P., Sandberg, R.: Full-length mRNA-Seq from single-cell levels of RNA and individual circulating tumor cells. *Nat Biotech* **30**(8), 777–782 (2012). doi:10.1038/nbt.2282
20. Hashimshony, T., Wagner, F., Sher, N., Yanai, I.: CEL-Seq: Single-Cell RNA-Seq by Multiplexed Linear Amplification. *Cell Reports* **2**(3), 666–673 (2012). doi:10.1016/j.celrep.2012.08.003
21. Picelli, S., Björklund, A.K., Faridani, O.R., Sagasser, S., Winberg, G., Sandberg, R.: Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Meth* **10**(11), 1096–1098 (2013). doi:10.1038/nmeth.2639
22. Sasagawa, Y., Nikaido, I., Hayashi, T., Danno, H., Uno, K., Imai, T., Ueda, H.: Quartz-Seq: a highly reproducible and sensitive single-cell RNA sequencing method, reveals non-genetic gene-expression heterogeneity. *Genome Biology* **14**(4), 31 (2013). doi:10.1186/gb-2013-14-4-r31
23. 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). doi:10.1126/science.1247651
24. Macosko, E.Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N., Martersteck, E.M., et al.: Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**(5), 1202–1214 (2015). doi:10.1016/j.cell.2015.05.002
25. 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**(5), 1187–1201. doi:10.1016/j.cell.2015.04.044
26. Amir, E.-a.D., Davis, K.L., Tadmor, M.D., Simonds, E.F., Levine, J.H., Bendall, S.C., Shenfeld, D.K., Krishnaswamy, S., Nolan, G.P., Pe'er, D.: viSNE enables visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of leukemia. *Nat Biotech* **31**(6), 545–552 (2013). doi:10.1038/nbt.2594
27. Mahfouz, A., van de Giessen, M., van der Maaten, L., Huisman, S., Reinders, M., Hawrylycz, M.J., Lelieveldt, B.P.F.: Visualizing the spatial gene expression organization in the brain through non-linear similarity embeddings. *Methods* **73**, 79–89 (2015). doi:10.1016/j.ymeth.2014.10.004
28. Shekhar, K., Brodin, P., Davis, M.M., Chakraborty, A.K.: Automatic Classification of Cellular Expression by Nonlinear Stochastic Embedding (ACCENSE). *Proceedings of the National Academy of Sciences of the United States of America* **111**(1), 202–207 (2014). doi:10.1073/pnas.1321405111
29. Islam, S., Zeisel, A., Joost, S., La Manno, G., Zajac, P., Kasper, M., Lonnerberg, P., Linnarsson, S.: Quantitative single-cell RNA-seq with unique molecular identifiers. *Nat Meth* **11**(2), 163–166 (2014). doi:10.1038/nmeth.2772
30. Tsafir, D., Tsafir, I., Ein-Dor, L., Zuk, O., Notterman, D.A., Domany, E.: Sorting points into neighborhoods (SPIN): data analysis and visualization by ordering distance matrices. *Bioinformatics* **21**(10), 2301–2308 (2005). doi:10.1093/bioinformatics
31. Qiu, P., Simonds, E.F., Bendall, S.C., Gibbs Jr, K.D., Bruggner, R.V., Linderman, M.D., Sachs, K., Nolan, G.P., Plevritis, S.K.: Extracting a cellular hierarchy from high-dimensional cytometry data with SPADE. *Nat Biotech* **29**(10), 886–891 (2011). doi:10.1038/nbt.1991
32. Levine, J.H., Simonds, E.F., Bendall, S.C., Davis, K.L., Amir, E.-a.D., Tadmor, M.D., Litvin, O., Fienberg, H.G., Jager, A., Zunder, E.R., Finck, R., Gedman, A.L., Radtke, I., Downing, J.R., Pe'er, D., Nolan, G.P.: Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. *Cell* **162**(1), 184–197 (2015). doi:10.1016/j.cell.2015.05.047
33. Marco, E., Karp, R.L., Guo, G., Robson, P., Hart, A.H., Trippa, L., Yuan, G.-C.: Bifurcation analysis of single-cell gene expression data reveals epigenetic landscape. *Proceedings of the National Academy of Sciences* **111**(52), 5643–5650 (2014). doi:10.1073/pnas.1408993111
34. 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.-I., Song, H.: Single-Cell RNA-Seq with Waterfall Reveals Molecular Cascades underlying Adult Neurogenesis. *Cell Stem Cell* **17**(3), 360–372. doi:10.1016/j.stem.2015.07.013
35. Haghverdi, L., Buettner, F., Theis, F.J.: Diffusion maps for high-dimensional single-cell analysis of differentiation data. *Bioinformatics* (2015). doi:10.1093/bioinformatics
36. Moignard, V., Woodhouse, S., Haghverdi, L., Lilly, A.J., Tanaka, Y., Wilkinson, A.C., Buettner, F., Macaulay, I.C., Jawaid, W., Diamanti, E., Nishikawa, S.-I., Piterman, N., Kouskoff, V., Theis, F.J., Fisher, J., Gottgens, B.: Decoding the regulatory network of early blood development from single-cell gene expression measurements. *Nat Biotech* **33**(3), 269–276 (2015). doi:10.1038/nbt.3154
37. Fan, J., Salathia, N., Liu, R., Kaeser, G., Yung, Y., Herman, J.L., Kaper, F., Fan, J.-B., Zhang, K., Chun, J., Kharchenko, P.: Characterizing transcriptional heterogeneity through pathway and gene set overdispersion analysis. *bioRxiv* (2015). doi:10.1101/026948
38. Saliba, A.-E., Westermann, A.J., Gorski, S.A., Vogel, J.: Single-cell RNA-seq: advances and future challenges. *Nucleic Acids Research* (2014). doi:10.1093/nar
39. Nicolae, M., Mangul, S., Mandoiu, I.I., Zelikovsky, A.: Estimation of alternative splicing isoform frequencies from RNA-Seq data. *Algorithms for Molecular Biology* **6**(1), 9 (2011)
40. Patro, R., Mount, S.M., Kingsford, C.: Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nature biotechnology* **32**(5), 462–464 (2014)
41. Bray, N., Pimentel, H., Melsted, P., Pachter, L.: Near-optimal probabilistic RNA-Seq quantification. *Nature Biotechnology*, in press (2016)
42. Pachter, L.: Models for transcript quantification from RNA-Seq. *arXiv preprint arXiv:1104.3889* (2011)
43. Gerstein, M.B., Bruce, C., Rozowsky, J.S., Zheng, D., Du, J., Korbel, J.O., Emanuelsson, O., Zhang, Z.D., Weissman, S., Snyder, M.: What is a gene, post-ENCODE? History and updated definition. *Genome research* **17**(6), 669–681 (2007)
44. Trapnell, C., Williams, B.A., Pertea, G., Mortazavi, A., Kwan, G., Van Baren, M.J., Salzberg, S.L., Wold, B.J.,

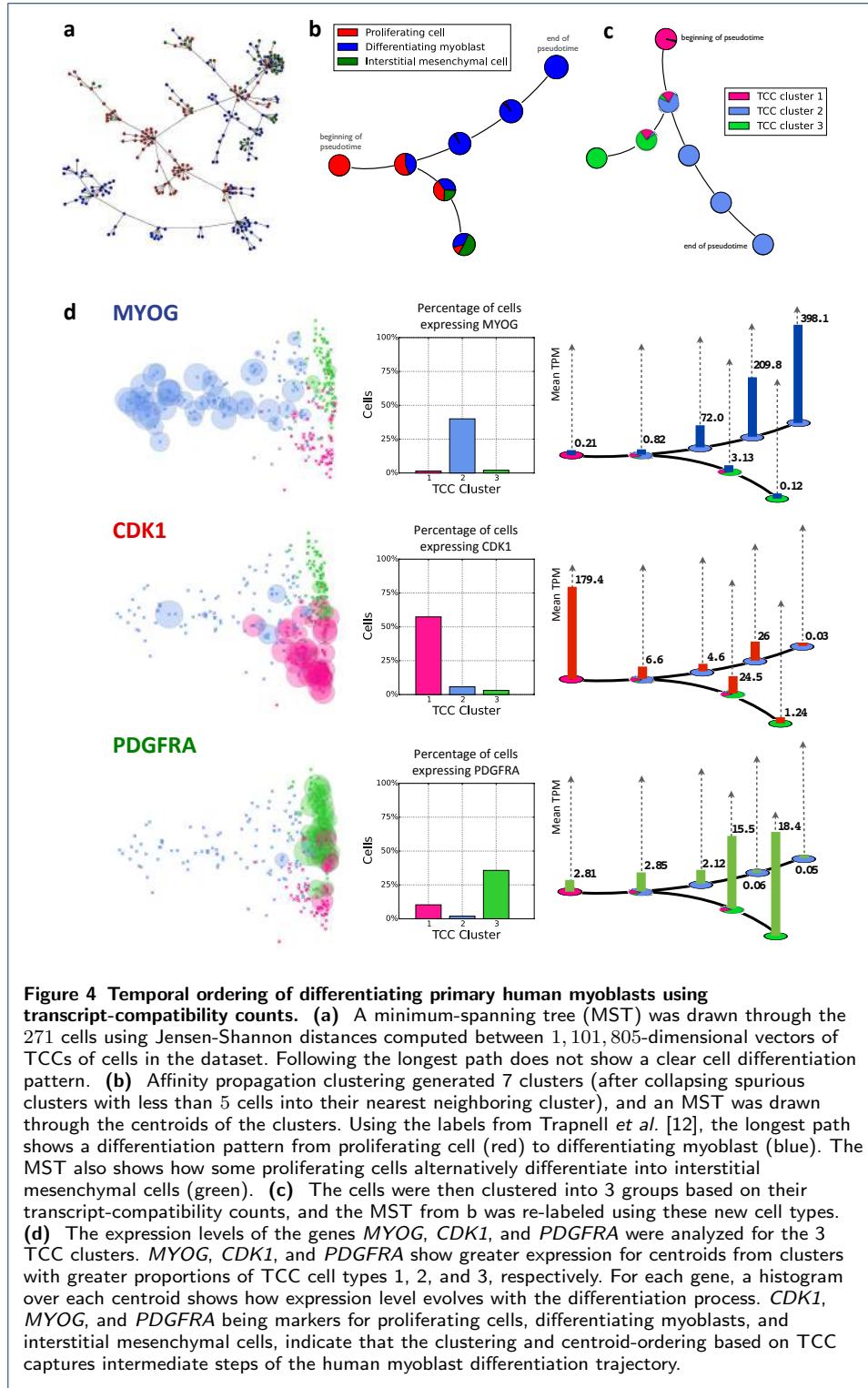
- Pachter, L.: Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nature biotechnology* **28**(5), 511–515 (2010)
- 45. Magwene, P.M., Lizardi, P., Kim, J.: Reconstructing the temporal ordering of biological samples using microarray data. *Bioinformatics* **19**(7), 842–850 (2003). doi:10.1093/bioinformatics
 - 46. Clustering on Transcript Compatibility Counts. Github repository (2016).
https://github.com/govinda-kamath/clustering_on_transcript_compatibility_counts
 - 47. Lin, J.: Divergence measures based on the Shannon entropy. *Information Theory, IEEE Transactions on* **37**(1), 145–151 (1991)
 - 48. Fuglede, B., Topsøe, F.: Jensen-Shannon divergence and Hilbert space embedding. In: *IEEE International Symposium on Information Theory*, pp. 31–31 (2004)
 - 49. Batu, T., Fortnow, L., Rubinfeld, R., Smith, W.D., White, P.: Testing that distributions are close. In: *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium On*, pp. 259–269 (2000). IEEE
 - 50. Frey, B.J., Dueck, D.: Clustering by passing messages between data points. *science* **315**(5814), 972–976 (2007)
 - 51. Van der Maaten, L., Hinton, G.: Visualizing data using t-SNE. *Journal of Machine Learning Research* **9**(2579–2605), 85 (2008)
 - 52. Mühlbacher, P.: A python implementation of the diffusion maps algorithm introduced by Lafon. GitHub (2015)
 - 53. Lafon, S.S.: Diffusion maps and geometric harmonics. PhD thesis, Yale University (2004)

Figures









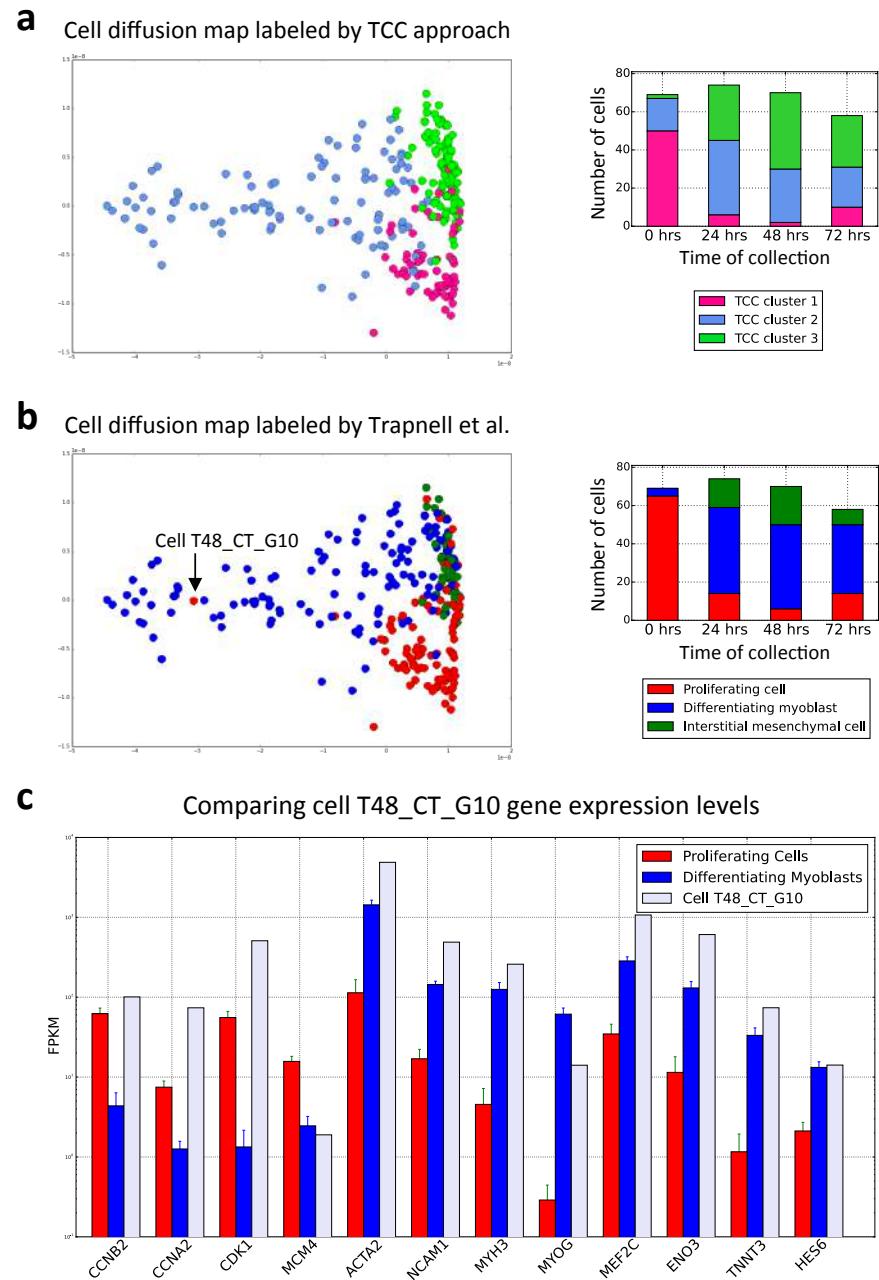
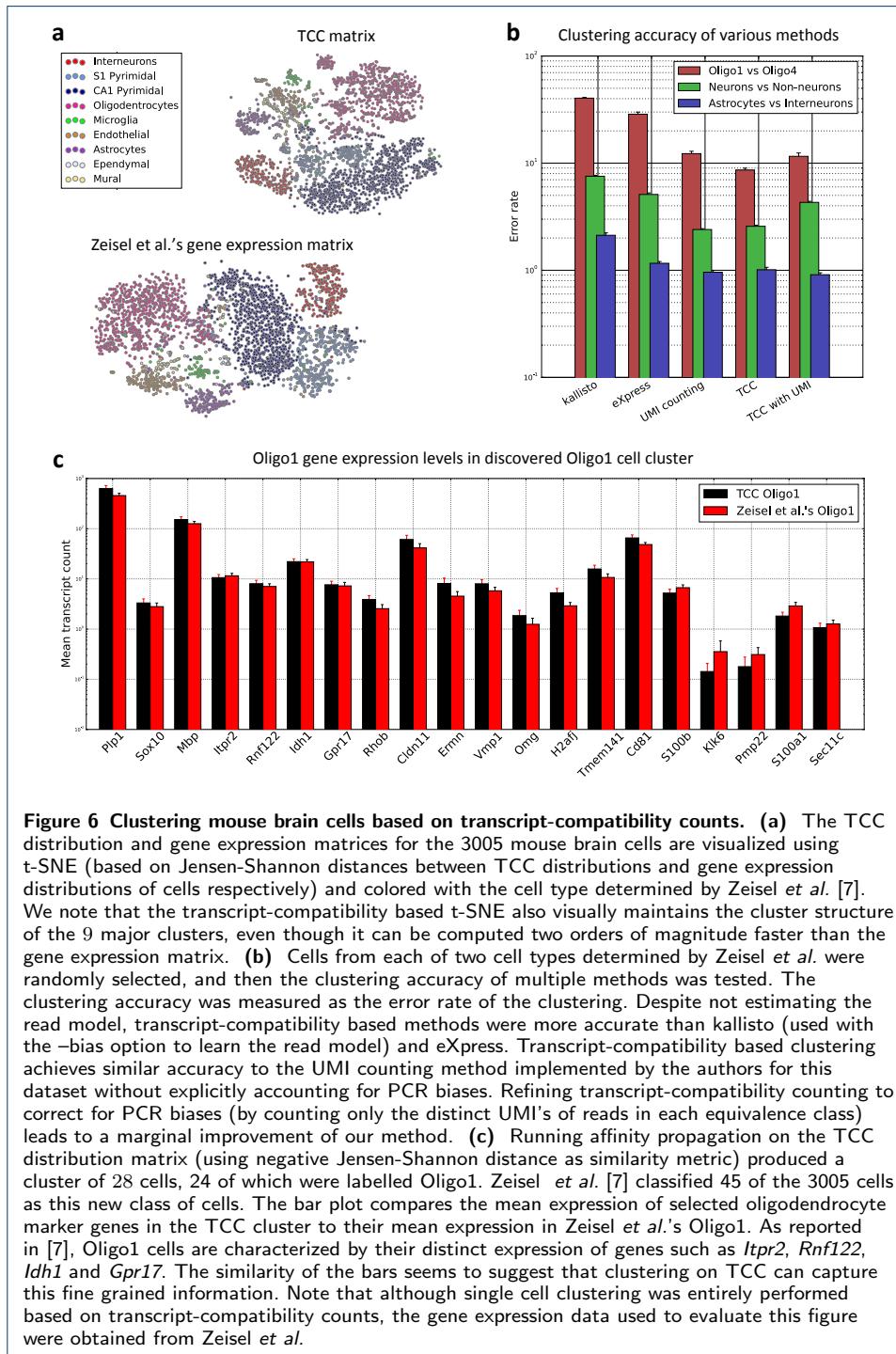
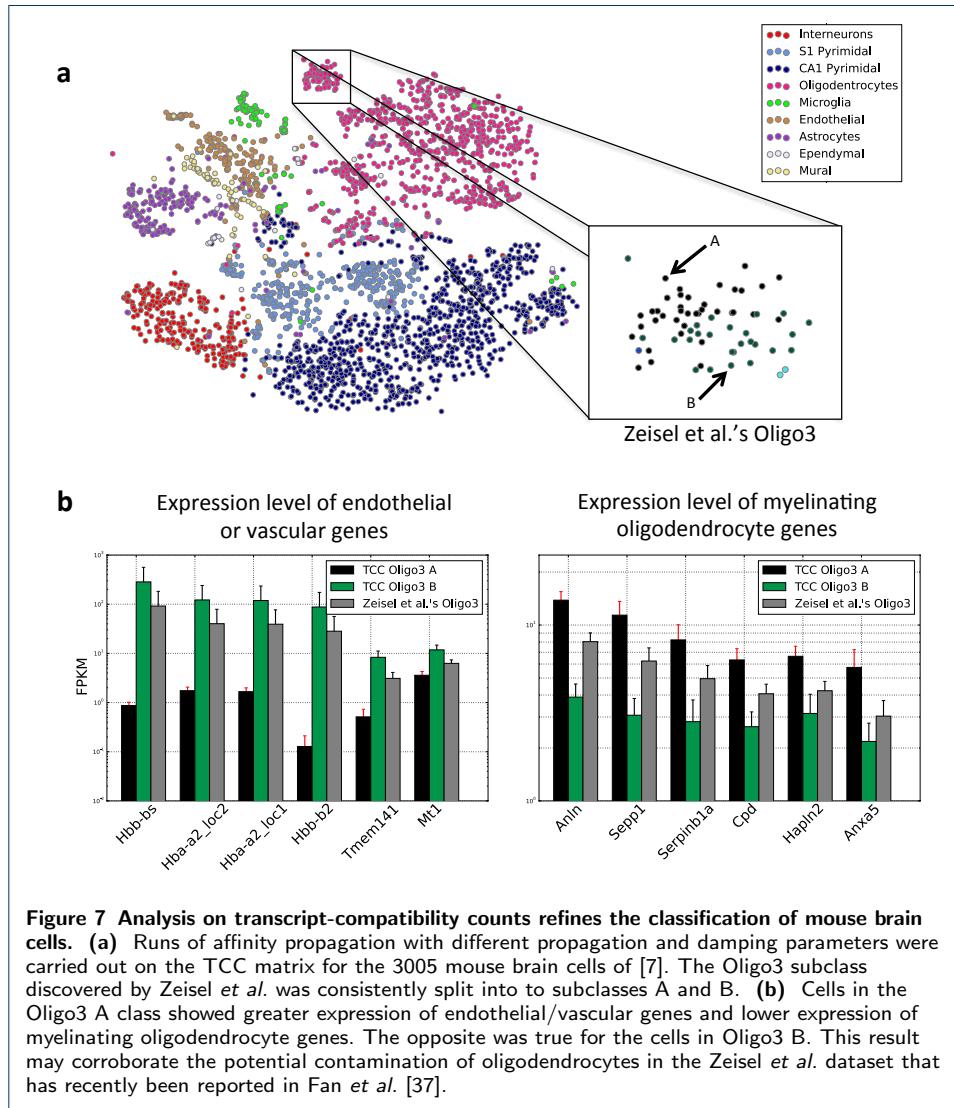
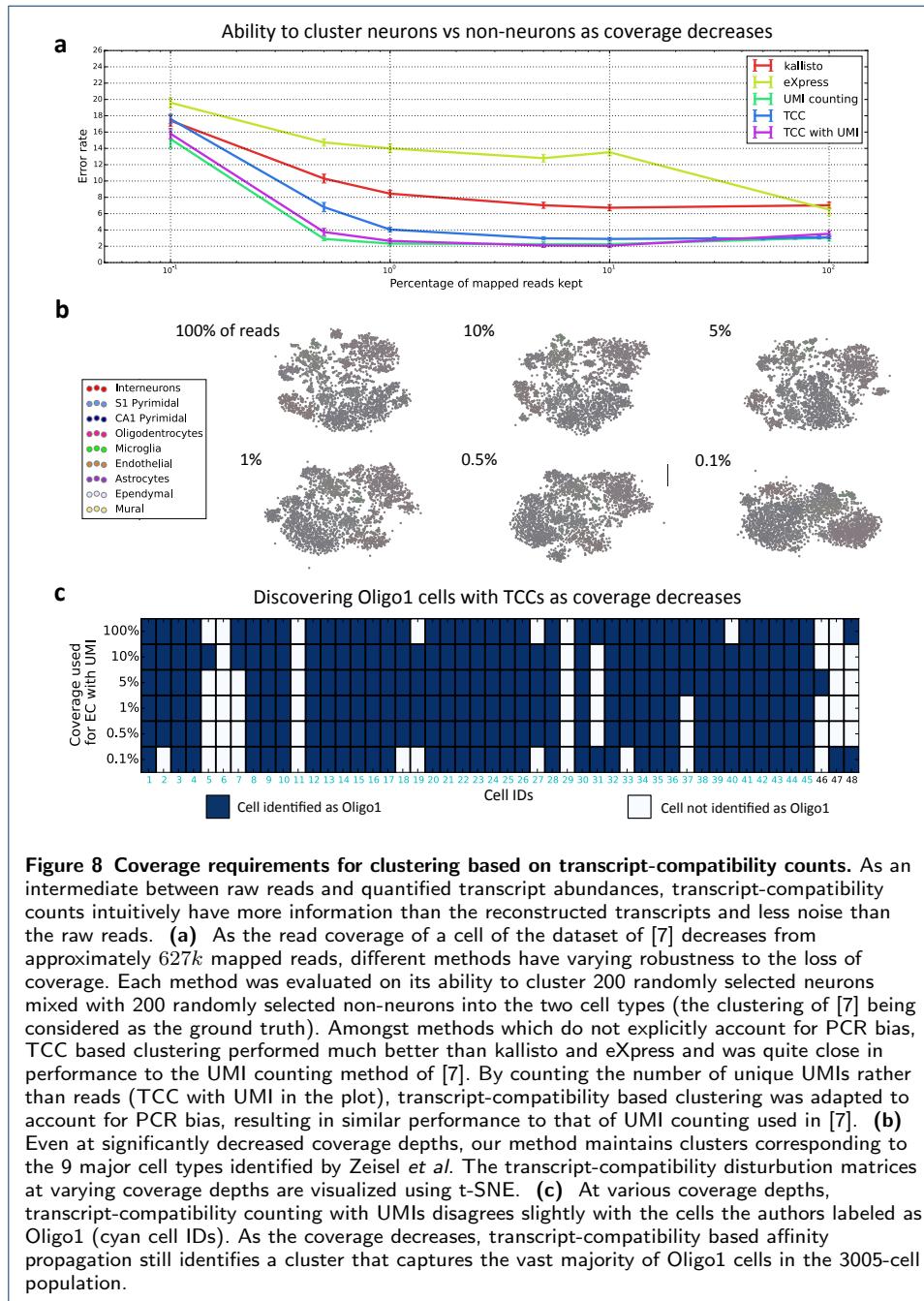


Figure 5 Clustering primary human myoblasts based on transcript-compatibility counts. (a) The transcript-compatibility counts matrix for 271 primary human myoblasts from [12] is visualized using a diffusion map. Three clusters obtained using affinity propagation are shown along with the distribution of these cells across the 4 cell-collection timepoints (0, 24, 48 and 72 hours). (b) The diffusion map obtained using transcript compatibility counts is relabeled using the cells reported by [12]. Clusters 1, 2, 3 generated by the transcript compatibility based method map to proliferating cells, differentiating myoblasts, and interstitial cells respectively. According to Trapnell et al.'s labels, the transcript compatibility based method seems to have severely misclassified cell T48_CT_G10 (SRR1033183) as a differentiating myoblast. (c) Comparing the expressions of 12 differentiating genes in T48_CT_G10 with those of the average proliferating cell and the average differentiating myoblast, 8 out of the 12 genes show expressions similar to what one would expect from a differentiating myoblast. MYOG seems to show an FPKM of 14, which while more than the mean expression of proliferating cells (around 0.28) is much less than the mean expression of differentiating myoblasts (around 61.33). We note that this cell has the highest expression of MYOG among all cells labelled by Trapnell et al. as proliferating cell (and the second highest cell has expression around 5.4). However there are 88 differentiating myoblasts with MYOG expression less than 15 FPKM. Hence it is reasonable to think that this MYOG expression is more typical of differentiating myoblasts than proliferating cells. Only genes CDK1 and CCNB2 show expressions close to what one would expect from a proliferating cell. Even though CDK1 is a highly specific marker for proliferating cells, the above gene profile indicates that classifying cell T48_CT_G10 as a differentiating myoblast seems reasonable.

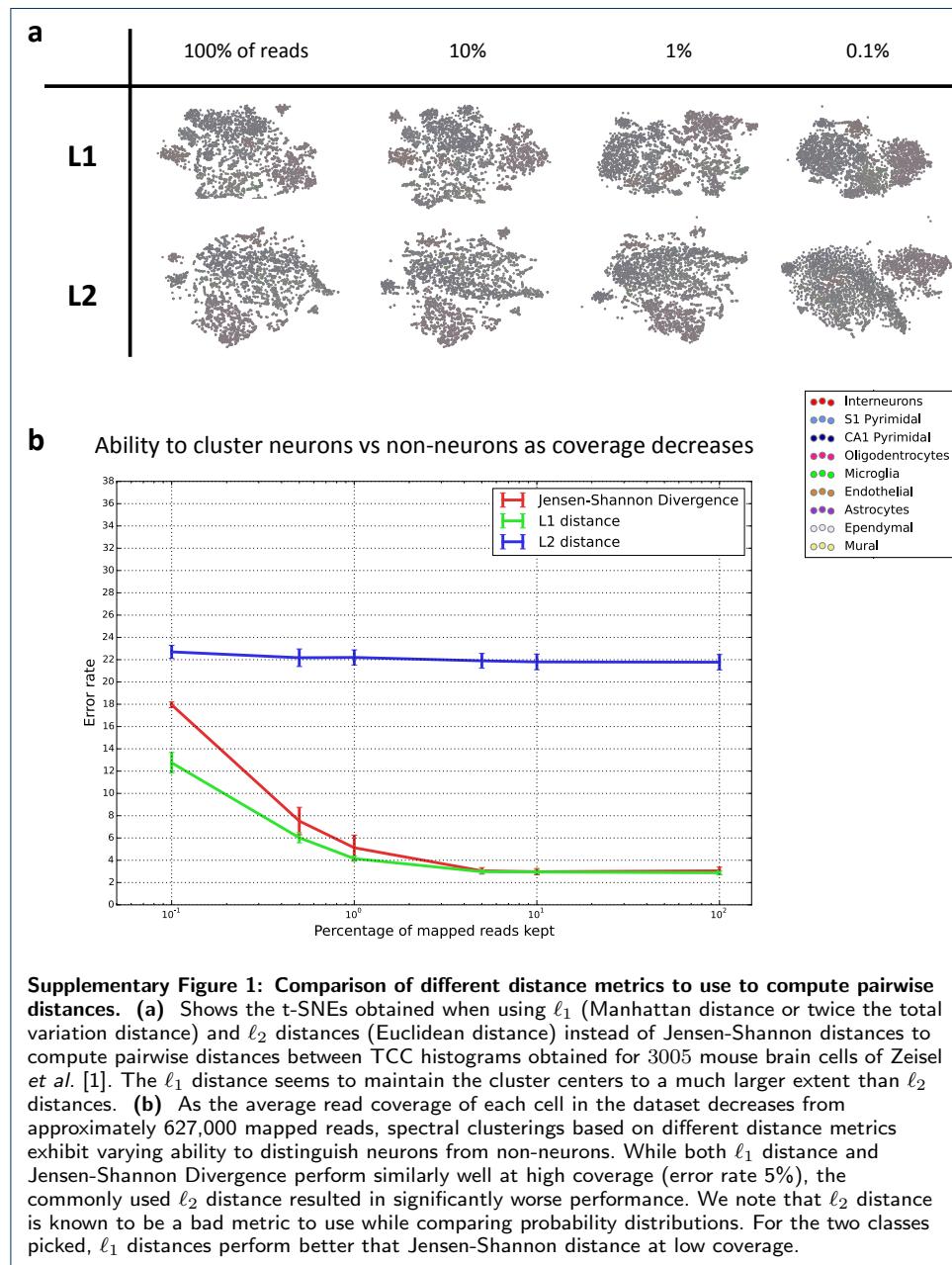


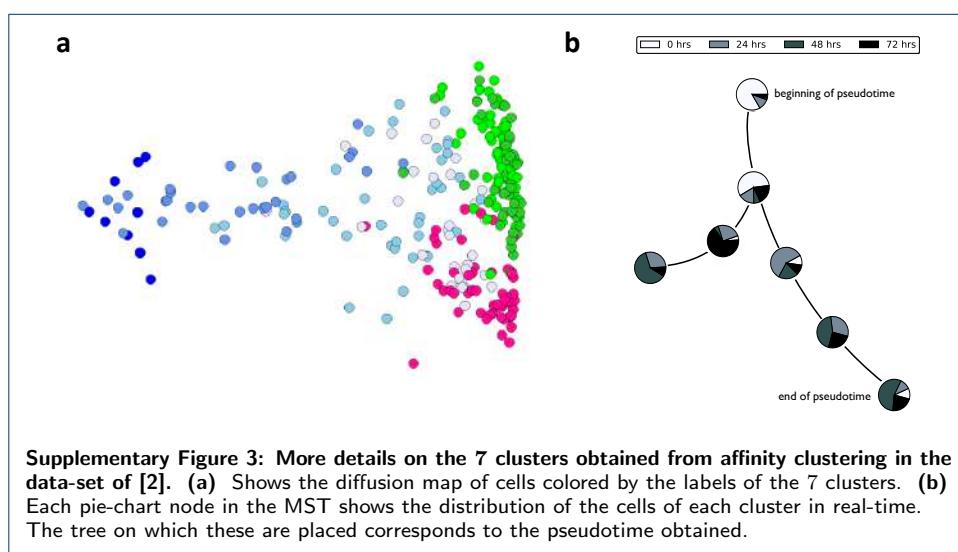
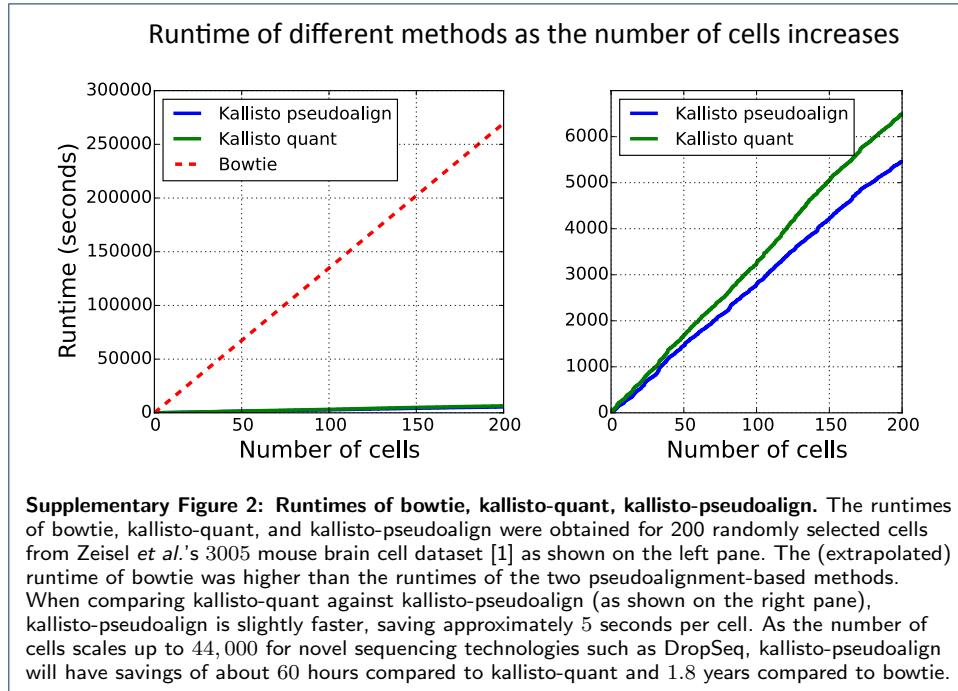


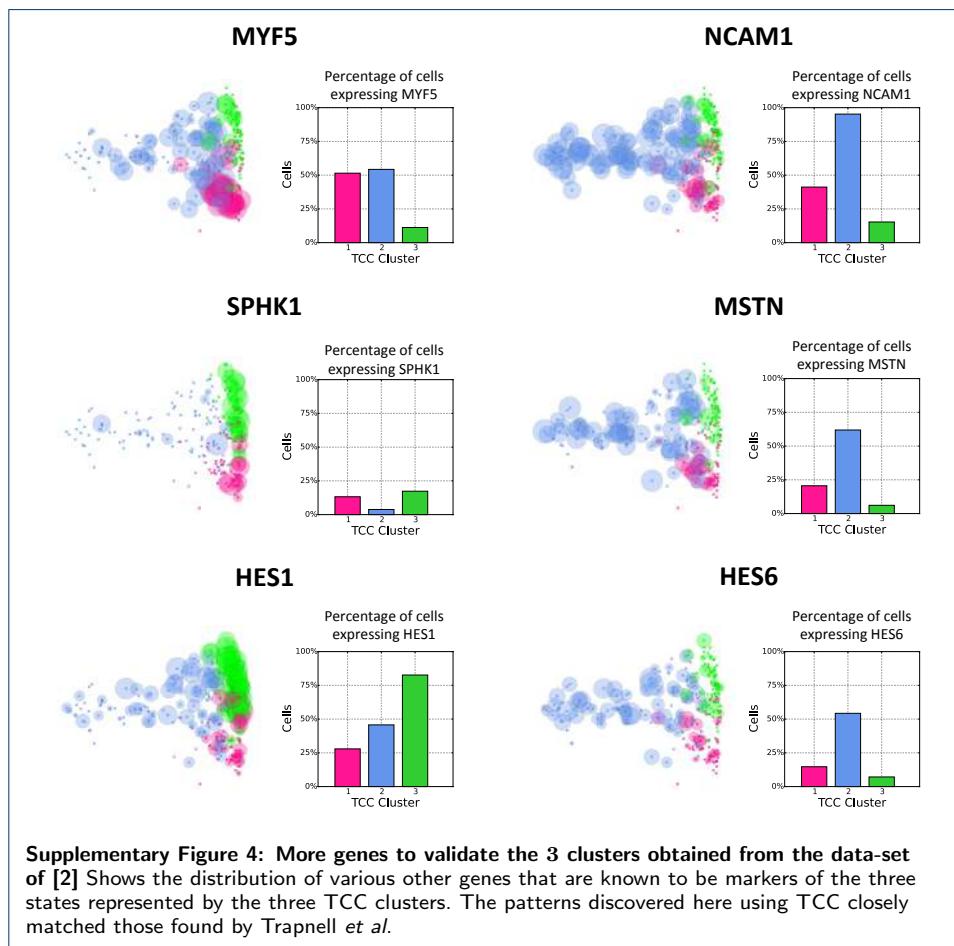


Supplementary Figures for

“Fast and accurate single-cell RNA-Seq analysis by clustering of transcript-compatibility counts”







References

1. Zeisel, A., Muñoz-Manchado, A.B., Codeluppi, S., Lönnerberg, P., La Manno, G., Jurèus, 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). doi:10.1126/science.aaa1934. <http://www.sciencemag.org/content/347/6226/1138.full.pdf>
2. Trapnell, C., Cacchiarelli, D., Grimsby, J., Pokharel, 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. *Nat Biotech* **32**(4), 381–386 (2014). doi:10.1038/nbt.2859