# Celda: a Bayesian model to perform co-clustering of genes into modules and cells into subpopulations using single-cell RNA-seq data

Zhe Wang [1,2,†], Shiyi Yang[2,†], Yusuke Koga[1,2], Sean E. Corbett[1,2], Conor V. Shea[1,2], W. Evan Johnson[1,2], Masanao Yajima[3] and Joshua D. Campbell[1,2,*]

[1]Bioinformatics Program, Boston University, Boston, MA, USA, [2]Division of Computational Biomedicine, Department of Medicine, Boston University School of Medicine, Boston, MA, USA and [3]Department of Mathematics and Statistics, Boston University, Boston, MA, USA

## ABSTRACT

**Single-cell RNA-seq (scRNA-seq) has emerged as a powerful technique to quantify gene expression in individual cells and to elucidate the molecular and cellular building blocks of complex tissues. We developed a novel Bayesian hierarchical model called Cellular Latent Dirichlet Allocation (Celda) to perform co-clustering of genes into transcriptional modules and cells into subpopulations. Celda can quantify the probabilistic contribution of each gene to each module, each module to each cell population and each cell population to each sample. In a peripheral blood mononuclear cell dataset, Celda identified a subpopulation of proliferating T cells and a plasma cell which were missed by two other common single-cell workflows. Celda also identified transcriptional modules that could be used to characterize unique and shared biological programs across cell types. Finally, Celda outperformed other approaches for clustering genes into modules on simulated data. Celda presents a novel method for characterizing transcriptional programs and cellular heterogeneity in scRNA-seq data.**

## INTRODUCTION

Complex biological systems can be conceptually defined into hierarchies where each level of the hierarchy is composed of different subunits which cooperate to perform distinct biological functions (1). For example, organisms can be subdivided into a collection of complex tissues: each complex tissue is composed of different cell types; each cell population is denoted by a unique combination of transcriptionally activated pathways (i.e. transcriptional modules); and each transcriptional module is composed of genes that are coordinately expressed to perform specific molecular functions. By identifying the 'building blocks' and their composition within each level of the hierarchy, we can more readily identify the patterns that define the behavior of these elements.

Single-cell RNA-seq (scRNA-seq) is a molecular assay that can quantify gene expression patterns in individual cells. In contrast to profiling of 'bulk' RNA from a sample, where only an average transcriptional signature across all the composite cells can be derived, scRNA-seq experiments can profile thousands of single-cell transcriptomes per sample and can thus offer an excellent opportunity to identify novel subpopulations of cells and to characterize transcriptional programs that define each subpopulation by examining co-varying patterns of gene expression across cells (2). However, analysis of scRNA-seq data presents several challenges. For example, the data tend to be sparse due to the difficulty in amplifying low amounts of RNA in individual cells. To combat noise from the amplification process, unique molecular identifiers (UMIs) are often incorporated to eliminate duplicate reads derived from the same mRNA molecule (3). The use of these UMIs enables the measurement of discrete counts of mRNA transcripts within each cell, making models constructed using discrete distributions a suitable approach for analyzing this type of data.

Discrete Bayesian hierarchical models have proven to be powerful tools for unsupervised modeling of discrete data types. In the text mining field, a plethora of models have been developed that can identify hidden topics across documents and/or cluster documents into distinct groups (4–8). These models generally treat each document as a 'bag-of-words' where each document is represented by a vector of counts or frequencies for each word in the vocabulary. Each document cluster (hidden topic) is represented by a Dirichlet distribution where words with higher probability are observed more frequently for the document cluster (5).

Given the success of topic models with sparse text data, and the discrete, sparse nature of transcriptional data generated by many scRNA-seq protocols, the application of such discrete Bayesian hierarchical models represents a promising approach to characterize structures in scRNA-seq data.

Various scRNA-seq tools have been developed to group cells into clusters, including ascend (9), BAMM-SC (10), CIDR (11), DESC (12), DIMM-SC (13), pcaReduce (14), SAFE-clustering (15), SAME-clustering (16), SC3 (17), scran (18), Seurat (19), SIMLR (20), TSCAN (21) and VPAC (22). Additionally, previous approaches for clustering genes have been developed for bulk RNA-seq and microarray data such as weighted gene co-expression network analysis (WGCNA) (23). However, methods that can cluster genes into modules based on co-expression patterns across cells using scRNA-seq data have not been reported. Other co-clustering methods such as QUBIC2 can identify blocks of co-expressed genes in a subset of samples (24). However, these methods are non-exhaustive and not strictly exclusive, meaning that not every gene or cell will get assigned into a block and individual genes or cells may be assigned into multiple blocks. When analyzing single-cell data, clustering of all genes into distinct, non-overlapping modules can be useful for characterizing the combinations of transcriptional programs that define unique, non-overlapping cell clusters.

Towards this end, we developed a model (Celda_CG) that performs exclusive and exhaustive co-clustering of cells into subpopulations and genes into transcriptional modules. In addition to clustering of genes and cells, Celda_CG also has the ability to describe the relationship between different layers of a biological hierarchy via probabilistic distributions. These distributions constitute dimensionally reduced representations of the data that can be used for downstream exploratory analysis. We demonstrate the utility of this approach by applying the Celda_CG model to a publicly available scRNA-seq dataset of peripheral blood mononuclear cells (PBMCs). Celda_CG identifies novel cell subpopulations missed by other approaches while characterizing transcriptional programs that are active to various degrees within and across major cell types.

## MATERIALS AND METHODS

### Celda_CG statistical model

The Celda_CG model uses sets of Dirichlet-multinomial distributions to model the hierarchies in the scRNA-seq data. The generative process for Celda_CG is outlined in Figure 1 and below, while the complete specification for the model can be found in the supplementary text. The generative process for Celda_CG is as follows :
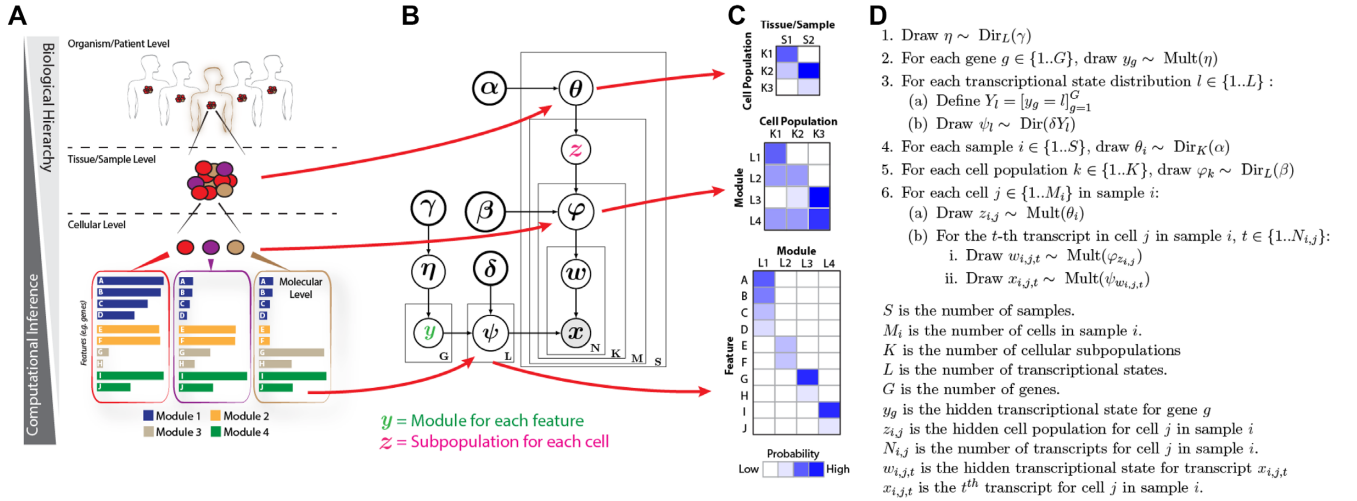
1. Draw $\eta \sim \mathrm{Dir}_L(\gamma)$
2. For each gene $g \in \{1, 2, \ldots, G\}$, draw $y_g \sim \mathrm{Mult}(\eta)$
3. For each transcriptional module distribution $l \in \{1, 2, \ldots, L\}$
   a. Define $Y_l = [y_g = l]_{g=1}^G$
   b. Draw $\psi_l \sim \mathrm{Dir}(\delta Y_l)$
4. For each sample $i \in \{1, 2, \ldots, S\}$, draw $\theta_i \sim \mathrm{Dir}_K(\alpha)$
5. For each cell population $k \in \{1, 2, \ldots, K\}$, draw $\varphi_k \sim \mathrm{Dir}_L(\beta)$

6. For each cell $j \in \{1, 2, \ldots, M_i\}$ in sample $i$
   a. Draw $z_{i,j} \sim \mathrm{Mult}(\theta_i)$
   b. For the $t$-th transcript in cell $j$ in sample $i$, $t \in \{1, 2, \ldots, N_{i,j}\}$
      i. Draw $w_{i,j,t} \sim \mathrm{Mult}(\varphi_{z_{i,j}})$
      ii. Draw $x_{i,j,t} \sim \mathrm{Mult}(\psi_{w_{i,j,t}})$

$\eta$ is from a Dirichlet distribution with symmetric concentration parameter $\gamma$ with length equal to the total number of transcriptional modules specified by $L$. $G$ is the number of genes. $y_g$ is the hidden transcriptional module label drawn from $\eta$ for gene $g$ and will return a value between 1 and $L$. $L$ is the number of transcriptional modules. '[]' refers to a Boolean operator and returns 1 when the expression within the bracket is true and 0 otherwise. We use this operator in step 3a to denote that the element corresponding to gene $g$ in $Y_i$ will be set to 1 if $y_g = 1$ and 0 otherwise. $Y_i$ will then be used as an indicator variable in step 3b to control the genes turned on in transcriptional module $l$. $\psi_i$ is from a Dirichlet distribution parameterized by $\delta Y_i$ where each element represents the probability of a gene in the module. If an element in $Y_i$ is zero, the parameter $\delta Y_i$ for the Dirichlet distribution will be zero along with the corresponding probability $\psi_i$ for that gene, thus turning off the expression of that gene in that module. The combination of these variables results in the 'hard-clustering' behavior by controlling the assignment of each gene to a single transcriptional module. $S$ is the number of samples. $\theta_i$ is from a Dirichlet distribution parameterized by the symmetric concentration parameter $\alpha$ that defines the probability of each cell population in each sample $i$. $K$ is the number of cellular subpopulations. Each cell population $k$ follows a Dirichlet distribution $\varphi_k$ parameterized by the symmetric concentration parameter $\beta$ where each element in $\varphi_k$ represents the probability of a transcriptional module in population $k$. $M_i$ is the number of cells in sample $i$. $z_{i,j}$ is the hidden cell population label for cell $j$ in sample $i$. $N_{i,j}$ is the number of transcripts for cell $j$ in sample $i$. $w_{i,j,t}$ is the hidden transcriptional module label for transcript $x_{i,j,t}$, and $x_{i,j,t}$ is the $t$-th transcript for cell $j$ in sample $i$. $z_{i,j}$ is drawn from $\theta_i$ and represents the hidden label denoting the population assignment for each cell. $w_{i,j,t}$ is the hidden label for transcript $t$ in cell $j$ drawn from $\varphi_{z_{i,j}}$ and represents the module assignment for that transcript. $x_{i,j,t}$ is the observed transcript which is drawn from $\psi_{w_{i,j,t}}$. Note that only the genes 'turned on' according to the indicators $Y_l$ will have a non-zero probability and will be selected from this draw.

The complete likelihood function of the Celda_CG model is then given as:

$$P(\eta, \psi, \theta, \varphi, Y, Z, W, X | \alpha, \beta, \gamma, \delta) = P(\eta|\gamma) \prod_{g=1}^{G} P(y_g|\eta)$$

$$\prod_{l=1}^{L} P(\psi_l|\delta, Y) \prod_{i=1}^{S} P(\theta_i|\alpha) \prod_{k=1}^{K} P(\varphi_k|\beta) \prod_{j=1}^{M_i} P(z_{i,j}|\theta_i)$$

$$\prod_{t=1}^{N_{i,j}} P(w_{i,j,t}|\varphi_{z_{i,j}}) P(x_{i,j,t}|\psi_{w_{i,j,t}}),$$

**Figure 1.** Celda identifies cell heterogeneity by clustering genes into modules and cells into subpopulations. (**A**) Example of a biological hierarchy. One way in which we try to understand complex biological systems is by organizing them into hierarchies. Individual organisms are composed of complex tissues. Each complex tissue is composed of different cellular populations with distinct functions; each cellular subpopulation contains a unique mixture of molecular pathways (i.e. modules); and each module is composed of groups of genes that are co-expressed across cells. (**B**) Plate diagram of the Celda_CG model. We developed a novel discrete Bayesian hierarchical model called Celda_CG to characterize the molecular and cellular hierarchies in biological systems. Celda_CG performs 'co-clustering' by assigning each gene to a module and each cell to a subpopulation. (**C**) In addition to clustering, Celda_CG also inherently performs a form of 'matrix factorization' by deriving three distinct probability matrices: (i) a cell population × sample matrix representing the probability that each population is present in each sample; (ii) a transcriptional module × cell population matrix representing the contribution of each transcriptional state to each cellular subpopulation; and (iii) a gene × module matrix representing the contribution of each gene to its module. (**D**) Generative process for the Celda_CG model.

where $\alpha$, $\beta$, $\gamma$, $\delta$ are the symmetric prior parameters in their corresponding Dirichlet distributions, and $Y$, $Z$, $W$, $X$ are the collections of $y_g$, $z_{i,j}$, $w_{i,j,t}$ and $x_{i,j,t}$, respectively.

### Estimation of model parameters

We use a heuristic hard Expectation Maximization (EM) procedure to estimate the cell population label $z_{i,j}$ for cell $j$ in sample $i$ and a collapsed Gibbs sampling procedure to estimate the hidden transcriptional module label $y_g$ for gene $g$ (Supplementary text). To estimate the hidden transcriptional module label for each gene, we integrate out $\psi$, $\varphi$ and $W$, and drop components related to $\theta$ that are invariant with respect to $Y$. The final formula after simplification is as follows:

$$P\left( y_g = l | Y_{-(g)}, Z, X, \alpha, \beta, \delta, \gamma \right) \propto \frac{\prod_{l=1}^{L} \Gamma\left(|V_l| + \gamma\right)}{\Gamma\left(\sum_{l=1}^{L}\left(|V_l| + \gamma\right)\right)} \times$$

$$\prod_{k=1}^{K}\left[\frac{\Gamma\left(n_{(\cdot),(k),(V_l)} + \beta\right)}{\Gamma\left(n_{(\cdot),(k),\left(V_l^{-(g)}\right)} + \beta\right)}\right] \times \left[\prod_{l=1}^{L}\frac{\Gamma\left(|V_l|\delta\right)}{\Gamma(\delta)^{|V_l|}}\right] \times$$

$$\left[\frac{\Gamma\left(\sum_{v \in V_l^{-(g)}}\left(n_{(\cdot),(\cdot),v} + \delta\right)\right)}{\Gamma\left(\sum_{v \in V_l}\left(n_{(\cdot),(\cdot),v} + \delta\right)\right)}\right],$$

where $L$ is the total number of modules, $K$ is the total number of cell populations, $|V_l|$ is the total number of genes in module $l$, $V_l^{-(g)}$ is the total number of genes in module $l$ leaving out the current gene $g$, $n_{(\cdot),(k),\left(V_l^{-(g)}\right)}$ is the total number of transcripts from genes in module $l$ across all the cells in

cluster $k$ leaving out those from gene $g$, $n_{(\cdot),(\cdot),v}$ is the total number of transcripts for gene $g$ across all cells and samples, $n_{(\cdot),(k),(V_l)}$ is the number of transcripts from all genes in module $l$ in population $k$, and $\Gamma$ is the gamma function. For estimating the hidden population label for each cell $z_{i,j}$, we relied on a heuristic 'hard' EM procedure to increase speed on large datasets with many cells. The collapsed Gibbs sampling equations for $z_{i,j}$ can also be found in the Supplementary text. First, we drop components related to $\psi$ that are invariant with respect to $Z$. The 'hard' EM obtains a point estimate of $z_{i,j}$ by maximizing the posterior with respect to point estimates of $\theta$ and $\varphi$ given the current configurations of $Z$ and $Y$:

$$\hat{z}_{i,j} = \text{argmax}_k \left\{ P\left( z_{i,j} = k | X_{i,j}, \hat{\theta}, \hat{\varphi} \right) \right\}$$

$$= \text{argmax}_k \left\{ \hat{\theta}_{i,k} \prod_{l=1}^{L} \hat{\varphi}_{k,l}^{n_{i,j,(V_l)}} \right\},$$

where $X_{i,j}$ is the collection of transcripts $x_{i,j,t}$ within cell $j$ in sample $i$, $n_{i,j,(V_l)}$ is the number of transcripts from all the genes that belong to module $l$ of cell $j$ in sample $i$. $\hat{\theta}_{i,k}$ is the point estimate of $\theta_{i,k}$ which is the probability of a cell belonging to cell population $k$ in sample $i$ and can be calculated as:

$$\hat{\theta}_{i,k} = \frac{m_{i,k} + \alpha}{M_i + K\alpha},$$

where $m_{i,k}$ is the total number of cells assigned to cluster $k$ in sample $i$ and $M_i$ is the total number of cells in sample $i$. $\hat{\varphi}_{k,l}$ is the point estimate of $\varphi_{k,l}$ which is the probability of

module $l$ in cell population $k$ and can be calculated as:

$$\hat{\varphi}_{k,l} = \frac{n_{(\cdot),(k),(V_l)} + \beta}{n_{(\cdot),(k),(\cdot)} + L\beta}.$$

Within each iteration of the optimization procedure, we apply the 'hard' EM procedure to estimate the cell population labels ($Z$) given a fixed set of transcriptional module labels ($Y$) and then apply the collapsed Gibbs sampling procedure to estimate the transcriptional module labels ($Y$) given a fixed set of cell population labels ($Z$). We generally run the model for a maximum number of iterations (200 by default) or until there has been no improvement in the log-likelihood for a pre-defined number of iterations (10 by default). The configurations of $Z$ and $Y$ that produced the highest likelihood are returned as the final solution.

In order to avoid a local optimum, we apply a heuristic cluster/module splitting procedure every 10 iterations. To apply the cell splitting procedure at a given iteration, we try to find a better configuration for $Z$ with a higher log-likelihood by splitting one population into two new clusters and removing another unsplit population. Let $K^*$ be the set of cell population clusters that have $>3$ cells and $|K^*|$ be the cardinality of set $K^*$. For one cell population $k^* \in K^*$, the cluster is split into two new clusters using the Celda_C model setting $K_c = 2$. Then parallelly for all other cell clusters $\{k' : k' \in \{1, 2, \ldots, K\} \wedge k' \neq k^*\}$, we redistribute all the cells in cluster $k'$ to their second most likely cluster according to EM probabilities of current $Z$ configuration. The log-likelihood is re-calculated for each of these $K - 1$ configurations. After repeating this procedure for all the $\{k^* : k^* \in K^*\}$, a total number of $|K^*| \times (K - 1)$ new possible configurations for $Z$ are obtained. The configuration that produced the highest likelihood will be set at the current solution. If none of the new configurations had a higher likelihood than the original configurations, then no splitting will be performed and the original configuration of $Z$ will be maintained. The module splitting procedure is similarly applied to the transcriptional modules to find a $Y$ that has a higher log-likelihood. One module $l^*$ is split using Celda_G with $L_G = 2$ and new likelihoods are calculated by redistributing the genes in each of the other modules. One potential limitation is that running Celda_G on all cells to split each module would result in a dramatic reduction in speed for large datasets. We therefore take each cell population cluster and split it up into 10 new clusters using Celda_C with $K = 10$ to produce a temporary configuration denoted $Z^*$. These temporary populations are used to potentially find a better configuration of module labels. Splitting each cell population into 10 temporary cell populations ensures that better splits of the modules can be obtained even if the current cluster labels $Z$ are suboptimal. Even though the modules are split with $Z^*$, the overall likelihood for all new splits of $Y$ is still calculated with the current configuration of $Z$ containing the $K$ subpopulations. As in the cell splitting approach, the module split with the best log-likelihood is chosen if it is higher than that from the current $Y$ configuration.

### Determining the number of cell populations and transcriptional modules

Perplexity has been commonly used in the topic models to measure how well a probabilistic model predicts observed samples (5). Here, we use perplexity to calculate the probability of observing expression counts given an estimated Celda_CG model. Rather than performing cross-validation which is computationally expensive, a series of test sets are created by sampling the counts from each cell according to a multinomial distribution defined by dividing the counts for each gene in the cell by the total number of counts for that cell. Perplexity is then calculated on each test set, with a lower perplexity indicating a better model fit (5). For a test set $x$, the perplexity of Celda_CG is given as

$$\text{Perplexity} = \exp\left\{ -\frac{\log\left(P\left(x\right)\right)}{\sum_{i=1}^{S} \sum_{j=1}^{M_i} N_{i,j}} \right\},$$

where $\log(P(x))$ is defined as:

$$\log\left(P(x)\right) = \sum_{i=1}^{S} \sum_{j=1}^{M_i} \log\left[ \sum_{k=1}^{K} \theta_{i,k} \prod_{g=1}^{G} \left(\sum_{l=1}^{L} \phi_{k,l} \psi_{l,g}\right)^{n_{i,j,g}} \right].$$

We compare perplexity values among different model settings and use rate of perplexity change (RPC) (25) to determine an appropriate number of cell populations and transcriptional modules. In particular, setting a fixed number of transcriptional modules, a series of Celda_CG models with a sequence of equally spaced $K$s arranged in ascending order are fitted. We then calculate the RPC along the course of increase of cell populations, and choose the smallest $K$ as the appropriate number of cell populations where the RPC is zero at a given precision. Similarly, setting a fixed number of cell populations, an appropriate number of transcriptional modules can be selected by calculating the RPC along a sequence of equally spaced $L$s.

### Data collection and pre-processing

PBMC_4k, 33k and 68k datasets were downloaded using R/Bioconductor package TENxPBMCData v1.8.0. They contain 4340 cells and 33 694 genes, 33 148 cells and 32 738 genes, and 68 597 cells and 32 738 genes, respectively. We applied DecontX (26) to remove inadvertent contamination using default settings. For the PBMC 4k dataset, 17 039 genes detected in fewer than three cells were excluded. We applied NormalizeData and FindVariableFeatures functions from Seurat v3.2.2 (19) using default settings and identified a set of the 2000 most variable genes for clustering by variance-stabilizing transformation (VST) (19). Principal component analysis (PCA) was performed on scaled normalized gene expression using the RunPCA function from Seurat v3.2.2 in default settings. For coloring of uniform manifold approximation and projections (UMAPs) and module heatmaps, the decontaminated counts were normalized by library size, square root-transformed, centered and scaled to unit variance. Values greater than 2 or less than –2 were trimmed.

## Selecting the number of transcriptional modules ($L$) and cell clusters ($K$)

We applied two stepwise splitting procedures as implemented in the recursiveSplitModule and recursiveSplitCell functions in Celda to determine the optimal $L$ and $K$. recursiveSplitModule uses the Celda_G model to cluster genes into modules for a range of possible $L$ values between 10 and 200. The module labels of the previous model with $L-1$ modules are used as the initial values in the current model with $L$ modules. The best split of an existing module, evaluated by best overall likelihood, is found to create the $L$th module. The RPC was calculated for each successive model generation. For the PBMC 4k dataset, we found that the model with 80 transcriptional modules had low RPC and included both known and novel gene programs. recursiveSplitCell uses the Celda_CG model to cluster cells into cell clusters for a range of possible $K$ values between 3 and 30. The module labels of genes from model $L=80$ was used to initialize the modules in recursiveSplitCell. We found that the model with 20 cell clusters had low RPC and included both known and novel cell populations (Supplementary Figure S1). The final Celda_CG model used in this analysis of the PBMC 4k dataset was extracted from the stepwise splitting results using the subsetCeldaList function.

## UMAP of PBMC cells based on Celda transcriptional modules

Dimensionality reduction for visualization by UMAP (27) is performed using the square root-transformed module probability (MP) matrix which contains the probability of each transcriptional module in each cell. Specifically, the MP matrix is defined as

$$\text{MP}_{i,j,l} = \frac{n_{i,j,(V_l)}}{N_{i,j}} \ ,$$

where $n_{i,j,(V_l)}$ denotes the sum of all counts belonging to genes in transcriptional module $l$, and $N_{i,j}$ is the total sum of counts for a cell. The square root transformation is applied as it can be applied to zero counts without the need to add a pseudocount as is required with log transformation. The umap function from the uwot R package was applied to the MP matrix using Euclidean distance to obtain two-dimensional coordinates for each cell with n_neighbors = 10, min_dist = 0.5 and default settings.

## Testing for differential expression

A hurdle model from MAST (28) was used for significance testing of differential expression between cell clusters. Benjamini and Hochberg false discovery rate (29) (FDR) adjusted $P$-values were used to reject the null hypotheses.

## Cell clustering and UMAP using Seurat

The same set of 2000 most variable genes in the decontaminated PBMC dataset were used for clustering by Seurat. PCA was performed on scaled normalized gene expressions using the RunPCA function from Seurat (19) v3.2.2 in default settings. The Shared Nearest Neighbor

(SNN) graph was constructed using the FindNeighbors function with the top 22 principal components (PCs). Clusters were identified by modularity optimization using FindClusters with default settings. UMAP was generated using the RunUMAP function and the top 22 PCs with n.neighbors = 10, min.dist = 0.5 and default settings.

## Cell clustering and UMAP using scran

The same set of 2000 most variable genes in the decontaminated PBMC dataset were used for clustering by scran. PCA was performed on normalized gene expressions using the runPCA function from scater (30) v1.18.3 using 2000 genes and default settings. The SNN graph was constructed using the buildSNNGraph function from scran (18) v1.18.3 with the top 28 PCs. Clusters were identified by random walks using the cluster_walktrap function from igraph (31) package v1.2.6 with default settings. UMAP was generated using the runUMAP function with n_neighbors = 5 and default settings.

## Biclustering of PBMCs using QUBIC2

The same set of 2000 most variable genes in the decontaminated PBMC dataset were used for biclustering by QUBIC2. The left-truncated mixture of Gaussian distribution option for data discretization was performed on counts per million (CPM) values with -F and -R flags. The 1.0 objective function and dual expansion biclustering option was performed on the discretized result with -d, -C, -N flags, and the number of biclusters set as 80 (-o 80).

## Clustering of mouse lung cells with Celda

The ExperimentHub package was used to download mouse single-cell data from The Tabula Muris Consortium (32) (ExperimentHub ID: EH1617). Cells were filtered to retain those from lung tissue that did not belong to a subtissue ($n = 2150$). Features with at least three counts in three cells were included in the Celda_CG analysis. The recursive splitting procedure was used to identify 125 modules ($L$) and 35 cell populations ($K$). Cells were displayed on a UMAP, and broad cell types were identified using marker genes.

## Pathway enrichment analysis of gene modules using Enrichr

The enrichment of pathways from the 'GO_Biological_Process_2021' database was determined for each of the 80 gene modules identified by Celda_CG from the PBMC 4k dataset using the enrichR (33) R package (v3.0). The Gene Ontology (GO) pathways with FDR $<0.05$ were considered significantly enriched. To assess the significance of the number of pathways enriched in Celda modules, a null distribution was created by performing 100 random permutations of module cluster labels and testing for pathway enrichment of the permuted labels with enrichR.

## Evaluation of module clustering accuracy with simulated data

Data were simulated based on the generative model of Celda_CG (see the Materials and Methods, and Supplementary text) using the simulateCells function in Celda.

Specifically, we set the model to 'celda_CG', S to 1, CRange between 4000 and 6000, NRange between 1000 and 10 000, G to 33 000, and K to 20. scRNA-seq count data were simulated using six combinations of concentration parameters $\beta$ and $\delta$ ranging from 1 to 40 representing six levels of clustering difficulties. For each combination of $\beta$ and $\delta$, a range of the number of transcriptional modules ($L$) from 10 to 200 was simulated with 10 replicates per $L$. After simulated data were generated, the 2000 most variable genes determined by VST were selected. Celda_CG clustering, k-means and PCA were applied to group the 2000 most variable genes to transcriptional modules whose numbers equal the true remaining number of modules for the 2000 genes. Adjusted Rand indices (ARIs) were calculated between the gene clustering results of Celda_CG, k-means or PCA and the true module labels of genes using the adjustedRandIndex function from mclust package v5.4.6. k-means was performed using the stats package v4.0.5 and default settings.

A heuristic approach was used to cluster genes based on the loadings from the PCA. After performing PCA to reduce the data to PCs whose number equals the number of true modules for the 2000 most variable genes, we first order the genes by the loadings in increasing order for each PC. For each PC, if the sum of the absolute values of the top 50 negative loadings is greater than the sum of the absolute values of the bottom 50 positive loadings, we rank the genes by loadings in this PC in increasing order. Otherwise, we rank the genes by loadings in this PC in decreasing order. This is to account for the bidirectionality of PCA loadings so that when genes are assigned to a PC, they are always in the same direction with respect to the orientation of the PC. After ranking the genes by loadings for each PC, we assign each gene to its highest ranking PC accordingly. If a gene has the same highest ranks in two or more PCs, this gene is not used for the calculation of ARI.

UMAPs of genes were generated to visualize the variability of genes and the clustering difficulties for the simulated data. For each combination of $\beta$ and $\delta$, UMAP was generated based on one of the 10 simulations at $L = 100$. Gene counts for each cell were collapsed to cell clusters before applying UMAP. Specifically, for each gene, the counts for each of the 20 true cell clusters were added and divided by total counts for this gene, so the number of features for UMAP were reduced from the total number of cells to 20. These cell cluster probabilities were then square root-transformed before applying UMAP. UMAP dimension reduction coordinates for cells were generated using the umap function from the uwot R package with n_neighbors = 10, min_dist = 0.5 and default settings.

## RESULTS

We developed a novel discrete Bayesian hierarchical model, called Cellular Latent Dirichlet Allocation (Celda), to perform exclusive and exhaustive co-clustering of genes into modules and cells into subpopulations (Figure 1; Supplementary text). Each level in the biological hierarchy is modelled as a mixture of components using Dirichlet distributions: sample $i$ is a mixture of cellular subpopulations ($\theta_i$), each cell subpopulation $k$ is a mixture ($\varphi_k$) of transcriptional modules, and each module l is a mixture ($\psi_l$) of fea-

tures such as genes. $\theta_{i,k}$ is the probability of cell population $k$ in sample $i$, $\varphi_{k,l}$ is the probability of module $l$ in population $k$ and $\psi_{l,g}$ is the probability of gene $g$ in module $l$ (Figure 1A, B). Each cell $j$ in sample $i$ has a hidden cluster label, $z_{i,j}$, denoting the population to which it belongs. Each transcript $x_{i,j,t}$ has a hidden label $w_{i,j,t}$ denoting the transcriptional module to which it belongs. A similarly structured topic model has previously been proposed called 'Latent Dirichlet Co-Clustering' (8). However, we add a unique and novel component to our model specifically geared towards gene expression analysis.
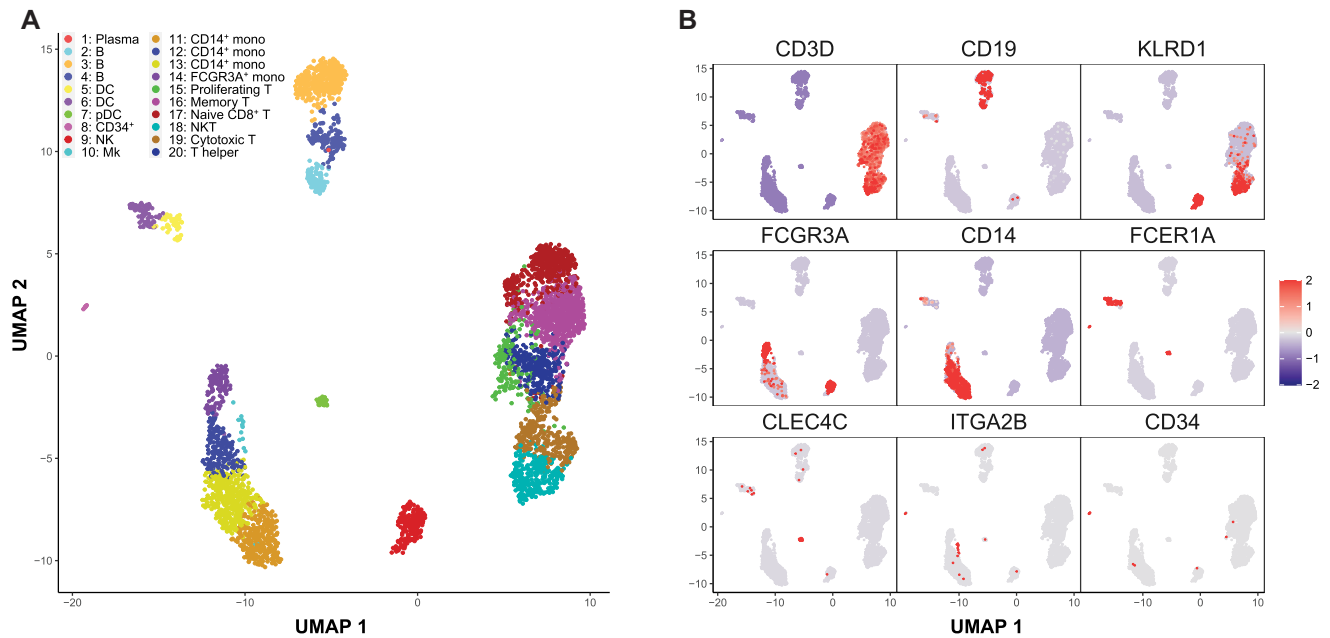
The goal of many gene expression clustering algorithms is to group genes into distinct, non-overlapping sets of genes (23,34–36) (i.e. hard-clustering of genes). The rationale for this type of clustering is that genes that co-vary across cells and samples are likely to be involved in the same biological processes and should be considered a single biological program (37). In order to enforce 'hard-clustering' of genes into modules, we modified an approach from Wang and Blei (6) regarding the sparse Topic Model (sparseTM) that has the capability to turn words 'on' or 'off' in different topics, by assigning a non-zero or zero probability to that word in each topic. In Celda_CG, we leverage this technique to turn off genes in all modules except one to enable the hard-clustering behavior.

While Celda can perform clustering, it also offers probabilistic distributions which describe the contribution of each 'building block' to each layer of the biological hierarchy (Figure 1C). These distributions can also be viewed as reduced dimensional representations of the data that can be used for downstream exploratory analyses. For example, the $\varphi$ matrix contains the probability of each module in each cell population and thus provides a high-level view of the structure of the dataset. The generative process for Celda_CG is shown in Figure 1D.

### Identification of cell populations in PBMCs

To assess Celda_CG's ability to identify biologically meaningful cell subpopulations in real-world scRNA-seq data, we applied it to a publicly available dataset provided by 10X Genomics. The dataset (PBMC 4k) contains 4340 PBMCs collected from a healthy donor. To determine the optimal number of transcriptional modules ($L$) and cell populations ($K$), we employed a step-wise splitting procedure first for the number of modules using a temporary cell-clustering solution and then for the number of cell populations using a fixed number of modules (see the Materials and methods and Supplementary Figure S1). The RPC (25) was measured at each split. An RPC closer to zero indicates that the addition of new modules or cell clusters is not substantially decreasing the perplexity. By observing the 'elbows' on these curves as a reference point in combination with manual review of module heatmaps and cell clusters, a solution of $L = 80$ transcriptional modules and $K = 20$ cell populations was chosen for further characterization.

A UMAP (27) dimension reduction representation was generated based on the estimated module probabilities for each cell, and the major subtypes of immune cells were identified by examining expression of known marker genes (Figure 2; Supplementary Table S1). Among the

**Figure 2.** Celda identifies immune cell subpopulations from PBMC scRNA-seq data. To demonstrate the utility of the Celda clustering model, we applied it to an scRNA-seq dataset of 4340 PBMCs generated using 10X Chromium platform and identified 80 transcriptional modules and 20 cell populations. (**A**) UMAP dimension reduction representation of 4340 PBMCs based on the transcriptional module probabilities. (**B**) Scaled normalized expression of representative gene markers shows clustering of cell subpopulations including T cells (CD3D), B cells (CD19), NKs (KLRD1), FCGR3A$^+$ monocytes (FCGR3A), CD14$^+$ monocytes (CD14), DCs (FCER1A), pDCs (CLEC4C), megakaryocytes (ITGA2B) and CD34$^+$ progenitor cells (CD34). Cell populations 1 (plasma cell) and 15 (proliferating T cells) are novel cell clusters identified by Celda, demonstrating Celda's ability to characterize additional cellular heterogeneity.

20 identified cell clusters in the PBMC sample, we identified major immune cell populations including CD19$^+$ B cells, FCER1A$^+$ dendritic cells (DCs), CLEC4C$^+$ plasmacytoid dendritic cells (pDCs), CD34$^+$ progenitor cells, KLRD1$^+$ natural killer cells (NKs), ITGA2B$^+$ megakaryocytes, CD14$^+$ monocytes, FCGR3A$^+$ monocytes, CCR7$^+$ memory T cells, CD8A$^+$CD8B$^+$ cytotoxic T cells and CD4$^+$ T helper cells. Cell subpopulations 15–20 show a consistently higher expression (FDRs <0.01) of the T-cell marker genes CD3D, CD3E and CD3G relative to all other clusters. Among these T-cell subpopulations, clusters 17, 18 and 19 show consistently higher expression (FDRs <0.01) of CD8A and CD8B (Supplementary Figure S2). Within these CD8A$^+$CD8B$^+$ T cells, cluster 17 has high expression (FDR <0.01) of naive T-cell marker CCR7, whereas cluster 18 has consistent high expression (FDRs <0.01) of NK markers GNLY, KLRG1 and granzyme genes GZMA and GZMH, so we classified them as naive CD8$^+$ T cells and NKT cells, respectively (38). Cell subpopulation 15 expressed T-cell markers as well as uniquely high levels of module 61, which contained genes associated with proliferation including MKI67, IL2RA, CENPM and CENPF (Supplementary Figure S2) commonly found in activated proliferating T cells (39,40).
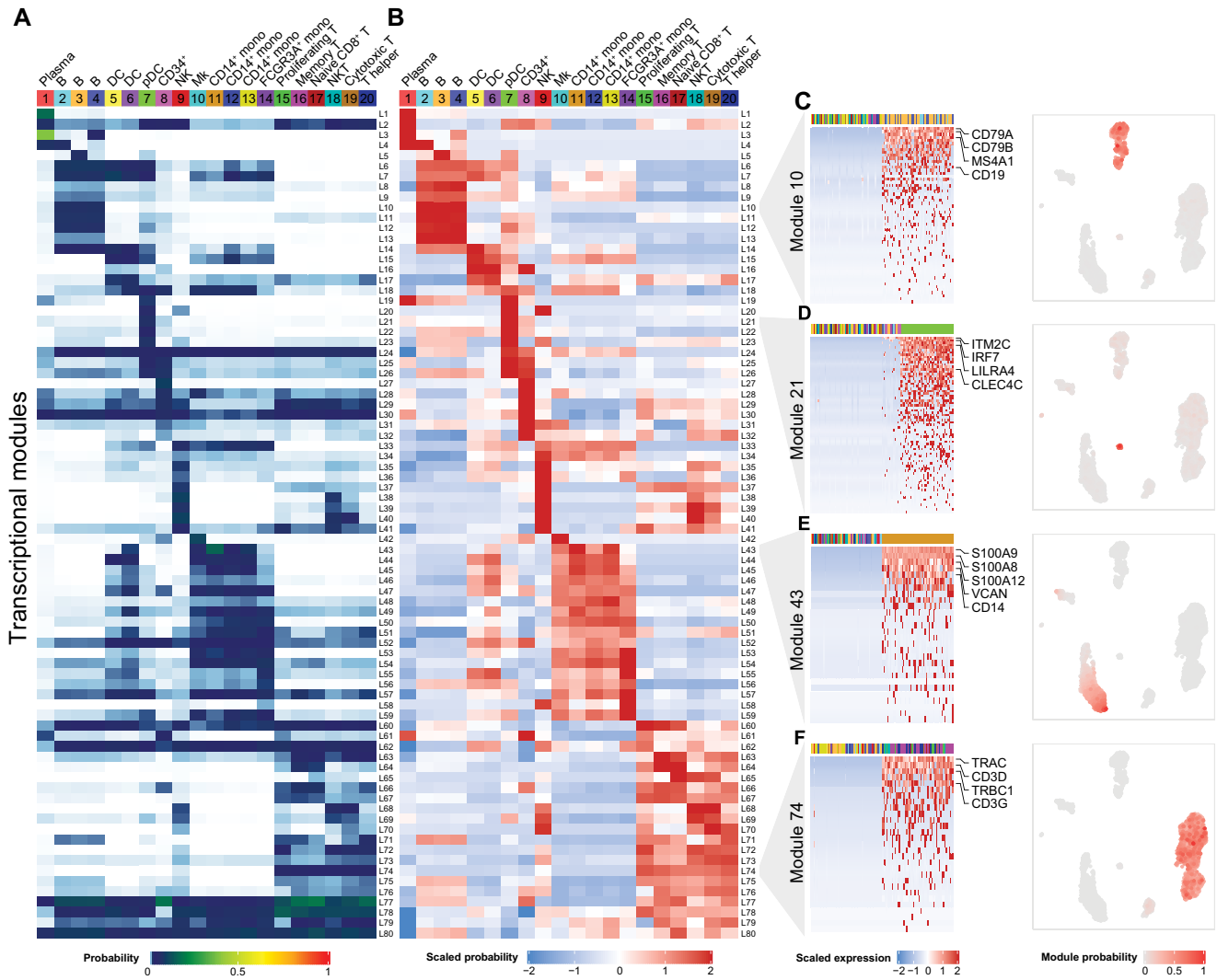
Cell subpopulation 1 contained a single cell which had the highest number of UMIs across the dataset (*n* = 48 443). This cell expressed several B-lineage markers such as CD79A, CD79B and CD19 but also contained a relatively high fraction (27%) of UMIs for immunoglobulin heavy chain and light chain genes IGHG1, IGHG3, IGLC2 and IGLC3. These genes were not observed in other cells and

suggest a plasma cell lineage (41) (Supplementary Figure S3). We also observed similar plasma cell subpopulations in PBMC 33k and 68k datasets (Supplementary Figures S4 and S5). The proportions of plasma cells were 0.02, 0.28 and 0.18% for PBMC 4k, 33k and 68k datasets respectively. Cell populations 1 and 15 were not identified by the analysis workflows and graph-based clustering methods used in Seurat (19) and scran (18) packages (Supplementary Figure S6), demonstrating Celda's ability to characterize additional cellular heterogeneity compared with other popular single-cell analysis workflows. We also applied QUBIC2 (42) to identify cell subpopulations through non-exclusive and non-exhaustive biclustering. All of the cells in each of the 80 biclusters identified by QUBIC2 were from the dendritic and/or monocyte subpopulations, and no biclusters contained B cells, T cells or NKs, pDCs or the novel subpopulations of proliferating T cells (Supplementary Figure S7).

Finally, we applied Celda to a dataset of mouse lung cells from The Tabula Muris Consortium (32) to demonstrate that it can find biologically relevant cell clusters and gene modules in another tissue type (Supplementary Figure S8).

**Identification of transcriptional modules with unique patterns of expression across cell populations**

Beyond assessment of individual marker genes, Celda has the ability to identify modules of co-expressed genes which can be further examined to characterize transcriptional programs active in one or more cell populations (Figure 3). An overview of the relationships between modules and

**Figure 3.** Celda produces a high-level overview of the relationships between transcriptional modules and cell populations. (**A**) The $\varphi$ matrix shows the probability of each of the 80 transcriptional modules (rows) in each of the 20 cellular subpopulations (column) and can be used to explore the relationship between modules within a cell population. (**B**) The row-scaled $\varphi$ matrix can be used to explore the relative probability of each module across cell populations. (**C–F**) Module heatmaps and UMAPs showing the gene expression profiles for cell type-specific transcriptional modules 10, 21, 43 and 74. The top annotation row indicates a total of 100 cells with the highest and lowest probabilities in the module and are colored by their cell cluster labels. Selected marker genes for B cells (CD79A, CD79B, MS4A1 and CD19), pDCs (ITM2C, IRF7, LILRA4 and CLEC4C), CD14$^+$ monocytes (S100A9, S100A8, S100A12, VCAN and CD14) and T cells (TRAC, CD3D, TRBC1 and CD3G) are highlighted on the right.

cell subpopulations can be explored with the $\varphi$ probability matrix which contains the probability of each module within each cell subpopulation (Figure 3A). This matrix gives insights into the absolute abundance of each module within the same cell subpopulation. For example, module 62 contains actin-related housekeeping genes such as ACTB and ARPC1B, and has higher expression than most other modules within each cell population. A relative probability heatmap can also be produced by taking the z-score of the module probabilities across cell subpopulations (Figure 3B). Examining the relative abundance of a transcriptional module among different cell populations can be useful for finding modules that exhibit specific patterns across cell populations even if they have an overall lower absolute probability compared with other modules. For example, module 65 contains CD8A and CD8B, and has an overall

lower abundance compared with other housekeeping modules such as module 62 within each cell population (as can be observed in Figure 3A). However, module 65 has higher relative expression in the T-cell populations 17, 18 and 19, and can be used to classify CD8$^+$ subpopulations (as can be observed in Figure 3B).

Traditional single-cell workflows such as those utilized in Seurat (19) and Scanpy (43) seek to identify genes that are specific to cell populations using differential expression between that population and all other cells. In Celda, several modules are specific to individual cell populations or cell types. For example, module 10 is expressed in clusters 2, 3 and 4 and contains the B-lymphocyte antigen receptor genes CD79A and CD79B, as well as the B-lymphocyte cell surface antigens MS4A1 and CD19 (Figure 3C). Module 21 contains pDC marker genes ITM2C, IRF7, LILRA4
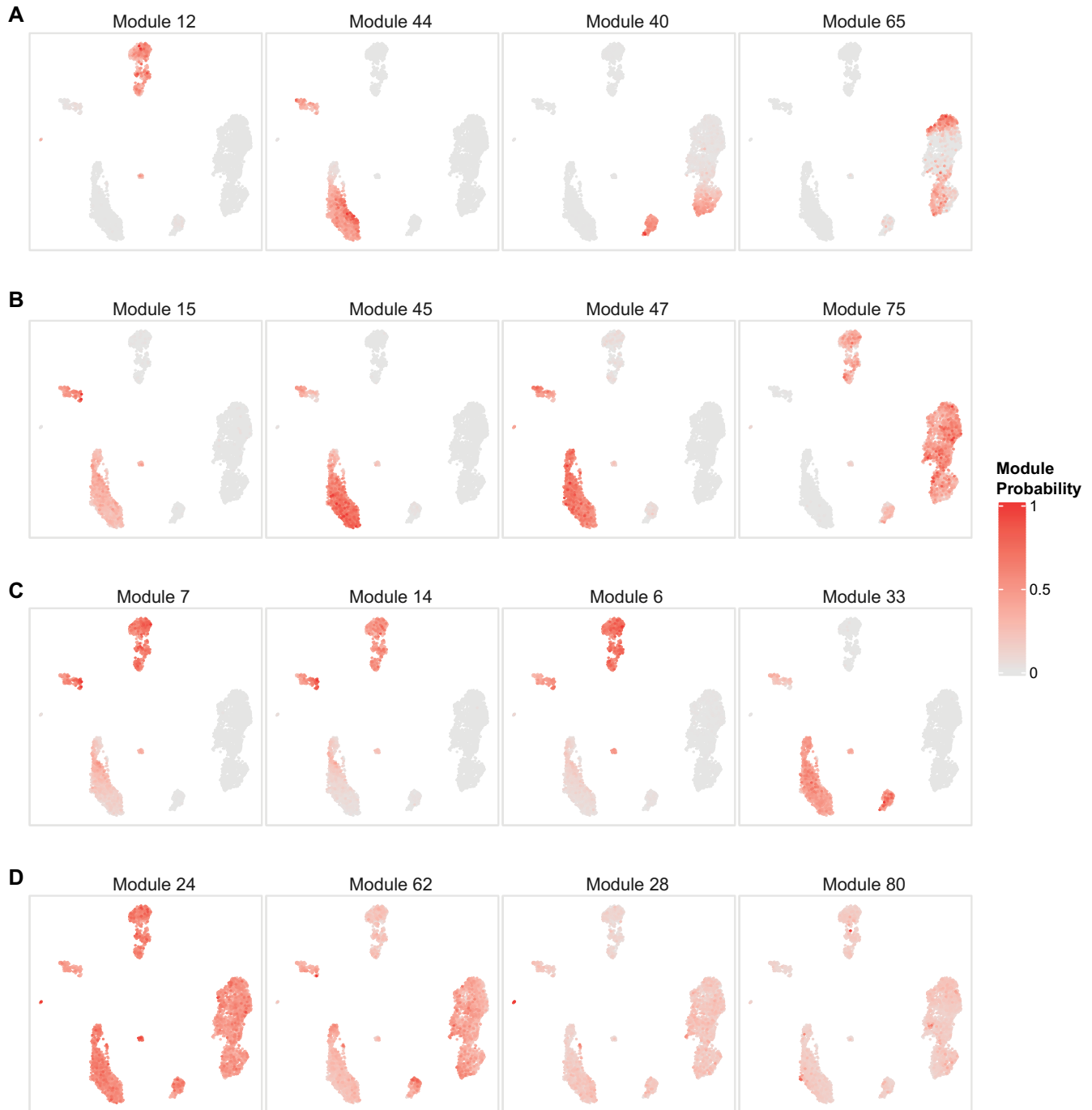
and CLEC4C, and has high probability in cell population 7 (Figure 3D). Module 43 contains monocyte cell markers S100A9, S100A8, S100A12, VCAN and CD14, and has high probabilities in cell populations 11, 12 and 13 (Figure 3E). Module 74 contains T-cell receptor genes TRAC, CD3D, TRBC1 and CD3G, and has high probabilities in cell populations 15–20 (Figure 3F). UMAPs colored by module probabilities can illustrate the patterns of transcriptional modules across cell populations.

In addition to the identification of co-expressed genes specific to a single cell type, Celda gene modules can also be used to identify transcriptional programs that are jointly expressed across multiple cell populations. For example, transcriptional modules 12, 44, 40 and 65 have high probability in at least two unique cell subpopulations (Figure 4A). Module 12 contains genes BANK1 and BLNK associated with B-cell activation, and genes FCGR2B and HLA-DOB associated with antigen processing and presentation, and have high probability in both B cells and pDCs (44–46). Module 44 contains genes including LYZ and SIRPA that are associated with both DCs and CD14$^+$ monocytes (47,48). Module 40 is present across NKs, cytotoxic T cells and NKT cells, and contains granzyme genes such as GZMA and GZMH important for cytolytic activity (49). Module 65 is expressed in both naive and cytotoxic T cells and contains genes for the CD8 receptor, CD8A and CD8B (50). Transcriptional modules 15, 45, 47 and 75 are present in at least three unique cell subpopulations (Figure 4B). Module 15 contains myeloid lineage genes CD33, CSF2RA and IL1R2 (51,52). Module 45 contains Toll-like receptor genes TLR2, TLR4 and TLR8. Module 47 contains C-type lectin domain family genes CLEC4A, CLEC7A, CLEC12A and CLEC4G, and leukocyte immunoglobulin-like receptor genes LILRA2 and LILRB3. These three modules all have high probabilities to varying degrees in DCs, pDCs and monocytes. Module 75 contains lymphoid lineage marker CD69 and has high expression in B, T and NK cells (53). Modules 7, 14, 6 and 33 span four unique cell subpopulations (Figure 4C). Modules 7, 14 and 6 have high probability in B cells, DCs, pDCs and monocytes. Modules 7 and 14 are predominated by major histocompatibilty complex (MHC) class II genes which are key determinants of antigen-presenting cells (54) (APCs). Module 6 contains CD74 which is an important chaperone that regulates antigen presentation (55). Module 33 contains genes such as transmembrane immune signaling adaptor TYROBP, IgE receptor gene FCER1G and macrophage inflammatory gene CCL3, and is expressed in DCs, pDCs, monocytes and NK cells (56). Modules 24, 28, 62 and 80 have high probability in almost all cell populations and contain many known housekeeping and essential genes (Figure 4D). Module 28 contains several common housekeeping genes such as GAPDH, HMGB2, HMGB3 and TUBA1C (57). Module 80 contains mitochondrial genes MT-CO1, MT-CO2 and MT-CO3. Although expressed to varying degrees in all cells, an extremely high proportion of these genes can indicate severe stress or poor quality within a cell (58,59). To assess the biological significance of the gene modules, we utilized Enrichr to identify the number of pathways enriched in each module (33). Using the GO Biological Process database, 62 of the 80 modules were enriched for at least
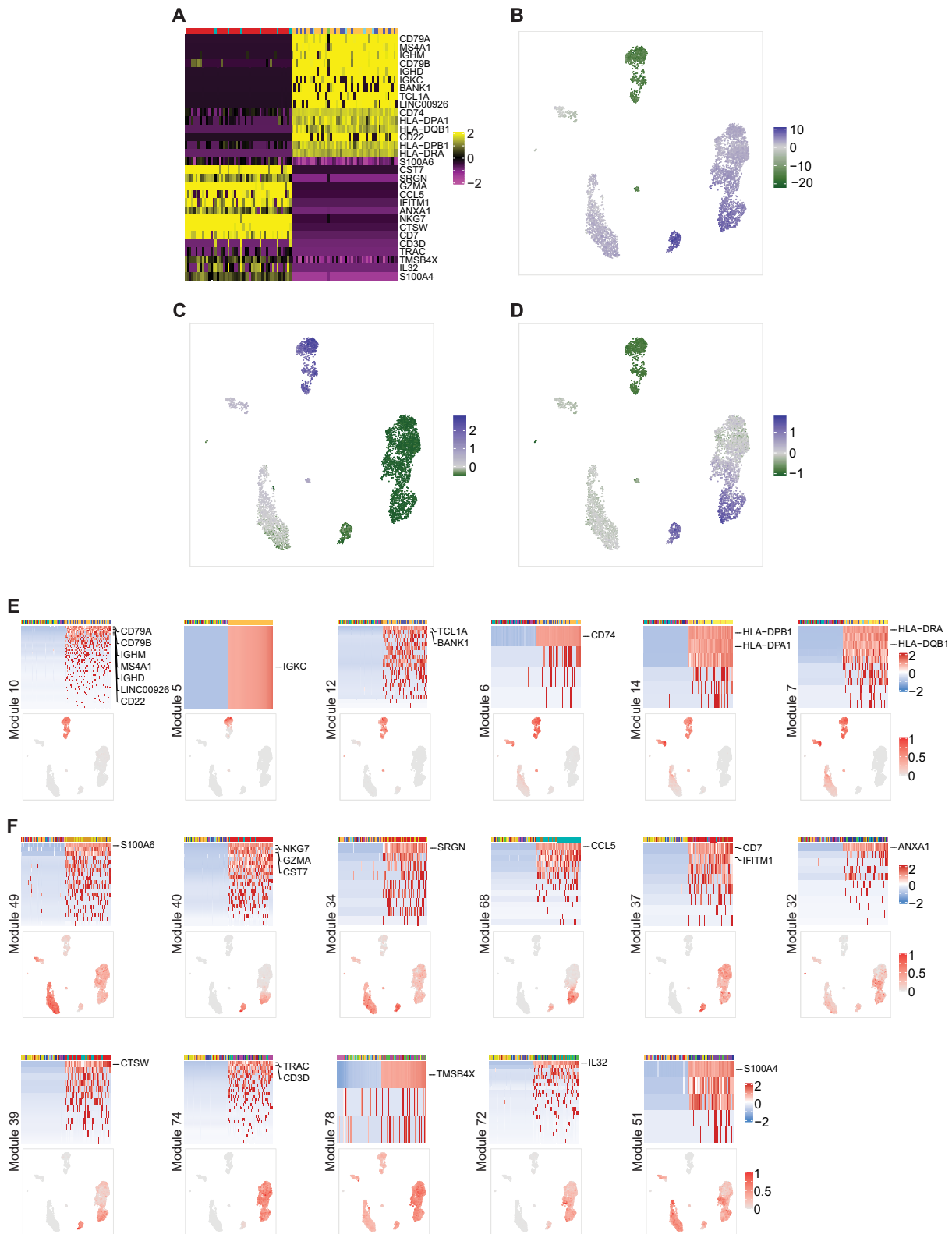
one term (FDR <0.05) and each module had 37 enriched terms on average. These numbers were higher than what was observed for 100 random permutations of the module cluster labels ($P < 0.01$; Supplementary Figure S9). Overall, these results suggest that Celda is able to cluster biologically related genes into modules.

**Qualitative comparison of Celda with principal components for module detection**

Many popular scRNA-seq clustering workflows, including ascend (9), Seurat (19) and TSCAN (21), perform ad hoc dimensionality reduction using PCA before cell clustering. The resulting PCs are used in downstream analyses such as clustering and 2-D embedding as they maintain the relative distance between cells while alleviating noise potentially present in genes expressed at a low level. Genes that have extreme loading scores (positive or negative) to a principal axis will be highly correlated with the corresponding PC and are often plotted together in a heatmap when assessing the quality of PCs or choosing the number of PCs (19,60). However, the biology of the genes associated with PCs is rarely assessed or utilized in downstream analyses. Advantages of Celda modules are that the per cell module probabilities can be used for dimensionality reductions similar to the PCs, and the biology of the co-expressed genes within each module can be used for discovery of biological programs. To qualitatively compare transcriptional modules from Celda with those that can be derived using PCA, we analyzed the same PBMC dataset using Seurat (Figure 5). There are three major issues when trying to define gene modules using PCA. The first issue is that biological programs from different cell types can be represented at each end of the PC. For example, when examining PC2 from the PCA generated by Seurat, the top 15 genes negatively correlated with PC2 contained B-cell markers CD79A, MS4A1 and CD79B, and MHC class II genes, while the top 15 genes positively correlated with PC2 contained T- and NK-cell marker genes including TRAC, CD3D, CD7, CTSW and NKG7 (Figure 5A). Similarly, the B-cell and pDC subpopulations were enriched with negative PC2 scores, while the NK cells and a subset of T cells were enriched with positive PC2 scores (Figure 5B). The average expressions of the top 15 genes negatively correlated with PC2 and the top 15 genes positively correlated with PC2 further confirmed enrichment of PC2-associated genes in different cell types (Figure 5C, D). The second issue is that transcriptional programs co-expressed in a subset of cell populations can be conflated within the same PC. For example, B-cell marker genes such as CD79A, MS4A1 and CD79B were negatively correlated with PC2 along with MHC class II genes such as HLA-DRA and HLA-DPA1. While the MHC class II genes are highly expressed in the B-cell populations, they are also highly expressed in the dendritic cell populations where B-cell marker genes are absent. The third issue is that a gene can be highly correlated with many PCs. For example, CST7, NKG7 and GZMA were among the top 15 genes in PCs 2, 3 and 4, while CD7 was among the top 15 genes in PCs 1, 2 and 8 (Supplementary Figure S10). Overall, these results illustrate that genes from different cell types and different biological programs can be associated with the same PC and a single gene can be asso-

**Figure 4.** Celda identifies transcriptional modules shared across cell populations. (**A**) Selected example UMAPs showing modules with high probabilities in at least two different cell types: module 12 in B cells and pDCs, module 44 in DCs and CD14$^+$ monocytes, module 40 in NK and NKT cells, and module 65 in naive cytotoxic T and cytotoxic T cells. (**B**) Selected UMAPs showing modules with high probabilities in at least three different cell types: modules 15, 45 and 47 in DCs, pDCs and monocytes, and module 75 in B, T and NK cells. (**C**) Selected UMAPs showing modules with high probabilities in at least four different cell types: modules 7, 14 and 6 in B cells, DCs, pDCs and monocytes, and module 33 in DCs, pDCs, monocytes and NK cells. (**D**) Selected UMAPs showing modules with high probabilities in all 20 cell clusters. Analyzing modules can reveal novel insights about biological programs active in one or more cell types.

**Figure 5.** Qualitative comparison of gene co-variation patterns derived from Celda and PCA. (**A**) The 15 genes with the most positive loadings for PC2 and 15 genes with the most negative loadings for PC2 are shown in rows of the heatmap. The 50 cells with the lowest PC2 scores and the 50 cells with the highest PC2 scores are shown in the columns of the heatmap. The top annotation row contains Celda cell subpopulation labels. (**B**) UMAP colored by scores for PC2. (**C**) UMAP colored by the average scaled expression of the top 15 genes negatively correlated with PC2. (**D**) UMAP colored by the average scaled expression of the top 15 genes positively correlated with PC2. (**E**) Heatmaps and UMAPs of six Celda modules containing the top 15 genes negatively correlated with PC2. (**F**) Heatmaps and UMAPs of 11 Celda modules containing the top 15 genes positively correlated with PC2. Overall, these results show that the genes most highly correlated with PC2 from PCA can have different patterns of expression across cell types. In contrast, Celda provided additional insight to gene co-variation by categorizing these top genes into more refined transcriptional modules.

ciated with multiple PCs, which can obscure the biological interpretation of these transcriptional programs.

Celda provided additional insight to gene co-variation by categorizing these top genes into more refined transcriptional modules. For example, among the top 15 genes negatively correlated with PC2, seven genes were in module 10 and expressed across all B-cell subpopulations (Figure 5E). However, other genes were clustered in five other modules and exhibited different patterns across cell populations. IGKC was found in module 5 by itself and was expressed only in one of the B-cell subpopulations (cell cluster 3). TCLA and BANK1 were in module 12 which was present in B-cell and pDC populations. Similarly, the MHC class II-associated genes were found in modules 6, 14 and 7. These three modules had high probability in B-cell, DC and pDC populations and moderate probability in different subsets of monocyte populations to varying degrees. Among the top 15 genes positively correlated with PC2, five were clustered in modules 40, 68 and 39 (Figure 5F). These modules showed enrichment in NK, NKT cell and cytotoxic T-cell populations which were enriched with positive PC2 scores. However, 10 remaining genes clustered in eight other modules showed patterns undetected by PC2. For example, S100 family genes S100A6 and S100A4 were found in modules 49 and 51, which were present in DCs, monocytes and subsets of T and NK cells. SRGN was in module 34 which was present in NKs, DCs, pDCs and monocytes, and had moderate probability in T cells. CD7 and IFITM1 were found in module 37 which had high probability in NKs and T cells. Annexin family gene ANXA1 was grouped in module 32 which was present in subsets of T cells, NKs, monocytes, DCs and CD34$^+$ cells to varying degrees. T-cell receptor genes TRAC and CD3D were grouped in module 74 which was present across all T-cell subpopulations. TMSB4X was grouped in module 78 which had high probability in T cells and moderate probability in all other cell populations. IL32 was in module 72 which had high probability in proliferating T cells and moderate probability in other T-cell populations. Overall, these results suggest that Celda can identify transcriptional programs representing unique biological processes with better clarity than what can be readily parsed by associating genes with PCs from PCA.
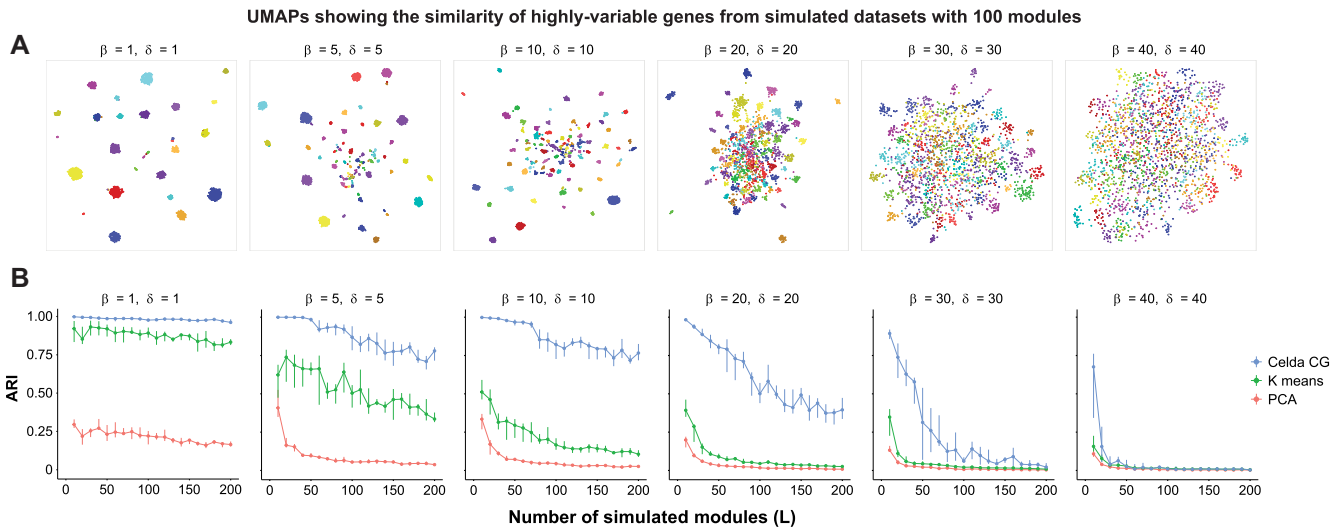
### Benchmarking of Celda clustering

To systematically benchmark Celda's ability to cluster genes into modules, we compared the performance of Celda_CG, PCA and k-means to accurately cluster genes into modules based on simulated data with true gene module labels. To create distinct, non-overlapping modules from PCA, each gene was assigned to a single PC based on the magnitude of its loading ranks across all PCs (see the Materials and Methods). Six datasets were simulated with increasing similarity between cell populations and modules (Figure 6A). Celda_CG outperformed PCs in clustering co-expressed genes into transcriptional modules for all simulated clustering difficulties (Figure 6B). Median ARIs for the six increasing clustering difficulties were 0.98, 0.88, 0.84, 0.57, 0.10 and 0.01 for Celda_CG, and 0.20, 0.06, 0.04, 0.02, 0.01 and

0.01 for PCs. These results demonstrate that Celda_CG was more accurate at identifying modules of co-expressed genes compared with a PCA and k-means-based approaches. For cell clustering, we utilized the 'DuoClustering2018' R package (61) and compared Celda's performance with that of 11 other algorithms across nine datasets using median ARI. Based on median ARI across all datasets, Celda ranked 6 out of the 12 algorithms tested. However, Celda performed within 0.1 median ARI of the top algorithm in 6 of the 9 datasets, suggesting that the cell clustering accuracy is relatively close to other tools (Supplementary Figure S11A). In the KohTCC and Zhengmix4eq datasets, performance of Celda increased with the addition of one extra cluster. These results suggest that cells can be clustered with Celda with accuracy comparable with other approaches (Supplementary Figure S11B). The run times for different tools used in Duo-Clustering benchmark are shown in Supplementary Figure S12. Finally, we observed that the speed of Celda_CG scales proportionately with dataset size. A median of 161.3, 218.8 and 312.9 s was used to generate the clustering results for 2000 variable genes, and 4340, 33 148 and 68 579 cells for PBMC 4k, 33k and 68k datasets, respectively (Supplementary Figure S13). A median of 122.7 and 344.2 s was used to generate the clustering results for the PBMC 4k dataset with 1000 and 4000 variable genes (Supplementary Figure S13).

## DISCUSSION

Celda is a novel discrete Bayesian hierarchical model for scRNA-seq data that can perform co-clustering of cells into subpopulations and genes into transcriptional modules. When applied to a well-characterized PBMC dataset, Celda revealed novel cell populations missed by other approaches and provided information about the combination of transcriptional programs that distinguished each population. Raw scRNA-seq count data are generally discrete and sparse after UMI corrections are applied. Many available workflows that perform cell clustering for scRNA-seq count data often require pre-processing the data before clustering. Seurat (19), ascend (9), TSCAN (21), SC3 (17), CIDR (11) and scran (18) all perform cell clustering based on dimensionality reduced data, which requires some of the pre-processing steps including normalization of total counts in each cell, logarithmic transformation and/or $z$-score standardization to center and scale the variables. Celda is based on hierarchical Dirichlet multinomial distributions which inherently work with sparse non-negative integer count data without prior normalization. Multinomial distributions have been shown to model UMI-corrected data without inflation better than conventional normalization strategies (62). For example, the single plasma B cell was identified by Celda because it had nearly twice the raw counts compared with any other cell in the PBMC dataset. We also used the top 2000 most variable genes determined by variance-stabilizing transformation (19) for clustering the PBMCs in this particular analysis. We note that this is not a requirement when running Celda. For example, we previously clustered this dataset with Celda by including 4529 genes with at least three counts across

**Figure 6.** Celda_CG achieves a higher accuracy for clustering of genes into modules compared with k-means and a PCA-based approach. Datasets were simulated according to the generative process of the Celda_CG model. Higher values of $\beta$ produced more similar transcriptional modules within each cell population and higher values of $\delta$ produced a more equal distribution of counts between genes within each module. A range of $L$ values from 10 to 200 was simulated for each combination of the parameters. (**A**) UMAPs of genes for one of 10 replicate simulations at $L = 100$ were generated to show the relationship between the 2000 most variable genes. Each point on the UMAP represents a single gene and is colored by its true module label. Genes closer together in the UMAP have more similar expression patterns across cells. (**B**) The ARI shows the similarity between the true module labels and the gene clustering results for Celda_CG, k-means or PCA. Points and vertical lines represent medians and interquartile ranges of 10 replicate simulations. Celda_CG achieved a higher median ARI compared with k-means and PCA for all $L$ values in five out of six datasets.

three cells (26). While the overall cluster solutions are similar, applying the variability filter in this analysis promoted the clustering of CD8A and CD8B into a unique module that helped to define the naive CD8$^+$ T-cell population. In general, limiting the analysis to variable genes can decrease the computational time and help identify modules of genes with lower overall counts, but will exclude some genes from being characterized in transcriptional modules.

Celda is a discrete Bayesian model that is based on solid statistical principles that borrow some ideas developed in the field of topic modeling. Popular topic models such as LDA and NMF can be conceived as a special case of Celda with a certain set of prior assumptions. Celda groups genes into modules which are co-expressed across all cells in the dataset, whereas topics from LDA or factors from NMF identify groups of genes that co-vary across a subset of cells. Furthermore, each gene is grouped into a single module in Celda whereas genes will be active in all topics or factors to varying degrees. We note that Celda can actually be used in conjunction with other factorization methods. Celda modules produce a reduced dimensional representation of the dataset which can in turn be used as input into other factorization methods. The factors identified by this procedure will represent combinations of modules that define continuous cell states active to different degrees within each cell. Other bi-clustering methods such as QUBIC2 identify blocks of co-expressed genes within a subset of samples. In the output of these tools, each gene and sample may belong to multiple blocks or not be assigned to a block at all. In contrast, Celda is a co-clustering method which assigns each cell to a single subpopulation and each gene to a single module. Given the differences in goals of the two types of

clustering approaches, we did not benchmark Celda against other bi-clustering methods.

One major challenge with clustering tools applied to any data type is determining the number of clusters. Statistical metrics to assess cluster stability can be used in conjunction with prior biological knowledge to settle on a solution that is robust and gives the most biological insight. Seurat implements modularity-based community detection where a resolution parameter is used to customize the granularity level at which community structures are detected, but does not provide inherent metrics for choosing the number of clusters (19,60,63). Ascend sets a supervised pruning window in the agglomerative hierarchical clustering procedure using Ward's minimum variance to determine the number of subpopulations (9,64). TSCAN uses Gaussian mixture modeling which relies on the Bayesian information criterion (BIC) to determine the number of clusters (21,65). In Celda, we use RPC (25) to assist in choosing the number of cell clusters ($K$) and transcriptional modules ($L$). We note that the elbows in the RPC plots can provide good starting point for choosing these numbers. However, further splitting of modules or cell populations by choosing higher $L$ or $K$ may be useful in some settings and can be performed after examining UMAPs and module heatmaps. Another limitation of our current model is that technical differences between batches of samples are not taken into account. In the future, we plan to develop distributions that can specifically model technical variation between groups of samples. Overall, Celda presents a novel model-based clustering approach towards simultaneously characterizing cellular and transcriptional heterogeneity in biological samples profiled with scRNA-seq assays.

## DATA AVAILABILITY

## SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

## ACKNOWLEDGEMENTS

## FUNDING

## REFERENCES

1. Kanehisa,M., Araki,M., Goto,S., Hattori,M., Hirakawa,M., Itoh,M., Katayama,T., Kawashima,S., Okuda,S., Tokimatsu,T. *et al.* (2008) KEGG for linking genomes to life and the environment. *Nucleic Acids Res.*, **36**, D480–D484.
2. Papalexi,E. and Satija,R. (2018) Single-cell RNA sequencing to explore immune cell heterogeneity. *Nat. Rev. Immunol.*, **18**, 35–45.
3. Potter,S.S. (2018) Single-cell RNA sequencing for the study of development, physiology and disease. *Nat. Rev. Nephrol.*, **14**, 479–492.
4. Blei,D.M. (2012) Probabilistic topic models. *Commun. ACM*, **55**, 77–84.
5. Blei,D.M., Ng,A.Y. and Jordan,M.I. (2003) Latent dirichlet allocation. *J. Mach. Learn. Res.*, **3**, 993–1022.
6. Wang,C. and Blei,D.M. (2009) In: *Proceedings of the 22nd International Conference on Neural Information Processing Systems*. Curran Associates Inc., Vancouver, British Columbia, Canada, pp. 1982–1989.
7. Yin,J. and Wang,J. (2014) In: *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, NY, pp. 233–242.
8. Shafiei,M.M. and Milios,E.E. (2006) Latent dirichlet co-clustering. *Sixth International Conference on Data Mining (ICDM'06)*. pp. 542–551.
9. Senabouth,A., Lukowski,S.W., Hernandez,J.A., Andersen,S.B., Mei,X., Nguyen,Q.H. and Powell,J.E. (2019) ascend: r package for analysis of single-cell RNA-seq data. *Gigascience*, **8**, giz087.
10. Sun,Z., Chen,L., Xin,H., Jiang,Y., Huang,Q., Cillo,A.R., Tabib,T., Kolls,J.K., Bruno,T.C., Lafyatis,R. *et al.* (2019) A bayesian mixture model for clustering droplet-based single-cell transcriptomic data from population studies. *Nat. Commun.*, **10**, 1649.
11. Lin,P., Troup,M. and Ho,J.W. (2017) CIDR: ultrafast and accurate clustering through imputation for single-cell RNA-seq data. *Genome Biol.*, **18**, 59.
12. Li,X., Wang,K., Lyu,Y., Pan,H., Zhang,J., Stambolian,D., Susztak,K., Reilly,M.P., Hu,G. and Li,M. (2020) Deep learning enables accurate clustering with batch effect removal in single-cell RNA-seq analysis. *Nat. Commun.*, **11**, 2338.
13. Sun,Z., Wang,T., Deng,K., Wang,X.F., Lafyatis,R., Ding,Y., Hu,M. and Chen,W. (2018) DIMM-SC: a dirichlet mixture model for clustering droplet-based single cell transcriptomic data. *Bioinformatics*, **34**, 139–146.
14. Zurauskiene,J. and Yau,C. (2016) pcaReduce: hierarchical clustering of single cell transcriptional profiles. *BMC Bioinf.*, **17**, 140.
15. Yang,Y., Huh,R., Culpepper,H.W., Lin,Y., Love,M.I. and Li,Y. (2019) SAFE-clustering: single-cell aggregated (from Ensemble) clustering for single-cell RNA-seq data. *Bioinformatics*, **35**, 1269–1277.
16. Huh,R., Yang,Y., Jiang,Y., Shen,Y. and Li,Y. (2020) SAME-clustering: single-cell aggregated clustering via mixture model Ensemble. *Nucleic Acids Res.*, **48**, 86–95.
17. Kiselev,V.Y., Kirschner,K., Schaub,M.T., Andrews,T., Yiu,A., Chandra,T., Natarajan,K.N., Reik,W., Barahona,M., Green,A.R. *et al.* (2017) SC3: consensus clustering of single-cell RNA-seq data. *Nat. Methods*, **14**, 483–486.
18. Lun,A.T., McCarthy,D.J. and Marioni,J.C. (2016) A step-by-step workflow for low-level analysis of single-cell RNA-seq data with bioconductor. *F1000Research*, **5**, 2122.
19. Stuart,T., Butler,A., Hoffman,P., Hafemeister,C., Papalexi,E., Mauck,W.M. 3rd, Hao,Y., Stoeckius,M., Smibert,P. and Satija,R. (2019) Comprehensive integration of single-cell data. *Cell*, **177**, 1888–1902.
20. Wang,B., Ramazzotti,D., De Sano,L., Zhu,J., Pierson,E. and Batzoglou,S. (2018) SIMLR: a tool for large-scale genomic analyses by multi-kernel learning. *Proteomics*, **18**, 1700232.
21. Ji,Z. and Ji,H. (2016) TSCAN: pseudo-time reconstruction and evaluation in single-cell RNA-seq analysis. *Nucleic Acids Res.*, **44**, e117.
22. Chen,S., Hua,K., Cui,H. and Jiang,R. (2019) VPAC: variational projection for accurate clustering of single-cell transcriptomic data. *BMC Bioinf.*, **20**, 139–151.
23. Langfelder,P. and Horvath,S. (2008) WGCNA: an R package for weighted correlation network analysis. *BMC Bioinf.*, **9**, 559.
24. Pontes,B., Giraldez,R. and Aguilar-Ruiz,J.S. (2015) Biclustering on expression data: a review. *J. Biomed. Inform.*, **57**, 163–180.
25. Zhao,W., Chen,J.J., Perkins,R., Liu,Z., Ge,W., Ding,Y. and Zou,W. (2015) A heuristic approach to determine an appropriate number of topics in topic modeling. *BMC Bioinf.*, **16**(Suppl. 13), S8.
26. Yang,S., Corbett,S.E., Koga,Y., Wang,Z., Johnson,W.E., Yajima,M. and Campbell,J.D. (2020) Decontamination of ambient RNA in single-cell RNA-seq with decontX. *Genome Biol.*, **21**, 57.
27. Becht,E., McInnes,L., Healy,J., Dutertre,C.A., Kwok,I.W.H., Ng,L.G., Ginhoux,F. and Newell,E.W. (2019) Dimensionality reduction for visualizing single-cell data using UMAP. *Nat. Biotechnol.*, **37**, 38–44.
28. Finak,G., McDavid,A., Yajima,M., Deng,J., Gersuk,V., Shalek,A.K., Slichter,C.K., Miller,H.W., McElrath,M.J., Prlic,M. *et al.* (2015) MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol.*, **16**, 278.
29. Benjamini,Y. and Hochberg,Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Statist. Soc. Ser. B (Methodol.)*, **57**, 289–300.
30. McCarthy,D.J., Campbell,K.R., Lun,A.T. and Wills,Q.F. (2017) Scater: pre-processing, quality control, normalization and visualization of single-cell RNA-seq data in R. *Bioinformatics*, **33**, 1179–1186.
31. Csardi,G. and Nepusz,T. (2006) The igraph software package for complex network research. *InterJournal, Complex Systems*, 1695.
32. The,Tabula Muris Consortium.2018. Single-cell transcriptomics of 20 mouse organs creates a *Tabula Muris. Nature*, **562**, 367–372.
33. Chen,E.Y., Tan,C.M., Kou,Y., Duan,Q., Wang,Z., Meirelles,G.V., Clark,N.R. and Ma'ayan,A. (2013) Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinf.*, **14**, 128.

34. Manning,C.D., Raghavan,P. and Schütze,H. (2008) *Introduction to information retrieval*. Cambridge University Press, Cambridge, England.

35. Sharma,A., Podolsky,R., Zhao,J. and McIndoe,R.A. (2009) A modified hyperplane clustering algorithm allows for efficient and accurate clustering of extremely large datasets. *Bioinformatics*, **25**, 1152–1157.

36. Thalamuthu,A., Mukhopadhyay,I., Zheng,X. and Tseng,G.C. (2006) Evaluation and comparison of gene clustering methods in microarray analysis. *Bioinformatics*, **22**, 2405–2412.

37. Pehkonen,P., Wong,G. and Toronen,P. (2005) Theme discovery from gene lists for identification and viewing of multiple functional groups. *BMC Bioinf.*, **6**, 162.

38. van der Leun,A.M., Thommen,D.S. and Schumacher,T.N. (2020) CD8$^+$ T cell states in human cancer: insights from single-cell analysis. *Nat. Rev. Cancer*, **20**, 218–232.

39. Soares,A., Govender,L., Hughes,J., Mavakla,W., de Kock,M., Barnard,C., Pienaar,B., Janse van Rensburg,E., Jacobs,G., Khomba,G. *et al.* (2010) Novel application of Ki67 to quantify antigen-specific in vitro lymphoproliferation. *J. Immunol. Methods*, **362**, 43–50.

40. Lindqvist,C.A., Christiansson,L.H., Simonsson,B., Enblad,G., Olsson-Stromberg,U. and Loskog,A.S. (2010) T regulatory cells control T-cell proliferation partly by the release of soluble CD25 in patients with B-cell malignancies. *Immunology*, **131**, 371–376.

41. Tellier,J. and Nutt,S.L. (2017) Standing out from the crowd: how to identify plasma cells. *Eur. J. Immunol.*, **47**, 1276–1279.

42. Xie,J., Ma,A., Zhang,Y., Liu,B., Cao,S., Wang,C., Xu,J., Zhang,C. and Ma,Q. (2020) QUBIC2: a novel and robust biclustering algorithm for analyses and interpretation of large-scale RNA-Seq data. *Bioinformatics*, **36**, 1143–1149.

43. Wolf,F.A., Angerer,P. and Theis,F.J. (2018) SCANPY: large-scale single-cell gene expression data analysis. *Genome Biol.*, **19**, 15.

44. Gomez Hernandez,G., Morell,M. and Alarcon-Riquelme,M.E. (2021) The role of BANK1 in B cell signaling and disease. *Cells*, **10**, 1184.

45. Fu,C., Turck,C.W., Kurosaki,T. and Chan,A.C. (1998) BLNK: a central linker protein in B cell activation. *Immunity*, **9**, 93–103.

46. Junker,F., Gordon,J. and Qureshi,O. (2020) Fc gamma receptors and their role in antigen uptake, presentation, and T cell activation. *Front. Immunol.*, **11**, 1393.

47. Collin,M. and Bigley,V. (2018) Human dendritic cell subsets: an update. *Immunology*, **154**, 3–20.

48. Kapellos,T.S., Bonaguro,L., Gemund,I., Reusch,N., Saglam,A., Hinkley,E.R. and Schultze,J.L. (2019) Human monocyte subsets and phenotypes in major chronic inflammatory diseases. *Front. Immunol.*, **10**, 2035.

49. Cullen,S.P., Brunet,M. and Martin,S.J. (2010) Granzymes in cancer and immunity. *Cell Death Differ.*, **17**, 616–623.

50. Philip,M. and Schietinger,A. (2021) CD8$^+$ T cell differentiation and dysfunction in cancer. *Nat. Rev. Immunol.*, **10**, 1184.

51. Borot,F., Wang,H., Ma,Y., Jafarov,T., Raza,A., Ali,A.M. and Mukherjee,S. (2019) Gene-edited stem cells enable CD33-directed immune therapy for myeloid malignancies. *Proc. Natl Acad. Sci. USA*, **116**, 11978–11987.

52. Autenshlyus,A., Arkhipov,S., Mikhailova,E., Marinkin,I., Arkhipova,V. and Varaksin,N. (2019) The relationship between cytokine production, CSF2RA, and IL1R2 expression in mammary adenocarcinoma, tumor histopathological parameters, and lymph node metastasis. *Technol. Cancer Res. Treat.*, **18**, 1533033819883626.

53. Lugthart,G., Melsen,J.E., Vervat,C., van Ostaijen-Ten Dam,M.M., Corver,W.E., Roelen,D.L., van Bergen,J., van Tol,M.J., Lankester,A.C. and Schilham,M.W. (2016) Human lymphoid tissues harbor a distinct CD69$^+$CXCR6$^+$ NK cell population. *J. Immunol.*, **197**, 78–84.

54. Roche,P.A. and Furuta,K. (2015) The ins and outs of MHC class II-mediated antigen processing and presentation. *Nat. Rev. Immunol.*, **15**, 203–216.

55. Leng,L., Metz,C.N., Fang,Y., Xu,J., Donnelly,S., Baugh,J., Delohery,T., Chen,Y., Mitchell,R.A. and Bucala,R. (2003) MIF signal transduction initiated by binding to CD74. *J. Exp. Med.*, **197**, 1467–1476.

56. Baba,T. and Mukaida,N. (2014) Role of macrophage inflammatory protein (MIP)-1alpha/CCL3 in leukemogenesis. *Mol. Cell Oncol.*, **1**, e29899.

57. Hounkpe,B.W., Chenou,F., de Lima,F. and De Paula,E.V. (2021) HRT atlas v1.0 database: redefining human and mouse housekeeping genes and candidate reference transcripts by mining massive RNA-seq datasets. *Nucleic Acids Res.*, **49**, D947–D955.

58. Osorio,D. and Cai,J.J. (2021) Systematic determination of the mitochondrial proportion in human and mice tissues for single-cell RNA sequencing data quality control. *Bioinformatics*, **37**, 963–967.

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

60. Butler,A., Hoffman,P., Smibert,P., Papalexi,E. and Satija,R. (2018) Integrating single-cell transcriptomic data across different conditions, technologies, and species. *Nat. Biotechnol.*, **36**, 411–420.

61. Duo,A., Robinson,M.D. and Soneson,C. (2018) A systematic performance evaluation of clustering methods for single-cell RNA-seq data. *F1000Research*, **7**, 1141.

62. Townes,F.W., Hicks,S.C., Aryee,M.J. and Irizarry,R.A. (2019) Feature selection and dimension reduction for single-cell RNA-Seq based on a multinomial model. *Genome Biol.*, **20**, 295.

63. Waltman,L. and van Eck,N.J. (2013) A smart local moving algorithm for large-scale modularity-based community detection. *Eur. Phys. J. B*, **86**, 471.

64. Nguyen,Q.H., Lukowski,S.W., Chiu,H.S., Senabouth,A., Bruxner,T.J.C., Christ,A.N., Palpant,N.J. and Powell,J.E. (2018) Single-cell RNA-seq of human induced pluripotent stem cells reveals cellular heterogeneity and cell state transitions between subpopulations. *Genome Res.*, **28**, 1053–1066.

65. Fraley,C. and Raftery,A.E. (2002) Model-based clustering, discriminant analysis, and density estimation. *J. Am. Statist. Assoc.*, **97**, 611–631.