

1 Scalable nonparametric clustering with unified marker gene 2 selection for single-cell RNA-seq data

3

4 Chibuikem Nwizu^{1,2}, Madeline Hughes³, Michelle L. Ramseier⁴⁻⁸, Andrew W. Navia⁴, Alex K. Shalek⁴⁻⁸,
5 Nicolo Fusi³, Srivatsan Raghavan^{4,9-11,§}, Peter S. Winter^{4,§}, Ava P. Amini^{3,§,†}, and Lorin Crawford^{3,§,†}

6 1 Center for Computational Molecular Biology, Brown University, Providence, RI, USA

7 2 Warren Alpert Medical School of Brown University, Providence, RI, USA

8 3 Microsoft Research, Cambridge, MA, USA

9 4 Broad Institute of MIT and Harvard, Cambridge, MA, USA

10 5 Institute for Medical Engineering and Science, Massachusetts Institute of Technology,
11 Cambridge, MA, USA

12 6 Koch Institute for Integrative Cancer Research, Massachusetts Institute of Technology,
13 Cambridge, MA, USA

14 7 Department of Chemistry, Massachusetts Institute of Technology, Cambridge, MA, USA

15 8 Ragon Institute of MGH, MIT, and Harvard, Cambridge, MA, USA

16 9 Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

17 10 Harvard Medical School, Boston, MA, USA

18 11 Department of Medicine, Brigham and Women's Hospital, Boston, MA, USA

19 § Authors Contributed Equally

20 † Corresponding E-mail: ava.amini@microsoft.com; lcrawford@microsoft.com

21 Abstract

22 Clustering is commonly used in single-cell RNA-sequencing (scRNA-seq) pipelines to characterize cellular
23 heterogeneity. However, current methods face two main limitations. First, they require user-specified
24 heuristics which add time and complexity to bioinformatic workflows; second, they rely on post-selective
25 differential expression analyses to identify marker genes driving cluster differences, which has been shown
26 to be subject to inflated false discovery rates. We address these challenges by introducing nonparamet-
27 ric clustering of single-cell populations (NCLUSION): an infinite mixture model that leverages Bayesian

sparse priors to identify marker genes while simultaneously performing clustering on single-cell expression data. NCLUSION uses a scalable variational inference algorithm to perform these analyses on datasets with up to millions of cells. Through simulations and analyses of publicly available scRNA-seq studies, we demonstrate that NCLUSION (i) matches the performance of other state-of-the-art clustering techniques with significantly reduced runtime and (ii) provides statistically robust and biologically relevant transcriptomic signatures for each of the clusters it identifies. Overall, NCLUSION represents a reliable hypothesis-generating tool for understanding patterns of expression variation present in single-cell populations.

Introduction

Recent advances in sequencing technologies have increased the throughput of genomic studies to millions of single cells, necessitating computational workflows to explore and analyze these data¹. In single-cell RNA sequencing (scRNA-seq), unsupervised clustering and marker gene selection are integral steps in the exploratory phase of analyses²⁻⁵. Clustering facilitates the identification of cell types, while marker gene selection enables the annotation of gene modules and cluster-specific biological programs. However, there has yet to be a consensus on the best approach to clustering cells and identifying the transcriptomic signatures that characterize them^{6,7}. This has resulted in a multitude of proposed clustering methods for single-cell data, many of which are reliant on various user-defined heuristics that prevent practitioners from performing an unbiased survey of data and limit each method’s “out-of-the-box” applicability when analyzing multiple studies.

Many current clustering approaches take a subset of highly variable genes as input, use dimensionality reduction techniques to simplify the representation of single-cell expression for these genes, and then perform clustering on top of this reduced representation. K-nearest neighbor (KNN) algorithms⁸, for example, generate nearest-neighbor (NN) graphs using transcriptomic similarity scores between cells and then perform Louvain clustering on the estimated graphs. Popular methods such as Seurat⁹ and scLCA¹⁰ use principal component analysis (PCA) and singular value decomposition to learn a lower-dimensional representation of cells, respectively. Ensemble approaches such as scCESS-SIMLR¹¹ learn over a mixture of kernels to generate a final cell-cell similarity matrix which is then used in a spectral clustering algorithm.

Notably, selecting an appropriate embedding for single-cell data is not always a straightforward task. Previous studies have shown that if the generated lower-dimensional representation does not accurately capture underlying biological relationships between cells, then both the quality of clustering and the generalizability of findings in downstream analyses can be compromised^{12,13}. Factors such as the number of highly variable genes retained during data preprocessing and the number of components used to define the lower-dimensional embedding can affect the ability of a clustering algorithm to identify fine-grained differences between cell types^{7,14}. Furthermore, nearly all state-of-the-art clustering methods require users to specify the number of clusters K to be used in the algorithm. Strategies such as consensus-finding^{11,15}, outlier detection¹⁶, and iterative cluster merging and splitting^{17–19} rely on human-in-the-loop interactive steps within their algorithms to determine an “optimal” choice for K . Generally, requiring users to make these additional decisions can add significant time and complexity when using clustering as a preliminary analysis in bioinformatic workflows.

Perhaps the biggest limitation of current single-cell clustering algorithms is that most do not directly identify top marker genes that are driving the inference of different biologically significant clusters; instead, they use post-selective inference to find genes that are differentially expressed between the inferred cell groups^{7,20}. Since the point of clustering algorithms is to separate dissimilar data into different groups, it is expected *a priori* that there are differences between the groups and any test statistics computed by comparing the groups are likely to be inflated due to “data double dipping”. Many studies have shown that performing this post-selective inference uncorrected can lead to inflated type I error rates^{20,21}. Though there has been work developed to correct for post-selection, these are still in nascency, and most do not yet scale to high-dimensional settings^{22–27}. Recently, others have proposed unified frameworks for simultaneous clustering and marker gene selection using hierarchical tree-based algorithms²⁸, regularized copula models²⁹, and *post hoc* sensitivity measures³⁰. However, these approaches rely on arbitrary thresholding to find “significant” marker genes and fail to theoretically test a well-defined null hypothesis, making them difficult to biologically interpret.

We present “Nonparametric CLustering of Single-cell populatiONs” (NCLUSION): a unified Bayesian nonparametric framework that simultaneously performs clustering and marker gene selection. NCLUSION works directly on normalized single-cell count data, bypassing the need to perform dimensionality reduction. By modeling the expression of each gene as a sparse hierarchical Dirichlet process normal mixture model^{31–37}, NCLUSION both learns the optimal number of clusters based on the variation observed

between cellular expression profiles and uses sparse prior distributions to identify genes that significantly influence cluster definitions. The key to our proposed integrative framework is that clustering and extracting marker genes concurrently is a more efficient approach to the exploratory analysis of scRNA-seq data, as it effectively allows each process to inform the other. Most importantly, our approach eliminates the need for human-in-the-loop decisions, significantly reducing the runtime and complexity of these analyses. Altogether, NCLUSION mitigates the need for heuristic choices (e.g., choosing specific lower-dimensional embeddings), avoids iterative hyper-parameter optimization, bridges the interpretability gap suffered by many unsupervised learning algorithms in single-cell applications, and scales to accommodate the growing sizes of emerging scRNA-seq datasets.

Results

NCLUSION simplifies traditional clustering workflows

Conventional scRNA-seq clustering approaches include numerous steps that require user heuristics or human-in-the-loop decisions which increase runtime and complexity (Fig 1A). These can range from deciding how to optimally embed high-dimensional expression data into a lower-dimensional space to selecting the number of clusters, K , to identify in the data. Furthermore, current methods require that marker gene selection is performed post-clustering, which can lead to inflated rates of false discovery^{20,21}. In this work, we aim to address these challenges using a new approach: NCLUSION.

NCLUSION leverages a Bayesian nonparametric mixture modeling framework to reduce the number of choices that users need to make while simultaneously performing variable selection to identify top cluster-specific marker genes for downstream analyses (Fig 1B; see Methods for details). There are three key components of our model formulation that distinguish it from traditional bioinformatic workflows. First, NCLUSION is fit directly on single-cell expression matrices and does not require the data to be embedded within a lower-dimensional space. Second, we implicitly assume *a priori* that cells can belong to one of infinitely many different clusters. This is captured by placing a Dirichlet process prior over the cluster assignment probabilities for each cell. By allowing the number of possible clusters $K = \infty$, we remove the need for users to iterate through different values until they find the optimal choice. Third, NCLUSION assumes that not all genes are important when assigning cells to a given cluster. To model this, we place a spike and slab prior on the mean expression of each gene within each cluster. This prior

shrinks the mean of genes that play an insignificant role in cluster formation towards zero.

To identify cluster-specific marker genes, we start by estimating posterior inclusion probabilities (PIPs), which represent our confidence that a gene’s mean expression within a cluster is nonzero. These PIPs act as a signature that can be used to distinguish clusters. Since NCLUSION fits each gene and cell independently, signatures learned between clusters can share the same subsets of genes. To select for unique cell type markers, we multiply each PIP with a usage weight, which is calculated by performing min-max normalization over the proportion of clusters in which a gene is determined to be statistically significant (i.e., using the median probability model threshold³⁸ $\text{PIP} \geq 0.5$). We use these adjusted inclusion probabilities with the effect size sign (ESS) and strictly standardized mean difference (SSMD) of each gene to filter for significant genes that are substantially up-regulated (indicated by positive ESS and large SSMD values). The genes remaining after filtering make up cluster-specific marker gene modules that provide insight into the biological features underlying cluster assignments.

To enable efficient posterior inference that can scale as the size of scRNA-seq datasets continues to grow, we train NCLUSION using variational expectation-maximization (EM). This algorithm leverages a “mean-field” assumption to approximate the true posterior distribution over model parameter estimates with a product of simpler distributions^{39–41}. In training, our objective is to minimize the Kullback–Leibler (KL) divergence between the variational posterior and the true posterior⁴². Optimization during model fitting occurs using a coordinate ascent procedure where parameters are sequentially updated based on their gradients (Methods and Supplementary Material). With this variational approach, NCLUSION is capable of scaling well up to 1 million cells without applying any dimensionality reduction to the input data.

We evaluated the runtime of NCLUSION against a set of state-of-the-art single-cell clustering methods using publicly available datasets. The methods we used for comparison include: Seurat⁹, scLCA¹⁰, K-nearest neighbors followed by the Leiden clustering algorithm (KNN+Leiden)⁴³, SOUP⁴⁴, and scCCESS-SIMLR¹¹. Each of these methods operates by first reducing the dimensionality of the input data and then performing clustering on the reduced representation. To the best of our knowledge, NCLUSION is the only method to date that clusters on the expression matrix directly while jointly identifying cluster-specific salient genes. NCLUSION does not incur additional runtime cost due to this model design choice. In fact, NCLUSION boasts a faster runtime compared to other methods, particularly as the number of cells in a dataset grows. We showcase this scalability on the BRAIN-LARGE⁴⁵ dataset (Fig 1C). To

144 facilitate comparison across all baseline methods, we follow Lopez et al.¹⁴ and limit the analysis to
145 720 genes, while subsampling the number of cells from 500 to 1 million cells. As a practical reference,
146 we use a grey dotted line to highlight the runtime for each method at the median scRNA-seq dataset
147 size as determined in 2020: 31,000 cells⁴⁶. We observed that only NCLUSION and KNN+Leiden were
148 able to scale past 100,000 cells, while NCLUSION was the only method able to run on 1 million cells.
149 Additionally, NCLUSION records competitive runtimes across varying numbers of genes (Fig S1). These
150 scalability results show the potential of NCLUSION’s utility for emerging large-scale single-cell studies.

151 **NCLUSION accurately identifies clusters and marker genes in simulations**

152 We used simulations to evaluate the performance of NCLUSION in controlled settings. Here, we used
153 **scDesign3**⁴⁷ to generate synthetic datasets consisting of 10,000 cells and 1,000 genes (due to computa-
154 tional constraints of the software) distributed over five clusters. We considered four different scenarios
155 (with 20 replicates per scenario), where we varied cluster size and marker gene composition (Methods).
156 Scenario I was the simplest, where we evenly distributed all cells across the five clusters, and each cluster
157 had 50 marker genes. In Scenarios II and III, all five clusters had 50 marker genes, but one cluster had
158 significantly fewer cells than the other four clusters. Lastly, in Scenario IV, each of the five clusters had
159 the same number of cells, but one cluster had a signature of only 20 marker genes, while the other four
160 clusters had 50 marker genes each.

161 We first compared the cluster identification accuracy of NCLUSION with the previously described al-
162 gorithms: Seurat⁹, scLCA¹⁰, K-nearest neighbors followed by the Leiden clustering algorithm (KNN+Leiden)⁴³,
163 SOUP⁴⁴, and scCCESS-SIMLR¹¹. The relatively smaller size of these simulated datasets also allowed
164 us to include three additional widely used methods: CIDR⁴⁸, SC3⁴⁹, and scDeepCluster⁵⁰. We should
165 highlight that all competing methods, except for NCLUSION, first perform dimensionality reduction
166 prior to clustering, allowing us to evaluate the impact of dimensionality reduction on cluster recovery.

167 The normalized mutual information (NMI) and adjusted Rand index (ARI) were calculated to quan-
168 titatively evaluate the clustering results given by each algorithm. In all scenarios, scDeepCluster had
169 the best performance, with NCLUSION, scLCA, and scCCESS-SIMLR rounding out the top four. The
170 reason for scDeepCluster’s top performance is that it first performs Leiden clustering on principal com-
171 ponents from the expression data to obtain the initial cluster assignments. Then, it uses a deep neural
172 network to refine the cluster assignments; therefore, performance is highly dependent on the initialization

of the cluster assignments and on the success of the initial clustering. NCLUSION, on the other hand, only performs clustering once to achieve comparable results. Notably, all methods performed relatively worse in Scenarios II and III than their performance in Scenarios I and IV (Fig S2; Table S1) due to the class imbalance in how the cells were distributed across clusters.

We also evaluated the performance of NCLUSION on marker gene detection using the same simulated datasets and compared it with three popular differential expression algorithms: DUBStepR⁵¹, singleCellHaystack²⁵, and FESETM⁵². Notably, of these methods, only NCLUSION is able to find cluster-specific marker genes. FESETM uses an alternative algorithm (via the Scott-Knott test) to assign marker genes to clusters found by an independent clustering method, while both DUBStepR and singleCellHaystack aim to identify differentially expressed genes. To that end, we evaluated the global marker gene detection of each model using true positive rate (TPR), false discovery rate (FDR), and false positive rate (FPR; computed as 1-Specificity for each method). We found no statistically significant differences when comparing the median power of NCLUSION to the other approaches (Kruskal–Wallis H -test $P > 0.99$; Fig S3 and Table S2). However, importantly, NCLUSION and FESETM were the only two methods to have high power while maintaining a low FDR and FPR. On the other hand, singleCellHaystack and DUBStepR achieved similar power but only because they incorrectly labeled many genes as being marker genes (resulting in markedly higher FDR and FPR).

NCLUSION achieves competitive clustering performance on PBMC data with less runtime

Next, we assessed the quality of clustering done by NCLUSION as compared to other baseline approaches on real data. Here, we analyzed scRNA-seq from FACS-purified peripheral blood mononuclear cells (PBMCs)⁵³. This dataset captures 10 cellular populations, including CD14+ monocytes, CD34+ cells, major lymphoid lineages (B and T cells), as well as other phenotypic lineages within the T cell population, including CD4+ helper T cells, CD8+ cytotoxic T cells, CD4+ regulatory T cells, and CD4+ memory T cells. After quality control (Methods), the final dataset contained 94,615 cells and 5,000 genes with the highest standard deviation (post log-normalization)^{20,54}. We evaluated the performance of NCLUSION, Seurat, scLCA, the KNN+Leiden algorithm, SOUP, and scCCESS-SIMLR by comparing inferred cluster assignments to the cell type annotations from the original study, which were obtained via a combination of FACS analysis and clustering with Seurat⁵³ (Fig 3A-B).

To qualitatively assess clustering performance, we used contingency heat maps to evaluate how well each method captured the unique cell types across clusters (Fig 3C). For a given method, each n -th row of the heat map represents an annotation from the original study and each k -th column represents an inferred cluster identified by the method. The color saturation of each (n, k) -th element in the heat map indicates the fraction of a n -th cell annotation that a given method assigned to the k -th inferred cluster. Overall, all baselines were able to distinguish the B cell population from other cells, and each approach uniquely clustered a majority of the CD14+ monocytes together (Fig 3C). Furthermore, all methods except scCCESS-SIMLR and KNN+Leiden were able to separate the natural killer (NK) cell population with an inferred cluster occupancy rate of greater than 90% (Table S3). NCLUSION's performance was most similar to that of Seurat (χ^2 -test $P = 0.99$ when assessing independence between their contingency tables). Both methods were able to identify major PBMC cell lineages and were the only approaches to divide the B cell, cytotoxic CD8+ T cells, and regulatory T cells into subpopulations⁵⁵ (Fig 3C).

As a quantitative assessment of each method's clustering performance, we used the FACS-derived experimental annotations as reference labels and computed NMI and ARI, each measuring how well the clustering algorithm's labels matched the reference labels (Methods). For both metrics, values closer to 1 indicate better clustering performance. To assess the robustness and consistency of each method, we ran them on five different randomly subsampled partitions containing 80% of the cells in the dataset. We report the mean metric score for each clustering algorithm across these partitions, along with corresponding 95% confidence intervals (Fig 3D). NCLUSION outperformed all competing approaches across both metrics. It obtained the highest mean NMI coefficient of 0.80 ($\pm 1.62 \times 10^{-2}$ standard deviation), with Seurat and scLCA each scoring lower values of 0.77 ($\pm 2.78 \times 10^{-3}$) and 0.62 ($\pm 3.04 \times 10^{-2}$), respectively. These differences were statistically significant, as determined by two-sided t-tests ($P = 1.08 \times 10^{-3}$ and $P = 1.90 \times 10^{-6}$, respectively). When comparing performance using the ARI, NCLUSION remained competitive and significantly outperformed other methods (Table S4). Specifically, NCLUSION achieved a higher mean ARI of 0.67 ($\pm 2.02 \times 10^{-2}$) when compared to 0.62 ($\pm 2.48 \times 10^{-3}$; two-sided t-tests $P = 1.17 \times 10^{-3}$) for Seurat and 0.50 ($\pm 2.52 \times 10^{-2}$; two-sided t-tests $P = 3.10 \times 10^{-6}$) for scLCA, respectively.

Lastly, NCLUSION recorded the shortest runtime for this analysis without any iterative processes or optimization of hyperparameters. It finished approximately 254 seconds (4.23 minutes) faster than Seurat, more than 12,900 seconds (215.00 mins or 3.58 hrs) faster than scLCA, and more than 41,331

seconds (688.85 mins or 11.48 hrs) faster than scCCESS-SIMLR.

NCLUSION is well-powered to identify PBMC-specific marker genes

The key distinguishing property of NCLUSION is its inherent ability to perform variable selection. NCLUSION thus provides users with cluster labels for each cell as well as unique gene signatures that define each cluster. The statistical model underlying NCLUSION selects cluster-specific marker genes based on two criteria: (i) an adjusted posterior inclusion probability, $PIP(j; k)$, which provides evidence that the j -th gene's mean expression is uniquely nonzero within the k -th cluster, (ii) the sign of the j -th gene's effect, $ESS(j; k)$, which is used to determine whether it is uniquely up-regulated or down-regulated within the k -th cluster, and (iii) the magnitude of the j -th gene's effect, $SSMD(j; k)$, on the definition of the k -th cluster (Methods). We assessed the marker genes identified by NCLUSION for each of the inferred clusters in the PBMC dataset to determine whether they provide insights into the biology of different cell types (Fig 4A-B).

Overall, NCLUSION successfully identified cluster-specific marker genes that are known to be associated with examined cell types (Fig 4C and Table S5). For example, in cluster 8 which has cells mapping back to the cytotoxic and naive cytotoxic T cell population, we observed that NCLUSION correctly identified marker genes known to play an important role in cytotoxic T cell biology, such as *CD8A* (adjusted $PIP = 0.90$), *CD8B* (adjusted $PIP = 0.70$)⁵⁶, and *CD27* (adjusted $PIP = 0.62$)⁵⁷. Furthermore, genes associated with cytotoxicity, such as *GZMM*, tended to be selected in clusters largely containing CD8+ T cells (Cluster 8; adjusted $PIP = 0.60$) and NK cells (Cluster 2; adjusted $PIP = 0.61$)—two cell types that have been shown to have functionally similar cytotoxic activity⁵⁸. In other clusters, where we observed genes associated with both B cells (e.g., in Cluster 1, *CD19*, adjusted $PIP = 1.00$; *LINC00926*, adjusted $PIP = 1.00$; *MS4A1*, $PIP = 0.90$)^{59–61} and myeloid lineages (e.g., in Cluster 3, *MS4A6A*, $PIP = 1.00$; *S100A8*, adjusted $PIP = 0.93$; *LYZ*, $PIP = 0.90$)⁶², NCLUSION distinguished genes known to play an important role in T cell biology as statistically significant. This observation suggests that NCLUSION accounted for the variance among T cell and T cell-like expression patterns when distinguishing cell types in the PBMC dataset. We observed that our criteria for cluster-specific marker genes, based on high PIPs and positive ESS scores, strongly agreed with the relative over-expression of each gene in its respective cluster (Fig 4C). Imposing a threshold on SSMD allowed NCLUSION to filter out less relevant genes, narrowing the list of cluster-specific over-expressed genes to those most salient.

To further evaluate marker gene quality, we computed gene module scores in order to compare the normalized expression for signature genes across clusters (Fig 4D and Table S6). Here, we find that each module exhibits the highest expression within its respective cluster, with the most definitive signatures occurring within the B (inferred cluster 1), NK (inferred cluster 2), monocytic (inferred cluster 3), and CD34+ cells (inferred cluster 4) (Fig 4E and Fig S4). In clusters that contained heterogeneous combinations of T cell subpopulations, we still see an increased relative expression among cluster-specific marker genes, although not as distinct as in the other cell types.

As an additional analysis, we compared the similarity between the marker genes identified by NCLUSION with the list of marker genes that are identified by using a *post hoc* differential expression analysis with Seurat. Here, we took the FACS-derived experimental annotations from Zheng et al.⁵³ and found differentially expressed genes by doing a one-versus-all Wilcoxon rank sum test for each cluster (mirroring the typical procedure in a conventional bioinformatic workflow). As expected, this post-selective inference procedure resulted in Seurat identifying a multitude of candidate marker genes for each cell type, even after Bonferroni correction. A direct comparison between NCLUSION and Seurat showed that the proposed Bayesian variable selection approach in the NCLUSION framework results in smaller and more refined transcriptomic signatures for downstream investigation (Fig 4F and Figs S5- S6). For example, NCLUSION identified 134 cluster-specific marker genes for NK cells, 97% of which were also included in the 1,780 marker genes selected *post hoc* by Seurat. In total, an average of 96% of the marker genes identified by NCLUSION were included in the much larger sets of differentially expressed genes selected by Seurat across each of the FACS-annotated cell types.

Over-representation analysis using gene product annotations in Gene Ontology (GO) further confirmed that the selective set of gene modules inferred by NCLUSION reflect known immune cell biology⁶³⁻⁶⁵ (Fig 4G-H and Table S7). For example, in the cluster containing predominantly B cells (inferred cluster 1), we observed an up-regulation of B cell receptor signaling (adjusted $P = 6.68 \times 10^{-7}$), B cell activation (adjusted $P = 2.3 \times 10^{-7}$), and B cell proliferation (adjusted $P = 2.51 \times 10^{-7}$)^{66,67}. Notably, some of the biologically relevant GO terms found when using NCLUSION were not statistically significant when applying the larger sets of marker genes provided *post hoc* by Seurat. For instance, in the cluster with NK cells (inferred cluster 2), significant gene sets from NCLUSION included the positive regulation of leukocyte chemotaxis (adjusted $P = 1.53 \times 10^{-4}$) and positive regulation of natural killer cell chemotaxis (adjusted $P = 4.88 \times 10^{-6}$)⁶⁸⁻⁷⁰. We also observed enrichment of the up-regulation of natural killer cell

mediated cytotoxicity (adjusted $P = 3.81 \times 10^{-5}$), consistent with the known highly cytotoxic behavior of NK cells^{68,71}. Each of these gene sets was insignificant when using differentially expressed genes from Seurat (Fig 4G). Lastly, in the monocyte-dominated cluster (inferred cluster 3), the NCLUSION-generated module was enriched for macrophage activation involved in immune response (adjusted $P = 2.41 \times 10^{-5}$) and antigen-presenting activity (e.g., MHC Class II antigen presentation, adjusted $P = 5.39 \times 10^{-4}$). Complete marker gene and GO analyses for all clusters inferred by NCLUSION and Seurat as a baseline in the PBMC dataset can be found in Table S7. Together, these results demonstrate that NCLUSION can identify cluster-specific gene signatures that reflect underlying cellular phenotypes.

NCLUSION's performance generalizes to the other single-cell datasets

Finally, we assessed the generalizability of NCLUSION by testing it on three additional large scRNA-seq datasets of various sample sizes: a pancreatic ductal adenocarcinoma (PDAC) dataset from Raghavan et al.⁷² with $N = 23,042$ cells; an acute myeloid leukemia (AML) dataset from van Galen et al.⁷³ with $N = 43,690$ cells; and a tissue immune (IMMUNE) atlas dataset from Domínguez Conde et al.⁷⁴ with $N = 88,057$ cells. These datasets represent a range of tissue and disease states to assess our method's performance in different use cases. Both the PDAC and AML datasets contain a mixture of malignant and non-malignant cells from different patient biopsies, while the IMMUNE dataset contains healthy white blood cells from different anatomical locations. After performing quality control (Methods), we had a total of 5,000 genes with the highest standard deviation (after log-normalization) for the analysis.

We observed similar scalability in the runtime of NCLUSION and competing baselines on all three datasets (Fig 5A). NCLUSION maintains its computational efficiency, now only being slightly outperformed by Seurat on the PDAC and AML datasets due to longer convergence time in its variational EM algorithm. We also found that NCLUSION continued to remain competitive in terms of clustering performance. When quantitatively evaluating the clustering ability of NCLUSION versus the competing baselines using the annotations provided by the original studies, NCLUSION was often statistically significantly better (as determined via a two-sided t-test, $P < 0.05$) according to ARI and NMI across all datasets (Fig 5B-C, Figs S7-S19, and Table S4).

We then analyzed the interpretability of the cluster-specific marker genes inferred by NCLUSION (Fig 5D-G, Figs S8-S24, and Tables S8-S19). For brevity, we highlight just notable results from the PDAC and IMMUNE datasets in the main text. Additional analyses for the AML dataset can be found

in the Supplementary Material (see Figs S12-S17 and Tables S16-S19).

To begin, we focused on evaluating the gene modules generated from the NCLUSION inferred clusters in the PDAC dataset (Fig 5D-G). As with the PBMC data, we observed higher module expression within the respective clusters (Fig 5F). NCLUSION appeared to use immune cell signatures as the primary axis for distinguishing malignant and non-malignant populations (Fig 5E and Table S10). For example, the inferred cluster 6 predominantly contained cells that were originally annotated as NK and T cells by Raghavan et al.⁷². This inferred cluster had immune cell type specific marker genes such as *CD2* (PIP = 0.92), *GZMB* (PIP = 0.85), *IL7R* (PIP = 0.85), and *NCAM1* (PIP = 1.00)^{75,76} (Fig 5D). An additional GO analysis of this cluster showed an enrichment of natural killer cell mediated cytotoxicity (adjusted $P = 6.60 \times 10^{-9}$) and T cell receptor signaling (adjusted $P = 1.17 \times 10^{-21}$) (Fig 5G).

Notably, the other clusters that primarily contained non-malignant cells (inferred clusters 4, 5, 11, 12, 13, and 14) also directly aligned with cell type labels originally annotated by Raghavan et al.⁷². For the clusters that primarily contained malignant and metastatic cells (i.e., inferred clusters 1, 2, 3, 9, and 10), a GO analysis revealed an enrichment of extracellular matrix (ECM) organization and cell migration processes (Fig 5G). Importantly, however, NCLUSION also had the power to divide these cells into more granular subpopulations based on their level of differentiation. For example, the inferred cluster 9 was enriched for both cell migration processes (e.g., MET-activated PTK2 signaling, adjusted $P = 2.97 \times 10^{-5}$; MET-promoted cell motility, adjusted $P = 4.02 \times 10^{-5}$) and fibroblast cell activity (pancreatic fibroblasts, adjusted $P = 2.1 \times 10^{-9}$; collagen formation, adjusted $P = 1.13 \times 10^{-7}$)^{77,78}.

Finally, we evaluated the granularity of NCLUSION's clustering on the IMMUNE dataset, which contained 33 manually annotated labels from experts⁷⁴. When analyzing this study, NCLUSION was able to delineate between multiple T cell sub-lineages, whereas methods like Seurat, KNN+Leiden, and scCCESS SIMLR merged these subpopulations into 1 or 2 clusters. For example, the inferred cluster 12 by NCLUSION was enriched (~95% occupancy rate) for CD8+ effector memory (T_{EM}) and effector memory cells re-expressing CD45RA (T_{EMRA}), while NCLUSION's inferred cluster 8 was enriched (~90% occupancy rate) for CD8+ tissue-resident memory (T_{RM}). These two populations have been shown to be functionally distinct subpopulations in the CD8+ T cell lineage^{79,80}. Likewise NCLUSION's inferred clusters 6 and 10 were enriched (~77% and ~72% occupancy rates, respectively) for functionally distinct subpopulations of CD4+ T cells, namely T follicular helper cells (T_{fh}) and CD4+ effector/effector memory T cells, respectively^{79,81,82} (Table S12).

The GO analysis of NCLUSION-generated gene modules also showed an enrichment of T cell phenotypes. The inferred cluster 12’s top ontology term was indeed “CD8+ Effector Memory T4” (adjusted $P = 4.18 \times 10^{-9}$), while the inferred cluster 10 had “CD4+ Central Memory T1” as a top enriched term (adjusted $P = 3.53 \times 10^{-3}$). Similar results showing how NCLUSION-generated gene modules granularly distinguish these cellular populations can be found in the Supplementary Material (Table S15).

Discussion

We present NCLUSION: a scalable Bayesian nonparametric framework designed to serve as an unbiased method for inferring phenotypic clusters and identifying cluster-specific marker genes in scRNA-seq experiments. We show how our approach simplifies traditional single-cell transcriptomic workflows, which often rely on the transformation of the data to a lower-dimensional representation to facilitate clustering and iteratively tune the number of clusters used to obtain optimal results. In contrast, NCLUSION operates on the full normalized gene expression matrix, eliminating the need for transformation to a lower-dimensional space; infers the optimal number of clusters without iterative user refinement; and simultaneously identifies the cluster-specific marker genes that significantly drive the clustering. By leveraging a variational inference algorithm, NCLUSION can scale to scRNA-seq studies with a million cells. Through the analysis of a collection of large-scale publicly available datasets, we show that NCLUSION not only achieves clustering performance comparable to state-of-the-art methods but also provides refined sets of gene candidates for downstream analyses. By unifying clustering and marker gene selection, NCLUSION provides a flexible and unified statistical framework for inferring complex differential gene expression patterns observed in heterogeneous tissue populations^{21,83,84}.

The current implementation of the NCLUSION framework offers many directions for future development and applications. First, NCLUSION assumes a normal mixture model for log-normalized gene expression data. We use this assumption both because log-based transformations have been shown to reduce the effects of sparsity in single-cell analyses^{20,85} and because the Gaussian-based specification offers computational advantages for scalable posterior inference. Still, future extensions of the NCLUSION framework should explore the utility of Poisson- and negative binomial-based likelihoods to deal with the zero-inflated nature of scRNA-seq studies in their raw form.

Second, the current formulation of NCLUSION models the gene expression of each cell indepen-

dently and does not consider, for example, the correlation between genes with similar functionality or co-expression patterns between genes within the same signaling pathway. One possible extension of NCLUSION would be to incorporate additional genomic information into the sparse prior distributions used for Bayesian variable selection. For example, previous studies have proposed an integrative approach where the importance of a variable also depends on an additional set of covariates^{40,86,87}. In the case of single-cell applications, we could assume that the prior probability of the j -th gene being a marker of the k -th cluster is also dependent upon its cellular pathway membership. Unlike the current spike-and-slab prior NCLUSION implements, this new prior would assume that biologically related pathways contain shared marker genes, essentially integrating the concept of gene set enrichment analysis into clustering. An alternative approach would be to extend NCLUSION to incorporate non-diagonal correlation structures by exploring sparse covariance models, which could provide a balance between the need for maintaining computational efficiency while representing a richer set of gene dependencies^{88,89}.

Third, although it helps NCLUSION scale to large datasets, variational expectation-maximization (EM) algorithms are known to both produce slightly miscalibrated parameter estimates and underestimate the total variation present within a dataset^{41,90,91}. While this does not greatly affect the performance of NCLUSION in the evaluations presented in this paper, this can be seen as a limitation depending on the application of interest. For example, in the PBMC dataset, NCLUSION is unable to resolve all the different T cell subtypes that were annotated by Zheng et al.⁵³ (Fig 3). This is most likely due to variational approximations being well-suited to describe the global variation across cells but at the cost of smoothing over local variation between smaller subpopulations. Considering other (equally scalable) ways to carry out approximate Bayesian inference may be relevant for future work⁹².

Lastly, a thrust of recent work in genomics has been to develop methods that identify spatially variable marker genes as a key step during analyses of spatially-resolved transcriptomics data⁹³. Future efforts could extend NCLUSION to this emerging modality by, for example, reformulating the method as a spatial Dirichlet process mixture model⁹⁴.

In sum, NCLUSION provides a unified framework for simultaneous clustering and marker gene selection in single-cell transcriptomic data, yielding improvements in computational efficiency, interpretability, and scalability. We envision that NCLUSION will accelerate key analytic steps universal to single-cell analysis across diverse applications.

Materials and methods

Overview of NCLUSION

We provide a brief overview of the probabilistic framework underlying the “Nonparametric CLustering of SIngle-cell populatiONs” (NCLUSION) model. Detailed derivations of the algorithm are provided in the Supplementary Material. Consider a study with single-cell RNA sequencing (scRNA-seq) expression data for $n = 1, \dots, N$ cells that each have measurements for $j = 1, \dots, J$ genes. Let this dataset be represented by the $N \times J$ matrix \mathbf{X} where the row-vector $\mathbf{x}_n = (x_{n1}, \dots, x_{nJ})$ denotes the expression profile for the n -th cell. We assume that the log-normalized gene expression for each cell follows a sparse hierarchical Dirichlet process normal mixture model^{31–33} of the form

$$x_{nj} \sim \sum_{k=1}^{\infty} \pi_k \mathcal{N}(\nu_j + \mu_{jk}, \sigma_j^2) \quad (1)$$

where π_k represents the marginal (unconditional) probability that a cell belongs to the k -th cluster, ν_j and σ_j^2 are the global means and variances for the j -th gene across all cells (i.e., not conditioned on cluster identity), and μ_{jk} is the mean shift of expression for the j -th gene within the k -th cluster. There are two key features in the model formulation of NCLUSION specified above. First, we assume that the formation of clusters is driven by a few important genes that have mean expression shifted away from a baseline gene-specific expression level, ν_j . To that end, we place a sparsity-inducing spike and slab prior distribution on the mean effect of each gene

$$\mu_{jk} \sim \eta \mathcal{N}(0, \lambda_{jk} \sigma_j^2) + (1 - \eta) \delta_0, \quad (2)$$

where δ_0 is a point mass at zero, λ_{jk} scales the global variance to form a cluster-specific “slab” distribution for each gene, and η is the prior probability that any given gene has a nonzero effect when assigning a cell to any cluster. In practice, there are many different ways to estimate η . Following previous work^{40,41,95–97}, one choice would be to assume a uniform prior over $\log \eta$ to reflect our lack of knowledge about the correct number of “marker” genes for each cell type that is present in the data. Instead, in this work, we assume $\eta \sim \text{Beta}(1, 1)$ to represent this uncertainty and learn its value during model inference. To facilitate posterior computation and interpretable inference, we introduce a binary indicator variable $\rho_{jk} \in \{0, 1\}$

where we implicitly assume *a priori* that $\Pr[\rho_{jk} = 1] = \eta$. Alternatively, we say that ρ_{jk} takes on a value of 1 when the effect of a gene μ_{jk} on cluster assignment is nonzero and deviates from the baseline gene expression level ν_j . As NCLUSION is trained, the posterior mean for unimportant genes will trend towards the global mean (i.e., $\mu_{jk} \rightarrow 0$) as the model attempts to identify subsets of marker genes that are relevant for each cluster. We then use posterior inclusion probabilities (PIPs) as general summaries of evidence that the j -th gene is statistically important in determining when a cell is assigned to the k -th cluster where

$$\text{PIP}(j; k) \equiv \Pr[\mu_{jk} \neq 0 \mid \mathbf{X}]. \quad (3)$$

The second key feature in Eq. (1) is that we do not assume to know the true number of clusters K . Instead, we take a nonparametric approach and attempt to learn K directly from the data. Once again, to facilitate posterior computation, we introduce a categorical latent variable ψ_n which indicates that the n -th cell is in the k -th cluster with prior probability π_k . Explicitly, we write this as $\Pr[\psi_n = k] = \pi_k$. Here, we implement the stick-breaking construction of the Dirichlet process³¹ where we say

$$\boldsymbol{\pi} \sim \text{Dir}(\alpha_0 \boldsymbol{\beta}), \quad \beta_k \sim \chi_k \prod_{l=1}^{k-1} (1 - \chi_l), \quad \chi_k \sim \text{Beta}(1, \gamma_0) \quad (4)$$

with $\boldsymbol{\pi} = (\pi_1, \dots, \pi_{>K})$ having mean $\boldsymbol{\beta}$ and variance determined by the concentration hyper-parameters α_0 and γ_0 ⁹⁸. The concentration hyper-parameters α_0 and γ_0 are both non-negative scalars that effectively help to determine the number of clusters used in the model^{31,39}. Larger values for these parameters increase the model's sensitivity to variation in the data and encourage the creation of a greater number of smaller clusters. Smaller values for these parameters, on the other hand, decrease the model's sensitivity to variation in the data and encourage the creation of fewer larger clusters. In this work, we encourage the creation of fewer clusters and fix α_0 and γ_0 to be less than or equal to 1 (Supplementary Material). After model training, we use the posterior distribution over the latent categorical indicators $\Pr[\psi_n = k \mid \mathbf{X}]$ to determine the cluster assignment for each cell. It is worth noting that, although the prior number of normal components is infinite in Eq. (1), the posterior number of components after model fitting will be finite. This truncation reflects the fact that not all infinite states are used when conditioning on finite data⁹⁸. Additionally, the algorithm used to estimate the parameters in the NCLUSION software

penalizes empty clusters (see Supplementary Material), and, as a result, the model has the flexibility to automatically adjust its complexity based on the inferred complexity of the data being analyzed. This helps to increase the utility and adaptability of NCLUSION across a wide range of single-cell applications.

Selection of cluster-specific marker genes

NCLUSION jointly performs clustering on single-cell populations while also learning cluster-specific gene signatures. To achieve this, we use the spike and slab prior distribution specified in Eq. (2) and the resulting PIPs defined in Eq. (3) to find the most salient genes per cluster. Since the model fits to each j -th gene expression for the n -th cell independently, signatures learned between clusters can share subsets of the same genes. Genes that are identified as “important” across many different clusters can effectively be seen as ubiquitous housekeeping variables rather than significant marker genes of unique cell types. Therefore, we down-weight the inclusion probabilities to proportionally penalize genes based on the number of clusters in which they appear

$$\text{PIP}^*(j; k) = w_j \times \text{PIP}(j; k), \quad w_j = \left(1 - \frac{S_j}{K^*}\right) / \left(1 - \frac{1}{K^*}\right) \quad (5)$$

where $K^* \leq K$ is the finite number of occupied clusters learned by the model, and S_j is the number of clusters that the j -th gene is significant in according to a given selection threshold. We set this threshold to be 0.5 which corresponds to the median probability criterion in Bayesian statistics³⁸.

While Eqs. (3) and (5) can be used to identify the genes that are differentially expressed in a given cluster, they do not indicate the direction or magnitude of this shift. Therefore, for each gene, we combine the adjusted posterior inclusion probabilities with effect size sign (ESS) and strictly standardized mean difference (SSMD) measures to find the most salient markers per cluster. Here, we obtain the effect size sign by taking the sign of Cohen’s d ⁹⁹ between the expression of the j -th gene for cells in the k -th cluster and cells not in the k -th cluster (denoted by k')

$$\text{ESS}(j; k) = \text{sgn} \left(\frac{\rho_{jk} \mu_{jk} - \bar{m}_{jk'}}{\sigma_j} \right) \quad (6)$$

where, in addition to previous notation, $\bar{m}_{jk'} = \sum_{k'} \rho_{jk'} \mu_{jk'} / (K^* - 1)$ is the average mean shift for the j -th gene in all clusters outside of the k -th. Here, $\text{sgn}(\cdot)$ is the piecewise sign function where $\text{sgn}(u) = +$

(i.e., positive) when $u > 0$, $\text{sgn}(u) = -$ (i.e., negative) when $u < 0$, and $\text{sgn}(u) = 0$ when $u = 0$.

The strictly standardized mean difference (SSMD) is a metric often used in high-throughput screenings to test for the significance of an effect size magnitude^{100–103}. It is computed as the following

$$\text{SSMD}(j; k) = \frac{\mu_{jk} - \bar{\mu}_{jk'}}{\sqrt{\sigma_j^2 [(N_k - 1)/N_k + (N_{k'} - 1)/N_{k'}]}} \quad (7)$$

where $\bar{\mu}_{jk'} = \sum_{k'} \mu_{jk'}/(K^* - 1)$ is the average global mean for the j -th gene in all clusters outside of the k -th. Asymptotically, the SSMD follows a normal distribution^{100,101,103}. To determine a significant value, we follow a previous procedure¹⁰³ by calculating a threshold $|\text{SSMD}(j; k)| \geq S^*(j; k)$ which controls for a predetermined false positive rate (FPR). Here, this threshold is given by

$$S^*(j; k) = \text{SSMD}_{\min} + \Phi^{-1} \left(1 - \frac{\text{FPR}}{2} \right) \varsigma_{jk} \quad (8)$$

where FPR is set to 0.05, $\Phi^{-1}(\cdot)$ is the inverse cumulative distribution function of a standard normal, and SSMD_{\min} is the minimum SSMD magnitude that one considers to be significant. In practice, this minimum value is often set between 0 and 0.25 in order to identify weak effect sizes. In the main text, we follow previous work^{103–105} and let $\text{SSMD}_{\min} = 0.15$. The parameter ς_{jk} is used to denote the asymptotic variance which is given by

$$\varsigma_{jk} = \frac{(N_k - 1)/N_k^2 + (N_{k'} - 1)/N_{k'}^2}{(N_k - 1)/N_k + (N_{k'} - 1)/N_{k'}} + \frac{(N_k - 1)^2/N_k^2 + (N_{k'} - 1)^2/N_{k'}^2}{2\sigma_j^2 [(N_k - 1)/N_k + (N_{k'} - 1)/N_{k'}]^3} (\mu_{jk} - \bar{\mu}_{jk'})^2. \quad (9)$$

In the main text, cluster-specific marker genes are selected as those that have a significant adjusted inclusion probability and are notably up-regulated in a given cluster meaning that they satisfy the following criteria: (1) $\text{PIP}^*(j; k) \geq 0.5$, (2) $\text{ESS}(j; k) = +$, and (3) $\text{SSMD}(j; k) \geq S^*(j; k)$, respectively.

Posterior inference via variational EM algorithm

We combine the likelihood in Eq. (1) and the prior distributions in Eqs. (2) and (4) to perform Bayesian inference. In current scRNA-seq datasets, it is less feasible to implement traditional Markov Chain Monte Carlo (MCMC) algorithms due to the large number of cells being studied. For model fitting, we instead use a variational expectation-maximization (EM) algorithm^{31,32,98}, which allows us to estimate

parameters within an optimization framework. The overall goal of variational inference is to approximate the true posterior distribution for model parameters using a set of approximating distributions. The EM algorithm optimizes parameters such that it minimizes the Kullback-Leibler divergence between the exact and approximate posterior distributions. To compute the variational approximations, we make the mean-field assumption that the true posterior can be “fully-factorized”¹⁰⁶. The algorithm then follows two general steps. In the first step, we iterate through a combination of hyper-parameter values and compute variational updates for the other parameters using coordinate ascent. In the second step, we empirically compute (approximate) posterior values for the main model parameters $\{\boldsymbol{\mu}, \boldsymbol{\rho}, \boldsymbol{\psi}\}$. Detailed steps in the variational EM algorithm, explicit coordinate ascent updates for the model parameters, pseudocode, and other derivations are given in the Supplementary Material. Parameters in the variational EM algorithm are initialized by taking random draws from their assumed prior distributions. Iterations in the algorithm are terminated when at least one of two stopping criteria are met: (i) the difference between the lower bound of two consecutive updates is within some small range (specified by argument ϵ), or (ii) a maximum number of iterations is reached. For the analyses run in this paper, we set $\epsilon = 1$ for the first criterion and used a maximum of 1×10^4 iterations for the second.

Simulation study design

Generating simulated datasets. To evaluate the robustness and sensitivity of NCLUSION under controlled conditions, we generated synthetic single-cell RNA-seq datasets using `scDesign3`⁴⁷ (v.1.4.0). The reference dataset we used was derived from the FACS-sorted peripheral blood mononuclear cell (PBMC) dataset produced by Zheng et al.⁵³. Initial preprocessing for this reference dataset included mitochondrial gene content assessment, ribosomal and hemoglobin gene filtering, and quality control to remove both low-quality cells and lowly expressed genes. Highly variable genes (HVGs) were identified using the `modelGeneVar` function in the `scrn` R package, and the top 1000 HVGs were retained for downstream simulation. We used five immune cell types (B cells, CD14+ monocytes, CD56+ natural killer (NK) cells, cytotoxic T cells, and regulatory T cells) for these analyses. To ensure balanced representation across cell types, we implemented a stratified subsampling scheme which selected an equal number of cells per type while enforcing non-zero gene expression across all selected cells and genes.

Each simulated dataset comprised of $N = 10,000$ cells across five clusters and 1,000 genes where we preserved realistic transcriptomic correlation structures through Gaussian copula modeling. Simulations

were conducted across four different scenarios (with 20 replicates per scenario), each varying in cluster size imbalance and marker gene composition.

- **Scenario I:** Balanced clusters of 2000 cells per cell type, each with 50 marker genes.
- **Scenario II:** Imbalanced cluster design where one small cluster had 200 cells and the other four larger clusters each had 2450 cells. All clusters contained 50 marker genes.
- **Scenario III:** Imbalanced cluster design where one cluster had 20 (rare) cells and the other four larger clusters each had 2495 cells each. All clusters contained 50 marker genes.
- **Scenario IV:** Balanced clusters of 2000 cells per cell type, but one cluster had only 20 marker genes while the other four clusters had 50 marker genes.

More specifically, synthetic datasets were generated using the `construct_data`, `fit_marginal`, `fit_copula`, `extract_para`, and `simu_new` functions within `scDesign3` to create gene expression vectors using a negative binomial distribution that is conditioned on cell type from the reference data. To introduce differentially expressed genes (DEGs), we first ranked genes by their cell type specific mean expression in the reference data and sampled a number of top-ranked genes to be markers. Then in the synthetic data, these DEGs were then artificially upregulated in one cluster while maintaining the baseline expression in others. This was done by apply a log-fold change factor sampled uniformly over the interval [1.5, 2.5]. This ensured that we maintained realistic variance but still had distinct signal between cell types.

Real datasets and preprocessing

Below we briefly describe all of the datasets and the preprocessing steps used in this work. Each of these datasets is relatively large (containing at least 20,000 cells) with unique molecular identifiers (UMI). The latter is important because prior research suggests that UMIs provide enough information to avoid overcounting issues due to amplification and zero-inflation^{14,107,108}. We use an asterisk by the BRAIN-LARGE dataset to indicate that it was exclusively to test the scalability of NCLUSION and competing methods; therefore, clustering performance was not recorded. For the other datasets, we use cell type annotations provided by the original study as “true” reference labels during our analyses. Cells were filtered for quality using a custom `scanpy`¹⁰⁹ (v.1.9.1) pipeline script (see Software availability). Unless otherwise stated, all data was preprocessed by taking the logarithm (to the base 2) of the counts, dividing

by a scaling factor of 10000, and then adding a pseudo-count of 1.0 for stability. Additionally, unless otherwise stated, all results were produced using the top 5000 highly variable genes (HVG), which were determined by sorting the standard deviation of the transformed counts³³.

BRAIN-LARGE*. This dataset originally contains 1.3 million mouse brain cells from 10x Genomics⁴⁵. During preprocessing, we subset the data to a collection of 720 genes following a procedure outlined by Lopez et al.¹⁴. Next, we further filtered by only keeping cells that had at least one of these genes expressed. This left a total of 64,071 cells. Since the original study did not provide cell labels, we exclusively used this dataset to compare runtime performance. To do so, we up-sampled by randomly selecting groups of 64,071 cells to create a synthetic dataset of 1 million cells. We report the runtime for each method on datasets with 500, 1K, 5K, 10K, 50K, 100K, 500K, and 1M cells.

PBMC. We took scRNA-seq data from fluorescence-activated cell sorted (FACS) populations of peripheral blood mononuclear cells (PBMCs) provided by Zheng et al.⁵³ and concatenated each population into one dataset. During preprocessing, we filtered out genes that were expressed in fewer than three cells. We also dropped cells with (i) fewer than 200 genes expressed, (ii) greater than 20% mitochondrial reads, and (iii) fewer than 5% ribosomal reads. This resulted in a final dataset with 94,615 high-quality cells representing 10 distinct cell types.

PDAC. We used scRNA-seq data from pancreatic ductal adenocarcinoma (PDAC) tissue obtained from 23 patients according to methods documented in Raghavan et al.⁷². This dataset contains 23,042 total cells made up of 15,302 non-malignant cells of 11 distinct cell types and 7,740 malignant cells.

AML. The scRNA-seq data obtained from van Galen et al.⁷³ contains 43,690 acute myeloid leukemia (AML) and non-malignant donor cells taken from 16 AML patients, 5 healthy donors, and 2 cell lines. It is comprised of 13,489 patient-derived malignant cells, 23,005 non-malignant donor cells, 6,018 cells from the MUTZ-3 AML cell line, and 1,178 cells from the OCI-AML3 cell line. To account for the biological differences between cell lines and donor cells of the same cell type annotation, we appended the cell line name onto cell type labels where applicable, producing 33 distinct cell types overall. To process the data, we filtered out all cells with “unclear” cell state labels, retaining only “malignant” or “non-malignant” cells.

IMMUNE. We obtained filtered scRNA-seq data from approximately 330,000 immune cells from 12 organ donors in Domínguez Conde et al.⁷⁴. To mitigate batch effects, we isolated 88,057 cells that were taken from a single organ donor (donor D496) with uniform chemistry annotations, containing 44 distinct immune cell types.

Other methods

We selected five additional methods to compare against the performance of NCLUSION in real data and in simulations: (1) a Louvain algorithm implemented using the `FindClusters` function in Seurat⁹ (v.4.3.0.1); (2) a spectral clustering method called scLCA¹⁰ (v.0.0.0.9000), which optimizes both intra- and inter-cluster similarity; (3) a combination of a K-Nearest Neighbor (KNN) classifier with the Louvain community detection algorithm to find clusters implemented via `scikit-learn`⁴³ (v.1.2.2) and `scanpy`¹⁰⁹ (v.1.9.1), respectively; (4) a semi-soft clustering algorithm called SOUP⁴⁴ (v.0.0.0.9000); and (5) an ensemble method called scCCESS-SIMLR (v.0.2.1), which leverages the spectral clustering approach SIMLR^{11,110}. In the simulation experiments, we also compared the clustering performance of NCLUSION against three additional methods: (6) a consensus clustering method, SC3⁴⁹ (v.1.34.0); (7) a deep-learning based method, scDeepCluster⁵⁰ (v.1.0.0); and (8) an imputation and dimensionality reduction method, CIDR⁴⁸ (v.0.1.5). Also in simulations, when assessing the ability of NCLUSION to perform robust marker gene selection, we compare it against: (9) a method that leverage differential correlation patterns in the local structure of a PCA-derived cell neighborhood graph, DUBStepR⁵¹ (v.1.2.0); (10) a feature selection via an EM algorithm, FESTEM⁵² (v.1.2.1); and (11) a divergence-based strategy with permutation tests, singleCellHaystack²⁵ (v.1.0.2). Additional details about each method are provided in the Supplementary Material.

Evaluation metrics

Below we describe the metrics and approaches used to compare performance across all methods. Our clustering evaluation procedure used extrinsic metrics that require reference labels to serve as the ground truth in our calculations.

Normalized mutual information (NMI). This metric is a normalized variant of mutual information (MI). It is an entropy-based metric that captures the amount of shared information between the inferred

label distribution and the reference label distribution. NMI ranges between $[0, 1]$ where 1 represents total information sharing between label sets and 0 represents no information sharing between label sets. NMI is calculated by

$$\text{NMI} = \frac{I(Q; R)}{\sqrt{\mathbb{H}(Q)\mathbb{H}(R)}}$$

where Q and R are the empirical label distributions from the inferred and reference labels, respectively. The function $I(\cdot)$ is the mutual information between the inferred labels distribution and reference labels distributions; $\mathbb{H}(\cdot)$ represents the Shannon entropy of a given label distribution^{110,111}.

Adjusted Rand index (ARI). This metric captures the similarity between labels inferred by a method and the reference labels. It is based on the Rand index (RI) but corrects for the measurement's sensitivity to chance. ARI ranges between $[-1, 1]$ where 1 represents perfect agreement between label sets, 0 represents random agreement, and -1 represents perfect disagreement. ARI is calculated by

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

where n_{ij} , a_i , and b_j are values obtained from a contingency table, and $n = \sum_{ij} n_{ij}$ ¹¹⁰⁻¹¹².

Metrics to evaluate marker gene selection in simulations. In the simulation studies, we evaluated the accuracy of marker gene detection of NCLUSION and competing methods by treating the task as a classification problem. In order to do so, we defined the confusion matrix defined below.

		(Inferred Label)		Total
		Marker gene	Non-marker gene	
(True Label)	Marker gene	TP	FN	b_1
	Non-marker gene	FP	TN	b_2
Total		a_1	a_2	n

Table 1. Confusion matrix showing true and inferred marker gene labels.

Here, TP represents the number of correctly identified marker genes (true positives), FN represents the number of incorrectly identified marker genes (false negatives), TN represents the number of correctly identified non-marker genes (true negatives), and FP represents the number of incorrectly identified non-

marker genes (false positives). In the table above, we let a_1 and a_2 be the total number of genes inferred as markers and the total number of genes not inferred as markers, respectively. Likewise, we let b_1 and b_2 be the total number of genes that are truly markers and non-markers, respectively. It follows that the total number of genes is defined as $n = a_1 + b_1 + a_2 + b_2$. From here, we can compute the following metrics.

- **True positive rate (TPR; also referred to as power)** captures the proportion of correctly identified marker genes using a given method. It is defined as $\text{TPR} = \text{TP}/(\text{TP} + \text{FP})$.
- **False discovery rate (FDR)** details the proportion of all identified marker genes that are actually non-marker genes. It is defined as $\text{FDR} = \text{FP}/(\text{TP} + \text{FP})$.
- **False positive rate (FPR)** captures the proportion of non-marker genes that will be incorrectly labeled marker genes. It is defined as $\text{FPR} = \text{FP}/(\text{TN} + \text{FP})$.

Note that the false positive rate can also be computed as $\text{FPR} = 1 - \text{Specificity}$.

Normalized module expression. Genes with significantly adjusted PIPs in Eq. (5), positive ESS in Eq. (6), and significant SSMD in Eq. (7) were used to generate modules (i.e., a collection of marker genes) for each cluster. We calculated a score for each cluster to assess the exclusivity of expression within each module. This was done using the `score_genes` function in `scanpy` (v.1.10.4). The violin plots were generated using the `violinplot` function in `matplotlib` (v.3.10.0).

Gene set over-enrichment analysis. We also performed gene set enrichment analysis on each of the learned gene modules across clusters. This was done via an over-enrichment analysis within the `GSEAPy` package¹¹³ (v.1.1.5) in Python (v.3.11.0). This method uses a hypergeometric test to calculate the enrichment of genes in a supplied module with respect to the gene sets within an ontology. In this work, we use the ontology labeled `GO_Biological_Process_2025`^{63–65,114}, `Tabula_Sapiens`^{115,116}, `Azimuth_Cell_Types_2021`¹¹⁷, `KEGG_2021_Human`^{118–120}, and `Reactome_2022`^{121–128}. The gene sets in this particular ontology represent a combination of biological processes, pathways, and phenotypes. In this analysis, we use q -values to determine the enrichment of a given gene set with a significance threshold set to 0.05. The q -value is the analog of a p -value that has been corrected for testing multiple hypotheses (i.e., an adjusted P).

Software availability

An open-source software implementation of NCLUSION is available on GitHub at <https://github.com/microsoft/Nclusion.jl>. Guided tutorials and all code needed to reproduce the results and figures in this work can be found at <https://microsoft.github.io/Nclusion.jl/>.

Data availability

All of the datasets analyzed in this paper are publicly available. The PDAC dataset from Raghavan et al.⁷² can be accessed at https://singlecell.broadinstitute.org/single_cell/study/SCP1644/microenvironment-drives-cell-state-plasticity-and-drug-response-in-pancreatic-cancer#/. The AML data from van Galen et al.⁷³ can be found at https://www.dropbox.com/s/399x045zc57fiut/Seurat_AML.rds?dl=0. The BRAIN-LARGE dataset can be accessed at <https://www.10xgenomics.com/datasets/1-3-million-brain-cells-from-e-18-mice-2-standard-1-3-0>. The individual PBMC data from Zheng et al.⁵³ can be downloaded directly from <https://www.10xgenomics.com/resources/datasets>. Lastly, the immune cell atlas dataset can be accessed at https://cellgeni.cog.sanger.ac.uk/pan-immune/CountAdded_PIP_global_object_for_cellxgene.h5ad.

Acknowledgements

We thank members of the Crawford, Raghavan, and Shalek Labs for insightful comments on earlier versions of this manuscript. This research was conducted by using a combination of computational resources and services provided by Microsoft Research and the Center for Computation and Visualization at Brown University. This research was also supported in part by an Alfred P. Sloan Research Fellowship and a David & Lucile Packard Fellowship for Science and Engineering awarded to LC. CN was a trainee supported under the Brown University Predoctoral Training Program in Biological Data Science (NIH T32GM128596). SR is supported by NCI K08 award 1K08CA260442, the Claudia Adams Barr Program in Innovative Basic Cancer Research, and the Dana-Farber Cancer Institute Hale Family Center for Pancreatic Cancer Research. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of any of the funders.

688 **Author contributions**

689 CN, APA, and LC conceived the study and developed the methods. CN and MH developed the algorithm
690 and software. CN, MH, and MR led the analyses. SR, PSW, APA, and LC provided resources, supervised
691 the project, and conducted secondary analyses. All authors interpreted the results and wrote and revised
692 the manuscript.

693 **Declaration of interests**

694 SR holds equity in Amgen. PSW reports compensation for consulting/speaking from Engine Ventures
695 and AbbVie unrelated to this work. AKS reports compensation for consulting and/or scientific advisory
696 board membership from Honeycomb Biotechnologies, Cellarity, Ochre Bio, Relation Therapeutics, Fog
697 Pharma, Bio-Rad Laboratories, IntrECate Biotherapeutics, Passkey Therapeutics and Dahlia Biosciences
698 unrelated to this work. SR and PSW receive research funding from Microsoft. MH, NF, APA, and LC are
699 employees of Microsoft and own equity in Microsoft. All other authors have declared that no competing
700 interests exist.

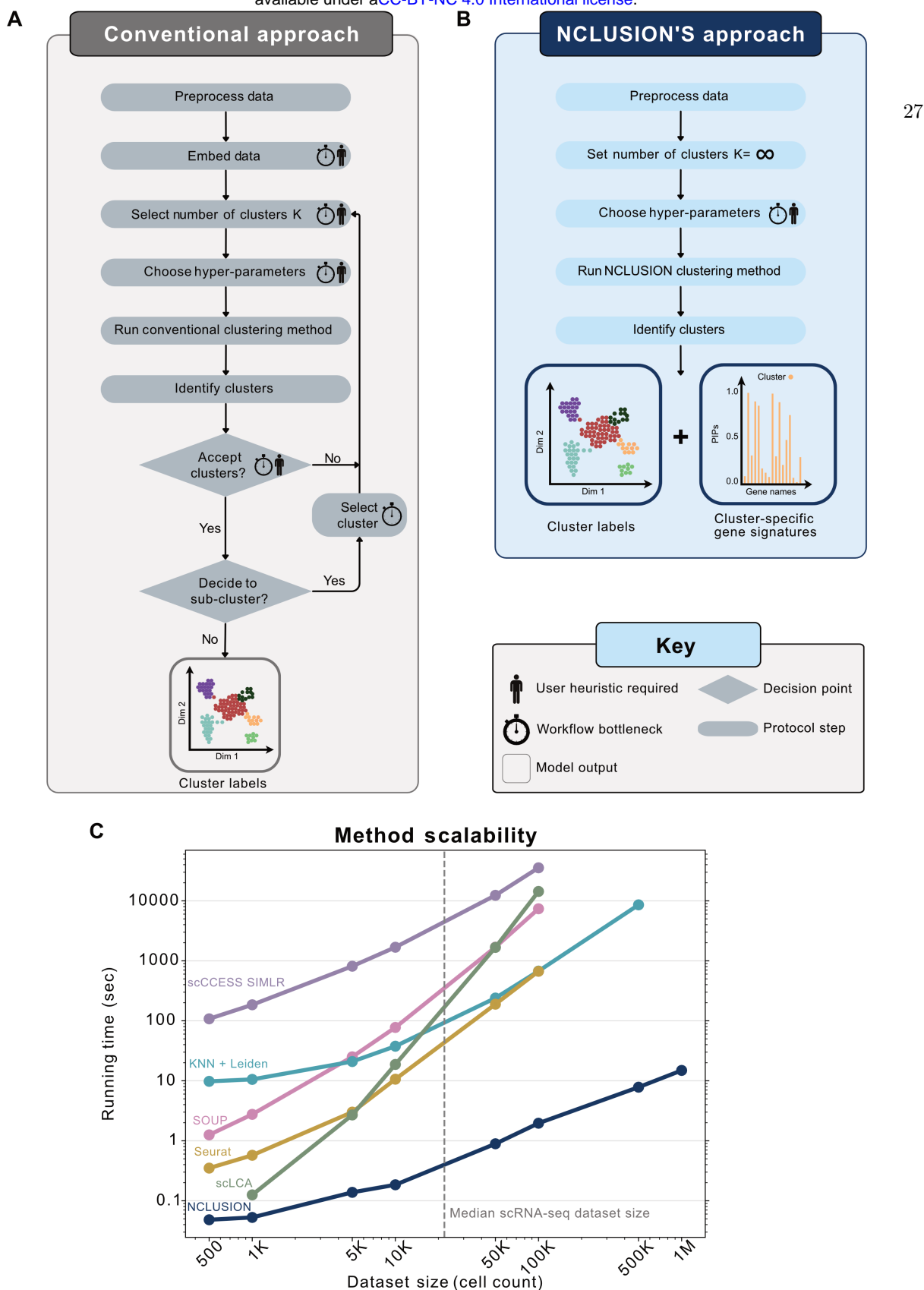


Fig 1. NCLUSION provides a scalable, unified workflow for both clustering and marker gene selection in single-cell analysis. (A) Conventional clustering algorithms require user heuristics and decision making steps that increase wall clock runtime (e.g., selection and human-in-the-loop refinement of the number of clusters K). (B) The nonparametric workflow of NCLUSION reduces the number of choices and heuristics that users have to make while also performing cluster-specific variable selection to identify top marker genes for downstream investigation. (C) Runtimes of NCLUSION and other baselines on the BRAIN-LARGE dataset with a fixed set of 720 genes and an increasing sample size ranging from $N = 500$ to 1 million cells.

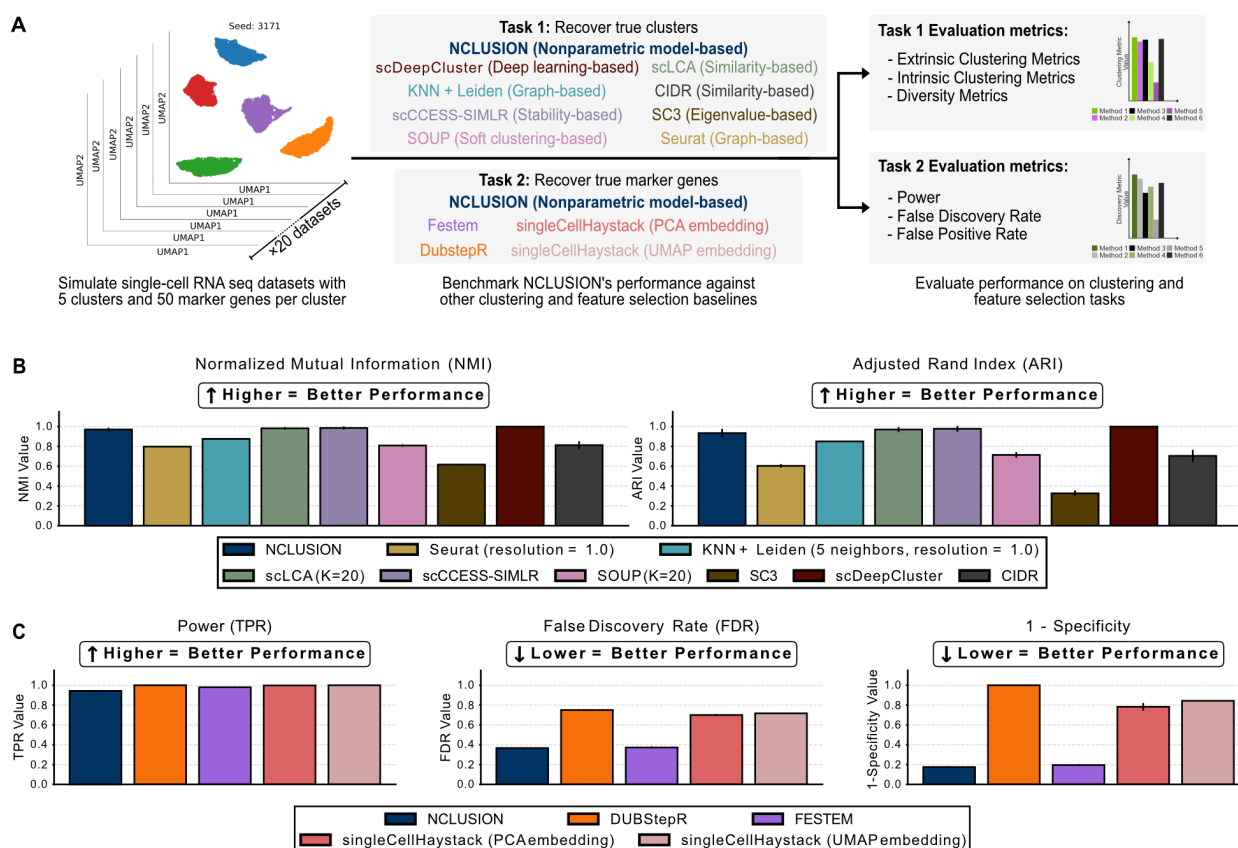


Fig 2. Comparing NCLUSION and competing algorithms on performing clustering and marker gene selection in a simulation study. Depicted are results for Scenario I where we evenly distributed all synthetically generated cells across five clusters and each cluster had a unique set of 50 marker genes. (A) Overview of the simulation framework used for evaluating the quality of clustering and marker gene selection for NCLUSION and each competing method. (B) Inferred cluster labels were compared to “true” annotations created during the simulation, where performance was measured according to (left) normalized mutual information (NMI) and (right) adjusted Rand index (ARI). (C) Assessment of marker gene selection was done on the global scale, where methods were evaluated on how well they could detect a “true” causal gene without taking cluster assignment into account. This was due to the limitation of competing methods not being able to identify cluster-specific genes. Evaluations were done by measuring the true positive rate (TPR; or power), false discovery rate (FDR), and false positive rate (FPR; computed as 1-Specificity) for each approach. Results for (B) and (C) are based on 20 simulations, with each bar plot representing the mean and the error bars covering a $\pm 95\%$ confidence interval.

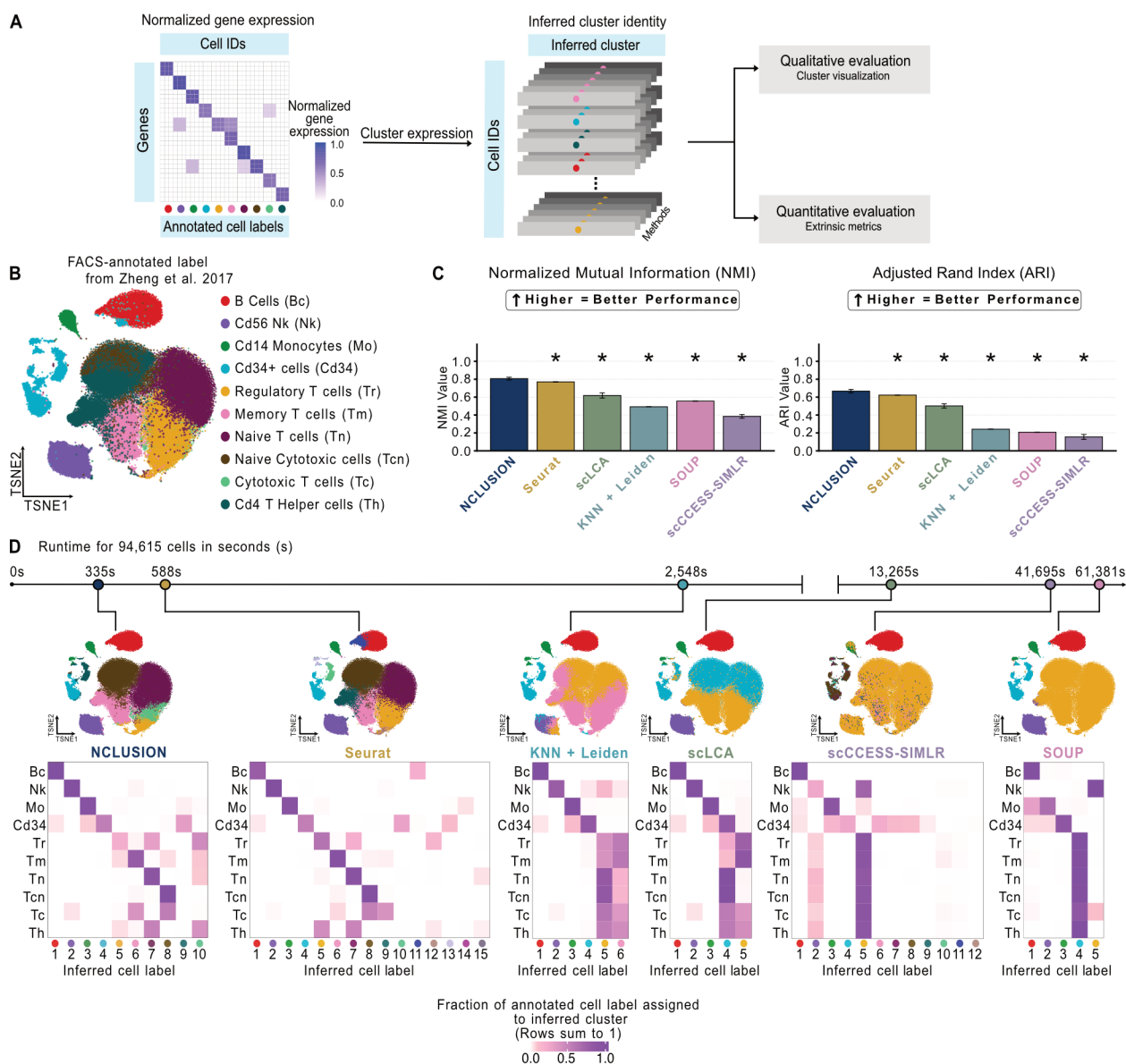


Fig 3. Clustering performance for NCLUSION and other baseline methods on the PBMC scRNA-seq dataset ($N = 94,615$ cells). (A) The framework used for evaluating the quality of clustering in each method. (B) Overview of FACS-based cell type annotations, visualized via t-distributed stochastic neighbor embedding (t-SNE), for the PBMC scRNA-seq dataset. These annotations serve as labels during the evaluation. (C) Assessment of the inferred cluster labels versus the experimental annotations, as quantified by two metrics: normalized mutual information (NMI) and adjusted Rand index (ARI) (for each method, we take five random 80% splits of the PBMC dataset; depicted in each bar plot is the mean \pm 95% confidence interval). Asterisks indicate that there is a statistically significant difference in performance between NCLUSION and a corresponding method (two-sided t-test $P < 0.05$). (D) Visualizing the structure of the inferred clusters across all baselines using t-SNEs and a contingency heat map showing the prevalence of each cell type within each cluster. Methods are ordered from fastest (left) to slowest (right) in terms of runtime. The same lower dimensional representation of the data is reused with relabeling of the plots according to the results from each clustering algorithm.

Fig 4. Evaluation of cluster-specific marker genes identified by NCLUSION on the PBMC dataset ($N = 94,615$ cells). (A) The framework used for assessing cluster-specific marker genes. (B) Embeddings of the experimental annotations for major cell types from the PBMC dataset compared to the clusters inferred by NCLUSION. (C) Heat maps of the adjusted posterior inclusion probabilities (PIPs) (left), effect size sign (ESS) (center), and strictly standardized mean difference (SSMD) (right) of significant genes in each cluster. Cluster-specific marker genes are selected as those that have a significant inclusion probability, are up-regulated in a given cluster, and have a large effect size magnitude such that $PIP \geq 0.5$, $ESS = +$, and $|SSMD(j; k)| \geq S^*(j; k)$, respectively. Here $S^*(j; k)$ is a threshold set to preserve a false positive rate of 0.05. (D) Highlighted location on t-SNEs of NCLUSION-inferred clusters that contain predominantly one cell type. (E) Violin plots comparing the normalized expression of cluster-specific marker genes in each of the inferred clusters. (F) Scatter plot comparing the marker genes identified using *post hoc* differential expression analysis with Seurat (yellow) versus the variable selection approach with NCLUSION (blue). Yellow points have $PIP \geq 0.5$ and $ESS = +$, while purple points have $PIP \geq 0.5$ and $ESS = -$, respectively. The vertical dashed line marks the median probability criterion³⁸, and the horizontal dashed line marks the Bonferroni-corrected threshold for significant q -values (i.e., an adjusted P). Genes in the top right quadrant are identified by both methods. (G) Scatter plot comparing gene ontology (GO) pathway enrichment analyses using cluster-specific marker genes from Seurat versus NCLUSION. The horizontal and vertical lines correspond to significant q -values being below 0.05. Pathways in the top right quadrant are selected by both approaches (red), while elements in the bottom right and top left quadrants are uniquely identified by NCLUSION (blue) and Seurat (orange), respectively. (H) Highlight of select top GO pathway enrichment analysis for the marker genes identified by NCLUSION. Plotted on the x-axis are the negative log-transformed q -values for each GO gene set. Gene sets with a q -value below 0.05 are deemed to be significant.

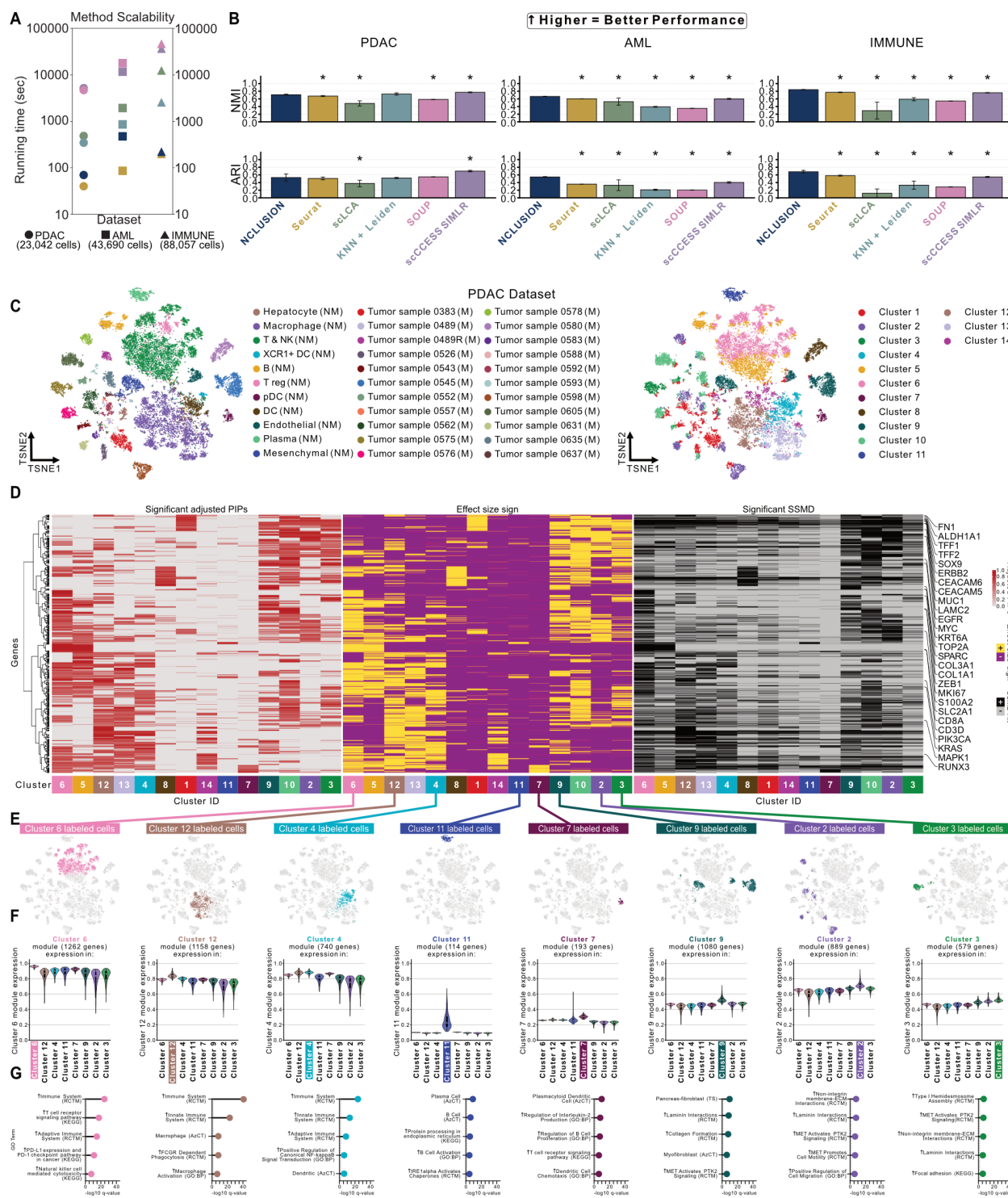


Fig 5. (Continued on the following page).

Fig 5. Scalability and generalizability of NCLUSION across diverse datasets. NCLUSION and baselines were applied to the following scRNA-seq datasets: PDAC ($N = 23,042$ cells)⁷², AML ($N = 43,690$ cells)⁷³, and IMMUNE ($N = 88,057$ cells)⁷⁴. **(A)** Runtimes for all methods when applied to each dataset. **(B)** Assessment of the inferred cluster labels from each method versus cell type annotations from the original studies. Evaluation is quantified by normalized mutual information (NMI) and adjusted Rand index (ARI). Asterisks indicate that there is a statistically significant difference in performance between NCLUSION and a corresponding method (two-sided t-test $P < 0.05$). Panels **(C)-(F)** depict results from running NCLUSION on the PDAC dataset. **(C)** Shown is a t-SNE visualization of the PDAC scRNA-seq dataset, annotated by the cell type labels from the PDAC study (top) compared to the clusters inferred by NCLUSION (bottom), where the “NM” labels indicate non-malignant cells and the “M” labels indicate malignant cells. **(D)** Heat maps of the adjusted posterior inclusion probabilities (PIPs) (left), effect size sign (ESS) (center), and strictly standardized mean difference (SSMD) (right) of the significant genes in each cluster. **(E)** Highlighted location on t-SNEs of NCLUSION-inferred clusters that contain predominantly one cell type. **(F)** Violin plots comparing the normalized expression of cluster-specific marker genes across clusters. **(G)** Gene ontology (GO) pathway enrichment analysis for the marker genes identified for each cluster. Gene sets with a q -value below 0.05 are deemed to be significant.

References

1. Zhen Miao, Benjamin D. Humphreys, Andrew P. McMahon, and Junhyong Kim. Multi-omics integration in the age of million single-cell data. *Nature Reviews Nephrology*, 17(11):710–724, November 2021. ISSN 1759-5061, 1759-507X. doi: 10.1038/s41581-021-00463-x.
2. F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. Scanpy: large-scale single-cell gene expression data analysis. *Genome Biology*, 19(1):15, February 2018. ISSN 1474-760X. doi: 10.1186/s13059-017-1382-0.
3. Minzhe Guo, Hui Wang, S. Steven Potter, Jeffrey A. Whitsett, and Yan Xu. Sincera: A pipeline for single-cell rna-seq profiling analysis. *PLOS Computational Biology*, 11(11):e1004575, November 2015. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1004575.
4. Aaron T.L. Lun, Davis J. McCarthy, and John C. Marioni. A step-by-step workflow for low-level analysis of single-cell rna-seq data with bioconductor. *F1000Research*, 5:2122, October 2016. ISSN 2046-1402. doi: 10.12688/f1000research.9501.2.
5. Ashraful Haque, Jessica Engel, Sarah A. Teichmann, and Tapio Lönnberg. A practical guide to single-cell rna-sequencing for biomedical research and clinical applications. *Genome Medicine*, 9(1):75, August 2017. ISSN 1756-994X. doi: 10.1186/s13073-017-0467-4.
6. Oliver Stegle, Sarah A. Teichmann, and John C. Marioni. Computational and analytical challenges in single-cell transcriptomics. *Nature Reviews Genetics*, 16(33):133–145, March 2015. ISSN 1471-0064. doi: 10.1038/nrg3833.
7. Vladimir Yu Kiselev, Tallulah S. Andrews, and Martin Hemberg. Challenges in unsupervised clustering of single-cell rna-seq data. *Nature Reviews Genetics*, 20(55):273–282, May 2019. ISSN 1471-0064. doi: 10.1038/s41576-018-0088-9.
8. Robert A. Amezquita, Aaron T. L. Lun, Etienne Becht, Vince J. Carey, Lindsay N. Carpp, Ludwig Geistlinger, Federico Marini, Kevin Rue-Albrecht, Davide Risso, Charlotte Soneson, Levi Waldron, Hervé Pagès, Mike L. Smith, Wolfgang Huber, Martin Morgan, Raphael Gottardo, and Stephanie C. Hicks. Orchestrating single-cell analysis with bioconductor. *Nature Methods*, 17(22):137–145, February 2020. ISSN 1548-7105. doi: 10.1038/s41592-019-0654-x.

9. Rahul Satija, Jeffrey A Farrell, David Gennert, Alexander F Schier, and Aviv Regev. Spatial reconstruction of single-cell gene expression data. Nature Biotechnology, 33:495–502, 2015. doi: 10.1038/nbt.3192. URL <https://doi.org/10.1038/nbt.3192>.
10. Changde Cheng, John Easton, Celeste Rosencrance, Yan Li, Bensheng Ju, Justin Williams, Heather L Mulder, Yakun Pang, Wenan Chen, and Xiang Chen. Latent cellular analysis robustly reveals subtle diversity in large-scale single-cell rna-seq data. Nucleic Acids Research, 47(22):e143, Dec 2019. ISSN 0305-1048. doi: 10.1093/nar/gkz826.
11. Bo Wang, Junjie Zhu, Emma Pierson, Daniele Ramazzotti, and Serafim Batzoglou. Visualization and analysis of single-cell rna-seq data by kernel-based similarity learning. Nature Methods, 14(44):414–416, Apr 2017. ISSN 1548-7105. doi: 10.1038/nmeth.4207.
12. Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. Journal of Machine Learning Research, 9(86):2579–2605, 2008. ISSN 1533-7928.
13. Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction, 2020. URL <http://arxiv.org/abs/1802.03426>. arXiv:1802.03426 [cs, stat].
14. Romain Lopez, Jeffrey Regier, Michael B. Cole, Michael I. Jordan, and Nir Yosef. Deep generative modeling for single-cell transcriptomics. Nature Methods, 15(12):1053–1058, Dec 2018. ISSN 1548-7091, 1548-7105. doi: 10.1038/s41592-018-0229-2.
15. Vladimir Yu Kiselev, Kristina Kirschner, Michael T. Schaub, Tallulah Andrews, Andrew Yiu, Tamir Chandra, Kedar N. Natarajan, Wolf Reik, Mauricio Barahona, Anthony R. Green, and Martin Hemberg. Sc3: consensus clustering of single-cell rna-seq data. Nature Methods, 14(55):483–486, May 2017. ISSN 1548-7105. doi: 10.1038/nmeth.4236.
16. Dominic Grün, Anna Lyubimova, Lennart Kester, Kay Wiebrands, Onur Basak, Nobuo Sasaki, Hans Clevers, and Alexander van Oudenaarden. Single-cell messenger rna sequencing reveals rare intestinal cell types. Nature, 525(75687568):251–255, September 2015. ISSN 1476-4687. doi: 10.1038/nature14966.

- 754 17. Bosiljka Tasic, Vilas Menon, Thuc Nghi Nguyen, Tae Kyung Kim, Tim Jarsky, Zizhen Yao, Boaz
755 Levi, Lucas T. Gray, Staci A. Sorensen, Tim Dolbeare, Darren Bertagnolli, Jeff Goldy, Nadiya
756 Shapovalova, Sheana Parry, Changkyu Lee, Kimberly Smith, Amy Bernard, Linda Madisen, Su-
757 san M. Sunkin, Michael Hawrylycz, Christof Koch, and Hongkui Zeng. Adult mouse cortical cell
758 taxonomy revealed by single cell transcriptomics. Nature Neuroscience, 19(22):335–346, February
759 2016. ISSN 1546-1726. doi: 10.1038/nn.4216.
- 760 18. Justina žurauskienė and Christopher Yau. pcareduce: hierarchical clustering of single cell
761 transcriptional profiles. BMC Bioinformatics, 17(1):140, March 2016. ISSN 1471-2105. doi:
762 10.1186/s12859-016-0984-y.
- 763 19. Amit Zeisel, Ana B. Muñoz-Manchado, Simone Codeluppi, Peter Lönnerberg, Gioele La Manno,
764 Anna Juréus, Sueli Marques, Hermany Munguba, Liqun He, Christer Betsholtz, Charlotte Rolny,
765 Gonçalo Castelo-Branco, Jens Hjerling-Leffler, and Sten Linnarsson. Cell types in the mouse
766 cortex and hippocampus revealed by single-cell rna-seq. Science, 347(6226):1138–1142, March
767 2015. doi: 10.1126/science.aaa1934.
- 768 20. Malte D Luecken and Fabian J Theis. Current best practices in single-cell rna-seq analysis: a
769 tutorial. Molecular Systems Biology, 15(6):e8746, Jun 2019. ISSN 1744-4292. doi: 10.15252/msb
770 .20188746.
- 771 21. David Lähnemann, Johannes Köster, Ewa Szczurek, Davis J. McCarthy, Stephanie C. Hicks,
772 Mark D. Robinson, Catalina A. Vallejos, Kieran R. Campbell, Niko Beerenwinkel, Ahmed Mah-
773 fouz, Luca Pinello, Pavel Skums, Alexandros Stamatakis, Camille Stephan-Otto Attolini, Samuel
774 Aparicio, Jasmijn Baaijens, Marleen Balvert, Buys de Barbanson, Antonio Cappuccio, Gia-
775 como Corleone, Bas E. Dutilh, Maria Florescu, Victor Guryev, Rens Holmer, Katharina Jahn,
776 Tamar Jessurun Lobo, Emma M. Keizer, Indu Khatri, Szymon M. Kielbasa, Jan O. Korbel,
777 Alexey M. Kozlov, Tzu-Hao Kuo, Boudewijn P.F. Lelieveldt, Ion I. Mandoiu, John C. Marioni,
778 Tobias Marschall, Felix Mölder, Amir Niknejad, Lukasz Raczkowski, Marcel Reinders, Jeroen de
779 Ridder, Antoine-Emmanuel Saliba, Antonios Somarakis, Oliver Stegle, Fabian J. Theis, Huan
780 Yang, Alex Zelikovsky, Alice C. McHardy, Benjamin J. Raphael, Sohrab P. Shah, and Alexander
781 Schönhuth. Eleven grand challenges in single-cell data science. Genome Biology, 21(1):31, Feb
782 2020. ISSN 1474-760X. doi: 10.1186/s13059-020-1926-6.

- 783 22. Lucy L. Gao, Jacob Bien, and Daniela Witten. Selective inference for hierarchical clustering.
784 Journal of the American Statistical Association, page 2116331, October 2022. ISSN 0162-1459,
785 1537-274X. doi: 10.1080/01621459.2022.2116331. Web of Science ID: WOS:000866074100001.
- 786 23. Anna Neufeld, Lucy L. Gao, Joshua Popp, Alexis Battle, and Daniela Witten. Inference af-
787 ter latent variable estimation for single-cell rna sequencing data. Biostatistics, December 2022.
788 ISSN 1465-4644, 1468-4357. doi: 10.1093/biostatistics/kxac047. Web of Science ID:
789 WOS:000896827000001.
- 790 24. Isabella N. Grabski, Kelly Street, and Rafael A. Irizarry. Significance analysis for clustering with
791 single-cell rna-sequencing data. Nature Methods, 20(88):1196–1202, Aug 2023. ISSN 1548-7105.
792 doi: 10.1038/s41592-023-01933-9.
- 793 25. Alexis Vandenberg and Diego Diez. A clustering-independent method for finding differentially
794 expressed genes in single-cell transcriptome data. Nature Communications, 11(11):4318, August
795 2020. ISSN 2041-1723. doi: 10.1038/s41467-020-17900-3.
- 796 26. Jesse M Zhang, Govinda M Kamath, and N Tse David. Valid post-clustering differential analysis
797 for single-cell rna-seq. Cell Systems, 9(4):383–392, 2019.
- 798 27. Alan DenAdel, Michelle L. Ramseier, Andrew W. Navia, Alex K. Shalek, Srivatsan Raghavan,
799 Peter S. Winter, Ava P. Amini, and Lorin Crawford. Artificial variables help to avoid over-
800 clustering in single-cell rna sequencing. The American Journal of Human Genetics, 0(0), March
801 2025. ISSN 0002-9297, 1537-6605. doi: 10.1016/j.ajhg.2025.02.014. URL [https://www.cell.c](https://www.cell.com/ajhg/abstract/S0002-9297(25)00061-8)
802 [om/ajhg/abstract/S0002-9297\(25\)00061-8](https://www.cell.com/ajhg/abstract/S0002-9297(25)00061-8).
- 803 28. Jesse M Zhang, Jue Fan, H Christina Fan, David Rosenfeld, and David N Tse. An interpretable
804 framework for clustering single-cell rna-seq datasets. BMC bioinformatics, 19(1):1–12, 2018.
- 805 29. Snehalika Lall, Sumanta Ray, and Sanghamitra Bandyopadhyay. Rgcop-a regularized copula
806 based method for gene selection in single-cell rna-seq data. PLOS Computational Biology, 17
807 (10):e1009464, 2021.
- 808 30. F William Townes, Stephanie C Hicks, Martin J Aryee, and Rafael A Irizarry. Feature selection

- and dimension reduction for single-cell rna-seq based on a multinomial model. Genome Biology,
20:1–16, 2019.
31. David M. Blei and Michael I. Jordan. Variational inference for dirichlet process mixtures. Bayesian
Analysis, 1(1):121–143, Mar 2006. ISSN 1936-0975, 1931-6690. doi: 10.1214/06-BA104.
32. Andrew Gelman, John B Carlin, Hal S Stern, David B Dunson, Aki Vehtari, and Donald B Rubin.
Bayesian Data Analysis, Third Edition. CRC Press, Hoboken, 2013. ISBN 978-1-4398-9820-8.
33. Sandhya Prabhakaran, Elham Azizi, Ambrose Carr, and Dana Pe’er. Dirichlet process mixture
model for correcting technical variation in single-cell gene expression data. JMLR Workshop and
Conference Proceedings, 48:1070–1079, 2016. ISSN 1938-7288.
34. Zhe Sun, Li Chen, Hongyi Xin, Yale Jiang, Qianhui Huang, Anthony R Cillo, Tracy Tabib, Jay K
Kolls, Tullia C Bruno, Robert Lafyatis, et al. A bayesian mixture model for clustering droplet-
based single-cell transcriptomic data from population studies. Nature communications, 10(1):
1649, 2019.
35. Tiehang Duan, José P Pinto, and Xiaohui Xie. Parallel clustering of single cell transcriptomic data
with split-merge sampling on dirichlet process mixtures. Bioinformatics, 35(6):953–961, 2019.
36. Zhe Sun, Ting Wang, Ke Deng, Xiao-Feng Wang, Robert Lafyatis, Ying Ding, Ming Hu, and Wei
Chen. Dimm-sc: a dirichlet mixture model for clustering droplet-based single cell transcriptomic
data. Bioinformatics, 34(1):139–146, 2018.
37. Nigatu A. Adossa, Kalle T. Rytkönen, and Laura L. Elo. Dirichlet process mixture models for
single-cell rna-seq clustering. Biology Open, 11(4):bio059001, April 2022. ISSN 2046-6390. doi:
10.1242/bio.059001. Citation Key: adossaDirichletProcessMixture2022.
38. Maria Maddalena Barbieri and James O. Berger. Optimal predictive model selection. The Annals
of Statistics, 32(3):870–897, Jun 2004. ISSN 0090-5364, 2168-8966. doi: 10.1214/00905360400000
0238.
39. Ping Zeng and Xiang Zhou. Non-parametric genetic prediction of complex traits with latent
dirichlet process regression models. Nature Communications, 8(1):456, Dec 2017. ISSN 2041-
1723. doi: 10.1038/s41467-017-00470-2.

- 836 40. Peter Carbonetto and Matthew Stephens. Scalable variational inference for bayesian variable
837 selection in regression, and its accuracy in genetic association studies. Bayesian Analysis, 7:
838 73–108, Mar 2012. ISSN 1936-0975, 1931-6690. doi: 10.1214/12-BA703.
- 839 41. Pinar Demetci, Wei Cheng, Gregory Darnell, Xiang Zhou, Sohini Ramachandran, and Lorin
840 Crawford. Multi-scale inference of genetic trait architecture using biologically annotated neural
841 networks. PLOS Genetics, 17(8):e1009754, 2021. ISSN 1553-7404. doi: 10.1371/journal.pgen.100
842 9754.
- 843 42. David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational inference: A review for statis-
844 ticians. Journal of the American Statistical Association, 112(518):859–877, April 2017. ISSN
845 0162-1459, 1537-274X. doi: 10.1080/01621459.2017.1285773. tex.ids: bleiVariationalInferenceRe-
846 view2017a arXiv: 1601.00670.
- 847 43. F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pret-
848 tenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Per-
849 rot, and E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning
850 Research, 12:2825–2830, 2011.
- 851 44. Lingxue Zhu, Jing Lei, Lambertus Klei, Bernie Devlin, and Kathryn Roeder. Semisoft clustering
852 of single-cell data. Proceedings of the National Academy of Sciences, 116(2):466–471, Jan 2019.
853 doi: 10.1073/pnas.1817715116.
- 854 45. 10x Genomics. Support: single cell gene expression datasets, 2023. URL [https://www.10xgen-](https://www.10xgenomics.com/resources/datasets)
855 [omics.com/resources/datasets](https://www.10xgenomics.com/resources/datasets).
- 856 46. Valentine Svensson, Eduardo da Veiga Beltrame, and Lior Pachter. A curated database reveals
857 trends in single-cell transcriptomics. Database: The Journal of Biological Databases and Curation,
858 2020:baaa073, Nov 2020. ISSN 1758-0463. doi: 10.1093/database/baaa073.
- 859 47. Dongyuan Song, Qingyang Wang, Guanao Yan, Tianyang Liu, Tianyi Sun, and Jingyi Jessica Li.
860 scdesign3 generates realistic in silico data for multimodal single-cell and spatial omics. Nature
861 Biotechnology, page 1–6, May 2023. ISSN 1546-1696. doi: 10.1038/s41587-023-01772-1.

- 862 48. Peijie Lin, Michael Troup, and Joshua W. K. Ho. Cidr: Ultrafast and accurate clustering through
863 imputation for single-cell rna-seq data. Genome Biology, 18(1):59, March 2017. ISSN 1474-760X.
864 doi: 10.1186/s13059-017-1188-0.
- 865 49. Vladimir Yu Kiselev, Kristina Kirschner, Michael T. Schaub, Tallulah Andrews, Andrew Yiu,
866 Tamir Chandra, Kedar N. Natarajan, Wolf Reik, Mauricio Barahona, Anthony R. Green, and
867 Martin Hemberg. Sc3: consensus clustering of single-cell rna-seq data. Nature Methods, 14(55):
868 483–486, May 2017. ISSN 1548-7105. doi: 10.1038/nmeth.4236.
- 869 50. Tian Tian, Ji Wan, Qi Song, and Zhi Wei. Clustering single-cell rna-seq data with a model-based
870 deep learning approach. Nature Machine Intelligence, 1(44):191–198, April 2019. ISSN 2522-5839.
871 doi: 10.1038/s42256-019-0037-0.
- 872 51. Bobby Ranjan, Wenjie Sun, Jinyu Park, Kunal Mishra, Florian Schmidt, Ronald Xie, Fatemeh
873 Alipour, Vipul Singhal, Ignasius Joanito, Mohammad Amin Honardoost, Jacy Mei Yun Yong,
874 Ee Tzun Koh, Khai Pang Leong, Nirmala Arul Rayan, Michelle Gek Liang Lim, and Shyam
875 Prabhakar. Dubstepr is a scalable correlation-based feature selection method for accurately clus-
876 tering single-cell data. Nature Communications, 12(1):5849, October 2021. ISSN 2041-1723. doi:
877 10.1038/s41467-021-26085-2. Citation Key: ranjanDUBStepRScalableCorrelationbased2021.
- 878 52. Zihao Chen, Changhu Wang, Siyuan Huang, Yang Shi, and Ruibin Xi. Directly selecting cell-type
879 marker genes for single-cell clustering analyses. Cell Reports Methods, 4(7):100810, July 2024.
880 ISSN 26672375. doi: 10.1016/j.crmeth.2024.100810. Citation Key: chenDirectlySelectingCell-
881 type2024.
- 882 53. Grace X. Y. Zheng, Jessica M. Terry, Phillip Belgrader, Paul Ryvkin, Zachary W. Bent, Ryan
883 Wilson, Solongo B. Ziraldo, Tobias D. Wheeler, Geoff P. McDermott, Junjie Zhu, Mark T. Gre-
884 gory, Joe Shuga, Luz Montesclaros, Jason G. Underwood, Donald A. Masquelier, Stefanie Y.
885 Nishimura, Michael Schnall-Levin, Paul W. Wyatt, Christopher M. Hindson, Rajiv Bharadwaj,
886 Alexander Wong, Kevin D. Ness, Lan W. Beppu, H. Joachim Deeg, Christopher McFarland,
887 Keith R. Loeb, William J. Valente, Nolan G. Ericson, Emily A. Stevens, Jerald P. Radich, Tarjei S.
888 Mikkelsen, Benjamin J. Hindson, and Jason H. Bielas. Massively parallel digital transcriptional

- 889 profiling of single cells. Nature Communications, 8(11):14049, Jan 2017. ISSN 2041-1723. doi:
890 10.1038/ncomms14049.
- 891 54. Alexander HS Vargo and Anna C Gilbert. A rank-based marker selection method for high through-
892 put scRNA-seq data. BMC bioinformatics, 21(1):1–51, 2020.
- 893 55. Lieke Michielsen, Marcel J. T. Reinders, and Ahmed Mahfouz. Hierarchical progressive learning
894 of cell identities in single-cell data. Nature Communications, 12(11):2799, May 2021. ISSN 2041-
895 1723. doi: 10.1038/s41467-021-23196-8.
- 896 56. Daniel M. Baume, Michael A. Caligiuri, Thomas J. Manley, John F. Daley, and Jerome Ritz.
897 Differential expression of $cd8\alpha$ and $cd8\beta$ associated with mhc-restricted and non-mhc-restricted
898 cytolytic effector cells. Cellular Immunology, 131(2):352–365, December 1990. ISSN 0008-8749.
899 doi: 10.1016/0008-8749(90)90260-X.
- 900 57. Jenny Hendriks, Loes A. Gravestein, Kiki Tesselaar, René A. W. van Lier, Ton N. M. Schumacher,
901 and Jannie Borst. Cd27 is required for generation and long-term maintenance of t cell immunity.
902 Nature Immunology, 1(5):433–440, November 2000. ISSN 1529-2916. doi: 10.1038/80877.
- 903 58. Suzanne Norris, Derek G. Doherty, Clive Collins, Gerry McEntee, Oscar Traynor, John E. Hegarty,
904 and Cliona O’Farrelly. Natural t cells in the human liver: cytotoxic lymphocytes with dual t cell
905 and natural killer cell phenotype and function are phenotypically heterogeneous and include $\nu\alpha 24$ -
906 $j\alpha q$ and $\gamma\delta$ t cell receptor bearing cells. Human Immunology, 60(1):20–31, Jan 1999. ISSN
907 0198-8859. doi: 10.1016/S0198-8859(98)00098-6.
- 908 59. Simon N. Willis, Julie Tellier, Yang Liao, Stephanie Trezise, Amanda Light, Kristy O’Donnell,
909 Lee Ann Garrett-Sinha, Wei Shi, David M. Tarlinton, and Stephen L. Nutt. Environmental sensing
910 by mature b cells is controlled by the transcription factors pu.1 and spib. Nature Communications,
911 8(11):1426, November 2017. ISSN 2041-1723. doi: 10.1038/s41467-017-01605-1.
- 912 60. Jonathan Zuccolo, Lili Deng, Tammy Unruh, Ratna Sanyal, Jeremy Bau, Jan Storek, Douglas
913 Demetrick, Joanne Luider, Iwona Auer-Grzesiak, Adnan Mansoor, and Julie Deans. Expression
914 of $ms4a$ and $tmem176$ genes in human b lymphocytes. Frontiers in Immunology, 4, 2013. ISSN
915 1664-3224. URL <https://www.frontiersin.org/articles/10.3389/fimmu.2013.00195>.

- 916 61. Zheng Chen, Mincheng Yu, Jiuliang Yan, Lei Guo, Bo Zhang, Shuang Liu, Jin Lei, Wen-
917 tao Zhang, Binghai Zhou, Jie Gao, Zhangfu Yang, Xiaoqiang Li, Jian Zhou, Jia Fan, Qing-
918 hai Ye, Hui Li, Yongfeng Xu, and Yongsheng Xiao. Pnoc expressed by b cells in cholan-
919 giocarcinoma was survival related and lair2 could be a t cell exhaustion biomarker in tu-
920 mor microenvironment: Characterization of immune microenvironment combining single-cell
921 and bulk sequencing technology. *Frontiers in Immunology*, 12, 2021. ISSN 1664-3224. URL
922 <https://www.frontiersin.org/articles/10.3389/fimmu.2021.647209>.
- 923 62. Elza Evren, Emma Ringqvist, Kumar Parijat Tripathi, Natalie Sleiers, Inés C6 Rives, Arlisa Al-
924 isjahbana, Yu Gao, Dhifaf Sarhan, Tor Halle, Chiara Sorini, Rico Lepzien, Nicole Marquardt,
925 Jakob Micha6lsson, Anna Smed-S6rensen, Johan Botling, Mikael C.I. Karlsson, Eduardo J. Vil-
926 lablanca, and Tim Willinger. Distinct developmental pathways from blood monocytes generate
927 human lung macrophage diversity. *Immunity*, 54(2):259–275.e7, February 2021. ISSN 10747613.
928 doi: 10.1016/j.immuni.2020.12.003.
- 929 63. Zhuoqing Fang, Xinyuan Liu, and Gary Peltz. Gseapy: a comprehensive package for performing
930 gene set enrichment analysis in python. *Bioinformatics*, 39(1):btac757, January 2023. ISSN 1367-
931 4811. doi: 10.1093/bioinformatics/btac757.
- 932 64. Steffen Durinck, Yves Moreau, Arek Kasprzyk, Sean Davis, Bart De Moor, Alvis Brazma, and
933 Wolfgang Huber. Biomart and bioconductor: a powerful link between biological databases and
934 microarray data analysis. *Bioinformatics*, 21(16):3439–3440, August 2005. ISSN 1367-4803. doi:
935 10.1093/bioinformatics/bti525.
- 936 65. Edward Y. Chen, Christopher M. Tan, Yan Kou, Qiaonan Duan, Zichen Wang, Gabriela Vaz
937 Meirelles, Neil R. Clark, and Avi Ma’ayan. Enrichr: interactive and collaborative html5 gene
938 list enrichment analysis tool. *BMC Bioinformatics*, 14:128, April 2013. ISSN 1471-2105. doi:
939 10.1186/1471-2105-14-128.
- 940 66. Michael Ashburner, Catherine A. Ball, Judith A. Blake, David Botstein, Heather Butler,
941 J. Michael Cherry, Allan P. Davis, Kara Dolinski, Selina S. Dwight, Janan T. Eppig, Midori A.
942 Harris, David P. Hill, Laurie Issel-Tarver, Andrew Kasarskis, Suzanna Lewis, John C. Matese,
943 Joel E. Richardson, Martin Ringwald, Gerald M. Rubin, and Gavin Sherlock. Gene ontology:

- 944 tool for the unification of biology. Nature Genetics, 25(11):25–29, May 2000. ISSN 1546-1718.
945 doi: 10.1038/75556.
- 946 67. Seth Carbon and Chris Mungall. Gene ontology data archive, Jul 2023. URL [https://zenodo](https://zenodo.org/record/8200914)
947 [.org/record/8200914](https://zenodo.org/record/8200914).
- 948 68. Eric Vivier, Elena Tomasello, Myriam Baratin, Thierry Walzer, and Sophie Ugolini. Functions
949 of natural killer cells. Nature Immunology, 9(55):503–510, May 2008. ISSN 1529-2916. doi:
950 10.1038/ni1582.
- 951 69. Sourav Paul and Girdhari Lal. The molecular mechanism of natural killer cells function and its
952 importance in cancer immunotherapy. Frontiers in Immunology, 8, 2017. ISSN 1664-3224. URL
953 <https://www.frontiersin.org/articles/10.3389/fimmu.2017.01124>.
- 954 70. Baptiste N. Jaeger, Jean Donadieu, Céline Cognet, Claire Bernat, Diana Ordoñez-Rueda, Vincent
955 Barlogis, Nizar Mahlaoui, Aurore Fenis, Emilie Narni-Mancinelli, Blandine Beaupain, Christine
956 Bellanné-Chantelot, Marc Bajénoff, Bernard Malissen, Marie Malissen, Eric Vivier, and Sophie
957 Ugolini. Neutrophil depletion impairs natural killer cell maturation, function, and homeostasis.
958 The Journal of Experimental Medicine, 209(3):565–580, March 2012. ISSN 0022-1007. doi: 10.1
959 084/jem.20111908.
- 960 71. Kun Jiang, Bin Zhong, Danielle L. Gilvary, Brian C. Corliss, Elizabeth Hong-Geller, Sheng Wei,
961 and Julie Y. Djeu. Pivotal role of phosphoinositide-3 kinase in regulation of cytotoxicity in
962 natural killer cells. Nature Immunology, 1(55):419–425, November 2000. ISSN 1529-2916. doi:
963 10.1038/80859.
- 964 72. Srivatsan Raghavan, Peter S. Winter, Andrew W. Navia, Hannah L. Williams, Alan DenAdel,
965 Kristen E. Lowder, Jennyfer Galvez-Reyes, Radha L. Kalekar, Nolawit Mulugeta, Kevin S. Kap-
966 ner, Manisha S. Raghavan, Ashir A. Borah, Nuo Liu, Sara A. Väyrynen, Andressa Dias Costa,
967 Raymond W. S. Ng, Junning Wang, Emma K. Hill, Dorisanne Y. Ragon, Lauren K. Brais, Alex M.
968 Jaeger, Liam F. Spurr, Yvonne Y. Li, Andrew D. Cherniack, Matthew A. Booker, Elizabeth F.
969 Cohen, Michael Y. Tolstorukov, Isaac Wakiro, Asaf Rotem, Bruce E. Johnson, James M. McFar-
970 land, Ewa T. Sicinska, Tyler E. Jacks, Ryan J. Sullivan, Geoffrey I. Shapiro, Thomas E. Clancy,
971 Kimberly Perez, Douglas A. Robinson, Kimmie Ng, James M. Cleary, Lorin Crawford, Scott R.

- Manalis, Jonathan A. Nowak, Brian M. Wolpin, William C. Hahn, Andrew J. Aguirre, and Alex K. Shalek. Microenvironment drives cell state, plasticity, and drug response in pancreatic cancer. Cell, 184(25):6119–6137.e26, Dec 2021. ISSN 0092-8674. doi: 10.1016/j.cell.2021.11.017.
73. Peter van Galen, Volker Hovestadt, Marc H. Wadsworth II, Travis K. Hughes, Gabriel K. Griffin, Sofia Battaglia, Julia A. Verga, Jason Stephansky, Timothy J. Pastika, Jennifer Lombardi Story, Geraldine S. Pinkus, Olga Pozdnyakova, Ilene Galinsky, Richard M. Stone, Timothy A. Graubert, Alex K. Shalek, Jon C. Aster, Andrew A. Lane, and Bradley E. Bernstein. Single-cell rna-seq reveals aml hierarchies relevant to disease progression and immunity. Cell, 176(6):1265–1281.e24, Mar 2019. ISSN 00928674. doi: 10.1016/j.cell.2019.01.031. tex.ids: vangalenSingleCellRNASeqReveals2019a.
74. C. Domínguez Conde, C. Xu, L. B. Jarvis, D. B. Rainbow, S. B. Wells, T. Gomes, S. K. Howlett, O. Suchanek, K. Polanski, H. W. King, L. Mamanova, N. Huang, P. A. Szabo, L. Richardson, L. Bolt, E. S. Fasouli, K. T. Mahbubani, M. Prete, L. Tuck, N. Richoz, Z. K. Tuong, L. Campos, H. S. Mousa, E. J. Needham, S. Pritchard, T. Li, R. Elmentaite, J. Park, E. Rahmani, D. Chen, D. K. Menon, O. A. Bayraktar, L. K. James, K. B. Meyer, N. Yosef, M. R. Clatworthy, P. A. Sims, D. L. Farber, K. Saeb-Parsy, J. L. Jones, and S. A. Teichmann. Cross-tissue immune cell analysis reveals tissue-specific features in humans. Science, 376(6594):eabl5197, May 2022. doi: 10.1126/science.abl5197.
75. Mazen Almeahmadi, Brian F. Flanagan, Naeem Khan, Suliman Alomar, and Stephen E. Christmas. Increased numbers and functional activity of cd56+ t cells in healthy cytomegalovirus positive subjects. Immunology, 142(2):258, April 2014. doi: 10.1111/imm.12250.
76. Emily R. Kansler and Ming O. Li. Innate lymphocytes—lineage, localization and timing of differentiation. Cellular and Molecular Immunology, 16(7):627–633, July 2019. ISSN 1672-7681. doi: 10.1038/s41423-019-0211-7.
77. L. Beviglia and R. H. Kramer. Hgf induces fak activation and integrin-mediated adhesion in mtln3 breast carcinoma cells. International Journal of Cancer, 83(5):640–649, November 1999. ISSN 0020-7136. doi: 10.1002/(sici)1097-0215(19991126)83:5(640::aid-ijc13)3.0.co;2-d.

- 999 78. Tobias Silzle, Gwendalyn J. Randolph, Marina Kreutz, and Leoni A. Kunz-Schughart. The fi-
1000 broblast: Sentinel cell and local immune modulator in tumor tissue. International Journal of
1001 Cancer, 108(2):173–180, 2004. ISSN 1097-0215. doi: 10.1002/ijc.11542.
- 1002 79. Serena Meraviglia, Paola Di Carlo, Diego Pampinella, Giuliana Guadagnino, Elena Lo Presti,
1003 Valentina Orlando, Giulia Marchetti, Francesco Dieli, and Consolato Sergi. T-cell subsets (tcm,
1004 tem, temra) and poly-functional immune response in patients with human immunodeficiency virus
1005 (hiv) infection and different t-cd4 cell response. Annals of Clinical & Laboratory Science, 49(4):
1006 519–528, 2019. ISSN 0091-7370, 1550-8080.
- 1007 80. Jason M. Schenkel and David Masopust. Tissue-resident memory t cells. Immunity, 41(6):886–897,
1008 December 2014. ISSN 1074-7613. doi: 10.1016/j.immuni.2014.12.007.
- 1009 81. Michael M. Opata, Samad A. Ibitokou, Victor H. Carpio, Karis M. Marshall, Brian E. Dillon,
1010 Jordan C. Carl, Kyle D. Wilson, Christine M. Arcari, and Robin Stephens. Protection by and
1011 maintenance of cd4 effector memory and effector t cell subsets in persistent malaria infection.
1012 PLOS Pathogens, 14(4):e1006960, 2018. ISSN 1553-7374. doi: 10.1371/journal.ppat.1006960.
- 1013 82. Shane Crotty. T follicular helper cell differentiation, function, and roles in disease. Immunity, 41
1014 (4):529–542, October 2014. ISSN 1074-7613. doi: 10.1016/j.immuni.2014.10.004.
- 1015 83. Dimitrios V. Vavoulis, Margherita Francescato, Peter Heutink, and Julian Gough. Dgeclust:
1016 differential expression analysis of clustered count data. Genome Biology, 16(1):39, February
1017 2015. ISSN 1465-6906. doi: 10.1186/s13059-015-0604-6.
- 1018 84. Charlotte Soneson and Mark D. Robinson. Bias, robustness and scalability in single-cell differential
1019 expression analysis. Nature Methods, 15(44):255–261, April 2018. ISSN 1548-7105. doi: 10.103
1020 8/nmeth.4612.
- 1021 85. Constantin Ahlmann-Eltze and Wolfgang Huber. Comparison of transformations for single-cell
1022 rna-seq data. Nature Methods, Apr 2023. ISSN 1548-7091, 1548-7105. doi: 10.1038/s41592-023
1023 -01814-1. URL <https://www.nature.com/articles/s41592-023-01814-1>.
- 1024 86. Xiang Zhu and Matthew Stephens. Large-scale genome-wide enrichment analyses identify new

- 1025 trait-associated genes and pathways across 31 human phenotypes. Nature Communications, 9(1):
1026 4361, 2018.
- 1027 87. Ying Ma, Shiquan Sun, Xuequn Shang, Evan T Keller, Mengjie Chen, and Xiang Zhou. Integrative
1028 differential expression and gene set enrichment analysis using summary statistics for scRNA-seq
1029 studies. Nature Communications, 11(1):1585, 2020.
- 1030 88. Jacob Bien and Robert J. Tibshirani. Sparse estimation of a covariance matrix. Biometrika, 98
1031 (4):807–820, December 2011. ISSN 0006-3444. doi: 10.1093/biomet/asr054.
- 1032 89. Youssef M Aboutaleb, Mazen Danaf, Yifei Xie, and Moshe E Ben-Akiva. Sparse covariance
1033 estimation in logit mixture models. The Econometrics Journal, 24(3):377–398, September 2021.
1034 ISSN 1368-4221. doi: 10.1093/ectj/utab008.
- 1035 90. David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statis-
1036 ticians. Journal of the American Statistical Association, 112(518):859–877, 2017.
- 1037 91. Ryan Giordano, Tamara Broderick, and Michael I Jordan. Covariances, robustness and variational
1038 bayes. Journal of Machine Learning Research, 19(51), 2018.
- 1039 92. Cheng Zhang, Babak Shahbaba, and Hongkai Zhao. Variational hamiltonian monte carlo via score
1040 matching. Bayesian Analysis, 13(2):485, 2018.
- 1041 93. Lukas M Weber, Arkajyoti Saha, Abhirup Datta, Kasper D Hansen, and Stephanie C Hicks.
1042 nnsVG for the scalable identification of spatially variable genes using nearest-neighbor gaussian
1043 processes. Nature Communications, 14(1):4059, 2023.
- 1044 94. Brian J Reich and Howard D Bondell. A spatial dirichlet process mixture model for clustering
1045 population genetics data. Biometrics, 67(2):381–390, 2011.
- 1046 95. Yongtao Guan and Matthew Stephens. Bayesian variable selection regression for genome-wide
1047 association studies and other large-scale problems. The Annals Applied Statistics, 5(3):1780–1815,
1048 2011. doi: 10.1214/11-AOAS455. URL [https://projecteuclid.org/443/euclid.aoas/1318](https://projecteuclid.org/443/euclid.aoas/1318514285)
1049 514285.
- 1050 96. Xiang Zhou, Peter Carbonetto, and Matthew Stephens. Polygenic modeling with Bayesian sparse
1051 linear mixed models. PLOS Genetics, 9(2):e1003264, 2013.

- 1052 97. Xiang Zhu and Matthew Stephens. Bayesian large-scale multiple regression with summary statis-
1053 tics from genome-wide association studies. The Annals Applied Statistics, 11(3):1561–1592, 2017.
1054 doi: 10.1214/17-AOAS1046. URL [https://projecteuclid.org:443/euclid.aoas/15071688](https://projecteuclid.org:443/euclid.aoas/1507168840)
1055 40.
- 1056 98. Michael C Hughes, Dae Il Kim, and Erik B Sudderth. Reliable and scalable variational inference
1057 for the hierarchical dirichlet process. Artificial Intelligence and Statistics, page 9, 2015.
- 1058 99. Jacob Cohen. Statistical Power Analysis for the Behavioral Sciences. Academic press, 2013.
- 1059 100. Xiaohua Douglas Zhang. A pair of new statistical parameters for quality control in rna interference
1060 high-throughput screening assays. Genomics, 89(4):552–561, April 2007. ISSN 0888-7543. doi:
1061 10.1016/j.ygeno.2006.12.014.
- 1062 101. Xiaohua Douglas Zhang. A new method with flexible and balanced control of false negatives
1063 and false positives for hit selection in rna interference high-throughput screening assays. SLAS
1064 Discovery, 12(5):645–655, August 2007. ISSN 24725552. doi: 10.1177/1087057107300645.
- 1065 102. Xiaohua Douglas Zhang. Strictly standardized mean difference, standardized mean difference and
1066 classical t-test for the comparison of two groups. Statistics in Biopharmaceutical Research, 2(2):
1067 292–299, May 2010. ISSN null. doi: 10.1198/sbr.2009.0074.
- 1068 103. Xiaohua Douglas Zhang. Optimal High-Throughput Screening: Practical Experimental Design
1069 and Data Analysis for Genome-Scale RNAi Research. Cambridge University Press, 1 edition,
1070 February 2011. ISBN 978-0-521-51771-3. doi: 10.1017/CBO9780511973888. URL [https:](https://www.cambridge.org/core/product/identifier/9780511973888/type/book)
1071 [//www.cambridge.org/core/product/identifier/9780511973888/type/book](https://www.cambridge.org/core/product/identifier/9780511973888/type/book).
- 1072 104. Xiaohua Douglas Zhang, Shane D. Marine, and Marc Ferrer. Error rates and powers in genome-
1073 scale rnai screens. SLAS Discovery, 14(3):230–238, March 2009. ISSN 2472-5552. doi: 10.1177/
1074 1087057109331475.
- 1075 105. Xiaohua Douglas Zhang, Raul Lacson, Ruojing Yang, Shane D. Marine, Alex McCampbell,
1076 Dawn M. Toolan, Tim R. Hare, Joleen Kajdas, Joel P. Berger, Daniel J. Holder, Joseph F.
1077 Heyse, and Marc Ferrer. The use of ssmd-based false discovery and false nondiscovery rates in

- genome-scale rnai screens. SLAS Discovery, 15(9):1123–1131, October 2010. ISSN 2472-5552. doi: 10.1177/1087057110381919.
106. Matthew P. Wand, John T. Ormerod, Simone A. Padoan, and Rudolf Fröhwrth. Mean field variational bayes for elaborate distributions. Bayesian Analysis, 6(4):847–900, 2011. doi: 10.1214/11-BA631. URL <https://doi.org/10.1214/11-BA631>.
107. Valentine Svensson. Droplet scrna-seq is not zero-inflated. Nature Biotechnology, 38(2):147–150, Feb 2020. ISSN 1546-1696. doi: 10.1038/s41587-019-0379-5.
108. Valentine Svensson. Reply to: Umi or not umi, that is the question for scrna-seq zero-inflation. Nature Biotechnology, 39(2):160–160, Feb 2021. ISSN 1546-1696. doi: 10.1038/s41587-020-00811-5.
109. F. Alexander Wolf, Philipp Angerer, and Fabian J. Theis. Scanpy: large-scale single-cell gene expression data analysis. Genome Biology, 19(1):15, Feb 2018. ISSN 1474-760X. doi: 10.1186/s13059-017-1382-0.
110. Lijia Yu, Yue Cao, Jean Y. H. Yang, and Pengyi Yang. Benchmarking clustering algorithms on estimating the number of cell types from single-cell rna-sequencing data. Genome Biology, 23(1):49, Feb 2022. ISSN 1474-760X. doi: 10.1186/s13059-022-02622-0.
111. Nguyen Xuan Vinh, Julien Epps, and James Bailey. Information theoretic measures for clusterings comparison: is a correction for chance necessary? In Proceedings of the 26th Annual International Conference on Machine Learning, pages 1073–1080, 2009.
112. Silke Wagner and Dorothea Wagner. Comparing clusterings - an overview, 2007. URL <https://publikationen.bibliothek.kit.edu/1000011477>.
113. Zhuoqing Fang, Xinyuan Liu, and Gary Peltz. Gseapy: a comprehensive package for performing gene set enrichment analysis in python. Bioinformatics, 39(1):btac757, Jan 2023. ISSN 1367-4811. doi: 10.1093/bioinformatics/btac757.
114. Suzi A Aleksander, James Balhoff, Seth Carbon, J Michael Cherry, Harold J Drabkin, Dustin Ebert, Marc Feuermann, Pascale Gaudet, Nomi L Harris, David P Hill, Raymond Lee, Huaiyu Mi, Sierra Moxon, Christopher J Mungall, Anushya Muruganugan, Tremayne Mushayahama, Paul W

- 1105 Sternberg, Paul D Thomas, Kimberly Van Auken, Jolene Ramsey, Deborah A Siegele, Rex L
1106 Chisholm, Petra Fey, Maria Cristina Aspromonte, Maria Victoria Nugnes, Federica Quaglia, Sil-
1107 vio Tosatto, Michelle Giglio, Suvarna Nadendla, Giulia Antonazzo, Helen Attrill, Gil dos Santos,
1108 Steven Marygold, Victor Strelets, Christopher J Tabone, Jim Thurmond, Pinglei Zhou, Saadul-
1109 lah H Ahmed, Praoparn Asanithong, Diana Luna Buitrago, Meltem N Erdol, Matthew C Gage,
1110 Mohamed Ali Kadhun, Kan Yan Chloe Li, Miao Long, Aleksandra Michalak, Angeline Pesala,
1111 Armalya Pritazahra, Shirin C C Saverimuttu, Renzhi Su, Kate E Thurlow, Ruth C Lovering, Colin
1112 Logie, Snezhana Oliferenko, Judith Blake, Karen Christie, Lori Corbani, Mary E Dolan, Harold J
1113 Drabkin, David P Hill, Li Ni, Dmitry Sitnikov, Cynthia Smith, Alayne Cuzick, James Seager,
1114 Laurel Cooper, Justin Elser, Pankaj Jaiswal, Parul Gupta, Pankaj Jaiswal, Sushma Naithani,
1115 Manuel Lera-Ramirez, Kim Rutherford, Valerie Wood, Jeffrey L De Pons, Melinda R Dwinell,
1116 G Thomas Hayman, Mary L Kaldunski, Anne E Kwitek, Stanley J F Lalederkind, Marek A
1117 Tutaj, Mahima VEDI, Shur-Jen Wang, Peter D'Eustachio, Lucila Aimo, Kristian Axelsen, Alan
1118 Bridge, Nevila Hyka-Nouspikel, Anne Morgat, Suzi A Aleksander, J Michael Cherry, Stacia R En-
1119 gel, Kalpana Karra, Stuart R Miyasato, Robert S Nash, Marek S Skrzypek, Shuai Weng, Edith D
1120 Wong, Erika Bakker, Tanya Z Berardini, Leonore Reiser, Andrea Auchincloss, Kristian Axelsen,
1121 Ghislaine Argoud-Puy, Marie-Claude Blatter, Emmanuel Boutet, Lionel Breuza, Alan Bridge,
1122 Cristina Casals-Casas, Elisabeth Coudert, Anne Estreicher, Maria Livia Famiglietti, Marc Feuer-
1123 mann, Arnaud Gos, Nadine Gruaz-Gumowski, Chantal Hulo, Nevila Hyka-Nouspikel, Florence
1124 Jungo, Philippe Le Mercier, Damien Lieberherr, Patrick Masson, Anne Morgat, Ivo Pedruzzi,
1125 Lucille Pourcel, Sylvain Poux, Catherine Rivoire, Shyamala Sundaram, Alex Bateman, Emily
1126 Bowler-Barnett, Hema Bye-A-Jee, Paul Denny, Alexandr Ignatchenko, Rizwan Ishtiaq, Antonia
1127 Lock, Yvonne Lussi, Michele Magrane, Maria J Martin, Sandra Orchard, Pedro Raposo, Elena
1128 Speretta, Nidhi Tyagi, Kate Warner, Rossana Zaru, Alexander D Diehl, Raymond Lee, Juancar-
1129 los Chan, Stavros Diamantakis, Daniela Raciti, Magdalena Zarowiecki, Malcolm Fisher, Christina
1130 James-Zorn, Virgilio Ponferrada, Aaron Zorn, Sridhar Ramachandran, Leyla Ruzicka, and Monte
1131 Westerfield. The gene ontology knowledgebase in 2023. *Genetics*, 224(1):iyad031, March 2023.
1132 ISSN 0016-6731. doi: 10.1093/genetics/iyad031.
- 1133 115. The Tabula Sapiens Consortium*, Robert C. Jones, Jim Karkanias, Mark A. Krasnow, An-
1134 gela Oliveira Pisco, Stephen R. Quake, Julia Salzman, Nir Yosef, Bryan Bulthaupt, Phillip Brown,

William Harper, Marisa Hemenez, Ravikumar Ponnusamy, Ahmad Salehi, Bhavani A. Sana-
gavarapu, Eileen Spallino, Ksenia A. Aaron, Waldo Concepcion, James M. Gardner, Burnett
Kelly, Nikole Neidlinger, Zifa Wang, Sheela Crasta, Saroja Kolluru, Maurizio Morri, Serena Y.
Tan, Kyle J. Travaglini, Chenling Xu, Marcela Alcántara-Hernández, Nicole Almanzar, Jane
Antony, Benjamin Beyersdorf, Deviana Burhan, Kruti Calcuttawala, Matthew M. Carter, Charles
K. F. Chan, Charles A. Chang, Stephen Chang, Alex Colville, Rebecca N. Culver, Ivana Cvijović,
Gaetano D'Amato, Camille Ezran, Francisco X. Galdos, Astrid Gillich, William R. Goodyer, Yan
Hang, Alyssa Hayashi, Sahar Houshdaran, Xianxi Huang, Juan C. Irwin, SoRi Jang, Julia Val-
Ive Juanico, Aaron M. Kershner, Soochi Kim, Bernhard Kiss, William Kong, Maya E. Kumar,
Angera H. Kuo, Rebecca Leylek, Baoxiang Li, Gabriel B. Loeb, Wan-Jin Lu, Sruthi Mantri,
Maxim Markovic, Patrick L. McAlpine, Antoine de Morree, Karim Mrouj, Shravani Mukherjee,
Tyler Muser, Patrick Neuhöfer, Thi D. Nguyen, Kimberly Perez, Ragini Phansalkar, Nazan Pu-
luca, Zhen Qi, Poorvi Rao, Hayley Raquer-McKay, Nicholas Schaum, Bronwyn Scott, Bobak
Seddighzadeh, Joe Segal, Sushmita Sen, Shaheen Sikandar, Sean P. Spencer, Lea C. Steffes,
Varun R. Subramaniam, Aditi Swarup, Michael Swift, Will Van Treuren, Emily Trimm, Stefan
Veizades, Sivakamasundari Vijayakumar, Kim Chi Vo, Sevahn K. Vorperian, Wanxin Wang, Han-
nah N. W. Weinstein, Juliane Winkler, Timothy T. H. Wu, Jamie Xie, Andrea R. Yung, Yue
Zhang, Angela M. Detweiler, Honey Mekonen, Norma F. Neff, Rene V. Sit, Michelle Tan, Jia
Yan, Gregory R. Bean, Vivek Charu, Erna Forgó, Brock A. Martin, Michael G. Ozawa, Oscar
Silva, Angus Toland, Venkata N. P. Vemuri, Shaked Afik, Kyle Awayan, Olga Borisovna Botvin-
nik, Ashley Byrne, Michelle Chen, Roozbeh Dehghannasiri, Adam Gayoso, Alejandro A. Grana-
dos, Qiqing Li, Gita Mahmoudabadi, Aaron McGeever, Julia Eve Olivieri, Madeline Park, Neha
Ravikumar, Geoff Stanley, Weilun Tan, Alexander J. Tarashansky, Rohan Vanheusden, Peter
Wang, Sheng Wang, Galen Xing, Les Dethlefsen, Po-Yi Ho, Shixuan Liu, Jonathan S. Maltz-
man, Ross J. Metzger, Koki Sasagawa, Rahul Sinha, Hanbing Song, Bruce Wang, Steven E.
Artandi, Philip A. Beachy, Michael F. Clarke, Linda C. Giudice, Franklin W. Huang, Ker-
wyn Casey Huang, Juliana Idoyaga, Seung K. Kim, Christin S. Kuo, Patricia Nguyen, Thomas A.
Rando, Kristy Red-Horse, Jeremy Reiter, David A. Relman, Justin L. Sonnenburg, Albert
Wu, Sean M. Wu, and Tony Wyss-Coray. The tabula sapiens: A multiple-organ, single-cell
transcriptomic atlas of humans. *Science*, May 2022. doi: 10.1126/science.abl4896. URL

- 1165 <https://www.science.org/doi/10.1126/science.abl4896>.
- 1166 116. Stephen R. Quake and The Tabula Sapiens Consortium. Tabula sapiens reveals transcription
1167 factor expression, senescence effects, and sex-specific features in cell types from 28 human organs
1168 and tissues. page 2024.12.03.626516, December 2024. doi: 10.1101/2024.12.03.626516. URL
1169 <https://www.biorxiv.org/content/10.1101/2024.12.03.626516v1>.
- 1170 117. Yuhao Hao, Stephanie Hao, Erica Andersen-Nissen, William M. Mauck, Shiwei Zheng, Andrew
1171 Butler, Maddie J. Lee, Aaron J. Wilk, Charlotte Darby, Michael Zager, Paul Hoffman, Mar-
1172 lon Stoeckius, Efthymia Papalexi, Eleni P. Mimitou, Jaison Jain, Avi Srivastava, Tim Stuart,
1173 Lamar M. Fleming, Bertrand Yeung, Angela J. Rogers, Juliana M. McElrath, Catherine A.
1174 Blish, Raphael Gottardo, Peter Smibert, and Rahul Satija. Integrated analysis of multimodal
1175 single-cell data. Cell, 184(13):3573–3587.e29, June 2021. ISSN 0092-8674, 1097-4172. doi:
1176 10.1016/j.cell.2021.04.048.
- 1177 118. M. Kanehisa and S. Goto. Kegg: kyoto encyclopedia of genes and genomes. Nucleic Acids
1178 Research, 28(1):27–30, January 2000. ISSN 0305-1048. doi: 10.1093/nar/28.1.27.
- 1179 119. Minoru Kanehisa. Toward understanding the origin and evolution of cellular organisms. Protein
1180 Science: A Publication of the Protein Society, 28(11):1947–1951, November 2019. ISSN 1469-
1181 896X. doi: 10.1002/pro.3715.
- 1182 120. Minoru Kanehisa, Miho Furumichi, Yoko Sato, Yuriko Matsuura, and Mari Ishiguro-Watanabe.
1183 Kegg: biological systems database as a model of the real world. Nucleic Acids Research, 53(D1):
1184 D672–D677, January 2025. ISSN 1362-4962. doi: 10.1093/nar/gkae909.
- 1185 121. Antonio Fabregat, Florian Korninger, Guilherme Viteri, Konstantinos Sidiropoulos, Pablo Marin-
1186 Garcia, Peipei Ping, Guanming Wu, Lincoln Stein, Peter D’Eustachio, and Henning Hermjakob.
1187 Reactome graph database: Efficient access to complex pathway data. PLoS computational biology,
1188 14(1):e1005968, January 2018. ISSN 1553-7358. doi: 10.1371/journal.pcbi.1005968.
- 1189 122. Antonio Fabregat, Konstantinos Sidiropoulos, Guilherme Viteri, Oscar Forner, Pablo Marin-
1190 Garcia, Vicente Arnau, Peter D’Eustachio, Lincoln Stein, and Henning Hermjakob. Reactome
1191 pathway analysis: a high-performance in-memory approach. BMC bioinformatics, 18(1):142,
1192 March 2017. ISSN 1471-2105. doi: 10.1186/s12859-017-1559-2.

- 1193 123. Antonio Fabregat, Konstantinos Sidiropoulos, Guilherme Viteri, Pablo Marin-Garcia, Peipei Ping,
1194 Lincoln Stein, Peter D'Eustachio, and Henning Hermjakob. Reactome diagram viewer: data struc-
1195 tures and strategies to boost performance. Bioinformatics (Oxford, England), 34(7):1208–1214,
1196 April 2018. ISSN 1367-4811. doi: 10.1093/bioinformatics/btx752.
- 1197 124. Johannes Griss, Guilherme Viteri, Konstantinos Sidiropoulos, Vy Nguyen, Antonio Fabregat,
1198 and Henning Hermjakob. Reactomegsa - efficient multi-omics comparative pathway analysis.
1199 Molecular & cellular proteomics: MCP, 19(12):2115–2125, December 2020. ISSN 1535-9484. doi:
1200 10.1074/mcp.TIR120.002155.
- 1201 125. Bijay Jassal, Lisa Matthews, Guilherme Viteri, Chuqiao Gong, Pascual Lorente, Antonio Fab-
1202 regat, Konstantinos Sidiropoulos, Justin Cook, Marc Gillespie, Robin Haw, Fred Loney, Bruce
1203 May, Marija Milacic, Karen Rothfels, Cristoffer Sevilla, Veronica Shamovsky, Solomon Shorser,
1204 Thawfeek Varusai, Joel Weiser, Guanming Wu, Lincoln Stein, Henning Hermjakob, and Peter
1205 D'Eustachio. The reactome pathway knowledgebase. Nucleic Acids Research, 48(D1):D498–D503,
1206 January 2020. ISSN 1362-4962. doi: 10.1093/nar/gkz1031.
- 1207 126. Marija Milacic, Deidre Beavers, Patrick Conley, Chuqiao Gong, Marc Gillespie, Johannes Griss,
1208 Robin Haw, Bijay Jassal, Lisa Matthews, Bruce May, Robert Petryszak, Eliot Ragueneau, Karen
1209 Rothfels, Cristoffer Sevilla, Veronica Shamovsky, Ralf Stephan, Krishna Tiwari, Thawfeek Varu-
1210 sai, Joel Weiser, Adam Wright, Guanming Wu, Lincoln Stein, Henning Hermjakob, and Pe-
1211 ter D'Eustachio. The reactome pathway knowledgebase 2024. Nucleic Acids Research, 52(D1):
1212 D672–D678, January 2024. ISSN 0305-1048. doi: 10.1093/nar/gkad1025.
- 1213 127. Konstantinos Sidiropoulos, Guilherme Viteri, Cristoffer Sevilla, Steve Jupe, Marissa Webber,
1214 Marija Orlic-Milacic, Bijay Jassal, Bruce May, Veronica Shamovsky, Corina Duenas, Karen
1215 Rothfels, Lisa Matthews, Heeyeon Song, Lincoln Stein, Robin Haw, Peter D'Eustachio, Peipei
1216 Ping, Henning Hermjakob, and Antonio Fabregat. Reactome enhanced pathway visualization.
1217 Bioinformatics (Oxford, England), 33(21):3461–3467, November 2017. ISSN 1367-4811. doi:
1218 10.1093/bioinformatics/btx441.
- 1219 128. Guanming Wu and Robin Haw. Functional interaction network construction and analysis for

1220 disease discovery. Methods in Molecular Biology (Clifton, N.J.), 1558:235–253, 2017. ISSN 1940-
1221 6029. doi: 10.1007/978-1-4939-6783-4_11.