


METHOD

Open Access



GoM DE: interpreting structure in sequence count data with differential expression analysis allowing for grades of membership

Peter Carbonetto^{1,2} , Kaixuan Luo¹, Abhishek Sarkar^{1,3}, Anthony Hung^{1,4}, Karl Tayeb^{1,5}, Sebastian Pott^{1,4} and Matthew Stephens^{1,6*}

*Correspondence:
mstephens@uchicago.edu

¹ Department of Human Genetics, University of Chicago, Chicago, IL, USA

² Research Computing Center, University of Chicago, Chicago, IL, USA

³ Vesalius Therapeutics, Cambridge, MA, USA

⁴ Section of Genetic Medicine, University of Chicago, Chicago, IL, USA

⁵ Committee on Genetics, Genomics and Systems Biology, University of Chicago, Chicago, IL, USA

⁶ Department of Statistics, University of Chicago, Chicago, IL, USA

Abstract

Parts-based representations, such as non-negative matrix factorization and topic modeling, have been used to identify structure from single-cell sequencing data sets, in particular structure that is not as well captured by clustering or other dimensionality reduction methods. However, interpreting the individual parts remains a challenge. To address this challenge, we extend methods for differential expression analysis by allowing cells to have partial membership to multiple groups. We call this grade of membership differential expression (GoM DE). We illustrate the benefits of GoM DE for annotating topics identified in several single-cell RNA-seq and ATAC-seq data sets.

Keywords: Gene expression, Single-cell RNA-seq, Single-cell ATAC-seq, Differential expression analysis, Dimensionality reduction, Parts-based representations, Matrix factorization, Topic modeling

Background

A key methodological aim in single-cell genomics is to learn structure from single-cell sequencing data in a systematic, data-driven way [1–3]. Clustering [4–7] and dimensionality reduction techniques such as PCA [8–10], *t*-SNE [11], or UMAP [12] are commonly used for this aim. Despite the fact that many of these techniques have been applied “out-of-the-box” (with some caveats [13–18]), they have been remarkably successful in revealing and visualizing biologically interesting substructures from single-cell data [7, 19–29].

Another class of dimensionality reduction approaches that have been used to identify structure from single-cell data are what are sometimes called *parts-based*



© The Author(s) 2023. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

representations—these approaches include non-negative matrix factorization (NMF) [30–44] and topic modeling [45–56], which also have formal connections [48, 57, 58]. Parts-based representations share some of the features of both a clustering and a dimensionality reduction: on the one hand, they learn a lower dimensional representation of the cells; on the other hand, the individual dimensions (the “parts”) of the reduced representation can identify discrete clusters or discrete subpopulations [59, 60]. However, parts-based representations are more flexible than clustering—the dimensions can also capture other features such as continuously varying cell states.

In this paper, we investigate the question of how to interpret the individual dimensions of a parts-based representation learned by fitting a topic model (in the topic model, the dimensions are also called “topics”). For topics that assign observations to discrete clusters, one could apply a standard method for differential expression analysis [61, 62] to compare expression between topics, then annotate these topics by the genes that are differentially expressed. The question, therefore, is what to do with topics that do not assign observations to discrete clusters. To tackle this question, we extend models that compare expression between groups by allowing observations to have *partial membership in multiple groups*. This more flexible differential expression analysis is implemented by taking an existing model and modifying it to allow for partial memberships to groups or topics. This modified model is a “grade of membership” model [63], so we call our new method *grade of membership differential expression* (GoM DE). The idea is that, by generalizing existing methods, we can continue to take advantage of existing elements of differential expression analysis but now apply them to learn about different types of cell features beyond discrete cell populations.

We describe the GoM DE approach more formally in the next section. Then, we evaluate the GoM DE approach in simulations, showing, in particular, that it recovers the same results as existing differential expression analysis methods when the cells can be grouped into discrete clusters. In case studies, we demonstrate how the GoM DE analysis can be used to uncover and interpret a variety of cell features from single-cell RNA-seq and ATAC-seq data sets.

Results

Methods overview and illustration

We begin by giving a brief overview of the topic model; then, we describe the new methods for annotating topics. To illustrate key concepts, we analyze a single-cell RNA-seq (scRNA-seq) data set obtained from peripheral blood mononuclear cells (PBMCs) [29] that has been used in several benchmarking studies (e.g., [4, 7, 8, 64, 65]). We refer to these data as the “PBMC data.”

Learning expression topics from single-cell RNA-seq data

The original aim of the topic model was to discover patterns from collections of text documents, in which text documents were represented as word counts [45, 50, 66–68]. By substituting genes for words and cells for documents, topic models can also be used to learn a reduced representation of cells by their membership in multiple “topics” [47].

When applied to scRNA-seq data generated using UMIs, the topic model assumes a multinomial distribution of the RNA molecule counts in a cell,

$$x_{i1}, \dots, x_{im} \sim \text{Multinomial}(s_i; \pi_{i1}, \dots, \pi_{im}). \quad (1)$$

where $s_i = x_{i1} + \dots + x_{im}$, and m is the number of genes, that is, the number of RNA molecules x_{ij} observed for gene j in cell i is a noisy observation of an underlying true expression level, π_{ij} [8, 69].

For n cells, the topic model is a *reduced representation* of the underlying expression,

$$\mathbf{\Pi} = \mathbf{L}\mathbf{F}^T, \quad (2)$$

where $\mathbf{\Pi}$, \mathbf{L} , \mathbf{F} are $n \times m$, $n \times K$, $m \times K$ matrices, respectively, with entries π_{ij}, l_{ik}, f_{jk} . Each cell i is represented by its “grade of membership” in K topics, a vector of proportions l_{i1}, \dots, l_{iK} , such that $l_{ik} \geq 0$, $\sum_{k=1}^K l_{ik} = 1$, and each “expression topic” is represented by a vector of (relative) expression levels f_{1k}, \dots, f_{mk} , $f_{jk} \geq 0$ (these are also constrained to sum to 1, which ensures that the π_{ij} s are multinomial probabilities). To efficiently fit the topic model to large single-cell data sets, we exploit the fact that the topic model is closely related to the Poisson NMF model [48].

The matrix \mathbf{L} in (2), which contains the membership proportions for all cells and topics, can be visualized using a “Structure plot.” Structure plots have been used to visualize the results of population genetics analyses (e.g., [70–72]) and, more recently, to visualize the topics learned from bulk and single-cell RNA-seq data [47].

A Structure plot visualizing the topic model fit to the PBMC data, with $K = 6$ topics, is given in Fig. 1. In this data set, the cells have been “sorted” into different cell types which provides a cell labeling to compare against. From the Structure plot, it is apparent

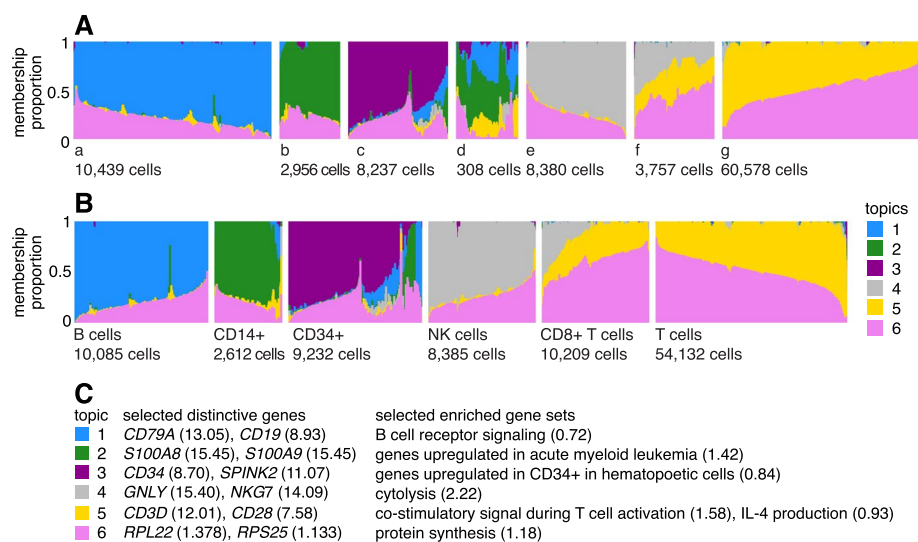


Fig. 1 **A** and **B** give two views of the topic model fit to the PBMC data [29] ($n = 94,655$ cells, $K = 6$ topics) using Structure plots [70, 71]. Cells are arranged horizontally; bar heights correspond to cell membership proportions. In **A**, the cells are arranged using the estimated membership proportions only. In **B**, the cells are grouped by the FACS labels (the “T cells” label combines all sorted T cell populations other than CD8+ cytotoxic T cells). In **C**, the topics are annotated by distinctive genes from the GoM DE analysis (Fig. 3) and by enriched gene sets. Numbers in parentheses next to genes give posterior mean i.e. LFCs, and for gene sets, they are enrichment coefficients. An enrichment coefficient is an estimate of the expected increase in the LFC for genes that belong to the gene set relative to genes that do not belong to the gene set. Note the groupings a–g in **A** are intended only to aid visualization. See also Additional file 1: Fig. S1 for an alternative visualization

that a subset of topics—topics 1, 2, and 3—correspond closely to the sorted subpopulations (B cells, CD14+ monocytes, CD34+ cells) (indeed, distinctive genes and enriched gene sets identified by the methods described below suggest these same subpopulations; Fig. 1C). Topics 4 and 5, on the other hand, are not confined to a single sorted cell type and instead appear to capture biological processes common to T cells and natural killer (NK) cells. CD8+ cytotoxic T cells have characteristics of both NK cells and T cells—these are T cells that sometimes become “NK-like” [73]—and this is captured in the topic model by assigning membership to both topics. Topic 6 also captures continuous structure, but, unlike topics 4 and 5, it is present in almost all cells, and therefore, its biological interpretation is not at all clear from the cell labeling. More generally, the topics, whether they capture largely discrete structure (topics 1–3) or more continuous structure (topics 4–6), can be thought of as a “soft” clustering [47].

Learning chromatin accessibility topics from single-cell ATAC-seq data

For single-cell ATAC-seq data, the observations x_{ij} denote the number of reads mapping to region j in cell i . However, it is common to “binarize” the read counts such that $x_{ij} = 1$ when at least one fragment in cell i maps to region j and $x_{ij} = 0$ otherwise.

Using the topic model to analyze (binarized) single-cell ATAC-seq data was first suggested by [49]. Therefore, they implicitly assumed a multinomial model (1) in which the x_{ij} s are binarized accessibility values instead of UMI counts. A binomial model for binarized accessibility data was proposed in [74]. As we explain in the “Methods” section, we view both models as approximations, and under reasonable assumptions the models are similar.

Differential expression analysis allowing for grades of membership

Having learned the topics, our aim now is to identify genes that are distinctive to each topic. In the simplest case, the topic is a distinct or nearly distinct cluster of cells, such as topic 1 or topic 2 in Fig. 1.

In the following, we describe methods for analyzing *differences in expression*, but they can also be understood as methods for analyzing *differences in chromatin accessibility*. Therefore, “expression,” “expressed,” and “gene” in the descriptions below may be substituted with “accessibility,” “accessible,” and “peak” (or “region”).

Consider a single gene, j . Provided unmodeled sources of variation are negligible relative to measurement error, a simple Poisson model of expression should suffice:

$$x_{ij} \sim \text{Poisson}(s_i \theta_{ij}). \quad (3)$$

In this model, θ_{ij} for gene j in cell i is controlled by the cell’s membership in the cluster: when cell i belongs to the cluster, $\theta_{ij} = p_{j1}$; otherwise, $\theta_{ij} = p_{j2}$. Under this model, differential expression (DE) analysis proceeds by estimating the log-fold change (LFC) in expression for each gene j ,

$$\text{LFC}(j) = \log_2 \frac{p_{j1}}{p_{j2}}. \quad (4)$$

Although simple, this Poisson model forms the basis for many DE analysis methods [75–80].

We now modify the Poisson model (3) in a simple way to analyze differential expression among topics. In a clustering, each cell belongs to a single cluster, whereas in the topic model, cells have *grades of membership* to the clusters [63] in which l_{ik} is the membership proportion for cluster or topic k . Therefore, we extend the model to allow for partial membership in the K topics:

$$\begin{aligned} x_{ij} &\sim \text{Poisson}(s_i \theta_{ij}) \\ \theta_{ij} &= \sum_{k=1}^K l_{ik} p_{jk}, \end{aligned} \quad (5)$$

in which the membership proportions l_{ik} are treated as known, and the unknowns p_{j1}, \dots, p_{jK} represent relative expression levels (a related model is used in C-SIDE [80] to model cell-type mixtures in DE analysis of spatial transcriptomics data). Note that p_{jk} will be similar to, but not the same as, f_{jk} in the topic model because the DE analysis is a gene-by-gene analysis, whereas the topic model considers all genes at once. The standard Poisson model (3) is recovered as a special case of (5) when $K = 2$ and all membership proportions l_{ik} are 0 or 1.

Recall, our aim is to identify genes that are *distinctive* to each topic. To this end, we estimate the *least extreme LFC* (i.e. LFC), which we define as

$$\begin{aligned} \text{LFC}_k^{\text{l.e.}}(j) &:= \text{LFC}_{k,l}(j) \\ &\text{such that } l = \underset{l' \neq k}{\text{argmin}} |\text{LFC}_{k,l'}(j)|, \end{aligned} \quad (6)$$

in which $\text{LFC}_{k,l}(j)$ is the *pairwise LFC*,

$$\text{LFC}_{k,l}(j) := \log_2 \frac{p_{jk}}{p_{jl}}. \quad (7)$$

In words, the l.e. LFC for topic k is the LFC comparing topics k and l , in which l is chosen to be topic that results in the smallest (“least extreme”) change. By this definition, a “distinctive gene” is one in which its expression is significantly different from its expression in *all other topics* (note the l.e. LFC reduces to the standard LFC (4) when $K = 2$). We then annotate topics by the distinctive genes. The estimation of l.e. LFCs and computation of related posterior statistics is described in the “[Methods](#)” section.

To illustrate what the least extreme LFC does and does not do, consider the following toy example with $K = 10$ topics (Fig. 2). Gene 1 has high expression in topic 1 and low expression in the other topics. Therefore, all the pairwise LFCs for topic 1 are large, $\text{LFC}_{1,k}(1) = \log_2(100)$, $k = 2, \dots, 10$, and this results in an l.e. LFC for topic 1 of $\log_2(100) \approx 6.6$. So gene 1 is a distinctive gene for topic 1. Next consider gene 2, which has high expression in topics 1 and 2 and low expression in the other topics. For gene 2, the pairwise LFCs for topic 1 are mostly large, $\text{LFC}_{1,k}(2) = \log_2(100)$, $k = 3, \dots, 10$, except for $\text{LFC}_{1,2}(2) = 0$. So, the l.e. LFC for topic 1 is zero and, as a result, gene 2, although potentially helpful for interpreting topic 1, is not a distinctive gene for topic 1.

Illustration of GoM DE analysis in PBMC data set

To illustrate, we applied the GoM DE analysis to the topic model shown in Fig. 1 and visualized the results in “volcano plots” (Fig. 3). We then used the GoM DE results

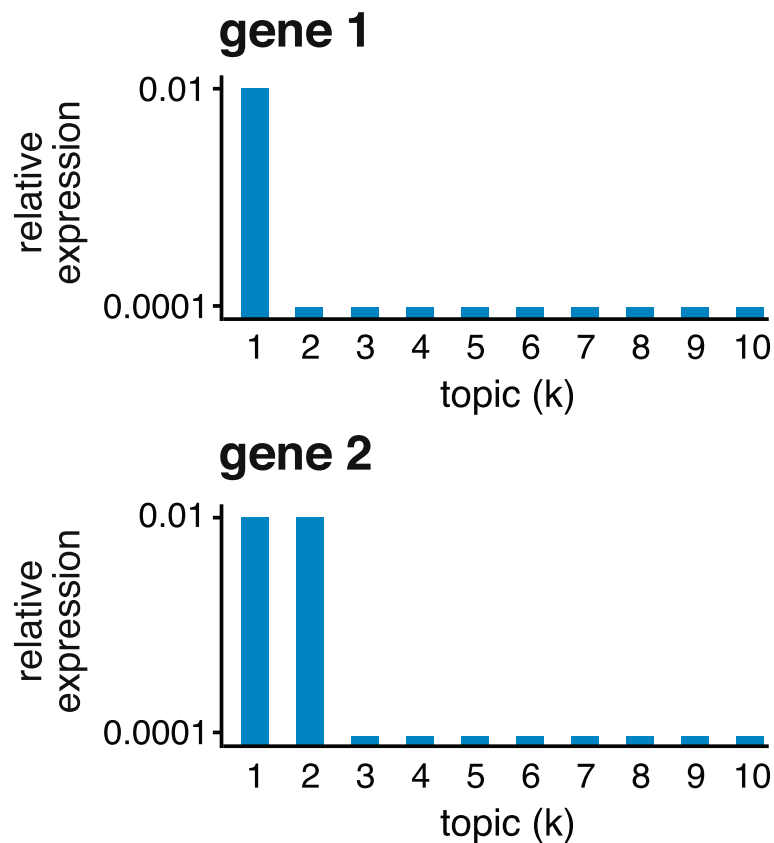


Fig. 2 Toy example illustrating the least extreme LFC. Gene 1 has high expression in topic 1 and low expression in the other topics; $p_{11} = 0.01, p_{1k} = 0.0001, k = 2, \dots, 10$. Gene 2 has high expression in topics 1 and 2 and low expression in the other topics; $p_{21} = p_{22} = 0.01, p_{2k} = 0.0001, k = 3, \dots, 10$

(Additional file 2: Table S1) to perform gene set enrichment analysis (Additional file 3: Tables S3, S4).

For the topics that closely correspond to cell types, the GoM DE analysis, as expected, identified genes and gene sets reflecting these cell types. For example, topic 1 corresponds to FACS B cells and is characterized by overexpression of *CD79A* (posterior mean l.e. LFC = 13.05) and enrichment of B cell receptor signaling genes (enrichment coefficient = 0.72). Topic 2 corresponds to myeloid cells and is characterized by overexpression of *S100A9* (l.e. LFC = 15.45) and enrichment of genes down-regulated in hematopoietic stem cells (enrichment coefficient = 0.90).

The close correspondence between topics 1 and 2 and FACS cell types (B cells, myeloid cells) provides an opportunity to contrast the GoM DE analysis with a standard DE analysis of the FACS cell types (Fig. 4). This is not a perfect comparison because the topics and FACS cell populations are not exactly the same, but the LFC estimates correlate well (Fig. 4A, B). This comparison illustrates two key differences:

1. Many more l.e. LFCs are driven toward zero in the GoM DE analysis (Fig. 4C), so the l.e. LFCs more effectively draw attention to the “distinctive genes” (Fig. 4A, B). This includes genes that are *distinctively underexpressed* such as *ID2* in B cells [81].

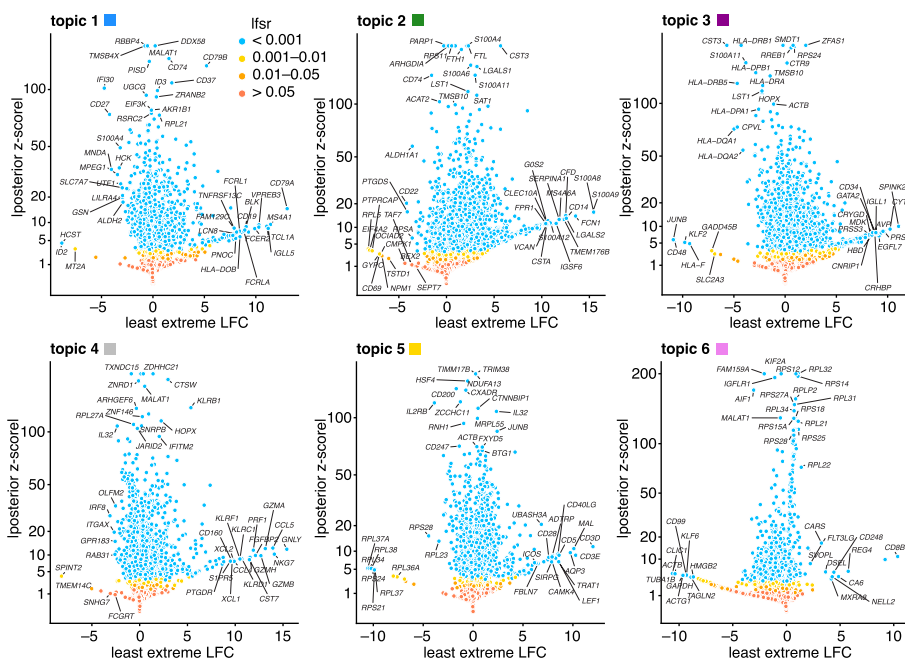


Fig. 3 GoM DE analysis of the PBMC data using the topic model shown in Fig. 1. The volcano plots show posterior mean estimates of the l.e. LFC vs. posterior z-scores for 17,055 genes. The posterior z-score is defined as the posterior mean l.e. LFC divided by the posterior standard error. Genes are colored according to the local false sign rate (*lfsr*) [82]. A few genes with extreme posterior z-scores are shown with smaller posterior z-scores so that they fit within the y-axis range. See also the detailed GoM DE results (Additional file 2: Table S1), detailed GSEA results (Additional file 3: Table S3, S4), and the interactive volcano plots (Additional file 4)

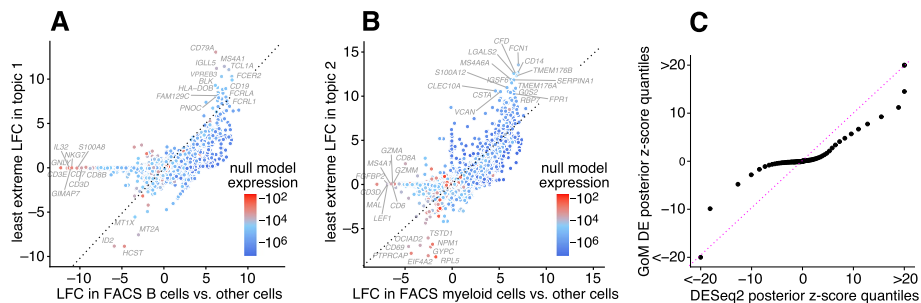


Fig. 4 GoM DE analysis vs. DESeq2 analysis in PBMC data. **A** and **B** compare differential expression in topics 1 and 2 (Fig. 1) with their closely corresponding FACS cell populations. Genes are only shown if the posterior z-score was greater than 2 in magnitude in at least one of the DE analyses. Genes are colored by the “null model” expression rate. The Q-Q plot (**C**) compares the overall distribution of posterior z-scores for B cells and myeloid cells (x-axis) and for topics 1 and 2 (y-axis). For better visualization of quantiles near zero, posterior z-scores larger than 20 in magnitude are shown as 20 or -20. Analysis of differential expression among the 6 FACS cell populations was performed using DESeq2 [79, 84]

2. The GoM DE analysis yields much larger LFC estimates of the cell-type-specific genes. This is because the topic model isolates the biological processes (topics 1 and 2) related to cell type while removing background biological processes (topic 6) that do not relate to cell type.

Other topics capture more continuous structure, such as topics 4 and 5 (Fig. 1). Although the GoM DE analysis of these topics is not comparable to a standard DE analysis, many of the distinctive genes and gene sets suggest NK and T cells, which are precisely the FACS-labeled cells with greatest membership to these topics: for example, for topic 4, overexpression of *NKG7* (posterior mean l.e. LFC = 14.09), enrichment of cytolysis genes (enrichment coefficient = 2.22); for topic 5, overexpression of *CD3D* (l.e. LFC = 12.01), enrichment co-stimulatory signaling during T-cell activation (enrichment coefficient = 1.58).

Topic 6 captures continuous structure and is present in almost all cells, so knowledge of the FACS cell types is not helpful for understanding this topic. Still, the GoM DE results for topic 6 show a striking enrichment of ribosome-associated genes (Fig. 3, Additional file 3: Tables S3, S4) (these ribosomal protein genes also account for a large fraction of the total expression in the cells [5]). This ability to annotate distinctly non-discrete structure is a distinguishing feature of the grade-of-membership approach, and below we will show more examples where this feature contributes to understanding of the cell populations.

Evaluation of DE analysis methods using simulated data

Having illustrated the features of this approach, we now evaluate the methods more systematically in simulated expression data sets. We began our evaluation by first considering the case of two groups in which there is no partial membership to these groups, that is, when the cells can be separated into two cell types. The GoM DE analysis should accommodate this special case and should compare well with existing DE analysis methods. We compared with DESeq2 [79] and MAST [83], both popular methods that have been shown to be competitive in benchmarking studies [61, 62, 85] (and are included in Seurat [25]).

To compare the ability of these methods to discover differentially expressed genes, we simulated RNA molecule count data for 10,000 genes and 200 cells in which 98% of cells were attributed to a single topic, with roughly the same number of cells assigned to each of the two topics (with membership proportions of 99% or greater). Note that although half the simulated genes had different expression levels in the two topics, most of these expression differences were small, and therefore the methods were not expected to identify most expression differences. This mimics the typical situation in gene expression studies whereby most expression differences are small. Molecule counts were simulated using a Poisson measurement model so that variation in expression across cells was due to either measurement error or true differences in expression levels between the two groups. For all DE analyses, we took group/topic assignments to be known so that incorrect assignment of cells to topics was not a source of error. Other aspects of the simulations were chosen to emulate molecule count data from scRNA-seq studies (see the “Methods” section). We repeated the simulations 20 times, and summarized the results of the DE analyses in Fig. 5 (also Additional file 1: Figs. S2, S3).

DESeq2 and the GoM DE analysis have several features in common: both are based on a Poisson model, and both use adaptive shrinkage [82, 84] to improve accuracy of the LFC estimates and test statistics. Therefore, we expected the GoM DE results to closely resemble DESeq2 in these simulations. Indeed, both methods produced nearly

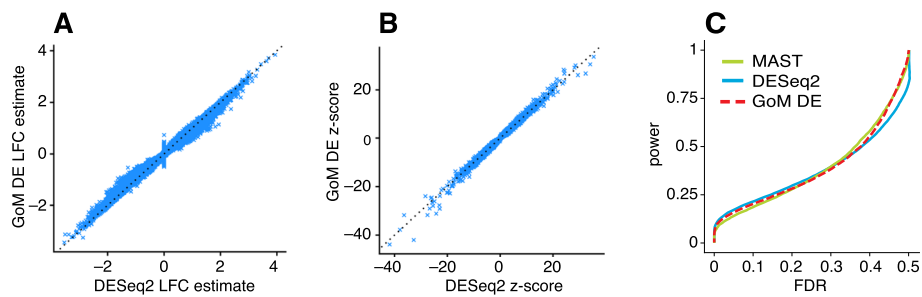


Fig. 5 Evaluation of DE analysis methods in single-cell expression data sets in which cells were simulated from two groups without partial membership to these groups. **A** and **B** compare posterior mean LFC estimates and posterior z-scores returned by DESeq2 [79] and GoM DE. Each plot shows 200,000 points for 10,000 genes \times 20 simulated data sets. **C** summarizes performance in identifying differentially expressed genes in all simulated data sets; it plots power and false discovery rates (FDR) for the three methods compared as the p -value (MAST [83]), s -value (DESeq2), or *lfsr* threshold (GoM DE) is varied from 0 to 1. Power and FDR are calculated from the number of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN) as $FDR = FP/(TP + FP)$ and $power = TP/(TP + FN)$. See also Additional file 1: Figs. S2, S3

identical posterior mean LFC estimates, posterior z-scores (Fig. 5A, B), and s -values (Additional file 1: Fig. S3) and achieved very similar performance (Fig. 5C). Although DESeq2 additionally estimates an overdispersion level for each gene, in these simulations, DESeq2 correctly determined that the level of overdispersion was small for genes with large expression differences, which explains the strong similarity of the LFC estimates and posterior z-scores. MAST, owing to an approach that is very different from DESeq2 and the GoM DE analysis, yielded estimates that were less similar (Additional file 1, Fig. S3), yet achieved comparable performance (Fig. 5C).

Next, we evaluated the GoM DE analysis methods in data sets in which the cells had varying degrees of membership to multiple topics. Since existing DE methods cannot handle the situation in which there are partial memberships to groups, we mainly sought to verify that the method behaves as expected in the ideal setting when data sets are simulated from the topic model (2). To provide some baseline for comparison, we also applied the method of Dey et al. [47], which is not strictly a DE analysis method but does provide a ranking of genes by their “distinctiveness” in each topic. This ranking is based on a simple Kullback-Leibler (K-L) divergence measure; large K-L divergences should signal large differences in expression, as well as high overall levels of expression, so large K-L divergences should correspond to small DE p -values. Since the K-L divergence is not a signed measure, we omitted tests for negative expression differences from the evaluations, which was roughly half of the total number of possible tests for differential expression.

We performed 20 simulations with $K = 2$ topics and $n = 200$ cells and another 20 simulations with $K = 6$ topics and $n = 1,000$ cells. To simplify evaluation, all genes either had the same rate of expression in all topics, or the rate was different in exactly one topic. As a result, the total number of expression differences in each data set was roughly the same regardless of the number of simulated topics. Other aspects of the simulations were kept the same as the first set of simulations (see Methods). Similar to before, we took the membership proportions to be known so that mis-estimation

of the membership proportions would not be source of error in the GoM DE analysis and in calculation of the K-L divergence scores.

The largest K-L divergence scores in the simulated data sets reliably recovered true expression differences (Fig. 6A, E). Therefore, the K-L divergence scores achieved good *true positives rates* (i.e., good power) at low *false positive rates*, $FPR = FP/(TN + FP)$ (see Fig. 5 for notation). However, for DE analysis, a more relevant performance measure is the *false discovery rate*, $FDR = FP/(TP + FP)$. Because the K-L divergence score does not fully account for uncertainty in the unknown gene expression differences, many genes with no expression differences among topics were also highly ranked, leading to poor FDR control (Fig. 6D, H). By contrast, the GoM DE analysis better accounted for uncertainty in the unknown expression levels. The GoM DE analysis also more accurately recovered true expression differences at small p -values or s -values (Fig. 6B, C, F, G) and therefore obtained much lower false discovery rates at corresponding levels of power (Fig. 6D, H). Comparing the GoM DE analysis with and without adaptive shrinkage, the adaptive shrinkage did not necessarily lead to better performance (Fig. 6D, H) but did provide more directly interpretable measures of significance (s -values or local false sign rates) by shrinking the LFC estimates and adapting the rate of shrinkage to the data; for example, the expression differences were shrunk more strongly in the $K = 6$ data sets, correctly reflecting the much smaller proportion of true expression differences (compare Fig. 6C and G).

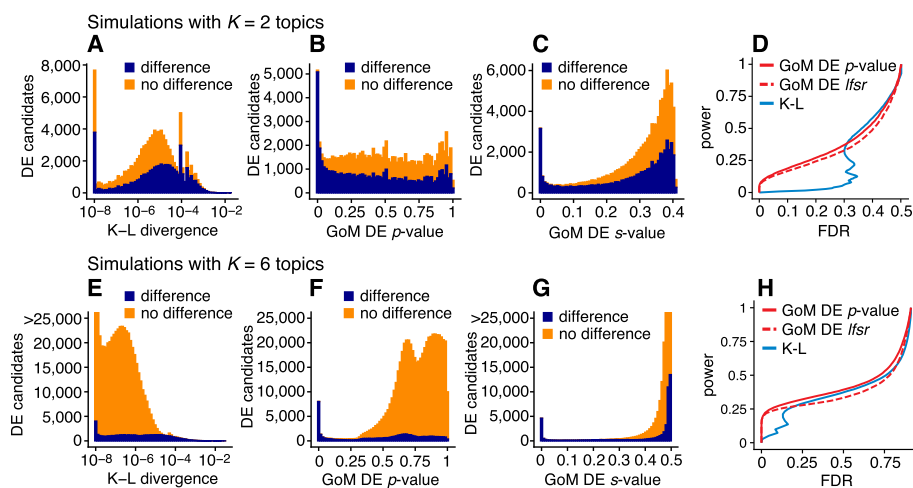


Fig. 6 Evaluation of methods for identifying expression differences in single-cell expression data sets in which cells were simulated with partial membership to 2 topics (A–D) or 6 topics (E–H). Methods compared are the Kullback-Leibler (K-L) divergence score of [47] and GoM DE with adaptive shrinkage (s -values, $lfsr$) and without adaptive shrinkage (p -values). The left-most panels (A, E) show the distribution of K-L divergence scores for all candidate expression differences (approximately half of 10,000 genes \times 2 or 6 topics \times 20 simulated data sets), shown separately for true expression differences (dark blue) and non-differences (orange). K-L divergence scores smaller than 10^{-8} are plotted as 10^{-8} . Similarly, B, C, F, and G show the distribution of GoM DE p -values or s -values with or without adaptive shrinkage, separately among differences and non-differences. D and H summarize performance in identifying expression differences; it shows power and FDR as the GoM DE p -value or $lfsr$ are varied from 0 to 1 or as the K-L divergence score is varied from large to small. Note that in E and G, some bar heights are actually larger than 25,000 but are cut off at 25,000 for better visualization

Case study: scRNA-seq epithelial airway data from Montoro et al. (2018)

We reanalyzed scRNA-seq data for $n = 7193$ single cells sampled from the tracheal epithelium in wild-type mice [86]. The original analysis [86] used a combination of methods, including t -SNE, community detection [87], diffusion maps [88], and partitioning around medoids (PAM) to identify 7 epithelial cell types: abundant basal and secretory (club) cells; rare, specialized epithelial cell types, including ciliated, neuroendocrine and tuft cells; a novel subpopulation of “ionocytes”; and a novel basal-to-club transitional cell type, “hillock” cells. Although not large in comparison to other modern single-cell data sets, this data set is challenging to analyze, with complex structure, and a mixture of abundant and rare cell types. In contrast to the PBMC data set, there are no existing cell annotations to interpret the topics, so we must rely on inferences made from the expression data alone to make sense of the results.

The topic model fit to the UMI counts with $K = 7$ topics is shown in Fig. 7A, and the results of the GoM DE analysis and subsequent GSEA are summarized in Fig. 7.

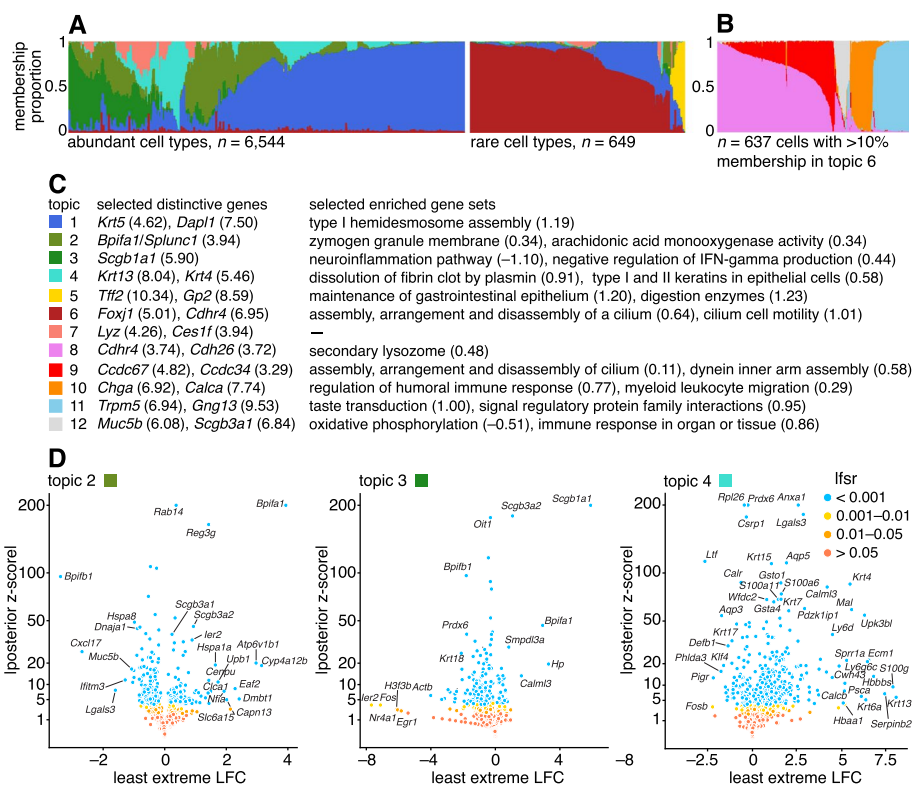


Fig. 7 Structure in mouse epithelial airway data ($n = 7193$ cells [86]) inferred from topic modeling (A, B), and GoM DE analysis (D) of selected topics using the membership proportions matrix L shown in A. In C, the topics are annotated by selected distinctive genes (numbers in parentheses are posterior mean i.e. LFCs) and selected enriched gene sets (numbers in parentheses are posterior mean estimates of the enrichment coefficients). In A, to better visualize the rare cell types, the cells were divided into two groups, “abundant” and “rare,” based on the estimated membership proportions, then the “abundant” cells were subsampled. The Structure plot in B was obtained by fitting another topic model, with $K = 5$ topics, to rare epithelial cell types (defined as the subset of 637 cells i with at least 10% membership to topic 6). The volcano plots show posterior estimates of i.e. LFC vs. posterior z-scores for 18,388 genes. A small number of genes with extreme posterior z-scores are shown with smaller posterior z-scores so that they fit within the y-axis range. See also the interactive volcano plots (Additional file 5: S6), GoM DE results (Additional file 2: Table S2, Additional file 1: Figs. S6; S7), and GSEA results (Additional file 3: Tables S5, S6)

Although we do not have cell labels to compare with, distinctive genes emerging from the GoM DE analysis help connect some of the topics to known cell types. For example, the most abundant topics correspond well with predominant epithelial cell types in the lung: topic 1 shows strong overexpression of basal cell marker gene *Krt5* [89] (posterior mean l.e. LFC = 4.62) and distinctive genes in topics 2 and 3 include key secretory genes in club cells such as *Bpifa1/Splunc1* [90] (l.e. LFC = 4.93) and *Scgb1a1* [91] (l.e. LFC = 5.90).

The “hillock” transitional cells, which were originally identified via a diffusion maps analysis [86], emerge as a single topic (topic 4, cyan), with *Krt13* (l.e. LFC = 8.04) and *Krt4* (l.e. LFC = 5.46) being among the most distinctive genes. The transitional nature of these cells is evoked by their mixed membership; only 237 out of the 7193 cells have > 90% membership to this topic.

Other less abundant epithelial cell types emerge as separate topics once a topic model is fit separately to the subpopulation of these rare cell types (Fig. 7B). These topics recover ciliated cells (topics 8, 9; *Ccdc153*, posterior l.e. LFC = 5.39), neuroendocrine cells (topic 10; *Chga*, l.e. LFC = 6.92), and tuft cells (topic 11; *Trpm5*, l.e. LFC = 6.94). Note that *Foxi1+* ionocytes were previously identified as a novel cell type from a small cluster of 26 cells [86], but our analysis failed to distinguish this very rare cell type from the neuroendocrine cells (Additional file 1: Figs. S4, S5).

The topics also capture biologically relevant *continuous substructure* in club cells (topics 2 and 3) and ciliated cells (topics 8 and 9) that was not discovered in the original analysis [86]. This continuous substructure may be reflective of finer scale cell differentiation or specialization of function. In particular, we interpret topic 3 as capturing “canonical” or “mature” (*Scgb1a1+*, l.e. LFC = 5.90) club cells [90], with negative regulation of inflammation, whereas cells with greater membership to topic 2 are “club-like” (*Bpifa1/Splunc1+*, l.e. LFC = 3.94) [89, 91]. Topic 9, similarly, appears to represent “canonical” ciliated cells, featuring upregulated genes such as such as *Ccdc67/Deup1* (l.e. LFC = 4.82) and *Ccdc34* (3.29) [89, 92, 93], and enrichment of Gene Ontology terms [94] such as cilium organization (GO:0044782) and axonemal dynein inner arm assembly (GO:0036159).

In summary, by taking a topic-model-based approach we identified and annotated well-characterized cell types such as basal cells, as well less distinct but potentially interesting substructures such as “Hillock” cells and club cell subtypes.

Case study: Mouse sci-ATAC-seq Atlas data from Cusanovich et al. (2018)

We reanalyzed data from the Mouse sci-ATAC-seq Atlas [97], comprising 81,173 single cells in 13 tissues. First, to provide an overview of the primary structure in the whole data set, we fit a topic model with $K = 13$ topics to these data. The topics correspond closely to the clusters identified in [97] (Additional file 1: Fig. S8), and several different tissues are distinguished by different topics (Fig. 8A). For the 4 tissues that have replicates, the replicates show a similar composition of the topics (Fig. 8A).

Next, we performed a more detailed analysis of just the kidney (6431 cells), fitting a topic model with $K = 10$ to just these cells. We focussed on the kidney cells because, as noted previously [97, 98], both expression and chromatin accessibility vary in relation to the spatial organization of the renal tubular cells, and we predicted that this spatial

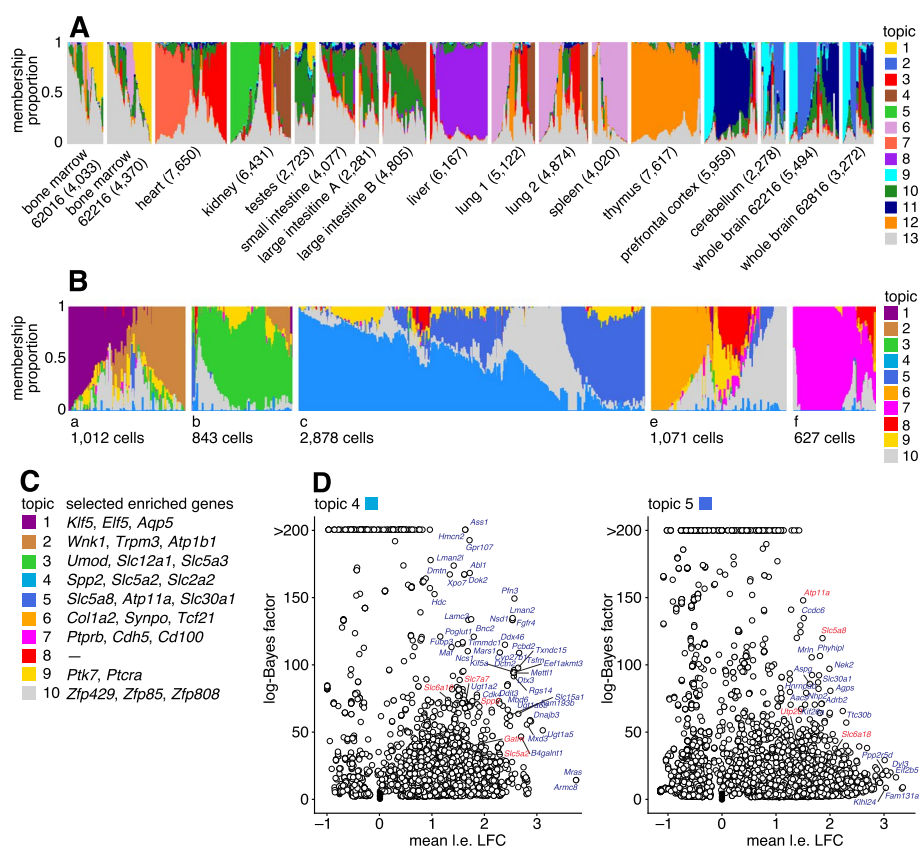


Fig. 8 **A** Structure in Mouse Atlas sci-ATAC-seq data ($n = 81,173$) inferred from topic modeling, with $K = 13$ topics. **B** Topic model fit to kidney cells ($n = 6,431$) with $K = 10$ topics. **C**, **D** Gene-based enrichment analysis of differentially accessible peaks for the kidney cell topics shown in **B**, in which peaks are linked to genes using Cicero [95]. In **A**, the cells are grouped by tissue, and replicates (for bone marrow, large intestine, lung and whole brain) are shown as separate tissues. Numbers in parentheses next to each tissue give the number of cells in that tissue. In **D**, marker genes for S1 (topic 4) and S3 (topic 5) proximal epithelial tubular cells are highlighted in red (see Table 1 of [96]). “Mean i.e. LFC” is the average i.e. LFC among all peaks connected to the gene, restricted to i.e. LFCs with $lfsr < 0.05$. Log-Bayes factors greater than 200 in the volcano plots. See Additional file 1: Fig. S9 and Additional file 6: Table S7 for more gene enrichment results. In **B**, the cells are subdivided into 5 groups (a–f) only to improve visualization. See also Additional file 1: Fig. S10 which compares the topics in **B** to cell-type predictions based on clustering [97]

structure could be better captured by topics rather than by traditional clustering methods. To interpret these topics obtained from chromatin accessibility data, we first used the GoM DE analysis to identify differentially accessible peaks for each topic; then, we used “co-accessibility” as predicted by Cicero [95, 97] to connect genes to peaks representing distal regulated sites. Finally, we performed a simple enrichment analysis to identify the “distinctive genes” for each topic, which we defined as the genes with many distal regulatory sites that were differentially accessible.

The results of these analyses are shown in Fig. 8. Many of the distinctive genes (Fig. 8, Additional file 1: Fig. S9, Additional file 6: Table S7) clearly relate topics to known kidney cell types. For example, topic 1 is enriched for genes *Klf5* and *Elf5* which relate to the collecting duct [98, 99]; topic 3 is enriched for genes *Umod* and *Slc12a1* associated with the loop of Henle [98, 100]; and topics 2, 6, and 7 are respectively enriched for

genes related to the distal convoluted tubule (*Wnk1*), podocytes (*Col1a2*) and glomerular endothelial cells (*Ptprb*).

Most interestingly, spatial organization of the proximal tubule is captured by two topics; topic 4 is enriched for *Slc5a2* (also known as *Sglt2*) and *Slc2a2* (also known as *Glut2*), associated with the S1 segment of the proximal tube [96, 101, 102], and topic 5 is enriched for *Slc5a8* (*Smct1*) and *Atp11a*, related to the S3 segment [96, 103]. This result illustrates the ability of the topic model to capture continuous variation in membership of two somewhat complementary processes, which traditional clustering methods are not designed for.

Case study: chromatin accessibility profiles of the hematopoietic system from Buenrostro et al. (2018)

Buenrostro et al. [104] studied 2034 single-cell ATAC-seq profiles of 10 cell populations isolated by FACS to characterize regulation of the human hematopoietic system. Both PCA and *t*-SNE showed, visually, the expected structure into the main developmental branches (Fig. 2 in [104]). However, neither PCA nor *t*-SNE isolated these branches as *individual dimensions* of the embedding. Identifying these branches may allow for more precise characterization of the underlying regulatory patterns. Here, by fitting a topic model to the data, the main developmental branches are identified as individual topics (Fig. 9A): topic 3, pDC; topic 4, erythroid (MEP); topic 5, lymphoid (CLP); and topic 6, myeloid (GMP and monocytes). Another topic captures the cells at the top of the developmental path (topic 1; HSC and MPP). Other cells at intermediate points in the developmental trajectory, such as CMP, GMP and LMPP cells, are more heterogeneous, and this is reflected by their high variation in topic membership.

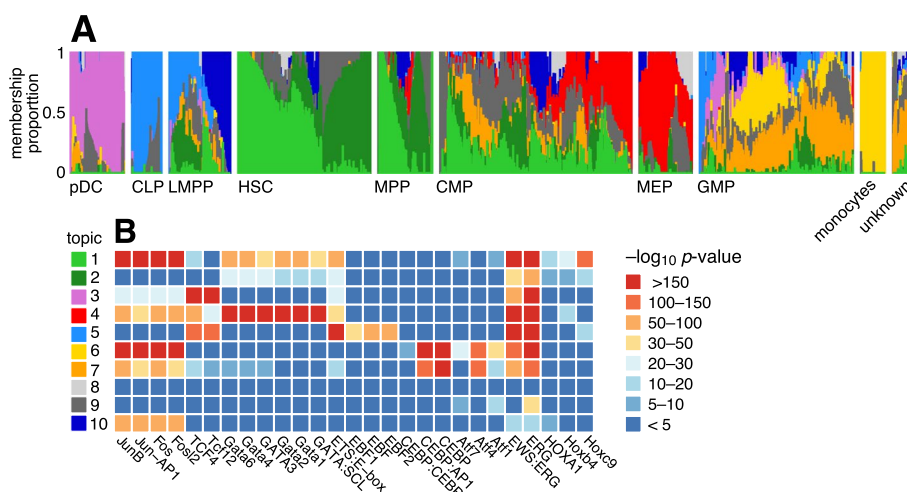


Fig. 9 Structure in human hematopoietic system data [104] ($n = 2034$ cells) inferred from the topic model with $K = 10$ topics (A) and HOMER motif enrichment analysis [105] applied to the results of the GoM DE analysis (B). In the Structure plot, the cells are grouped by FACS, as well as an unknown population from human bone marrow [104]. B shows HOMER enrichment results for selected motifs (for the full results, see Additional file 7: Table S8). Acronyms used: common lymphoid progenitor (CLP); common myeloid progenitor (CMP); granulocyte-macrophage progenitor (GMP); hematopoietic stem cell (HSC); lymphoid-primed multipotent progenitor (LMPP); megakaryocytic-erythroid progenitor (MEP); multi-potent progenitor (MPP); plasmacytoid dendritic cells (pDC)

To better interpret the regulatory patterns behind each topic, we identified transcription factor (TF) motifs that were enriched for differentially accessible regions in each topic (Fig. 9B, Additional file 7: Table S8). Many of the top TF motifs (as ranked by HOMER p -values [105]) point toward regulation of the main developmental trajectories, such as EBF motifs in topic 5 (lymphoid), CEBP motifs in topics 6 and 7 (myeloid), and Hox motifs in topic 1 (HSC and MPP cells). A few topics (topics 8–10) are much less abundant and do not align well with the FACS cell types, and their motif enrichment results were correspondingly more difficult to interpret.

A complication that arose in analyzing these data, which was also noted in [104], is that the cells were obtained from different sources, and this shows up as systematic variation in the chromatin accessibility. This donor effect is captured by topics 1 and 2 in HSC and MPP cells and, to a lesser extent, in CMP and LMPP cells (Additional file 1: Fig. S11). Topic 1 is enriched for Jun and Fos TF motifs, similar to what was found in [104].

Discussion

The GoM DE analysis is part of a *topic-model-based pipeline* for analysis of single-cell RNA-seq [47] or ATAC-seq data [49]. This pipeline includes the following steps: (1) fit a topic model to the data; (2) visualize the structure inferred by the topic model; (3) run the GoM DE analysis with the estimated topics; and, optionally, (4) perform other downstream analyses using the results of the GoM DE analysis, e.g., gene set enrichment analysis (for RNA-seq data) or motif enrichment analysis (for ATAC-seq data). Unlike most analysis pipelines for clustering and dimensionality reduction (e.g., [4, 19, 23, 26, 27]), the topic-model-based pipeline is directly applied to the “raw” count data and therefore does not require an initial step to transform and normalize the data which can lead to downstream issues in the statistical analysis [8, 106–108]. We presented several case studies illustrating the use of the topic-model-based pipeline to analyze single-cell RNA-seq and ATAC-seq data sets. From these case studies, we have drawn a few lessons on the practical challenges that may arise in applying topic modeling approaches to single-cell data, and we share these lessons here (see also [47, 49] for related discussion).

One practical question is how to choose K , the number of topics. Many papers have suggested different criteria for determining K . Our view, following [47], is that there is no single “best” K , and we recognize the advantages of learning topics at multiple settings of K ; in some data sets, different K s can reveal structure at different levels of granularity (for example, increasing the number of topics in the Mouse sci-ATAC-seq Atlas data revealed more structure within tissues; see <https://tinyurl.com/2p99swdk>). We have found that it is often helpful to start with a smaller K to elucidate the less granular structure, which is often easier to interpret, then rerun the topic modeling with larger K to identify finer structure.

We proposed annotating topics by distinctive genes identified using the l.e. LFC. One drawback is that this does not reveal the commonalities that may exist among multiple topics, for example, topics corresponding to subpopulations within a common class of cells. A simple alternative to the l.e. LFC, which is also implemented in the `fast-Topics` R package, is to compare against expression under the “null model” (see the

“Methods”). We view this as a complementary LFC metric that may reveal additional insights into the topics.

Donor, batch or other technical effects in the single-cell RNA-seq or ATAC-seq data can complicate the analysis and interpretation of the topics if these effects are not small. Since these effects are usually not known, usually we must assess their impact indirectly [109]. For example, the Mouse sci-ATAC-seq Atlas data included several replicates, but the replicate effects appeared to be small judging by the fact that the replicates showed a similar composition of topics. By contrast, the donor effects in the human hematopoietic system data were much larger, and in the topic model, these donor effects were at least partially captured by individual topics. The broader question of how to deal with non-ignorable donor or batch effects—in particular, how to separate technical effects from biological effects of interest—remains a question of considerable debate and continued investigation [25, 39, 109–116]. In particular, it has been noted that attempting to “correct” for effects can sometimes remove differences that we would like to learn about such as differences in cell-type proportions among the batches.

For modeling UMI counts, an open question is whether the Poisson or multinomial model (1) is sufficient or whether more flexible models are needed (this question was investigated in [69] for single-gene models, but not for multi-gene models). Alternative models such as the negative binomial [117] or Poisson log-normal [80, 118], which can capture additional random variation (“overdispersion”) in underlying expression or measurement error, may result in more robust estimation of the topics.

In single-cell ATAC-seq data, the GoM DE analysis identifies differentially accessible peaks or regions. Usually, these peak-level results need to be translated into biological units that are more useful for annotating the topics (e.g., genes, gene sets, transcription factors). In the analysis of the hematopoietic system single-cell ATAC-seq data, we used HOMER [105] to identify TF motifs enriched for differentially accessible peaks. In the analysis of the Mouse sci-ATAC-seq Atlas data, we identified genes enriched for differentially accessible distal regularity sites. Clearly, the quality of the gene enrichment results will depend on our ability to accurately associate peaks with genes. For this, we used the scores computed in [97] using Cicero [95]. However, there are now several alternatives to Cicero that may be preferred [19, 27, 28, 119–122], and in principle any of these approaches could be combined with the peak-level GoM DE results to identify relevant genes.

Recently developed technologies profile both transcription and chromatin accessibility in single cells [123, 124]. For such data, one could fit two topic models, one to the RNA-seq data and another to the ATAC-seq data. With a careful initialization of the topic model fitting algorithm, the topics may be more consistent across the two modalities. But it would be preferable to analyze the multimodal data jointly for improved accuracy [125–130]. Potentially, the strategy used in MOFA [131, 132] could be adapted for topic modeling—that is, the transcripts and accessibility profiles would share the same membership proportions, L , but each modality would have a different F . However, it remains to be seen how well this strategy works in practice.

Conclusions

To summarize, we have described a new method that aids in annotating and interpreting the “parts” of cells learned by fitting a topic model to scRNA-seq data or single-cell ATAC-seq data. Our method, GoM DE (differential expression analysis allowing for grades of membership), can be viewed as an extension of existing differential expression methods that allows for mixed membership to multiple groups or topics.

Methods

Models for single-cell ATAC-seq data

In single-cell ATAC-seq data, x_{ij} is the number of unique reads mapping to peak or region j in cell i . Although x_{ij} can take non-negative integer values, it is common to “binarize” the accessibility data (e.g., [19, 74, 133–135]), meaning that $x_{ij} = 1$ when at least one read in cell i maps to region j and $x_{ij} = 0$ otherwise. For this reason, one might prefer to model the binarized accessibility values as binomial (Bernoulli) random variables. A multinomial model, on the other hand, should better capture the sampling process for reads mapping to regions but does not account for the truncation of read counts above 1. Therefore, we view both the binomial and multinomial models as approximations. As we explain next, under reasonable assumptions the binomial and multinomial models are similar to each other so it may not matter which model one chooses.

The multinomial topic model for analyzing single-cell ATAC-seq data was suggested by [49]. They assumed the multinomial model (1) in which the x_{ij} s are binarized accessibility values instead of UMI counts.

A binomial model was proposed in [74],

$$x_{ij} \sim \text{Binom}(1, t_i r_j \theta_{ij}), \quad (8)$$

where $t_i > 0$ is a cell-specific factor that depends on sequencing coverage and other properties (e.g., amplification, read post-processing [136]), $r_j > 0$ is a region-specific factor (say, proportional to the size of the region), and the θ_{ij} s capture additional variation in accessibility across cells and regions. Moving forward, we make the simplifying assumption that the regions are all approximately the same size; that is, $r_j = 1$ for all $j = 1, \dots, m$.

The binomial model (8) is closely related to a multinomial model. To make the connection, we first note that the binomial model with $r_j = 1$ for all j can be approximated by a Poisson model,

$$x_{ij} \sim \text{Pois}(t_i \theta_{ij}). \quad (9)$$

This will be a good approximation when the θ_{ij} s are small and the cell-specific factors t_i are large, which is usually the case in single-cell ATAC-seq data. Next, we note that the Poisson model (9) and multinomial model (1) are closely related if we choose the size factors to be $t_i = s_i$ [69, 137]; this implies $\Theta \approx \Pi$, where Θ is the $n \times m$ matrix with entries θ_{ij} . By these arguments, the binomial model (8) (also the model used in [74]) and the multinomial model (1) (also the model used in [49]) are similar, and connecting the two models clarifies the assumptions made by each of the models. In particular, the multinomial topic model (1–2) used here and in [49] assumes a low-rank structure in the θ_{ij} s across cells and regions; *i.e.*, $\Theta \approx \mathbf{L}\mathbf{F}^T$.

Differential expression analysis allowing for grades of membership

Derivation of GoM DE model

In the “Methods overview,” we motivated the GoM DE model (3) as extending a basic Poisson model expression to allow for partial membership to K groups or topics. The GoM DE model can also be motivated from an approximation to the topic model. Recall, the topic model, is a multinomial model (1) in which the multinomial probabilities π_{ij} are given by affine combinations of the expression levels f_{jk} in the K topics, $\pi_{ij} = \sum_{k=1}^K l_{ik} f_{jk}$. The non-negativity constraints $l_{ik} \geq 0, f_{jk} \geq 0$ and sum-to-one constraints $\sum_{k=1}^K l_{ik} = 1, \sum_{j=1}^p f_{jk} = 1$ ensure that the π_{ij} s are multinomial probabilities. From a basic identity relating the multinomial and Poisson distributions [138, 139], the multinomial likelihood for the topic model can be replaced with a likelihood formed by a simple product of independent Poissons, that is,

$$\text{Multinomial}(\mathbf{x}_i; s_i, \boldsymbol{\pi}_i) \propto \prod_{j=1}^m \text{Pois}(x_{ij}; s_i \pi_{ij}), \quad (10)$$

where $\mathbf{x}_i = (x_{i1}, \dots, x_{im})$ and $\boldsymbol{\pi}_i = (\pi_{i1}, \dots, \pi_{im})$. The approximation then comes from no longer requiring the π_{ij} s to be multinomial probabilities by removing the constraint that $f_{1k} + \dots + f_{mk} = 1$. This allows us to analyze the genes $j = 1, \dots, m$ independently. This is a good approximation so long as s_i is large and the f_{jk} s are small (a similar approximation was used for GLM-PCA [8]). To be explicit about this approximation, we say $\pi_{ij} \approx \theta_{ij}$ (which are no longer guaranteed to be multinomial probabilities) and $f_{jk} \approx p_{jk}$ (which are no longer guaranteed to sum to one), resulting in the GoM DE model, which for convenience we restate here:

$$\begin{aligned} x_{ij} &\sim \text{Poisson}(s_i \theta_{ij}), \\ \theta_{ij} &= \sum_{k=1}^K l_{ik} p_{jk}. \end{aligned} \quad (11)$$

“Null” model

The simplest Poisson model of the form (3) is one in which θ_{ij} is the same across all cells i , that is, $\theta_{ij} = p_{j0}$ for all $i = 1, \dots, n$. We treat this a “null” model, which can be used to make certain comparisons, e.g., to estimate changes in expression in relative to expression in all cells. The maximum-likelihood estimate (MLE) of p_{j0} under the null model is

$$\hat{p}_{j0} = \frac{\sum_{i=1}^n x_{ij}}{\sum_{i=1}^n s_i}. \quad (12)$$

Estimation of log-fold change

In practice, we have found the l.e. LFC to work well, so in our results we use the l.e. LFC. But the l.e. LFC may not be appropriate in all circumstances, and for this reason, we note that the GoM DE analysis framework is quite general and accommodates alternatives to the l.e. LFC. Two alternatives are implemented in the software. One alternative is to compare with the “null” model,

$$\text{LFC}_k^{\text{null}}(j) := \log_2 \frac{p_{jk}}{p_{j0}}. \quad (13)$$

Another treats one topic l as a “reference topic”, and compares all other topics $k \neq l$ to l using (4).

Maximum-likelihood estimation

A convenience of the Poisson model allowing for grades of membership is that we can reuse Poisson NMF computations (described below and in more detail in [48]) to compute MLEs of the unknowns p_{jk} : if we consider all genes $j = 1, \dots, m$ simultaneously, we recover a Poisson NMF model, $x_{ij} \sim \text{Poisson}(\lambda_{ij})$, $\lambda_{ij} = \sum_{k=1}^K h_{ik} w_{jk}$, by setting $h_{ik} = s_i l_{ik}$, $w_{jk} = p_{jk}$. Therefore, we can reuse the Poisson NMF algorithms to compute MLEs of the unknowns p_{jk} .

Maximum a posteriori estimation

To improve numerical stability in the parameter estimation, we compute *maximum a posteriori* (MAP) estimates of p_{j1}, \dots, p_{jK} in which each p_{jk} is assigned a gamma prior, $p_{jk} \sim \text{Gamma}(\alpha, \beta)$, with $\alpha = 1 + \varepsilon$, $\beta = 1$, and $\varepsilon > 0$. Typically, ε will be some small, positive number, e.g., $\varepsilon = 0.1$. Here, we use the parameterization of the gamma distribution from [140] in which α is the shape parameter and β is the inverse scale parameter; under this parameterization, the mean is α/β and the variance is α/β^2 . The maximum-likelihood computations can be reused for MAP estimation with this gamma prior by adding “pseudocounts” to the data; specifically, MAP estimation of p_{j1}, \dots, p_{jK} given counts x_{1j}, \dots, x_{nj} and membership proportions \mathbf{L} and is equivalent to maximum-likelihood estimation of p_{j1}, \dots, p_{jK} given counts $x_{1j}, \dots, x_{nj}, \varepsilon, \dots, \varepsilon$ and membership proportions matrix $\begin{bmatrix} \mathbf{L} \\ \mathbf{I}_K \end{bmatrix}$, where \mathbf{I}_K is the $K \times K$ identity matrix. Unless otherwise stated, we added $\varepsilon = 0.1$ pseudocounts to the data.

Quantifying uncertainty and stabilizing LFC estimates

We implemented a simple Markov chain Monte Carlo (MCMC) algorithm [141, 142] to quantify uncertainty in the LFC estimates. Although normal approximations to likelihoods are typically used by DE methods to quickly obtain analytical measures of uncertainty (e.g., standard errors, confidence intervals) for LFCs, we found that normal approximations to the likelihoods from (5) were sometimes poorly behaved, particularly for lowly expressed genes. Another consideration was that the analytical solutions provide confidence intervals for the unknowns p_{jk} , but ultimately we are interested in quantifying uncertainty in the i.e. LFCs (6) which do not have a simple linear relationship to the p_{jk} s. Therefore, it is unclear whether the standard analytical solutions can be applied to the i.e. LFCs without making further approximations or simplifications.

MCMC is typically computationally intensive, but with careful implementation (e.g., use of sparse matrix operations and multithreaded computations) the MCMC algorithm is quite fast. Other benefits of using MCMC is that the algorithm can straightforwardly

accommodate different choices of LFC statistics and no normality assumptions are needed.

The basic idea behind the MCMC algorithm is as follows: for a given gene j , simulate the posterior distribution of the LFC statistic by performing a “random walk” on $\mathbf{g}_j = (g_{j1}, \dots, g_{jK})$, where $g_{jk} := \log p_{jk}$, $k = 1, \dots, K$. The random walk generates a sequence of states $\mathbf{g}_j^1, \dots, \mathbf{g}_j^{(n_s)}$, in which n_s denotes the pre-specified length of the simulated Markov chain. After choosing an initial state $\mathbf{g}_j^{(0)}$, each new state $\mathbf{g}_j^{(s+1)}$ is generated from the current state $\mathbf{g}_j^{(s)}$ by the following procedure: first, a topic $k \in \{1, \dots, K\}$ is chosen uniformly at random; next, a proposed state \mathbf{g}_j^* is generated as $g_{jk}^* = g_{jk} + \delta$, $\delta \sim N(0, \sigma^2)$, with $g_{jk'}^* = g_{jk'}$ for all $k' \neq k$. Assuming an (improper) uniform prior for the unknowns, $\Pr(p_{jk}) \propto 1$, the proposed state is accepted into the Markov chain with probability

$$A(\mathbf{g}_j^{(s)}, \mathbf{g}_j^*) = \left\{ 1, \frac{\Pr(\mathbf{x}_j | \mathbf{p}_j^*)}{\Pr(\mathbf{x}_j | \mathbf{p}_j^{(s)})} \times \frac{p_{jk}^*}{p_{jk}^{(s)}} \right\}, \quad (14)$$

in which \mathbf{x}_j is the j th column of the counts matrix \mathbf{X} , $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})$, and $\Pr(\mathbf{x}_j | \mathbf{p}_j)$ is the likelihood at \mathbf{p}_j , $\Pr(\mathbf{x}_j | \mathbf{p}_j) = \prod_{i=1}^n \text{Poisson}(x_i; s_i \theta_i)$ (note that \mathbf{x}_j may include pseudocounts). The standard deviation of the Gaussian proposal distribution, σ , is a tuning parameter (unless otherwise stated, we used $\sigma = 0.3$) The additional $p_{jk}^*/p_{jk}^{(s)}$ term in the acceptance probability is needed to account for the fact that we are simulating the log-transformed parameters \mathbf{g}_j , not \mathbf{p}_j ([143], p. 11). When the proposal is not accepted, the new state is simply copied from the previous state, $\mathbf{g}_j^{(s+1)} = \mathbf{g}_j^{(s)}$.

Most of the effort in running the MCMC goes into computing the acceptance probabilities (14), so we have carefully optimized these computations. For example, we have taken advantage of the fact that the count vectors \mathbf{x}_j are typically very sparse. Additionally, these computations can be performed in parallel since the Markov chains are simulated independently for each gene j .

Once Monte Carlo samples $\mathbf{g}_j^{(s)}$, for $s = 1, \dots, n_s$, have been simulated by this random-walk MCMC, we compute posterior mean LFC estimates and quantify uncertainty in the LFC estimates. For example, expressing the l.e. LFC for gene j and topic k as a function of the unknowns, $\text{LFC}_k^{\text{l.e.}}(\mathbf{p}_j)$, the posterior mean l.e. LFCs are calculated as $E[\text{LFC}_k^{\text{l.e.}}(\mathbf{p}_j)] \approx \sum_{s=1}^{n_s} \text{LFC}_k^{\text{l.e.}}(\mathbf{p}_j^{(s)})/n_s$.

The final step in the GoM DE analysis is to perform adaptive shrinkage [82] to stabilize the posterior mean estimates. To implement this step, we used the `ash` function from the `ashr` R package [144]. We used the same settings as DESeq2 to replicate as closely as possible the performance of DESeq2 with adaptive shrinkage. DESeq2 calls `ash` with `method = "shrink"`, which sets the prior to be a mixture of uniforms without a point-mass at zero.

The adaptive shrinkage method takes as input a collection of effect estimates $\hat{\beta}_1, \dots, \hat{\beta}_m$ and associated standard errors $\hat{s}_1, \dots, \hat{s}_m$. In this setting, it is not immediately obvious what are the standard errors, in part because the posterior distribution of the unknowns is not always symmetric about the mean or median. To provide a reasonable substitute summarizing uncertainty in the estimates, we computed Monte Carlo estimates of highest posterior density (HPD) intervals. A $(1 - \alpha)$ HPD interval is the smallest interval that

contains $100(1 - \alpha)\%$ of the probability mass [145, 146]. Specifically, let $[a_{jk}, b_{jk}]$ denote the $(1 - \alpha)$ HPD interval for the LFC estimate of gene j in topic k , and let $\hat{\beta}_{jk}$ denote the posterior mean. We defined the standard error as $\hat{s}_{jk} = b_{jk} - \hat{\beta}_{jk}$ when $\hat{\beta}_{jk} < 0$; otherwise, $\hat{s}_{jk} = \hat{\beta}_{jk} - a_{jk}$. Defining the standard errors in this way prevented overshrinking of estimates that were uncertain but had little overlap with zero. We set the size of the HPD intervals to $1 - \alpha = 0.68$ so that the \hat{s}_{jk} would recover conventional standard error calculations when the posterior distribution is well approximated by the normal distribution. The revised posterior means and standard errors returned by the adaptive shrinkage method were then used by `ashr` to calculate test statistics including posterior z -scores (defined as the posterior mean divided by the posterior standard error [147]), local false sign rates (*lfsr*), and s -values.

An important question is the choice of n_s . One heuristic way to assess whether n_s is large enough is to perform two independent MCMC runs initialized with different pseudorandom number generator states (“seeds”) and check consistency of the posterior estimates from the two runs (we checked consistency of the posterior estimates after stabilizing the estimates using adaptive shrinkage, as described above). In simulated data sets (below), comparison of two independent MCMC runs suggested that $n_s = 10,000$ was sufficient to obtain reasonably accurate estimates of posterior means and posterior z -scores for all genes (Additional file 1: Fig. S2). Therefore, we performed initial MCMC simulations for all single-cell data sets using $n_s = 10,000$. The runtimes for performing these MCMC simulations on the single-cell data sets (described below), with $n_s = 10,000$, are given in Table 1.

Although this consistency check suggested that running a simulation with $n_s = 10,000$ would be “good enough,” to provide additional assurance we performed another consistency assessment in the PBMC data set. We found that even better consistency was achieved with $n_s = 100,000$ (Additional file 1: Fig. S12). Therefore, to provide more reliable results, the final GoM DE results were generated with $n_s = 100,000$.

The GoM DE analysis methods are implemented in the `de_analysis` function in the `fastTopics` package [148].

Single-cell data sets

All data sets analyzed were stored as sparse $n \times m$ matrices \mathbf{X} , where n was the number of cells and m was the number of genes or regions. The data sets are summarized in Table 1.

Table 1 GoM DE simulation running times for the single-cell data sets with $n_s = 10,000$ simulation states; n is the number of cells, m is the number of genes or accessibility peaks analyzed, and K is the number of topics. See “Computing environment” for more details

Data set	n	m	K	runtime
PBMC [29]	94,655	21,952	6	5.1 h
Epithelial airway [86]	7193	18,388	7	1.0 h
Mouse Atlas, kidney only [97]	6431	270,864	10	3.2 h
Hematopoietic system [104]	2034	126,719	10	1.1 h

Preparation of scRNA-seq data

Since the topic model is a multinomial model of count data, no log-normalization or other transformation of the scRNA-seq molecule counts was needed. Furthermore, we kept all genes other than those with no variation in the data set (this is done in part to demonstrate that our methods are robust to including genes with little variation). Also note that due to the use of sparse matrix techniques in our software implementations, including genes with low variation did not greatly increase computational effort.

Preparation of single-cell ATAC-seq data

As previously suggested [19, 133–135]), we “binarized” the single-cell ATAC-seq data, that is, we assigned $x_{ij} = 1$ (“accessible”) when least one fragment in cell i mapped to peak or region j ; otherwise, $x_{ij} = 0$ (“inaccessible”). There are at least a couple reasons for doing this. For small peaks (say, < 5 kb), read counts do not provide a reliable quantitative measure of accessibility in single cells. This is because the (random) first insertion restricts the space for subsequent insertions. Additionally, insertions could occur within the same site on the same allele or on each of the two alleles, complicating interpretation of the read counts.

Like the RNA molecule count data (see above), we kept all regions except those that showed no variation.

PBMC data from Zheng et al. (2017)

We combined reference transcriptome profiles generated from 10 bead-enriched subpopulations of PBMCs (donor A) processed using Cell Ranger 1.1.0 [29, 149]. We downloaded the “Gene/cell matrix (filtered)” `tar.gz` file from the 10x Genomics website for each of the following 10 FACS-purified data sets: CD14+ monocytes, CD19+ B cells, CD34+ cells, CD4+ helper T Cells, CD4+/CD25+ regulatory T Cells, CD4+/CD45RA+/CD25- naive T cells, CD4+/CD45RO+ memory T Cells, CD56+ natural killer cells, CD8+ cytotoxic T cells, and CD8+/CD45RA+ naive cytotoxic T cells. After combining these 10 data sets, then filtering out unexpressed genes, the combined data set contained molecule counts for 94,655 cells and 21,952 genes; 97.1% of the molecule counts were zero.

In Fig. 1, the 54,132 cells from these data sets were labeled as “T cells”: CD4+ helper T Cells, CD4+/CD25+ regulatory T Cells, CD4+/CD45RA+/CD25- naive T cells, CD4+/CD45RO+ memory T Cells, and CD8+/CD45RA+ naive cytotoxic T cells.

Epithelial airway data from Montoro et al. (2018)

We analyzed a mouse epithelial airway data set from [86, 150]. These were gene expression profiles of trachea epithelial cells in C57BL/6 mice obtained using droplet-based 3' scRNA-seq, processed using the GemCode Single Cell Platform. We downloaded file `GSE103354_Trachea_droplet_UMIcounts.txt.gz`. This file also contained the cluster assignments that we compared with (in [86], the samples were subdivided into 7 clusters using a community detection algorithm). After removing genes that were not expressed in any of the cells, the data set contained molecule counts for 7193 cells and 18,388 genes (90.7% of counts were zero).

Mouse Atlas data from Cusanovich et al. (2018)

Cusanovich et al. [97] profiled chromatin accessibility by single-cell combinatorial indexing ATAC-seq (sci-ATAC-seq) [151, 152] in nuclei from 13 distinct tissues of a 8-week-old male C57BL/6J mouse. Replicates for 4 of the 13 tissues were obtained by profiling chromatin accessibility in a second mouse. We downloaded the (sparse) binarized peak \times cell matrix in RDS format, `atac_matrix.binary.qc_filtered.rds`, from the Mouse sci-ATAC-seq Atlas website [153]. We also downloaded `cell_metadata.txt` which included cell types estimated by a clustering of the cells (see Table S1 in [97]). The full data set used in our analysis (13 tissues, including 4 replicated tissues) consisted of the binary accessibility values for 81,173 cells and 436,206 peaks (1.2% overall rate of accessibility). Note that all peaks had fragments mapping to at least 40 cells, so no extra step was taken to filter out peaks.

Separately, we analyzed the sci-ATAC-seq data from kidney only, in which peaks with fragments mapping to fewer than 20 kidney cells were removed, resulting in data set containing binary accessibility values for 6431 cells and 270,864 peaks. Base-pair positions of the peaks were based on Mouse Genome Assembly mm9 (NCBI and Mouse Genome Sequencing Consortium, Build 37, July 2007).

From the Mouse sci-ATAC-seq website, we also downloaded the file `master_cicero_conns.rds` containing the Cicero co-accessibility predictions [95, 153], which we used to link chromatin accessibility peaks to genes. For the kidney data, we connected a peak given in the “Peak2” column of the Cicero co-accessibility data table to a gene given in the “peak1.tss.gene_id” column if the “cluster” column was 11, 18, 22, or 25 (these four clusters were the main kidney-related clusters identified in [97]). This extracted, for each gene, the distal and proximal sites connected to the gene associated with Peak1 (specifically, a gene in which the transcription start site overlaps with Peak1). Among the 22,194 genes associated with at least one peak, the median number of peaks connected to a gene was 19, and the largest number of peaks was 179 (for *Bahcc1* on chromosome 11). Among the 270,864 peaks included in the topic modeling analysis, 113,489 (42%) were connected to at least one gene, 95% of peaks were connected to 10 genes or fewer, and the largest number of connected genes was 60.

Human hematopoietic system data from Buenrostro et al. (2018)

Buenrostro et al. [104] used FACS to isolate 10 hematopoietic cell populations from human bone marrow and blood; then, the cells were assayed using single-cell ATAC-seq. The processed single-cell ATAC-seq data were downloaded from [154], specifically file `GSE96769_scATACseq_counts.txt.gz` containing the fragment counts and file `GSE96769_PeakFile_20160207.bed.gz` containing peaks obtained from bulk ATAC-seq data [104]. Although there may be benefits to calling peaks using aggregated single-cell data instead [155], we used the original accessibility data based on the bulk ATAC-seq peaks so that our analysis was more directly comparable to the analysis of [104].

Following [104, 155], we extracted the 2034 samples passing quality control filters; then, we “binarized” the counts. The list of 2034 cells considered “high quality” was obtained from file `metadata.tsv` included in the online benchmarking repository [155]. After removing peaks with fragments mapping to fewer than 20 cells, the final

data set used in our analysis consisted of binary accessibility values for 2034 cells and 126,719 peaks (4.6% overall rate of accessibility). Base-pair positions of the peaks were based on human genome assembly 19 (Genome Reference Consortium Human Build 37, February 2009).

In [104], a large, patient-specific batch effect was identified in the accessibility profiles for the HSC cells, and therefore, steps were taken in [104] to normalize the accessibility data before performing PCA. We instead fit the topic model to the unnormalized binary accessibility values, in part to find out how well the topic model can cope with the complication of a batch effect. In agreement with [104], this batch effect is at least partly captured by the topics, although in our analysis, the batch effect also appeared in MPP cells and, to a lesser extent, in CMP cells (Additional file 1: Fig. S11).

Fitting the topic models

In brief, we took the following steps to fit a topic model. All these steps are implemented in the R package `fastTopics`.

First, we fit a Poisson NMF model [37, 156],

$$\begin{aligned} x_{ij} &\sim \text{Poisson}(\lambda_{ij}) \\ \Lambda &= \mathbf{H}\mathbf{W}^T, \end{aligned} \tag{15}$$

where $\Lambda \in \mathbf{R}^{n \times m}$ is a matrix of the same dimension as \mathbf{X} with entries $\lambda_{ij} \geq 0$ giving the Poisson rates for the counts x_{ij} . The parameters of the Poisson NMF model are stored as two matrices, $\mathbf{H} \in \mathbf{R}^{n \times K}$, $\mathbf{W} \in \mathbf{R}^{m \times K}$, with non-negative entries h_{ik}, w_{jk} . `fastTopics` has efficient implementations of algorithms for computing maximum-likelihood estimates (MLEs) of \mathbf{W}, \mathbf{H} [48].

Second, we recovered MLEs of \mathbf{F}, \mathbf{L} from MLEs of \mathbf{W}, \mathbf{H} by a simple reparameterization [48].

In an empirical comparison of Poisson NMF algorithms with count data sets, including scRNA-seq data sets [48], we found that a simple co-ordinate descent (CD) algorithm [157, 158], when accelerated with the extrapolation method of Ang and Gillis [159], almost always produced the best Poisson NMF (and topic model) fits, and in the least amount of time. To confirm this, we compared topic model fits obtained by running the same four algorithms that were compared in [48]—EM and CD, with and without extrapolation—on the PBMC data set and assessed the quality of the fits. We evaluated the model fits in two ways: using the likelihood and using the residuals of the Karush-Kuhn-Tucker (KKT) first-order conditions (the residuals of the KKT system should vanish as the algorithm approaches maximum-likelihood estimates of \mathbf{W}, \mathbf{H}). Following [48], to reduce the possibility that multiple optimizations converge to different local maxima of the likelihood, which could complicate these comparisons, we first ran 1000 EM updates; then, we examined the performance of the algorithms after this initialization phase (Additional file 1: Figs. S13, S14). Consistent with [48], the extrapolated CD updates always produced the best fit, or at the very least a fit that was no worse than the other algorithms, and almost always converged on a solution more quickly than the other algorithms. Therefore, subsequently we used the extrapolated CD updates to fit the topic models.

In more detail, the pipeline for fitting topic models consisted of the following steps: (1) initialize \mathbf{W} using Topic-SCORE [160], (2) perform 10 CD updates of \mathbf{H} , with \mathbf{W} fixed, (3) perform 1000 EM updates (without extrapolation) to get close to a solution (“prefitting phase”), (4) run an additional 1000 extrapolated CD updates to improve the fit (“refinement phase”), and (5) recover \mathbf{F} , \mathbf{L} from \mathbf{W} , \mathbf{H} by a simple transformation. The prefitting phase was implemented by calling `fit_poisson_nmf` from `fastTopics` with these settings: `numiter = 1000, method = "em", control = list(numiter = 4)`. The refinement phase was implemented with a second call to `fit_poisson_nmf`, with `numiter = 1000, method = "scd", control = list(numiter = 4, extrapolate = TRUE)`, in which the model fit was initialized using the fit from the prefitting phase. The topic model fit was recovered by calling `poisson2multinom` in `fastTopics`. Note that only the estimates of \mathbf{L} were used in the GoM DE analysis.

For each data set, we fit topic models with different choices of K and compared the fits for each K by comparing their likelihoods (Additional file 1: Fig. S5).

Visualizing the membership proportions

The membership proportions matrix \mathbf{L} can be viewed as an embedding of the cells $i = 1, \dots, n$ in a continuous space with $K - 1$ dimensions [50] (it is $K - 1$ dimensions because of the constraint that the membership proportions for each cell must add up to 1). A simple way to visualize this embedding in 2-d is to apply a nonlinear dimensionality reduction technique such as t -SNE [11, 161] or UMAP [12] to \mathbf{L} ([49] used t -SNE). We have also found that plotting principal components (PCs) of the membership proportions can be an effective way to explore the structure inferred by the topic model (Additional file 1: Figs. S1, S4). However, we view these visualization techniques as primarily for exploration, and a more powerful approach is to visualize all $K - 1$ dimensions simultaneously using a Structure plot [70, 71]. Here, we describe some improvements to the Structure plot for better visualization. These improvements are implemented in the `structure_plot` function in `fastTopics`.

When cells were labeled, we compared topics against labels by grouping the cells by these labels in the Structure plot. We then applied t -SNE to the \mathbf{L} matrix, separately for each group, to arrange the cells on a line within each group. For this, we used the R package `Rtsne` [162] (in `fastTopics`, we also implemented options to arrange the cells in each group using UMAP or PCA, but in our experience we found that UMAP and PCA produced “noisier” visualizations).

Arranging the cells by 1-d t -SNE worked best for smaller groups of cells with less complex structure. For large groups of cells, or for unlabeled single-cell data sets, we randomly subsampled the cells to reduce t -SNE runtime (when cells number in the thousands, it is nearly impossible to distinguish individual cells in the Structure plot anyhow). Even with this subsampling, the Structure plot sometimes did not show fine-scale substructures or rare cell types. Therefore, in more complex cases, we first subdivided the cells into smaller groups based on the membership proportions, then ran t -SNE on these smaller groups. These groups were either identified visually from PCs of \mathbf{L} or in a more automated way by running k -means on PCs of \mathbf{L} (see [163]).

Gene enrichment analysis based on differential accessibility of peaks connected to genes

Here, we describe a simple approach to obtain gene-level statistics from the results of a differential accessibility analysis. This approach was applied in the topic modeling analysis of the Mouse Atlas kidney cells.

Cusanovich et al. [97] used the Cicero co-accessibility predictions and the binarized single-cell ATAC-seq data to compute a “gene activity score” R_{ki} for each gene k and cell i . Here, we have a related but different goal: we would like to use the results of the differential accessibility analysis, which generates differential accessibility estimates and related statistics for each peak and each topic, to rank genes according to their importance to a given topic. A difficulty, however, with ranking the genes is that the Cicero co-accessibility predictions are uncertain, and they are only partially informative about which peaks are relevant to a gene. In aggregate, however, the expectation is that the “most interesting” genes will be genes that are (a) predicted by Cicero to be connected many peaks that are differentially accessible and (b) the differences in accessibility are mainly in the same direction. This suggests an enrichment analysis in which, for each gene, we test for enrichment of differential accessibility among the peaks connected to that gene. Here, we describe a simple enrichment analysis for (a) and (b).

For (a), we computed a Bayes factor [164] measuring the support for the hypothesis that at least one of the peaks is differentially accessible (the LFC is not zero) against the null hypothesis that none of the peaks are differentially accessible. For (b), we computed the *average LFC* among all differentially accessible peaks (that is, peaks with nonzero LFC according to some significance criterion).

We implemented this gene enrichment analysis by running adaptive shrinkage [82] separately for each gene and topic. This had the benefit of adapting the shrinkage separately to each gene in each topic. In particular, in comparison to the usual adaptive shrinkage step for a GoM DE analysis (see above), it avoided overshrinking differences for genes exhibiting strong patterns of differential accessibility. We took the following steps to implement this adaptive shrinkage analysis. First, we ran function `ash` from the `ashr` package [144] once on the posterior mean i.e. LFC estimates $\hat{\beta}_{jk}$ and their standard errors $\hat{\sigma}_{jk}$ for all topics k and all peaks j , with settings `mixcompdist = "normal", method = "shrink"`. This was done only to determine the variances in the mixture prior and to get a “default” model fit to be used in the subsequent adaptive shrinkage analyses.

Next, we ran `ash` separately for gene and each topic k using the i.e. LFC estimates $\hat{\beta}_{jk}$ and standard errors $\hat{\sigma}_{jk}$ from the peaks j connected to that gene. We set the variances in the mixture prior to the variances determined from all the i.e. LFC estimates, and used `ash` settings `mixcompdist = "normal"` and `pointmass = FALSE`. One issue with running adaptive shrinkage using only the i.e. LFC estimates for the peaks connected to a gene is that some genes have few Cicero connections, leading to potentially unstable fits and unreliable posterior estimates. We addressed this issue by encouraging the fits toward the “default model” that was fitted to all genes and all topics; specifically, we set the Dirichlet prior on the mixture proportions to be $\text{Dirichlet}(\alpha_1, \dots, \alpha_K)$ with prior sample sizes $\alpha_k = 1.01 + n_0 \hat{\pi}_k^{\text{default}}$, where here K denotes the number of components of the prior mixture (not the number of

topics), $\hat{\pi}_k^{\text{default}}$ denotes the k th mixture proportion in the adaptive shrinkage prior for the fitted “default” model, and $n_0 = 20$. This stabilized the fits for genes with few Cicero connections while still allowing some ability to adapt to genes with many connections.

Finally, we used the `logLR` output from `ash` as a measure of support for enrichment (this is the Bayes factor on the log-scale), and we computed the mean l.e. LFC as the average of the posterior mean estimates of the l.e. LFCs taken over all peaks j connected to the gene and with posterior $lfsr < 0.05$.

Motif enrichment analysis for differentially accessible regions

We used HOMER [105] to identify transcription factor (TF) motifs enriched for differentially accessible regions, separately for each topic estimated from the single-cell ATAC-seq data. For each topic $k = 1, \dots, K$, we applied the HOMER Motif Analysis tool `findMotifsGenome.pl` to estimate motif enrichment in differentially accessible regions; specifically, we took “differentially accessible regions” to be those with p -value less than 0.05 in the GoM DE analysis (Additional file 1: Fig. S16). These differentially accessible regions were stored in a BED file `positions.bed`. The exact call from the command-line shell was `findMotifsGenome.pl positions.bed hg19 homer -len 8,10,12 -size 200 -mis 2 -S 25 -p 4 -h`.

Note that the adaptive shrinkage step was skipped in the GoM DE analysis, so these are the p -values for the unmoderated l.e. LFC estimates. The reason for skipping the adaptive shrinkage step is that the shrinkage is performed uniformly for the LFC estimates for all regions, and since the vast majority of regions have l.e. LFC estimates that are indistinguishable from zero, the result is that very few differentially accessible regions remain shrinkage.

Gene sets

Human and mouse gene sets for the gene set enrichment analyses (GSEA) were compiled from the following gene set databases: NCBI BioSystems [165], Pathway Commons [166, 167], and MSigDB [168–170], which includes Gene Ontology (GO) gene sets [94, 171]. Specifically, we downloaded `bsid2info.gz` and `biosystems_gene.gz` from the NCBI FTP site (<https://ftp.ncbi.nih.gov/gene>) on March 22, 2020; `PathwayCommons12.All.hgnc.gmt.gz` from the Pathway Commons website (<https://www.pathwaycommons.org>) on March 20, 2020; and `msigdb_v7.2.xml.gz` from the MSigDB website (<https://www.gsea-msigdb.org>) on October 15, 2020. For the gene set enrichment analyses, we also downloaded human and mouse gene information (“gene info”) files `Homo_sapiens_gene_info.gz` and `Mus_musculus_gene_info.gz` from the NCBI FTP site on October 15, 2020. Put together, we obtained 37,856 human gene sets and 33,380 mouse gene sets. In practice, we filtered gene sets based on certain criteria before running the GSEA. To facilitate integration of these gene sets into our analyses, we have compiled these gene sets into an R package [172].

Gene set enrichment analysis

We took a simple multiple linear regression approach to the gene set enrichment analysis (GSEA), in which we modeled the l.e. LFC estimate for gene i in a given topic, here denoted by y_i , as $y_i = \mu_i + \sum_{j=1}^n x_{ij}b_j + e_i$, $e_i \sim N(0, \sigma^2)$, in which $x_{ij} \in \{0, 1\}$ indicates gene set membership; $x_{ij} = 1$ if gene i belongs to gene set j , otherwise $x_{ij} = 0$ (we represented the gene-set membership as a sparse matrix since most x_{ij} s are zero). Here, n denotes the number of candidate gene sets, and σ^2 is the residual variance to be estimated. The idea behind this simple approach was that the most relevant gene sets are those that best explain the log-fold changes y_i , and therefore, in the multiple regression, we sought to identify these gene sets by finding coefficients b_j that were nonzero with high probability. See [173, 174] for similar ideas using logistic regression. Additionally, since many genes were typically differentially expressed in a given topic, modeling LFCs helped distinguish among DE genes that showed only a slight increase in expression versus those that were highly overexpressed [175, 176]. Of course, this simple multiple linear approach ignores uncertainty in the LFC estimates y_i , which is accounted for in most gene set enrichment analyses. We addressed this issue by shrinking the l.e. LFC estimates *prior to running the GSEA*, that is, we took y_i to be the the posterior mean LFC estimate after applying adaptive shrinkage, as described above (see the “[Quantifying uncertainty and stabilizing LFC estimates](#)” section). The result was that genes that we were more uncertain about had have an l.e. LFC estimate y_i that was zero or near zero.

We implemented this multiple linear regression approach using SuSiE (`susieR` version 0.12.10) [177]. A benefit to using SuSiE is that it automatically organized similar or redundant gene sets into “credible sets” (CSs), making it easier to quickly recognize complementary gene sets; see [178–183] for related ideas.

In detail, the GSEA was performed as follows. We performed a separate GSEA for each topic, $k = 1, \dots, K$. Specifically, for each topic k , we ran the `susieR` function `susie` with the following options: `L = 10`, `intercept = TRUE`, `standardize = FALSE`, `estimate_residual_variance = TRUE`, `refine = FALSE`, `compute_univariate_zscore = FALSE` and `min_abs_corr = 0`. We set `L = 10` so that SuSiE returned at most 10 credible sets. For a given topic k , we reported a gene set as being enriched if it was included in at least one CS. We organized the enriched gene sets by (95%) credible sets. We also recorded the Bayes factor for each CS, which gives a measure of the level of support for that CS. For each gene set included in a CS, we reported the posterior inclusion probability (PIP) and the posterior mean estimate of the regression coefficient b_j . In the results, we refer to b_j as the “enrichment coefficient” for gene set j since it is an estimate of the expected increase in the l.e. LFC for genes that belong to gene set j relative to genes that do not belong to the gene set.

Table 2 Number of gene sets included in each GSEA

Data set	All gene sets	Curated only
PBMC	23,193	12,225
Mouse epithelial airway	20,917	9946
—rare epithelial cell types only	20,288	9450

Often, a CS contained only one gene set, in which case the PIP for that gene set was close to 1. In several other cases, the CS contained multiple similar gene sets; in these cases, the smaller PIPs indicated that it was difficult to choose among the gene sets because they are similar to each other (note that the sum of the PIPs in a 95% CS should always be above 0.95 and less than 1). Occasionally, SuSiE returned a CS with a small Bayes factor containing a very large number of gene sets. We excluded such CSs from the results.

We repeated these gene set enrichment analyses with two collections of gene sets: (1) all gene sets other than the MSigDB collections C1, C3, C4, and C6 and “archived” MSigDB gene sets and (2) only gene sets from curated pathway databases, specifically Pathway Commons, NCBI BioSystems and “canonical pathways” (CP) in the MSigDB C2 collection, and Gene Ontology (GO) gene sets in the MSigDB C5 collection. In all cases, we removed gene sets with fewer than 10 genes and with more than 400 genes. Table 2 gives the exact number of gene sets included in each GSEA.

Simulations

For evaluating the DE analysis methods, we generated matrices of UMI counts $\mathbf{X} \in \mathbf{R}^{n \times m}$ for $m = 10,000$ genes and $n = 200$ or $n = 1000$ cells. We simulated the UMI counts x_{ij} from a Poisson NMF model (15) in which \mathbf{W} and \mathbf{H} were chosen to emulate UMI counts from scRNA-seq experiments.

The matrices \mathbf{W} and \mathbf{H} were generated as follows. First, for each cell i , we generated membership proportions l_{i1}, \dots, l_{iK} then set $h_{ik} = s_i l_{ik}$, for $k = 1, \dots, K$, where s_i is the total UMI count. To simulate the wide range of total UMI counts often seen in scRNA-seq data sets, total UMI counts s_i were normally distributed on the log-scale, $s_i = 10^{u_i}$, $u_i \sim N(0, 1/5)$, where $N(\mu, \sigma)$ denotes the univariate normal distribution with mean μ and standard deviation σ .

Membership proportions l_{ik} for each cell i were generated so as to obtain a wide range of mixed memberships, according to the following procedure: the number of nonzero proportions was set to $K' \in \{1, \dots, K\}$ with probability $2^{-K'}$; the K' selected topics $t_1, \dots, t_{K'} \subseteq \{1, \dots, K\}$ were drawn uniformly at random (without replacement) from $1, \dots, K$; then, the membership proportions for the selected topics were set to 1 when $K' = 1$, or, when $K' > 1$, they were drawn from the Dirichlet distribution with shape parameters $\alpha_{t_1}, \dots, \alpha_{t_{K'}}$.

Expression rates w_{jk} were generated so as to emulate the wide distribution of gene expression levels observed in single-cell data sets and to allow for differences in expression rates among topics. The procedure for generating the expression rates for each gene j was as follows: with probability 0.5, the expression rates were the same across all topics and were generated as $f_{j1} = \dots = f_{jK} = 2^{v_j}$, $v_j \sim N(-4, 2)$. Otherwise, with probability 0.5, the expression rates were the same in all topics except for one topic. The differing topic k' was chosen uniformly at random from $1, \dots, K$; then, the expression rate for topic k' was set to $f_{jk'} = 2^{v_j + e_j}$, $e_j \sim N(0, 1)$. As a result, the expression rates were roughly normally distributed on the log-scale, and the expression differences were also normally distributed on the log-scale. About half of genes had an expression difference among the topics.

Using this simulation procedure, we generated three collections of data sets. The simulation settings were altered slightly for each collection. In the first, data sets were simulated with $K = 2$, $\alpha = (1/100, 1/100)$, $n = 200$ so that most membership proportions were equal or very close to 0 or 1. In the second, we used $K = 2$, $\alpha = (1, 1)$, $n = 200$ to allow for a range of mixed memberships. In the third, we generated data sets with $K = 6$, $\alpha = (1, \dots, 1)$, $n = 1,000$.

For the data sets simulated with $K = 2$, $\alpha = (1/100, 1/100)$, the cells could essentially be subdivided into two groups. Therefore, we ran MAST [83, 184] and DESeq2 [78, 84] to test for genes that were differentially expressed between the two groups. MAST (R package version 1.20.0) was called via the FindMarkers interface in Seurat [25] (Seurat 4.0.3, SeuratObject 4.0.2) with the following settings: `ident.1 = "2", ident.2 = NULL, test.use = "MAST", logfc.threshold = 0, min.pct = 0`. DESeq was called from the DESeq2 R package (version 1.34.0) using settings recommended in the package vignette: `test = "LRT", reduced = ~1, useT = TRUE, minmu = 1e-6, minReplicatesForReplace = Inf`. Size factors were calculated using the `calculateSumFactors` method from scran version 1.22.1 [23]. The LFC estimates returned by DESeq were subsequently revised using adaptive shrinkage [82] by calling `lfcShrink` in DESeq2 with `type = "ashr", svalue = TRUE` (as in the GoM DE analysis, the DESeq2 posterior z -scores were defined as the posterior means divided by the posterior standard errors returned by the adaptive shrinkage).

To perform the GoM DE analysis in each of the simulations, we first fit a Poisson NMF model to the simulated counts \mathbf{X} using `fit_poisson_nmf` from the `fastTopics` R package [48, 148] (version 0.6-97). The loadings matrix \mathbf{H} was fixed to the matrix used to simulate the data, and \mathbf{W} was estimated by running 40 co-ordinate ascent updates on \mathbf{W} alone (`update.loadings = NULL, method = "scd", numiter = 40`). The equivalent topic model fit was then recovered. Three GoM DE analyses were performed using the `de_analysis` function from the `fastTopics` R package, with the topic model fit provided as input: one analysis without adaptive shrinkage (`shrink.method = "none"`), and two analyses with adaptive shrinkage (`shrink.method = "ashr", ashr` version 2.2-51 [144]) in which the MCMC was initialized with different pseudorandom number generator states. In all three runs, posterior calculations were performed with $n_s = 10,000$, $\varepsilon = 0.01$. Comparison of the two MCMC runs (with adaptive shrinkage) suggested that $n_s = 10,000$ was sufficient to obtain reasonably accurate posterior estimates in these simulations (Additional file 1: Fig. S2).

Computing environment

Most computations on real data sets were run in R 3.5.1 [185], linked to the OpenBLAS 0.2.19 optimized numerical libraries, on Linux machines (Scientific Linux 7.4) with Intel Xeon E5-2680v4 ("Broadwell") processors. For performing the Poisson NMF optimization, which included some multithreaded computations, as many as 8 CPUs and 16 GB of memory were used. The DESeq2 analysis of the PBMC data was performed in R 4.1.0, using 4 CPUs and 264 GB of memory. The evaluation of the DE analysis methods in simulated data sets was performed in R 4.1.0, using as many as 8 CPUs as 24 GB of memory. More details about the computing environment, including the R packages used, are recorded in the workflow pages in the companion code repositories [186, 187].

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s13059-023-03067-9>.

Additional file 1: Supplementary Figures. Contains **Figures S1–S16**.

Additional file 2: Tables S1, S2. Tables giving detailed statistics about DE genes identified in GoM DE analyses of PBMC (**Table S1**) and epithelial airway data (**Table S2**), with *lfsr* less than 0.01. Columns from left to right are: topic; ensembl id of gene (**Table S1** only); gene symbol; posterior mean estimate of the l.e. LFC; posterior z-score; *lfsr*.

Additional file 3: Tables S3–S6. Tables summarizing results of gene set enrichment analyses for PBMC (**Tables S3, S4**) and mouse epithelial airway data (**Tables S5, S6**). Columns from left to right are: topic; credible set (CS); log-Bayes factor (lbf); posterior inclusion probability (pip); SuSiE posterior mean estimate of the enrichment coefficient (coef); gene set name; gene set id; database; accession, when applicable; sub-category code, when applicable; organism; brief description of the gene set; and a list of the top 10 genes, defined as the members of the gene set with the largest l.e. LFC (by magnitude).

Additional file 4. Interactive volcano plots for PBMC data. Interactive volcano plots for browsing the results from the GoM DE analyses of the PBMC data. Detailed l.e. LFC statistics are displayed on mouseover: lower and upper limit of HPD interval; posterior mean estimate; posterior z-score; and *lfsr*. The maximum-likelihood estimate of the expression rate, p_{j0} , in the “null” expression model (12) is also given. Note that the lower and upper HPD intervals were not updated in the adaptive shrinkage step and therefore should be ignored.

Additional file 5. Interactive volcano plots for epithelial airway data. Interactive volcano plots for browsing the results from the GoM DE analyses of the mouse epithelial airway data. Detailed l.e. LFC statistics are displayed on mouseover: lower and upper limit of HPD interval; posterior mean estimate; posterior z-score; and *lfsr*. The maximum-likelihood estimate of the expression rate, p_{j0} , in the “null” expression model (12) is also given. Note that the lower and upper HPD intervals were not updated in the adaptive shrinkage step and therefore should be ignored.

Additional file 6: Table S7. Table giving detailed statistics about genes identified in the gene enrichment analysis for the Mouse Atlas kidney cells. All genes with log-Bayes factor > 17 are included in this table (here, we use the natural logarithm). Assuming (conservatively) that 1 out of the 22,142 tested genes is enriched, this Bayes factor corresponds to a posterior odds (PO) of $PO = e^{17}/22,142 \approx 1,000$, or a posterior probability of about 0.999. Columns from left to right are: topic; gene symbol; Ensembl gene id; log-Bayes factor; average l.e. LFC from differentially accessible peaks connected to the gene (*lfsr* < 0.05).

Additional file 7: Table S8. Table giving results of the HOMER motif enrichment analysis for the human hematopoietic system data. Columns from left to right are: motif; consensus sequence; (base-10 logarithm of) the enrichment *p*-value for the 10 topics.

Additional file 8. Review history. Document with the review history.

Acknowledgements

We thank Mihai Anitescu, Nicolas Chevrier, Michihiro Takahama, Kushal Dey, Adam Gruenbaum, Youngseok Kim, John Novembre, Alan Selewa, Katie Rhodes, Yusha Liu, Dongyue Xie, Zepeng Mu, and Jason Willwerscheid for their valuable input. We also thank the staff at the Research Computing Center for providing the high-performance computing resources used to implement the numerical experiments.

Peer review information

Andrew Cosgrove was the primary editor of this article and managed its editorial process and peer review in collaboration with the rest of the editorial team.

Review history

The review history is available as Additional file 8.

Authors' contributions

PC implemented the methods based on mathematical derivations by PC and MS, with contributions from KL, AS, and SP. PC performed empirical evaluations of the methods using simulated data sets. PC and KL analyzed the single-cell data sets. PC wrote the draft manuscript, with contributions, revisions, and discussion by all authors. All authors approved the final manuscript.

Funding

This work was supported by the NHGRI at the National Institutes of Health under award number 5R01HG002585.

Availability of data and materials

The fastTopics R package is available on GitHub (<https://github.com/stephenslab/fastTopics>) and CRAN (<https://cran.r-project.org/package=fastTopics>). A Seurat wrapper for fastTopics is available from <https://github.com/stephenslab/seurat-wrappers>. The data sets supporting the conclusions of this article are available in Zenodo repositories [186, 187]. These Zenodo repositories also include the source code implementing the analyses and workflow websites [188] for browsing the code and results. Permission to use the source code in these repositories is granted under the MIT license. Numerical implementations of the contributed statistical methods, including tools for visualizing the results generated by these methods, are available from the fastTopics R package [48, 148] under the MIT license. The gene sets used in the GSEA were compiled into an R package [172], also distributed under the MIT license. All data sets used in the study were obtained from public sources [149, 150, 153, 154]. A description of how these data sets were used is provided in the “Methods” section.

Declarations

Ethics approval and consent to participate

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Received: 3 March 2023 Accepted: 20 September 2023

Published online: 19 October 2023

References

1. Kharchenko PV. The triumphs and limitations of computational methods for scRNA-seq. *Nat Methods*. 2021;18(7):723–32.
2. Stegle O, Teichmann SA, Marioni JC. Computational and analytical challenges in single-cell transcriptomics. *Nat Rev Genet*. 2015;16(3):133–45.
3. Wagner A, Regev A, Yosef N. Revealing the vectors of cellular identity with single-cell genomics. *Nat Biotechnol*. 2016;34(11):1145–60.
4. Duò A, Robinson MD, Sonesson C. A systematic performance evaluation of clustering methods for single-cell RNA-seq data. *F1000Research*. 2018;7:1141.
5. Freytag S, Tian L, Lönnstedt I, Ng M, Bahlo M. Comparison of clustering tools in R for medium-sized 10x Genomics single-cell RNA-sequencing data. *F1000Research*. 2018;7:1297.
6. Kiselev VY, Andrews TS, Hemberg M. Challenges in unsupervised clustering of single-cell RNA-seq data. *Nat Rev Genet*. 2019;20(5):273–82.
7. Sun S, Zhu J, Ma Y, Zhou X. Accuracy, robustness and scalability of dimensionality reduction methods for single-cell RNA-seq analysis. *Genome Biol*. 2019;20:269.
8. Townes FW, Hicks SC, Aryee MJ, Irizarry RA. Feature selection and dimension reduction for single-cell RNA-Seq based on a multinomial model. *Genome Biol*. 2019;20:295.
9. Hotelling H. Analysis of a complex of statistical variables into principal components. *J Educ Psychol*. 1933;24(6):417–41.
10. Risso D, Perraudeau F, Gribkova S, Dudoit S, Vert JP. A general and flexible method for signal extraction from single-cell RNA-seq data. *Nat Commun*. 2018;9:284.
11. van der Maaten L, Hinton G. Visualizing data using t-SNE. *J Mach Learn Res*. 2008;9(86):2579–605.
12. McInnes L, Healy J, Saul N, Großberger L. UMAP: uniform manifold approximation and projection. *J Open Source Softw*. 2018;3(29):861.
13. Becht E, McInnes L, Healy J, Dutertre CA, Kwok IWH, Ng LG, et al. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat Biotechnol*. 2019;37:38–44.
14. Cooley SM, Hamilton T, Aragonés SD, Ray JCJ, Deeds EJ. A novel metric reveals previously unrecognized distortion in dimensionality reduction of scRNA-seq data. *bioRxiv*. 2022. <https://doi.org/10.1101/689851>.
15. Chari T, Pachter L. The specious art of single-cell genomics. *PLoS Comput Biol*. 2023;19(8):1011288.
16. Kobak D, Berens P. The art of using t-SNE for single-cell transcriptomics. *Nat Commun*. 2019;10(1):5416.
17. Kobak D, Linderman GC. Initialization is critical for preserving global data structure in both t-SNE and UMAP. *Nat Biotechnol*. 2021;39(2):156–7.
18. Wattenberg M, Viégas F, Johnson I. How to use t-SNE effectively. *Distill*. 2016. <https://doi.org/10.23915/distill.00002>.
19. Granja JM, Corces MR, Pierce SE, Bagdatli ST, Choudhry H, Chang HY, et al. ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nat Genet*. 2021;53(3):403–11.
20. Heiser CN, Lau KS. A quantitative framework for evaluating single-cell data structure preservation by dimensionality reduction techniques. *Cell Rep*. 2020;31(5):107576.
21. Linderman GC, Rachh M, Hoskins JG, Steinerberger S, Kluger Y. Fast interpolation-based t-SNE for improved visualization of single-cell RNA-seq data. *Nat Methods*. 2019;16(3):243–5.
22. Linderman GC, Steinerberger S. Clustering with t-SNE. *Provably SIAM J Math Data Sci*. 2019;1(2):313–32.
23. Lun ATL, 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.
24. McCarthy DJ, Campbell KR, Lun ATL, Wills QF. Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics*. 2017;33(8):1179–86.
25. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM, et al. Comprehensive integration of single-cell data. *Cell*. 2019;177(7):1888–902.
26. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat Biotechnol*. 2018;36(5):411–20.
27. Stuart T, Srivastava A, Madad S, Lareau CA, Satija R. Single-cell chromatin state analysis with Signac. *Nat Methods*. 2021;18(11):1333–41.
28. Fang R, Preissl S, Li Y, Hou X, Lucero J, Wang X, et al. Comprehensive analysis of single cell ATAC-seq data with SnapATAC. *Nat Commun*. 2021;12:1337.
29. Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. *Nat Commun*. 2017;8:14049.

30. Brunet JP, Tamayo P, Golub TR, Mesirov JP. Metagenes and molecular pattern discovery using matrix factorization. *Proc Natl Acad Sci*. 2004;101(12):4164–9.
31. Donoho D, Stodden V. When does non-negative matrix factorization give a correct decomposition into parts? In: *Proceedings of the 16th International Conference on Neural Information Processing Systems*. Cambridge: MIT Press; 2003. p. 1141–1148.
32. Durif G, Modolo L, Mold JE, Lambert-Lacroix S, Picard F. Probabilistic count matrix factorization for single cell expression data analysis. *Bioinformatics*. 2019;35(20):4011–9.
33. Gong W, Rasmussen TL, Singh BN, Koyano-Nakagawa N, Pan W, Garry DJ. Dpath software reveals hierarchical haemato-endothelial lineages of Etv2 progenitors based on single-cell transcriptome analysis. *Nat Commun*. 2017;8:14362.
34. Elyanow R, Dumitrascu B, Engelhardt BE, Raphael BJ. netNMF-sc: leveraging gene-gene interactions for imputation and dimensionality reduction in single-cell expression analysis. *Genome Res*. 2020;30(2):195–204.
35. Ho YJ, Anaparthi N, Molik D, Mathew G, Aicher T, Patel A, et al. Single-cell RNA-seq analysis identifies markers of resistance to targeted BRAF inhibitors in melanoma cell populations. *Genome Res*. 2018;28(9):1353–63.
36. Kotliar D, Veres A, Nagy MA, Tabrizi S, Hodis E, Melton DA, et al. Identifying gene expression programs of cell-type identity and cellular activity with single-cell RNA-seq. *eLife*. 2019;8:43803.
37. Lee DD, Seung HS. Learning the parts of objects by non-negative matrix factorization. *Nature*. 1999;401(6755):788–91.
38. Levitin HM, Yuan J, Cheng YL, Ruiz FJ, Bush EC, Bruce JN, et al. De novo gene signature identification from single-cell RNA-seq with hierarchical Poisson factorization. *Mol Syst Biol*. 2019;15:8557.
39. Welch JD, Kozareva V, Ferreira A, Vanderburg C, Martin C, Macosko EZ. Single-cell multi-omic integration compares and contrasts features of brain cell identity. *Cell*. 2019;177(7):1873–1887.e17.
40. Shao C, Höfer T. Robust classification of single-cell transcriptome data by nonnegative matrix factorization. *Bioinformatics*. 2016;33(2):235–42.
41. Sun S, Chen Y, Liu Y, Shang X. A fast and efficient count-based matrix factorization method for detecting cell types from single-cell RNAseq data. *BMC Syst Biol*. 2019;13:28.
42. Venkatasubramanian M, Chetal K, Schnell DJ, Atluri G, Salomonis N. Resolving single-cell heterogeneity from hundreds of thousands of cells through sequential hybrid clustering and NMF. *Bioinformatics*. 2020;36(12):3773–80.
43. Zhang S, Yang L, Yang J, Lin Z, Ng MK. Dimensionality reduction for single cell RNA sequencing data using constrained robust non-negative matrix factorization. *NAR Genomics Bioinforma*. 2020;2(3):lqaa064. <https://doi.org/10.1093/nargab/lqaa064>.
44. Pierson E, Yau C. ZIFA: Dimensionality reduction for zero-inflated single-cell gene expression analysis. *Genome Biol*. 2015;16:241.
45. Blei DM, Ng AY, Jordan MI. Latent Dirichlet allocation. *J Mach Learn Res*. 2003;3:993–1022.
46. DuVerle DA, Yotsukura S, Nomura S, Aburatani H, Tsuda K. Cell Tree: an R/bioconductor package to infer the hierarchical structure of cell populations from single-cell RNA-seq data. *BMC Bioinformatics*. 2016;17:363.
47. Dey KK, Hsiao CJ, Stephens M. Visualizing the structure of RNA-seq expression data using grade of membership models. *PLoS Genet*. 2017;13(3):1006599.
48. Carbonetto P, Sarkar A, Wang Z, Stephens M. Non-negative matrix factorization algorithms greatly improve topic model fits. 2021. [arXiv preprint arXiv:2105.13440](https://arxiv.org/abs/2105.13440).
49. González-Blas C, Minnoye L, Pappasokrati D, Aibar S, Hulselmans G, Christiaens V, et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data. *Nat Methods*. 2019;16(5):397–400.
50. Hofmann T. Probabilistic latent semantic indexing. In: *Proceedings of the 22nd Annual International ACM SIGIR Conference*. New York: Association for Computing Machinery; 1999. p. 50–57.
51. Bielecki P, Riesenfeld SJ, Hütter JC, Torlai Triglia E, Kowalczyk MS, Ricardo-Gonzalez RR, et al. Skin-resident innate lymphoid cells converge on a pathogenic effector state. *Nature*. 2021;592:128–32.
52. Housman G, Briscoe E, Gilad Y. Evolutionary insights into primate skeletal gene regulation using a comparative cell culture model. *PLoS Genet*. 2022;18(3):1010073.
53. Hung A, Housman G, Briscoe EA, Cuevas C, Gilad Y. Characterizing gene expression in an in vitro biomechanical strain model of joint health. *F1000Research*. 2022;11:296.
54. Rhodes K, Barr KA, Popp JM, Strober BJ, Battle A, Gilad Y. Human embryoid bodies as a novel system for genomic studies of functionally diverse cell types. *eLife*. 2022;11:71361.
55. Schenkel JM, Herbst RH, Canner D, Li A, Hillman M, Shanahan SL, et al. Conventional type I dendritic cells maintain a reservoir of proliferative tumor-antigen specific TCF-1+ CD8+ T cells in tumor-draining lymph nodes. *Immunity*. 2021;54(10):2338–2353.e6.
56. Xu H, Ding J, Porter CBM, Wallrapp A, Tabaka M, Ma S, et al. Transcriptional atlas of intestinal immune cells reveals that neuropeptide α -CGRP modulates group 2 innate lymphoid cell responses. *Immunity*. 2019;51(4):696–708.
57. Ding C, Li T, Peng W. On the equivalence between non-negative matrix factorization and probabilistic latent semantic indexing. *Comput Stat Data Anal*. 2008;52(8):3913–27.
58. Gaussier E, Goutte C. Relation between PLSA and NMF and implications. In: *Proceedings of the 28th Annual International ACM SIGIR Conference*. New York: Association for Computing Machinery; 2005. p. 601–602.
59. Gillis N. *Nonnegative matrix factorization*. Philadelphia: Society for Industrial and Applied Mathematics; 2021.
60. Kim J, Park H. *Sparse nonnegative matrix factorization for clustering*. Georgia Institute of Technology; 2008.
61. Sonesson C, Robinson MD. Bias, robustness and scalability in single-cell differential expression analysis. *Nat Methods*. 2018;15(4):255–61.
62. Wang T, Li B, Nelson CE, Nabavi S. Comparative analysis of differential gene expression analysis tools for single-cell RNA sequencing data. *BMC Bioinformatics*. 2019;20:40.
63. Erosheva E, Fienberg S, Lafferty J. Mixed-membership models of scientific publications. *Proc Natl Acad Sci*. 2004;101(Supplement 1):5220–7.
64. Abdelaal T, Michielsen L, Cats D, Hoogduin D, Mei H, Reinders MJT, et al. A comparison of automatic cell identification methods for single-cell RNA sequencing data. *Genome Biol*. 2019;20:194.

65. Diaz-Mejia JJ, Meng EC, Pico AR, MacParland SA, Ketela T, Pugh TJ, et al. Evaluation of methods to assign cell type labels to cell clusters from single-cell RNA-sequencing data. *F1000Research*. 2019;8:296.
66. Blei DM, Lafferty JD. Topic models. In: Srivastava AN, Sahami M, editors. *Text mining: classification, clustering, and applications*. Boca Raton: Chapman and Hall/CRC; 2009. p. 71–94.
67. Griffiths TL, Steyvers M. Finding scientific topics. *Proc Natl Acad Sci*. 2004;101(Supplement 1):5228–35.
68. Hofmann T. Unsupervised learning by probabilistic latent semantic analysis. *Mach Learn*. 2001;42(1):177–96.
69. Sarkar A, Stephens M. Separating measurement and expression models clarifies confusion in single-cell RNA sequencing analysis. *Nat Genet*. 2021;53(6):770–7.
70. Rosenberg NA. Genetic structure of human populations. *Science*. 2002;298(5602):2381–5.
71. Rosenberg NA. *distrupt*: a program for the graphical display of population structure. *Mol Ecol Notes*. 2004;4(1):137–8.
72. Haak W, Lazaridis I, Patterson N, Rohland N, Mallick S, Llamas B, et al. Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature*. 2015;522(7555):207–11.
73. Pereira BI, De Maeyer RPH, Covre LP, Nehar-Belaid D, Lanna A, Ward S, et al. Sestrins induce natural killer function in senescent-like CD8+ T cells. *Nat Immunol*. 2020;21(6):684–94.
74. Ashuach T, Reidenbach DA, Gayoso A, Yosef N. PeakVI: a deep generative model for single-cell chromatin accessibility analysis. *Cell Rep Methods*. 2022;2(3): 100182.
75. Marioni JC, Mason CE, Mane SM, Stephens M, Gilad Y. RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome Res*. 2008;18(9):1509–17.
76. Robinson MD, Smyth GK. Moderated statistical tests for assessing differences in tag abundance. *Bioinformatics*. 2007;23(21):2881–7.
77. Robinson MD, McCarthy DJ, Smyth GK. *edgeR*: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2009;26(1):139–40.
78. Anders S, Huber W. Differential expression analysis for sequence count data. *Genome Biol*. 2010;11:106.
79. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014;15(12):550.
80. Cable DM, Murray E, Shanmugam V, Zhang S, Zou LS, Diao M, et al. Cell type-specific inference of differential expression in spatial transcriptomics. *Nat Methods*. 2022;19(9):1076–87.
81. Becker-Herman S, Lantner F, Shachar I. Id2 negatively regulates B cell differentiation in the spleen. *J Immunol*. 2002;168(11):5507–13.
82. Stephens M. False discovery rates: a new deal. *Biostatistics*. 2016;18(2):275–94.
83. Finak G, McDavid A, Yajima M, Deng J, Gersuk V, Shalek AK, 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.
84. Zhu A, Ibrahim JG, Love MI. Heavy-tailed prior distributions for sequence count data: removing the noise and preserving large differences. *Bioinformatics*. 2019;35(12):2084–92.
85. Squair JW, Gautier M, Kathe C, Anderson MA, James ND, Hutson TH, et al. Confronting false discoveries in single-cell differential expression. *Nat Commun*. 2021;12:5692.
86. Montoro DT, Haber AL, Biton M, Vinarsky V, Lin B, Birket SE, et al. A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. *Nature*. 2018;560(7718):319–24.
87. Rosvall M, Bergstrom CT. Maps of random walks on complex networks reveal community structure. *Proc Natl Acad Sci*. 2008;105(4):1118–23.
88. Coifman RR, Lafon S, Lee AB, Maggioni M, Nadler B, Warner F, et al. Geometric diffusions as a tool for harmonic analysis and structure definition of data: diffusion maps. *Proc Natl Acad Sci*. 2005;102(21):7426–31.
89. Ruiz Garcia S, Deprez M, Lebrigand K, Cavard A, Paquet A, Arguel MJ, et al. Novel dynamics of human mucociliary differentiation revealed by single-cell RNA sequencing of nasal epithelial cultures. *Development*. 2019;146(20):dev177428. <https://doi.org/10.1242/dev.177428>.
90. Barkauskas CE, Chung MI, Fioret B, Gao X, Katsura H, Hogan BLM. Lung organoids: current uses and future promise. *Development*. 2017;144(6):986–97.
91. Rawlins EL, Okubo T, Xue Y, Brass DM, Auten RL, Hasegawa H, et al. The role of Scgb1a1+ clara cells in the long-term maintenance and repair of lung airway, but not alveolar. *Epithelium Cell Stem Cell*. 2009;4(6):525–34.
92. Spassky N, Meunier A. The development and functions of multiciliated epithelia. *Nat Rev Mol Cell Biol*. 2017;18(7):423–36.
93. Zhao H, Zhu L, Zhu Y, Cao J, Li S, Huang Q, et al. The cep63 paralogue *deup1* enables massive de novo centriole biogenesis for vertebrate multiciliogenesis. *Nat Cell Biol*. 2013;15(12):1434–44.
94. The Gene Ontology Consortium. The Gene Ontology resource: enriching a GOld mine. *Nucleic Acids Res*. 2020 12;49(D1):325–34.
95. Pliner HA, Packer JS, McFaline-Figueroa JL, Cusanovich DA, Daza RM, Aghamirzaie D, et al. Cicero predicts cis-regulatory DNA interactions from single-cell chromatin accessibility data. *Mol Cell*. 2018;71(5):858–8718.
96. Wu C, Tao Y, Li N, Fei J, Wang Y, Wu J, et al. Prediction of cellular targets in diabetic kidney diseases with single-cell transcriptomic analysis of db/db mouse kidneys. *J Cell Commun Signal*. 2023;17:169–88.
97. Cusanovich DA, Hill AJ, Aghamirzaie D, Daza RM, Pliner HA, Berletch JB, et al. A single-cell atlas of in vivo mammalian chromatin accessibility. *Cell*. 2018;174(5):1309–24.
98. Der E, Ranabothu S, Suryawanshi H, Akat KM, Clancy R, Morozov P, et al. Single cell RNA sequencing to dissect the molecular heterogeneity in lupus nephritis. *JCI Insight*. 2017;2(9):93009.
99. Grassmeyer J, Mukherjee M, DeRiso J, Hettinger C, Bailey M, Sinha S, et al. E1f5 is a principal cell lineage specific transcription factor in the kidney that contributes to *Aqp 2* and *Avpr 2* gene expression. *Dev Biol*. 2017;424(1):77–89.
100. Park J, Shrestha R, Qiu C, Kondo A, Huang S, Werth M, et al. Single-cell transcriptomics of the mouse kidney reveals potential cellular targets of kidney disease. *Science*. 2018;360(6390):758–63.

101. Ghezzi C, Loo DDF, Wright EM. Physiology of renal glucose handling via SGLT1, SGLT2 and GLUT2. *Diabetologia*. 2018;61(10):2087–97.
102. Thiagarajan RD, Georgas KM, Rumballe BA, Lesieur E, Chiu HS, Taylor D, et al. Identification of anchor genes during kidney development defines ontological relationships, molecular subcompartments and regulatory pathways. *PLoS ONE*. 2011;6(2):17286.
103. Gopal E, Umapathy NS, Martin PM, Ananth S, Gnana-Prakasam JP, Becker H, et al. Cloning and functional characterization of human SMCT2 (SLC5A12) and expression pattern of the transporter in kidney. *Biochim Biophys Acta Biomembr*. 2007;1768(11):2690–7.
104. Buenostro JD, Corces MR, Lareau CA, Wu B, Schep AN, Aryee MJ, et al. Integrated single-cell analysis maps the continuous regulatory landscape of human hematopoietic differentiation. *Cell*. 2018;173(6):1535–48.
105. Heinz S, Benner C, Spann N, Bertolino E, Lin YC, Laslo P, et al. Simple combinations of lineage-determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol Cell*. 2010;38(4):576–89.
106. Boeshaghi AS, Pachter L. Normalization of single-cell RNA-seq counts by $\log(x + 1)$ or $\log(1 + x)$. *Bioinformatics*. 2021;37(15):2223–4.
107. Lun A. Overcoming systematic errors caused by log-transformation of normalized single-cell RNA sequencing data. *bioRxiv*. 2018. <https://doi.org/10.1101/404962>.
108. Warton DI. Why you cannot transform your way out of trouble for small counts. *Biometrics*. 2018;74:362–8.
109. Amezquita RA, Lun ATL, Becht E, Carey VJ, Carpp LN, Geistlinger L, et al. Orchestrating single-cell analysis with Bioconductor. *Nat Methods*. 2020;17(2):137–45.
110. Barkas N, Petukhov V, Nikolaeva D, Lozinsky Y, Demharter S, Khodosevich K, et al. Joint analysis of heterogeneous single-cell RNA-seq dataset collections. *Nat Methods*. 2019;16(8):695–8.
111. Korsunsky I, Millard N, Fan J, Slowikowski K, Zhang F, Wei K, et al. Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat Methods*. 2019;16(12):1289–96.
112. Haghverdi L, Lun AT, Morgan MD, Marioni JC. Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nat Biotechnol*. 2018;36(5):421–7.
113. Parker HS, Leek JT, Favorov AV, Consideine M, Xia X, Chavan S, et al. Preserving biological heterogeneity with a permuted surrogate variable analysis for genomics batch correction. *Bioinformatics*. 2014;30(19):2757–63.
114. Tran HTN, Ang KS, Chevrier M, Zhang X, Lee NYS, Goh M, et al. A benchmark of batch-effect correction methods for single-cell RNA sequencing data. *Genome Biol*. 2020;21:12.
115. Richards LM, Riverin M, Mohanraj S, Ayyadhury S, Croucher DC, Díaz-Mejía JJ, et al. A comparison of data integration methods for single-cell RNA sequencing of cancer samples. *bioRxiv*. 2021. <https://doi.org/10.1101/2021.08.04.453579>.
116. Fan J, Slowikowski K, Zhang F. Single-cell transcriptomics in cancer: computational challenges and opportunities. *Exp Mol Med*. 2020;52(9):1452–65.
117. Gouvert O, Oberlin T, Févotte C. Negative binomial matrix factorization for recommender systems. 2018. *arXiv preprint arXiv:1801.01708*.
118. Gu J, Wang X, Halakivi-Clarke L, Clarke R, Xuan J. BADGE: a novel Bayesian model for accurate abundance quantification and differential analysis of RNA-Seq data. *BMC Bioinformatics*. 2014;15(S9):S6.
119. Wang C, Sun D, Huang X, Wan C, Li Z, Han Y, et al. Integrative analyses of single-cell transcriptome and regulome using MAESTRO. *Genome Biol*. 2020;21:198.
120. Lareau CA, Duarte FM, Chew JG, Kartha VK, Burkett ZD, Kohlway AS, et al. Droplet-based combinatorial indexing for massive-scale single-cell chromatin accessibility. *Nat Biotechnol*. 2019;37(8):916–24.
121. Kartha VK, Duarte FM, Hu Y, Ma S, Chew JG, Lareau CA, et al. Functional inference of gene regulation using single-cell multi-omics. *Cell Genomics*. 2022;2(9): 100166.
122. Bravo González-Blas C, De Winter S, Hulselmans G, Hecker N, Matetovici I, Christiaens V, et al. SCENIC+: single-cell multiomic inference of enhancers and gene regulatory networks. *Nat Methods*. 2023;20(9):1355–67.
123. Ma S, Zhang B, LaFave LM, Earl AS, Chiang Z, Hu Y, et al. Chromatin potential identified by shared single-cell profiling of RNA and chromatin. *Cell*. 2020;183(4):1103–1116.e20.
124. Cao J, Cusanovich DA, Ramani V, Aghamirzaie D, Pliner HA, Hill AJ, et al. Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science*. 2018;361(6409):1380–5.
125. Stuart T, Satija R. Integrative single-cell analysis. *Nat Rev Genet*. 2019;20(5):257–72.
126. Shiga M, Seno S, Onizuka M, Matsuda H. SC-JNMF: single-cell clustering integrating multiple quantification methods based on joint non-negative matrix factorization. *PeerJ*. 2021;9:12087.
127. Zhang S, Liu CC, Li W, Shen H, Laird PW, Zhou XJ. Discovery of multi-dimensional modules by integrative analysis of cancer genomic data. *Nucleic Acids Res*. 2012;40(19):9379–91.
128. Yang Z, Michailidis G. A non-negative matrix factorization method for detecting modules in heterogeneous omics multi-modal data. *Bioinformatics*. 2015;32(1):1–8.
129. Jin S, Zhang L, Nie Q. scAl: an unsupervised approach for the integrative analysis of parallel single-cell transcriptomic and epigenomic profiles. *Genome Biol*. 2020;21:25.
130. Hao Y, Stuart T, Kowalski MH, Choudhary S, Hoffman P, Hartman A, et al. Dictionary learning for integrative, multi-modal and scalable single-cell analysis. *Nat Biotechnol*. 2023. <https://doi.org/10.1038/s41587-023-01767-y>.
131. Argelaguet R, Velten B, Arnol D, Dietrich S, Zenz T, Marioni JC, et al. Multi-omics factor analysis—a framework for unsupervised integration of multi-omics data sets. *Mol Syst Biol*. 2018;14:8124.
132. Argelaguet R, Arnol D, Bredikhin D, Deloro Y, Velten B, Marioni JC, et al. MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data. *Genome Biol*. 2020;21:111.
133. Baker SM, Rogerson C, Hayes A, Sharrocks AD, Rattray M. Classifying cells with Scasat, a single-cell ATAC-seq analysis tool. *Nucleic Acids Res*. 2019;47(2):10.
134. Nair S, Kim DS, Perricone J, Kundaje A. Integrating regulatory DNA sequence and gene expression to predict genome-wide chromatin accessibility across cellular contexts. *Bioinformatics*. 2019;35:108–16.
135. Pott S, Lieb JD. Single-cell ATAC-seq: strength in numbers. *Genome Biol*. 2015;16:172.

136. Yan F, Powell DR, Curtis DJ, Wong NC. From reads to insight: a hitchhiker's guide to ATAC-seq data analysis. *Genome Biol.* 2020;21:22.
137. Taddy M. Distributed multinomial regression. *Ann Appl Stat.* 2015;9(3):1394–414.
138. Fisher RA. On the interpretation of χ^2 from contingency tables, and the calculation of P. *J R Stat Soc.* 1922;85(1):87–94.
139. Good IJ. Some statistical applications of Poisson's work. *Stat Sci.* 1986;1(2):157–70.
140. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian data analysis*. 3rd ed. Boca Raton: CRC Press; 2013.
141. Andrieu C, de Freitas N, Doucet A, Jordan MI. An introduction to MCMC for machine learning. *Mach Learn.* 2003;50:5–43.
142. Robert CP. *Monte Carlo statistical methods*. 2nd ed. New York: Springer; 2004.
143. Devroye L. *Non-uniform random variate generation*. New York: Springer-Verlag; 1986.
144. Stephens M, Carbonetto P, Gerard D, Lu M, Sun L, Willwerscheid J, et al. ashR: methods for adaptive shrinkage, using empirical Bayes. 2020. R package version 2.2-51. <https://github.com/stephens999/ashr>. Accessed 5 Mar 2023.
145. Chen MH, Shao QM. Monte Carlo estimation of Bayesian credible and HPD intervals. *J Comput Graph Stat.* 1999;8(1):69–92.
146. Box GEP, Tiao GC. *Bayesian inference in statistical analysis*. Reading: Addison-Wesley; 1992.
147. Gelman A, Hill J, Yajima M. Why we (usually) don't have to worry about multiple comparisons. *J Res Educ Eff.* 2012;5(2):189–211.
148. Carbonetto P, Luo K, Dey K, Hsiao J, Sarkar A, Hung A, et al. fastTopics: fast algorithms for fitting topic models and non-negative matrix factorizations to count data. 2022. R package version 0.6-142. <https://github.com/stephens999/fastTopics>.
149. Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, et al. Massively parallel digital transcriptional profiling of single cells. *Data sets. 10x Genomics.* 2017. <https://www.10xgenomics.com/support/single-cell-gene-expression>. Accessed 5 Mar 2023.
150. Montoro DT, Haber AL, Biton M, Vinarsky V, Chen S, Villoria J, et al. A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. *Data sets. Gene Expression Omnibus; 2018.* <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE103354>. Accessed 5 Mar 2023.
151. Cusanovich DA, Daza R, Adey A, Pliner HA, Christiansen L, Gunderson KL, et al. Multiplex single-cell profiling of chromatin accessibility by combinatorial cellular indexing. *Science.* 2015;348(6237):910–4.
152. Cusanovich DA, Reddington JP, Garfield DA, Daza RM, Aghamirzaie D, Marco-Ferreres R, et al. The cis-regulatory dynamics of embryonic development at single-cell resolution. *Nature.* 2018;555(7697):538–42.
153. Cusanovich DA, Hill AJ, Aghamirzaie D, Daza RM, Pliner HA, Berletch JB, et al. A single-cell atlas of *in vivo* mammalian chromatin accessibility. *Data sets. Mouse sci-ATAC-seq Atlas; 2018.* <https://shendurelab.github.io/mouse-atac/>. Accessed 5 Mar 2023.
154. Buenostro JD, Corces MR, Lareau CA, Wu B, Schep AN, Aryee MJ, et al. Single-cell epigenomics maps the continuous regulatory landscape of human hematopoietic differentiation. *Data sets. Gene Expression Omnibus; 2018.* <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE96772>. Accessed 5 Mar 2023.
155. Chen H, Lareau C, Andreani T, Vinyard ME, Garcia SP, Clement K, et al. Assessment of computational methods for the analysis of single-cell ATAC-seq data. *Genome Biol.* 2019;20:241.
156. Hien LTK, Gillis N. Algorithms for nonnegative matrix factorization with the Kullback-Leibler divergence. *J Sci Comput.* 2021;87(3):93.
157. Hsieh CJ, Dhillon IS. Fast coordinate descent methods with variable selection for non-negative matrix factorization. In: *Proceedings of the 17th ACM SIGKDD International Conference. New York: Association for Computing Machinery; 2011.* p. 1064–1072.
158. Lin X, Boutros PC. Optimization and expansion of non-negative matrix factorization. *BMC Bioinformatics.* 2020;21:7.
159. Ang AMS, Gillis N. Accelerating nonnegative matrix factorization algorithms using extrapolation. *Neural Comput.* 2019;31(2):417–39.
160. Ke ZT, Wang M. A new SVD approach to optimal topic estimation. 2019. arXiv preprint [arXiv:1704.07016](https://arxiv.org/abs/1704.07016).
161. van der Maaten L. Accelerating t-SNE using tree-based algorithms. *J Mach Learn Res.* 2014;15(93):3221–45.
162. Krijthe JH. Rtsne: t-distributed stochastic neighbor embedding using Barnes-Hut implementation. 2015. R package version 0.15. <https://github.com/krijthe/Rtsne>. Accessed 5 Mar 2023.
163. Ding C, He X. K-means clustering via principal component analysis. In: *21st International Conference on Machine Learning. New York: Association for Computing Machinery; 2004.*
164. Kass RE, Raftery AE. Bayes factors. *J Am Stat Assoc.* 1995;90(430):773–95.
165. Geer LY, Marchler-Bauer A, Geer RC, Han L, He J, He S, et al. The NCBI BioSystems database. *Nucleic Acids Res.* 2009;38(supplement-1):492–6.
166. Cerami EG, Gross BE, Demir E, Rodchenkov I, Babur Ö, Anwar N, et al. Pathway Commons, a web resource for biological pathway data. *Nucleic Acids Res.* 2010;39(supplement 1):685–90.
167. Rodchenkov I, Babur Ö, Luna A, Aksoy BA, Wong JV, Fong D, et al. Pathway Commons 2019 Update: integration, analysis and exploration of pathway data. *Nucleic Acids Res.* 2019;48(D1):489–97.
168. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci.* 2005;102(43):15545–50.
169. Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The molecular signatures database hallmark gene set collection. *Cell Syst.* 2015;1(6):417–25.
170. Liberzon A, Subramanian A, Pinchback R, Thorvaldsdóttir H, Tamayo P, Mesirov JP. Molecular signatures database (MSigDB) 3.0. *Bioinformatics.* 2011;27(12):1739–40.

171. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology: tool for the unification of biology. *Nat Genet.* 2000;25(1):25–9.
172. Carbonetto P, Stephens M. pathways: gene set enrichment analysis using human and mouse gene sets. 2021. R package version 0.1–20. <https://github.com/stephenslab/pathways>. Accessed 5 Mar 2023.
173. Fang T, Davydov I, Marbach D, Zhang JD. Gene-set enrichment with regularized regression bioRxiv. 2019. <https://doi.org/10.1101/659920>.
174. Sartor MA, Leikauf GD, Medvedovic M. LRpath: a logistic regression approach for identifying enriched biological groups in gene expression data. *Bioinformatics.* 2009;25(2):211–7.
175. Harrison PF, Pattison AD, Powell DR, Beilharz TH. Topconfects: a package for confident effect sizes in differential expression analysis provides a more biologically useful ranked gene list. *Genome Biol.* 2019;20:67.
176. McCarthy DJ, Smyth GK. Testing significance relative to a fold-change threshold is a TREAT. *Bioinformatics.* 2009;25(6):765–71.
177. Wang G, Sarkar A, Carbonetto P, Stephens M. A simple new approach to variable selection in regression, with application to genetic fine mapping. *J R Stat Soc Ser B.* 2020;82(5):1273–300.
178. Bauer S, Gagneur J, Robinson PN. GOing Bayesian: model-based gene set Analysis of Genome-Scale Data. *Nucleic Acids Res.* 2020;38(11):3523–32.
179. Ebrahimipoor M, Spitali P, Hettne K, Tsonaka R, Goeman J. Simultaneous enrichment analysis of all possible gene-sets: unifying self-contained and competitive methods. *Brief Bioinform.* 2019;21(4):1302–12.
180. Fontanillo C, Nogales-Cadenas R, Pascual-Montano A, Rivas JDL. Functional analysis beyond enrichment: non-redundant reciprocal linkage of genes and biological terms. *PLoS ONE.* 2011;6(9):24289.
181. Lu Y, Rosenfeld R, Simon I, Nau GJ, Bar-Joseph Z. A probabilistic generative model for GO enrichment analysis. *Nucleic Acids Res.* 2008;36(17):109.
182. Simillion C, Liechti R, Lischer HEL, Ioannidis V, Bruggmann R. Avoiding the pitfalls of gene set enrichment analysis with SetRank. *BMC Bioinformatics.* 2017;18:151.
183. Vivar JC, Pemu P, McPherson R, Ghosh S. Redundancy control in pathway databases (ReCiPa): an application for improving gene-set enrichment analysis in omics studies and “big data” biology. *Omics.* 2013;17(8):414–22.
184. McDavid A, Finak G, Yajima M. MAST: model-based analysis of single cell transcriptomics. 2021. R package version 1.20.0. <https://github.com/RGLab/MAST>. Accessed 5 Mar 2023.
185. R Core Team. R: a language and environment for statistical computing. Vienna; 2018. R Found Stat Comput. <https://www.R-project.org>. Accessed 5 Mar 2023.
186. Carbonetto P, Lao K, Sarkar A, Hung A, Tayeb K, Pott S, et al. Analysis of single-cell RNA-seq data sets for this manuscript. Zenodo. 2023. <https://doi.org/10.5281/zenodo.7962782>.
187. Carbonetto P, Lao K, Sarkar A, Hung A, Tayeb K, Pott S, et al. Analysis of single-cell ATAC-seq data sets for this manuscript. Zenodo. 2023. <https://doi.org/10.5281/zenodo.7962831>.
188. Blischak JD, Carbonetto P, Stephens M. Creating and sharing reproducible research code the workflow way. *F1000Research.* 2019;8:1749.

Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

