1 **Information-theory-based benchmarking and feature selection algorithm improve cell type**

2 **annotation and reproducibility of single cell RNA-seq data analysis pipelines**

3 Ziyou Ren[1], Martin Gerlach[2], Hanyu Shi[2], GR Scott Budinger[1], Luís A. Nunes Amaral[1,2,3,4, 5#]

4

5 [1]Department of Medicine, Division of Pulmonary and Critical Care Medicine, Northwestern

6 University, Chicago, IL, 60611, USA.

7 [2]Department of Chemical and Biological Engineering, Northwestern University, Evanston, IL

8 60208, USA.

9 [3]Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL 60208,

10 USA

11 [4] Department of Physics and Astronomy, Northwestern University, Evanston, IL 60208, USA
12
13 [5] Department of Molecular Biosciences, Northwestern University, Evanston, IL 60208, USA
14
15
16
17 #Corresponding author:
18
19 Luís A. Nunes Amaral, PhD

20 Chemical & Biological Engineering

21 Northwestern University

22 2145 Sheridan Road (Room E136)

23 Evanston, IL 60208, US

24 Phone: (847) 491-7850

25 Email: amaral@northwestern.edu

26

1   **Abstract:**

2   Single cell RNA sequencing (scRNA-seq) data are now routinely generated in experimental

3   practice because of their promise to enable the quantitative study of biological processes at the

4   single cell level. However, cell type and cell state annotations remain an important

5   computational challenge in analyzing scRNA-seq data. Here, we report on the development of a

6   benchmark dataset where reference annotations are generated independently from transcriptomic

7   measurements. We used this benchmark to systematically investigate the impact on labelling

8   accuracy of different approaches to feature selection, of different clustering algorithms, and of

9   different sets of parameter values. We show that an approach grounded on information theory

10  can provide a general, reliable, and accurate process for discarding uninformative features and to

11  optimize cluster resolution in single cell RNA-seq data analysis.

12

1    **Introduction**

2    The application of single cell RNA sequencing (scRNA-seq) as a hypothesis generation

3    step is becoming increasingly popular in biology and medical research because of the promise to

4    (i) efficiently identify novel cell types from transcriptomic data (Patel et al. 2014; Villani et al.

5    2017);  (ii) identify heterogeneity within defined cell populations during homeostasis and disease

6    (Reyfman et al. 2019; Lawson et al. 2015); (iii) follow cellular lineage trajectories (Treutlein et

7    al. 2014; Rizvi et al. 2017; Collin et al. 2019); or (iv) identify novel genetic markers for disease

8    (Miyamoto et al. 2015; Min et al. 2015; Lindström et al. 2019).

9    A key step in scRNA-seq analysis is cell type identification using computational

10    clustering algorithms. While supervised clustering approaches, i.e. identifying cell types based

11    on reference dataset, have been developed to use the data produced by large international efforts

12    such as the Human Primary Cell Atlas (HPCA) and the Encyclopedia of DNA Elements

13    (ENCODE) (Aran et al. 2019; Mabbott et al. 2013; Davis et al. 2018), many studies still rely on

14    unsupervised clustering approaches i.e. identifying cell types based on expression patterns.

15    Among all clustering algorithms, Louvain method (Blondel et al. 2008) is the most popular

16    method in scRNA-seq because of its high efficiency in analyzing high dimensional data. It is the

17    default clustering algorithm in one widely used analysis package, Seurat (Satija et al. 2015).

18    Besides Louvain methods, these integrated packages implement a number of normalizations,

19    feature selection and other cell clustering algorithms. Nevertheless, all unsupervised clustering

20    algorithms try to quantify the similarity of every pairwise cell expression profile in order to

21    identify "clusters" of cells of the same type. Clusters are then assigned to known cell types via

22    the manual identification of genes expressed by cells in the cluster based on *a priori* biological

23    knowledge. Visual inspection is then used to "validate" the quality of the clustering, or the

1    clustering produced by a new algorithm is compared to the clustering reported using a previously

2    published algorithm.

3         As single cell data are increasingly used to define differences in cell state, such as might

4    develop during disease, in addition to cell type, trustworthy quantitative tools are necessary to

5    assess whether the partitioning of cell populations by an algorithm is justified. An issue with

6    many currents approaches is that they have not been validating against label assignment that are

7    independent of the transcriptomic data used by the clustering algorithm. We focus here first the

8    development of a benchmarking data set that does not rely on transcriptomic data for generating

9    cell annotations. Specifically, we take surface protein measurements and use a decision tree

10   model to generate cell type labels. We use the normalized mutual information to quantify the

11   agreement between pairs of cell clustering partitions and systemically evaluate the impact of

12   sparsity and genes on the accuracy and reproducibility of current clustering algorithms against

13   the surface-protein-based reference labels.

14        Additionally, we validate an entropy based, non-parametric feature selection algorithm to

15   evaluate the information content for genes in scRNA-seq data. Using only the genes with high

16   information content, both the clustering accuracy and reproducibility were improved for several

17   clustering algorithms.

18   **Results**

19        To avoid circularity in comparing the performance of algorithms for clustering single cell

20   RNA-seq data, we must generate a reference dataset that does not rely on annotations generated

21   making use of the same RNA-seq data (Fig. 1A). A publicly available peripheral blood

22   mononuclear cells (PBMC) dataset compiled by 10x Genomics enables us to accomplish this

23   goal. This dataset provides protein staining with TotalSeq-B antibodies for every cell and single

4

1    cell transcriptomic profiling using the 10x platform. The protein staining data allow us to assign

2    cell type labels following the approaches used traditionally for cell classification in

3    immunohistochemistry and flow cytometry and independently from the RNA-seq data (Fig. 1B).

4    This benchmarking data set then enables it to systematically evaluate the impact of different

5    feature selection algorithms and different clustering algorithms on cell labelling accuracy.

6        Following the traditional approach in cell classification using surface protein markers

7    with reproducible workflow, we produce a decision tree for cell type annotation (Fig. 1C). It is

8    visually apparent that the distribution of the surface protein levels is consistent in all cases with a

9    bimodal distribution. Therefore, we model the distribution using a Gaussian mixture model with

10   two peaks and determine a binary cutoff by identifying the value for which the likelihoods of a

11   cell belonging to each of the two groups are identical. We list the estimated parameters for the

12   decision tree in Supplementary Table 1. We use this approach to classify every one of the 713

13   cells in the dataset into one of ten cell types. These cell types cover all of the well-known cell

14   types in PBMC (Supplementary Table 2).

15       In order to test the robustness of our classification and to acknowledge uncertainty in a

16   cell's classification, we also consider the case in which we set cells with protein levels that yield

17   odds ratios smaller than 5:1 between estimated Gaussian functions for two peaks as "unclassified

18   cells". Using this procedure, we are still able to classify 696 of 713 cells (Fig. S1 and

19   Supplementary Table 3). Importantly, our results and conclusions do not change when

20   considering this alternative labelling approach.

21       Next, we investigate whether our benchmarking approach enables us to quantify the

22   impacts of data sparsity and potentially uninformative genes on the performance of clustering

23   algorithms. For concreteness, we focus first on the performance of build-in Louvain methods in

5

1   Seurat, which we will simply refer to as Seurat in the remainder of the manuscript. Seurat is one

2   of the most widely adopted tools in scRNA-seq analysis. To test the hypothesis that high sparsity

3   and overwhelming numbers of uninformative genes will induce poor classification accuracy, we

4   generate synthetic datasets that allow us to vary data sparsity and the number of potentially

5   uninformative genes.

6       As a first step, we apply Seurat to the surface protein level dataset in which sparsity is

7   negligible and all features are by construction meaningful. Using the full data for the seventeen

8   cell-surface proteins, we are able to classify cell types with very high precision (Fig. 2A-B and

9   Fig. S2). Indeed, the only misclassifications occur in separating cells types that are relatively

10  infrequent in the dataset. Thus, as one would expect, Seurat performs very well when applied to

11  data with negligible data sparsity.

12      To systemically investigate how data sparsity affects clustering accuracy without the

13  complications due to uninformative features, we model data sparsity by setting a model for the

14  probability that a specific surface protein's level is too low to be detected (refer to Methods for

15  details). We reason that, compared with an available method for producing sparsity where counts

16  are drawn from the negative binomial distribution (Zappia, Phipson, and Oshlack 2017), or a

17  naïve model where counts are randomly changed to 0, this method better captures the physical

18  process occurring in real experiments where the presence of a value of zero counts will be

19  correlated with the expression level.

20      We systematically change the value of the threshold level $t$, and quantify the clustering

21  accuracy using the normalized mutual information (NMI) between the reference labels and the

22  labels generated by the clustering algorithm on the synthetic dataset. NMI is a commonly used

23  metric to quantify the overlap between different partition. An NMI equal to 1 indicates perfect

1    overlap, whereas a value of 0 indicates random overlap (Witten and Frank 2002). Figure 2C

2    shows that our benchmarking approach correctly reflects that steady decrease in clustering

3    accuracy as data sparsity increases, that is, as the detectability threshold increases to higher

4    levels.

5         Next, we examine how a large number of potentially uninformative features interacts

6    with data sparsity to affect the accuracy of clustering algorithm as another important part of

7    benchmarking. To this end, we add the data from 16,000 randomly selected genes from the

8    transcriptomics dataset to the synthetic data for protein levels with distinct detectability threshold

9    generated in the previous step (Fig. 2D). As expected, the large number of likely uninformative

10   features leads to a decrease in clustering accuracy. Surprisingly, this decrease is even more

11   dramatic than one would expect as it completely overshadows the impact of data sparsity.

12        In order to quantify the specific effect of uninformative features, we add to the full

13   surface protein level data different numbers of randomly selected genes from the transcriptomic

14   dataset. It is remarkable that transcriptomic data for even as few as 1,000 randomly selected

15   genes reduces the algorithm's NMI by 25% (from 0.8 to 0.6, Fig. 2E). This decrease clearly

16   demonstrates that more data are not always helpful. In fact, an increase in the number of

17   uninformative features dramatically decreases performance.

18        These results demonstrate that the inclusion of features that are uninformative for the

19   purpose of cell labeling may have an unexpected impact on clustering accuracy (Kiselev,

20   Andrews, and Hemberg 2019).  Before proceeding, we must emphasize that important genes for

21   cell specific functions are not necessarily informative of cell type if they are expressed in

22   multiple cell populations.  For example, *MALAT1* and *HLA-A*, genes which are highly conserved

23   or constitutively expressed, are critical for cell function (Hutchinson et al. 2007; Kovats et al.

7

1    1990), however, their broad expression across cell populations means they are uninformative for

2    cellular classification.  Due to the importance of the degree to which a gene may provide

3    *information* for cell classification, we pursue an approach grounded on information theory that

4    we have recently demonstrated can dramatically improve the performance of algorithms for the

5    classification of texts (Gerlach, Shi, and Amaral 2019). Unlike other currently pursued

6    approaches, such as filtering out genes with low coefficient of variation, ours has a solid

7    theoretical grounding and does not require any assumptions about the distribution of counts of

8    unique molecular identifiers (UMIs).

9        Concisely, conditional entropy $H(g|S)$ is calculated for each gene given the UMI

10   frequency. We also created null models by random distribution to calculate average conditional

11   entropy over different realizations, $\langle \widetilde{H}(g|C) \rangle$ (details refer to method). The information content

12   $I(g)$ is defined as the difference between $H(g|S)$ and $\langle \widetilde{H}(g|C) \rangle$. Higher $I(g)$ indicates that the

13   gene has high expression in a subset of cells, which will be useful to determine cell types.

14       To test the robustness of our filtering approach, we calculate $H(g|S)$, $\langle \widetilde{H}(g|C) \rangle$, and

15   $I(g)$ for three PBMC datasets publicly available on the 10x Genomics website. As predicted, we

16   observe a very strong dependency of $\langle \widetilde{H}(g|C) \rangle$ on gene total expression level (Fig. 3A). We

17   also find that most genes have $I(g) \approx 0$, indicating that those genes cannot possibly provide

18   meaningful information concerning cell identities. Surprisingly, 15 out of 17 of the genes coding

19   the surface protein markers we use in generating labels are not informative for determining cell

20   type.

21       A few genes such as *GNL*Y, which codes for a T-cell activation protein, and *PTGDS*,

22   which codes for a protein involved in platelet aggregation, are informative about cell type across

23   the three datasets (Fig. 3A, B). Importantly, we observe a large and consistent number of

1    overlaps for the set of top 5%, 10% and 20 % most informative genes for cell type classification

2    across the three datasets (Fig. 3D and Fig. S3B-D). We also observe a large overlap of

3    informative genes using the default filtering algorithm in Seurat and our information-based

4    filtering but not when we consider the coefficient of variation filtering approach (Fig. 3E). These

5    findings suggest that the genes informative of cell type will be generally conserved across

6    biological replicates.  Importantly, because this approach requires no prior biological knowledge

7    and contains no fitted parameters, it can be applied without adjustment to samples from different

8    tissues or from different organisms (Fig. S3A).

9         We next compare the impact of using three feature filtering algorithms on the clustering

10   accuracy of Seurat: coefficient of variation, default filtering algorithm in Seurat, and our

11   information-based algorithm (Fig. 4A). Our analysis shows that using the coefficient of variation

12   for filtering features performs poorly at detecting uninformative features.  In contrast, the default

13   Seurat filtering and our information-based filtering are truly able to select informative features. A

14   point of distinction between the filtering algorithm implemented with Seurat and our proposed

15   information-based filtering is their theoretical foundations.  The filtering algorithm in Seurat

16   algorithm models the mean-variance relationship within the data, as such is limited by the scope

17   of its assumption and by the need to select parameter values.  Information theory provides the

18   best-grounded non-parametric approach to the problem of selecting informative features.

19   Moreover, it requires the fewest assumptions and provides a roadmap for improvement and

20   generalization. To our best knowledge, it is unlikely that any *ad hoc* approach being proposed

21   now and not grounded on information theory would be superior to that provided by information

22   theory.

1        We next evaluate the impact of using our filtering algorithm on clustering accuracy of

2    three algorithms: SC3 (Kiselev et al. 2017), Seurat (Satija et al. 2015), and Topic Mapping

3    (Lancichinetti et al. 2015). Topic Mapping is a highly reproducible high-accuracy algorithm for

4    clustering documents. We implement a straightforward analogy: words to genes and documents

5    to cells.

6        As the fraction of filtered genes increases, the classification accuracy of all three

7    algorithms also increases (Fig. 4B and Fig. S4A). Interestingly, the greatest gain in accuracy

8    occurs once we remove the 50% genes with the lowest $I(g)$, demonstrating that most genes

9    measured in scRNA-seq truly contain no useful information for cell classification. Indeed,

10    clustering accuracy increases up to the removal of the bottom 90% of genes by information

11    content.

12        Filtering out uninformative genes greatly improves classification stability (Fig. 4C-D and

13    Fig. S4B). Specifically, the number of clusters returned by Seurat and Topic Mapping become

14    more stable as the percentage of genes filtered increases. Additionally, the confidence with

15    which cluster assignment are made in Topic Mapping also increase with percentage of filtered

16    genes. Importantly, the genes that are most informative in our analysis are dissimilar to those

17    used for cellular classification based on protein measurements (Fig. 3A). This suggests that

18    classification schema based on genes coding for proteins targeted in flow cytometry,

19    immunohistochemistry or other protein-based measures (Bhattacharya et al. 2014), will not

20    necessarily extract the optimal amount of information from scRNA-seq data.

21        Our results also suggest that "noise" from uninformative genes leads to ambiguous

22    assignment of cells to clusters. After removing uninformative genes, the three algorithms largely

23    correctly assign major cell groups such as B cells or monocytes (Figs. 5A, C, E). Seurat and SC3,

1    but not by Topic Mapping, also correctly identify NK cells. However, all three have trouble

2    distinguishing between other cell types. CD4, CD8 and NKT cells are placed into the same

3    clusters by the three methods, but while SC3 and Seurat arbitrarily break the cells from those

4    types into two clusters, Topic Mapping puts them all in a single cluster. The confusion matrices

5    are similar when we compare the same clustering results against a supervised annotation

6    approach using SingleR and ENCDOE or HPCA as reference dataset (Figs. 5 B, D, F and Fig.

7    S5C). While the supervised annotation approach identifies boarder cell groups, our protein-based

8    annotation significantly overlaps with their results (Figs. S5A, B).

9

10   **Study Limitations**

11        There are limitations in our study. Our benchmarking method is limited to CITE-seq

12   platform where both surface protein and mRNA are quantified at the same time. The number of

13   CITE-seq datasets is limited in current literature. However, CITE-seq analysis is becoming

14   increasing popular in medicine (Saigusa and Ley 2020; Bandyopadhyay et al. 2019; Kotliarov et

15   al. 2020) and our method can be readily applied to these new datasets. On the other side, our

16   filtering approach is not limited to CITE-seq platform and can be applied in any scRNA-seq data

17   without any prior biology knowledge. Secondly, we recognize that researchers perform some

18   filtering procedures to remove the undetected genes before clustering algorithms as best practice.

19   However, most of current filtering rely on parametric tests, there could be potential biases in the

20   assumptions. Our method, on the other hand, does not reply on parametric test and has minimal

21   assumptions. Lastly, there exist supervised cell annotation tools such as SingleR (Aran et al.

22   2019). We demonstrate our annotation method largely overlaps with the results based on

1  ENCODE and HPCA reference but with more details. Furthermore, our labels can be used as a

2  reference dataset in those supervised approach.

3

4  **Discussion**

5  Data sparsity due to low RNA capture rates and uninformative genes are frequently

6  acknowledged as challenges in the analysis of scRNA-seq data (Kharchenko, Silberstein, and

7  Scadden 2014). Depending on the sequencing platform, experimental setup, and sample

8  preparation, on the order of 80% of the scRNA-seq measurements for each of the approximately

9  20,000 human genes yield zero counts (Angerer et al. 2017). Confirming prior research, we find

10  that high sparsity and the large number of uninformative genes significantly reduces the

11  performance of clustering algorithms. We show that this reduction in clustering accuracy and

12  reproducibility can be addressed using a quantitative approach based on information theory. We

13  show this approach provides a quantitative assessment of the information content for cell

14  classification of any gene in a specific dataset without the need for any prior biological

15  knowledge. When applied to a benchmarking data set generated independently of RNA-seq data,

16  this method provides a significant improvement in classification accuracy when compared to

17  other widely used tools (Kiselev, Andrews, and Hemberg 2019).

18  Our study highlights the importance of independent validation for benchmarking scRNA-

19  seq analysis tools to avoid circularity and lack of reproducibility. The entanglement of cell type

20  assignment with assessment of clustering quality carries an inherent risk of confirmation bias.

21  This situation is exacerbated by the widespread practice of using visual inspection to confirm

22  clustering quality, in which the authors use the specificity of gene assignment to a given cluster

23  (feature plots) to support the accuracy of the clustering algorithm. Information theory based

12

1  approaches offer an objective alternative by providing a quantitative measure of the information

2  conferred by each gene for determining cellular annotations.  We are hopeful this approach will

3  be applied by others to additional benchmark datasets with a goal of increasing reproducibility

4  and accuracy of scRNA-seq analysis.

5

1 **Methods**

2 **Benchmarking CITE-seq protein data with decision tree modeling**

3 We model the distribution using a Gaussian mixture model with two peaks and determine a

4 binary cutoff by identifying the value for which the likelihoods of a cell belonging to each of the

5 two groups are identical. We list the estimated parameters for the decision tree in the

6 supplementary information.

7    In order to test the robustness of our classification and to acknowledge uncertainty in a

8 cell's classification, we also consider the case in which we set cells with protein levels that yield

9 odds ratios smaller than 5:1 between estimated Gaussian functions for two peaks as "unclassified

10 cells".

11

12 **Clustering PBMC scRNA-seq data**

13 **Seurat:** We download version 3.2.1 of the Seurat package from CRAN R project and use default

14 defaults settings illustrated by the online tutorial (https://satijalab.org/seurat/vignettes.html).

15 Different parameters in resolution are sampled to reflect the change in clustering accuracy.

16 **SC3:** We download version 1.16.0 of the SC3 package from Bioconductor project and use

17 default settings except for the selection of the number of clusters. Several cluster numbers are

18 used to examine the clustering accuracy.

19 **Topic Mapping:** We obtained the Topic Mapping package from the authors. We implement a

20 straightforward analogy: words to genes and documents to cells. The TopicMapping algorithm

21 outputs two conditional probability vectors: p(gene|topic), which is used to identify marker

22 genes for each cluster, and p(topic|cell), which is used to identify cell clusters based on

23 probability distribution over topics. Each topic is treated as a cell type and the cells are assigned

1    to the topic with the highest probability. There are only two parameters in determining the

2    clustering results $P$, which specifies the statistical significance threshold for acceptance of new

3    clusters, and $t$, which specifies which required minimal manipulation. We use default values ($P =$

4    0.05 and $t = 10$) in our analyses.

5

**6    Generating synthetic datasets to model sparsity and uninformative genes**

7    We model data sparsity by setting a model for the probability that a specific surface protein's

8    level is too low to be detected. We use a sigmoidal function to generate the probability that

9    protein $j$ is detected by the assay for an individual cell $c$:

$$p(x_{c,j}; x0) = \frac{1}{1+e^{\ln(x0/x_{c,j})/\sigma}}$$

11    where $x$ is the actual protein level, ln is the natural logarithm function, $x0$ is the selected

12    detectability threshold, and $\sigma$ controls the stochasticity of the detection process around the

13    threshold. For simplicity, we set $\sigma = 1$ for all simulations.

14          To model the impact of uninformative genes on clustering accuracy, we add the data

15    from 500, 1000, 2000, 8000, 16,000 randomly selected genes from the transcriptomics dataset to

16    the synthetic data for protein levels with distinct detectability threshold generated in the previous

17    step.

18

**19    Calculating the information content of all genes for a sample**

20    Our method is generalized from that reported by (Gerlach, Shi, and Amaral 2019). Consider a

21    sample $S$ comprised of $L$ cells and denote by $n(g, c)$ the number of UMIs for gene $g$ and cell $c$.

22    The total reading depth $N$ can be expressed as the sum over all genes $g$ and all cells $c$ of $n(g, c)$.

23    The probability that a UMI for gene $g$ occurs in cell $c \in S$ is:

15

1

$$p(c|g) = \frac{p(g,c)}{p(g)} = \frac{n(g,c)}{n(g)}$$

2     Where $p(g) = N(g)/N$ is the relative frequency of UMIs for gene $g$ across the entire sample.

3     We can use Shannon entropy to quantify the heterogeneity of this distribution:

4

$$H(g|S) = -\sum_{c \epsilon S} p(c|g) \ln [p(c|g)].$$

5     The entropy $H(g|S)$ thus describes the information we can hope to extract from the distribution

6     of number of UMIs for gene $g$. A gene specific to a certain cell type will have less uncertainty in

7     its expression pattern across the sample and thus will have a lower entropy than a 'house-

8     keeping' gene which will not be specific to a certain cell type. The maximum possible value of

9     the conditional entropy occurs for a uniform distribution of number of UMIs,

10

$$H_{max} = -\sum_{c \in S} \frac{1}{L} \ln \left(\frac{1}{L}\right) = \ln(L).$$

11     Due to high sparsity nature of scRNA-seq dataset, most genes will occur at very low frequencies

12     and the limit expressed by $H_{max}$ will not be accurate. Therefore, we estimate a null model for the

13     $n(g)$ in order to calculate the expected value of the entropy under the null hypothesis. To this

14     end, we generate random distributions of UMI counts across genes and cells while preserving the

15     marginal counts $N(g)$ and $N(c) = \sum_g n(g,c)$. By calculating the expected entropy from many

16     realizations of this null model, we define the "true" information content of a gene, $I(g)$, as:

17

$$I(g) = \langle \tilde{H}(g|S) \rangle - H(g|S)$$

18     where $\langle \tilde{H}(g|C) \rangle$ denotes the average conditional entropy over different realizations of the null

19     model, and higher $I(g)$ indicates higher information-content. In contrast, genes with $I(g) \approx 0$

20     are uninformative and should be excluded from the list of features during the clustering analysis.

21

16

1    **Comparison of protein-based cell type annotation with supervised methods.**

2    We compare our protein-based cell type annotation with the supervised cell type annotation

3    method using SingleR with both Human Primary Cell Atlas (HPCA) and Encyclopedia of DNA

4    Elements (ENCODE) dataset. SingleR package is obtained through CRAN R project and default

5    parameters are used in cell type annotation for single cells.

6

7    **Data and code availability**

8    We use a publicly available peripheral blood mononuclear cells (PBMC) dataset compiled by

9    10x Genomics: PBMC dataset 1 (https://support.10xgenomics.com/single-cell-gene-

10    expression/datasets/3.0.0/pbmc_1k_protein_v3). PBMC dataset 2

11    (https://support.10xgenomics.com/single-cell-gene-

12    expression/datasets/3.0.0/pbmc_10k_protein_v3) and PBMC dataset 3

13    https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k). Human

14    and mouse pancreas cells data were downloaded on 04/03/2018 from https://hemberg-

15    lab.github.io/scRNA.seq.datasets/human/pancreas/.

16    The source code is publicly available at https://github.com/amarallab/Benchmark_scRNA_seq.

17

18    **Statistical analysis**

19    All statistical analyses are performed using R (Team 2014) and figures are produced using the

20    package ggplot2 (Wickham). The error bar and confident intervals are calculated using

21    bootstrapping. P-value is calculated using student's t-test with two side for figure 2 and

22    hypergeometric test for Venn diagram in figure 3 and figure S3 using build-in function in R.

23
24    **Acknowledgements**

17

5

6    **Conflicts of interests**

7    The authors declare no conflicts of interests.

8

1

**Figure 1. Limitations of current framework for optimizing scRNA-seq cell type**

**classification algorithms and development of an externally validated dataset.**

A. Current implicit frameworks for optimizing scRNA-seq classification algorithms assume that

some algorithm, typically Louvain method from Seurat, yields a ground-truth classification

against which the accuracy of other algorithms is then determined.

B. A reproducible, objective framework would make use of an independently-obtained, robust,

and reproducible independently-generated dataset against which the accuracy of scRNA-seq

19

1    classification algorithms can be objectively determined. In order to avoid circularity,  reference

2    labeling should be based on independent approaches such as surface protein expression or

3    immunostaining.

4    C. Creation of externally validated labelling for peripheral blood mononuclear cells (PBMCs)

5    from a healthy donor released by 10x Genomics (https://support.10xgenomics.com/single-cell-

6    gene-expression/datasets/3.0.0/pbmc_1k_protein_v3).  The dataset includes simultaneous single

7    cell levels for 17 surface proteins measured using ToptalSeq-B antibodies, as well as

8    transcriptomic profiles for 713 cells.  Surface protein level from the cells is consistent with a

9    bimodal distribution. We model empirical distribution as a mixture of two Gaussian peaks and

10   detect a threshold for binary classification of protein level, which can be used for classification

11   with a biologically-grounded decision tree enabling us to classify every cell into one of ten cell

12   types.

13

1

**Figure S1. Creation of externally validated labelling for peripheral blood mononuclear cells (PBMCs) with uncertainty thresholds.** We model the empirical distribution as a mixture of two Gaussian peaks, which can be used for classification with a biologically-grounded decision tree enabling us to classify every cell into one of ten cell types.

6

**Supplementary Table 1. The estimated parameters for each decision tree using Gaussian mixture model.** (Data shown in separated file)

**Supplementary Table 2. Cell type classification based on binary cutoffs.** (Data shown in separated file)

**Supplementary Table 3. Cell type classification based on uncertainty regions.** (Data shown in separated file)

**Figure 2. Clustering accuracy of Seurat decreases dramatically with increasing levels of data sparsity and increasing number of uninformative variables.**

A. Sankey diagram of Seurat classification based solely on surface protein labels. The left side shows the clusters identified by Seurat with the parameter set to the default value (0.8). The right side shows the cell types we identify using surface protein data and a decision tree.

B. Confusion matrix between Seurat classification and surface protein label in A. The most likely cell type is highlighted in red on the diagonal line.

C. Effect of data sparsity on accuracy of Seurat algorithm for four values of the parameter, *resolution*, ranging from 0.3 to 3.0. We model data sparsity by setting a threshold level for the detectability of the log protein counts (see text for details). As we increase data sparsity, by

22

1  moving the mid-point of sigmoidal function, we find a steady decrease in accuracy. Notably, this

2  decrease does not depend on the value of *resolution*. Sparsity measured by percentage of zero

3  counts (black line) with the increasing threshold of log protein counts. About 90% of the counts

4  become zero for $t \geq 8$.

5  D. Effect of number of possibly uninformative variables and data sparsity on accuracy. We

6  repeat the procedure described above for surface protein levels, but now include in the data the

7  measured levels of 16,000 randomly selected mRNAs. The large number of uninformative

8  variables included in the transcriptomic data decreases clustering accuracy by nearly 40% for

9  low levels of data sparsity. Interestingly, it yields best accuracy with *resolution* >1.6, not with

10  the default value.

11  E.  Effect of number of possibly uninformative variables on accuracy. We generate 10 replicates

12  by adding to the protein-levels data randomly-selected sets with data from varying numbers of

13  mRNAs in order to estimate the confidence intervals for accuracy of classification. If none of the

14  possibly uninformative mRNA is included, then classification accuracy measured is very high.

15  As we include data from an increasing number of mRNAs, we find a steady decrease in accuracy

16  across all parameter settings.

17

1

**Figure S2. Sankey diagram of classification based on surface protein levels using**

**uncertainty region.**

For the first two panels, the right side shows the cell types we identify using binary cutoffs of

surface proteins. The left side shows the clusters identified by A. SC3 algorithm with the user

input of correct cluster number (10); B. Topic Mapping algorithm.

For the last three panels, the right side shows the cell types we identify using binary cutoffs of

surface proteins. The left side shows the clusters identified by C. Seurat algorithm at default

*resolution* (0.8); D. SC3 algorithm with the user input of correct cluster number (10);

Topic Mapping algorithm.

1

1 **Figure 3. Context-specific identification of genes that are uninformative for cell type**

2 **classification.** We consider 3 PMBC datasets published by 10x Genomics. PBMC dataset 1 is

3 the dataset we consider in the previously presented figures. PBMC dataset 2 and 3 comprise

4 2,700 and 7865 PBMCs from human blood samples (https://support.10xgenomics.com/single-

5 cell-gene-expression/datasets/3.0.0/pbmc_10k_protein_v3 and

6 https://support.10xgenomics.com/single-cell-gene-expression/datasets/1.1.0/pbmc3k).

7 A. The conditional entropy across cells enables the identification of genes that cannot provide

8 information for a classification task. The red lines show the expected entropy for genes that are

9 randomly distributed across cells. The conditional entropies are normalized by $H_{max}$. Most genes

10 have distributions with entropies very close to the random expectation, implying that those genes

11 cannot provide any information concerning cell classification. Two such uninformative genes are

12 *MALAT1* and *HLA-A*, which code for highly-conserved proteins and are constitutively expressed

13 in almost all cell types. In contrast, *GNLY*, which codes for a T-cell activation protein, and

14 *PTGDS*, which codes for a protein involved in platelet aggregation, are quite informative about

15 blood cell type. Most of the surface protein coding genes such as CD16 or CD19 are close to the

16 null expectation and thus provides no additional information for cell type assignment.

17 B. Distribution of information content for genes with different UMI frequency. Each dot

18 represents a gene from PBMC dataset 1 (left), dataset 2 (middle), dataset 3 (right). The

19 information content is calculated by the difference between randomly distributed null model and

20 actual conditional entropy.

21 C. Distribution of UMI counts for *GNLY* (orange) and *MALAT1* (purple) gene across cells in

22 PBMC dataset 1 (left), dataset 2 (middle), dataset 3 (right). *MALAT1* is highly expressed in most

26

1    cells and has a uniform distribution. *GNLY* is highly expressed in a small fraction of the cells and

2    has a bimodal distribution.

3    D. Overlap of the most informative genes from three datasets of peripheral blood mononuclear

4    cells. The expected overlap across the three sets if we had simply selected 630 mRNAs at

5    random from 12,600 would have been only about 1.2 mRNA on average, whereas overlap

6    between two sets would have been above for most informative genes.  Suggesting the robustness

7    of the approach, and its biological significance, we find for both most and least informative

8    mRNAs a statistically significant excess in the number of mRNAs that overlap across the three

9    sets (p-values were calculated using non-parametric test with10,000 iterations).

10   E. Overlap of most informative genes using Seurat default filtering approach, CV based filtering

11   approach and our entropy-based approach using PBMC dataset 1. We select the same number of
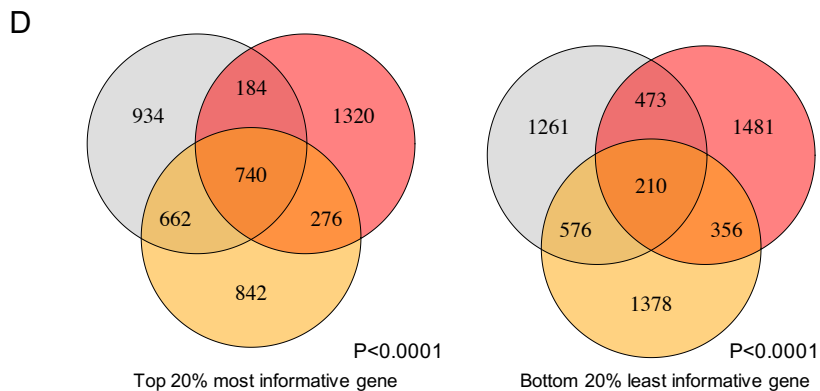
12   top informative genes from three filtering outputs.

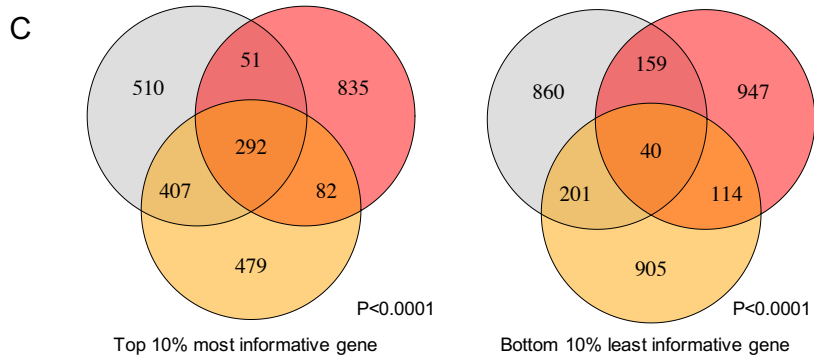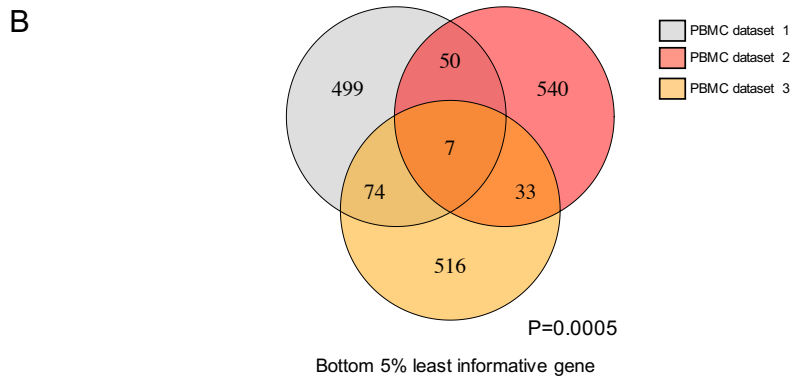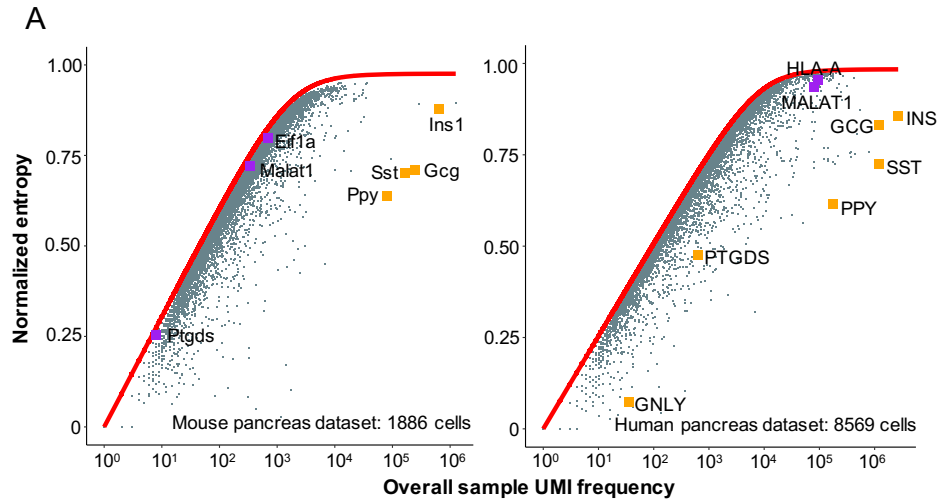1 **Figure S3. Context-specific identification of genes that are uninformative for cell type**

2 **classification.**

3 A. The conditional entropy across cells does not require *a priori* biological knowledge in

4 determining uninformative genes. We consider a published dataset of pancreas cells from either

5 human or mouse and calculate the information content for each gene. Genes such as INS that are

6 critical for pancreas cell types identification are with high information content. *MALAT1* or *HLA-*

7 *A* are uninformative in both human and mouse datasets like what we find in PBMC datasets.

8 Some genes such as *PTGDS* are informative for human cell types but not for mouse.

9 B-D: Overlap of most informative and least informative genes (B: 5%; C:10%; D: 20%) using

10 information contents from three datasets of peripheral blood mononuclear cells.

1

**Figure 4. Filtering uninformative genes yields dramatic improvements in classification**

**accuracy.**

A. Effect of filtering non-informative genes on accuracy using three different filtering

approaches. We filter the genes using Seurat default filtering, our info filtering and coefficient of

variance filtering and test the effects of classification accuracy on Seurat algorithm.

B. Effect of filtering non-informative genes on accuracy of two standard algorithms (Seurat,

SC3) and a new algorithm (Topic Mapping) for optimal parameter settings. We rank genes by

the information content of the distribution and filter out an increasing percentage of the least

informative genes. We generate 15 replicates using bootstrapping (sampling with replacement) in

order to estimate confidence intervals for the accuracy of the different classification algorithms.

We observe that there is no significant increase in accuracy until we remove at least 50%

mRNAs, indicating that most mRNA measured in scRNA-seq screens contain no useful

information for cell classification. As we then increase number of non-informative gene filtered

1 out, we reveal a steady increase in accuracy across all testing algorithms regardless of parameter

2 settings (see Fig. S4A). We also plot the results using Seurat default filtering algorithm. We also

3 plot the results using Seurat default filtering algorithm.

4 C. Filtering out non-informative genes does not alter the estimated number of detected clusters.

5 D. Filtering out non-informative genes reduces the percentage of unclassified cells in Topic

6 Mapping. The number of unclassified cells reduce significantly after we remove 50% of mRNAs

7 for Topic Mapping. For this dataset there are no unclassified cells in SC3. Seurat always

8 classifies every cell.

9

1

**Figure S4. Filtering uninformative genes yields dramatic improvements in classification accuracy.**

A. Effect of filtering non-informative genes on accuracy of two standard (Seurat, SC3) and a new (Topic Mapping) algorithms for different parameter settings.

B. Filtering out non-informative genes reduces the percentage of unclassified cells in Topic Mapping for different parameter settings.

8

**A — Seurat vs Protein**

| Seurat | B cells | Monocytes | Non-classical monocytes | CD4 naive | CD4 effector memory | CD8 naive | CD8 effector memory | NK cells | DC | NKT |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 75 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| 2 | 0 | 226 | 0 | 1 | 2 | 0 | 0 | 0 | 7 | 0 |
| 3 | 0 | 5 | 15 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 4 | 2 | 1 | 0 | 80 | 49 | 48 | 18 | 7 | 0 | 10 |
| 5 | 3 | 12 | 0 | 46 | 27 | 31 | 14 | 2 | 0 | 6 |
| 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 23 | 0 | 0 |

NMI: 0.58

**B — Seurat vs ENCODE**

| Seurat | B cells | Monocytes | CD4 T cells | CD8 T cells | NK cells | HSC | Erythrocytes | Eosinophils |
|---|---|---|---|---|---|---|---|---|
| 1 | 77 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 3 | 232 | 0 | 0 | 0 | 1 | 0 | 0 |
| 3 | 0 | 21 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 2 | 0 | 106 | 95 | 9 | 3 | 0 | 0 |
| 5 | 3 | 11 | 69 | 48 | 5 | 1 | 2 | 2 |
| 6 | 0 | 0 | 0 | 0 | 23 | 0 | 0 | 0 |

NMI: 0.63

**C — SC3 vs Protein**

| SC3 | B cells | Monocytes | Non-classical monocytes | CD4 naive | CD4 effector memory | CD8 naive | CD8 effector memory | NK cells | DC | NKT |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 78 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 222 | 0 | 0 | 2 | 0 | 0 | 0 | 4 | 0 |
| 3 | 0 | 15 | 15 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 4 | 0 | 0 | 0 | 73 | 46 | 43 | 16 | 0 | 0 | 3 |
| 5 | 1 | 6 | 0 | 47 | 28 | 32 | 14 | 0 | 0 | 6 |
| 6 | 0 | 0 | 0 | 5 | 1 | 1 | 0 | 0 | 0 | 2 |
| 7 | 0 | 1 | 0 | 2 | 1 | 2 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 32 | 0 | 3 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 |
| 10 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 2 |

NMI: 0.63

**D — SC3 vs ENCODE**

| SC3 | B cells | Monocytes | CD4 T cells | CD8 T cells | NK cells | HSC | Erythrocytes | Eosinophils |
|---|---|---|---|---|---|---|---|---|
| 1 | 78 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 0 | 228 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 31 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 0 | 97 | 84 | 0 | 0 | 0 | 0 |
| 5 | 1 | 5 | 71 | 49 | 3 | 1 | 2 | 2 |
| 6 | 0 | 0 | 5 | 4 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 2 | 0 | 0 | 4 | 0 | 0 |
| 8 | 0 | 0 | 0 | 3 | 34 | 0 | 0 | 0 |
| 9 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 1 | 0 | 0 | 3 | 0 | 0 | 0 | 0 |

NMI: 0.68

**E — Topic Mapping vs Protein**

| Topic Mapping | B cells | Monocytes | Non-classical monocytes | CD4 naive | CD4 effector memory | CD8 naive | CD8 effector memory | NK cells | DC | NKT |
|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 77 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| 2 | 0 | 235 | 15 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| 3 | 0 | 0 | 0 | 125 | 76 | 73 | 32 | 28 | 0 | 15 |
| 4 | 0 | 2 | 0 | 0 | 0 | 3 | 0 | 0 | 5 | 0 |
| Unclassified | 3 | 7 | 0 | 2 | 2 | 2 | 0 | 4 | 3 | 1 |

NMI: 0.67

**F — Topic Mapping vs ENCODE**

| Topic Mapping | B cells | Monocytes | CD4 T cells | CD8 T cells | NK cells | HSC | Erythrocytes | Eosinophils |
|---|---|---|---|---|---|---|---|---|
| 1 | 77 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 2 | 0 | 252 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 0 | 174 | 141 | 33 | 0 | 1 | 0 |
| 4 | 5 | 1 | 0 | 0 | 0 | 3 | 0 | 1 |
| Unclassified | 3 | 11 | 1 | 2 | 4 | 2 | 0 | 1 |

NMI: 0.73

**Figure 5. Confusion matrices for Seurat, SC3 and Topic Mapping classification of mRNA with 80% filtering against surface-protein based classification (A, C, E) and ENCODE based annotation (B, D, F), respectively.** For Topic Mapping, unclassified cells are listed in a separate row and excluded from NMI calculation.

1



A. ENCODE vs Protein

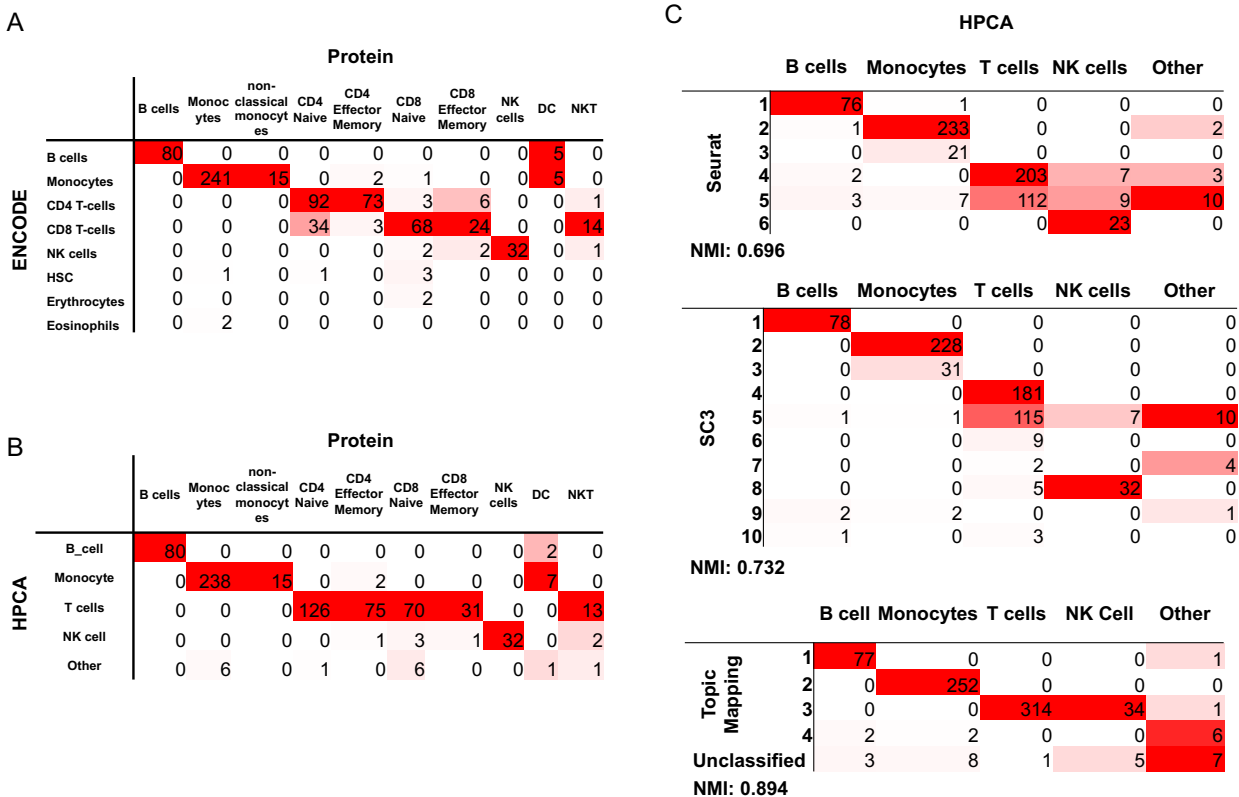| ENCODE | B cells | Monocytes | non-classical monocytes | CD4 Naive | CD4 Effector Memory | CD8 Naive | CD8 Effector Memory | NK cells | DC | NKT |
|---|---|---|---|---|---|---|---|---|---|---|
| B cells | 80 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 |
| Monocytes | 0 | 241 | 15 | 0 | 2 | 1 | 0 | 0 | 5 | 0 |
| CD4 T-cells | 0 | 0 | 0 | 92 | 73 | 3 | 6 | 0 | 0 | 1 |
| CD8 T-cells | 0 | 0 | 0 | 34 | 3 | 68 | 24 | 0 | 0 | 14 |
| NK cells | 0 | 0 | 0 | 0 | 0 | 2 | 2 | 32 | 0 | 1 |
| HSC | 0 | 1 | 0 | 1 | 0 | 3 | 0 | 0 | 0 | 0 |
| Erythrocytes | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 |
| Eosinophils | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

B. HPCA vs Protein

| HPCA | B cells | Monocytes | non-classical monocytes | CD4 Naive | CD4 Effector Memory | CD8 Naive | CD8 Effector Memory | NK cells | DC | NKT |
|---|---|---|---|---|---|---|---|---|---|---|
| B_cell | 80 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| Monocyte | 0 | 238 | 15 | 0 | 2 | 0 | 0 | 0 | 7 | 0 |
| T cells | 0 | 0 | 0 | 126 | 75 | 70 | 31 | 0 | 0 | 13 |
| NK cell | 0 | 0 | 0 | 0 | 1 | 3 | 1 | 32 | 0 | 2 |
| Other | 0 | 6 | 0 | 1 | 0 | 6 | 0 | 0 | 1 | 1 |

C. HPCA

Seurat

| | B cells | Monocytes | T cells | NK cells | Other |
|---|---|---|---|---|---|
| 1 | 76 | 1 | 0 | 0 | 0 |
| 2 | 1 | 233 | 0 | 0 | 2 |
| 3 | 0 | 21 | 0 | 0 | 0 |
| 4 | 2 | 0 | 203 | 7 | 3 |
| 5 | 3 | 7 | 112 | 9 | 10 |
| 6 | 0 | 0 | 0 | 23 | 0 |

NMI: 0.696

SC3

| | B cells | Monocytes | T cells | NK cells | Other |
|---|---|---|---|---|---|
| 1 | 78 | 0 | 0 | 0 | 0 |
| 2 | 0 | 228 | 0 | 0 | 0 |
| 3 | 0 | 31 | 0 | 0 | 0 |
| 4 | 0 | 0 | 181 | 0 | 0 |
| 5 | 1 | 1 | 115 | 7 | 10 |
| 6 | 0 | 0 | 9 | 0 | 0 |
| 7 | 0 | 0 | 2 | 0 | 4 |
| 8 | 0 | 0 | 5 | 32 | 0 |
| 9 | 2 | 2 | 0 | 0 | 1 |
| 10 | 1 | 0 | 3 | 0 | 0 |

NMI: 0.732

Topic Mapping

| | B cell | Monocytes | T cells | NK Cell | Other |
|---|---|---|---|---|---|
| 1 | 77 | 0 | 0 | 0 | 1 |
| 2 | 0 | 252 | 0 | 0 | 0 |
| 3 | 0 | 0 | 314 | 34 | 1 |
| 4 | 2 | 2 | 0 | 0 | 6 |
| Unclassified | 3 | 8 | 1 | 5 | 7 |

NMI: 0.894

2
3
4
5 **Figure S5. Supervised annotation methods highly overlap with our protein-based**

6 **annotation.**

7 A. Confusion matrix for supervised annotation method using ENCODE SingleR and our protein-

8 based annotation.

9 B. Confusion matrix for supervised annotation method using HPCA SingleR and our protein-

10 based annotation.

11 C. Confusion matrix for Seurat (top), SC3 (middle) and Topic Mapping (bottom) classification

12 of mRNA with 80% filtering against HPCA based annotation.

13

14

**References:**

Aizawa, Akiko. 2003. 'An information-theoretic perspective of tf–idf measures', *Information Processing & Management*, 39: 45-65.

Amaral, Luis AN, and Julio M Ottino. 2004. 'Complex networks', *The European Physical Journal B*, 38: 147-62.

Angerer, Philipp, Lukas Simon, Sophie Tritschler, F Alexander Wolf, David Fischer, and Fabian J Theis. 2017. 'Single cells make big data: New challenges and opportunities in transcriptomics', *Current Opinion in Systems Biology*, 4: 85-91.

Aran, D., A. P. Looney, L. Liu, E. Wu, V. Fong, A. Hsu, S. Chak, R. P. Naikawadi, P. J. Wolters, A. R. Abate, A. J. Butte, and M. Bhattacharya. 2019. 'Reference-based analysis of lung single-cell sequencing reveals a transitional profibrotic macrophage', *Nat Immunol*, 20: 163-72.

Bandyopadhyay, Gautam, Jacquelyn Lillis, Ravi S Misra, Jason R Myers, John M Ashton, Heidie L Huyck, Daria Krenitsky, Stephen T Romas, Cory J Poole, and Jeanne Holden-Wiltse. 2019. 'Identification and Characterization of Cellular Heterogeneity within Human Late Developmental Stage Dissociated Lung by CITE-Seq', *The FASEB Journal*, 33: 847.5-47.5.

Bhattacharya, Subarna, Paul W Burridge, Erin M Kropp, Sandra L Chuppa, Wai-Meng Kwok, Joseph C Wu, Kenneth R Boheler, and Rebekah L Gundry. 2014. 'High efficiency differentiation of human pluripotent stem cells to cardiomyocytes and characterization by flow cytometry', *JoVE (Journal of Visualized Experiments)*: e52010.

Blondel, Vincent D, Jean-Loup Guillaume, Renaud Lambiotte, and Etienne Lefebvre. 2008. 'Fast unfolding of communities in large networks', *Journal of statistical mechanics: theory and experiment*, 2008: P10008.

Cheng, Jie, Russell Greiner, Jonathan Kelly, David Bell, and Weiru Liu. 2002. 'Learning Bayesian networks from data: An information-theory based approach', *Artificial intelligence*, 137: 43-90.

Collin, J., R. Queen, D. Zerti, B. Dorgau, R. Hussain, J. Coxhead, S. Cockell, and M. Lako. 2019. 'Deconstructing Retinal Organoids: Single Cell RNA-Seq Reveals the Cellular Components of Human Pluripotent Stem Cell-Derived Retina', *Stem Cells*, 37: 593-98.

Davis, Carrie A, Benjamin C Hitz, Cricket A Sloan, Esther T Chan, Jean M Davidson, Idan Gabdank, Jason A Hilton, Kriti Jain, Ulugbek K Baymuradov, and Aditi K Narayanan. 2018. 'The Encyclopedia of DNA elements (ENCODE): data portal update', *Nucleic Acids Res*, 46: D794-D801.

Gerlach, Martin, Hanyu Shi, and Luís A Nunes Amaral. 2019. 'A universal information theoretic approach to the identification of stopwords', *Nature Machine Intelligence*, 1: 606-12.

Hutchinson, J. N., A. W. Ensminger, C. M. Clemson, C. R. Lynch, J. B. Lawrence, and A. Chess. 2007. 'A screen for nuclear transcripts identifies two linked noncoding RNAs associated with SC35 splicing domains', *BMC Genomics*, 8: 39.

Kharchenko, P. V., L. Silberstein, and D. T. Scadden. 2014. 'Bayesian approach to single-cell differential expression analysis', *Nat Methods*, 11: 740-2.

Kiselev, V. Y., T. S. Andrews, and M. Hemberg. 2019. 'Challenges in unsupervised clustering of single-cell RNA-seq data', *Nat Rev Genet*, 20: 273-82.

Kiselev, V. Y., K. Kirschner, M. T. Schaub, T. Andrews, A. Yiu, T. Chandra, K. N. Natarajan, W. Reik, M. Barahona, A. R. Green, and M. Hemberg. 2017. 'SC3: consensus clustering of single-cell RNA-seq data', *Nat Methods*, 14: 483-86.

Kotliarov, Yuri, Rachel Sparks, Andrew J Martins, Matthew P Mulè, Yong Lu, Meghali Goswami, Lela Kardava, Romain Banchereau, Virginia Pascual, and Angélique Biancotto. 2020. 'Broad immune activation underlies shared set point signatures for vaccine responsiveness in healthy individuals and disease activity in patients with lupus', *Nature Medicine*, 26: 618-29.

Kovats, S., E. K. Main, C. Librach, M. Stubblebine, S. J. Fisher, and R. DeMars. 1990. 'A class I antigen, HLA-G, expressed in human trophoblasts', *Science*, 248: 220-3.

Lancichinetti, Andrea, M Irmak Sirer, Jane X Wang, Daniel Acuna, Konrad Körding, and Luís A Nunes Amaral. 2015. 'High-reproducibility and high-accuracy method for automated topic classification', *Physical Review X*, 5: 011007.

Lawson, D. A., N. R. Bhakta, K. Kessenbrock, K. D. Prummel, Y. Yu, K. Takai, A. Zhou, H. Eyob, S. Balakrishnan, C. Y. Wang, P. Yaswen, A. Goga, and Z. Werb. 2015. 'Single-cell analysis reveals a stem-cell program in human metastatic breast cancer cells', *Nature*, 526: 131-5.

Lindström, N. O., G. De Sena Brandine, A. Ransick, and A. P. McMahon. 2019. 'Single-Cell RNA Sequencing of the Adult Mouse Kidney: From Molecular Cataloging of Cell Types to Disease-Associated Predictions', *Am J Kidney Dis*, 73: 140-42.

Mabbott, N. A., J. K. Baillie, H. Brown, T. C. Freeman, and D. A. Hume. 2013. 'An expression atlas of human primary cells: inference of gene function from coexpression networks', *BMC Genomics*, 14: 632.

Min, J. W., W. J. Kim, J. A. Han, Y. J. Jung, K. T. Kim, W. Y. Park, H. O. Lee, and S. S. Choi. 2015. 'Identification of Distinct Tumor Subpopulations in Lung Adenocarcinoma via Single-Cell RNA-seq', *PLoS One*, 10: e0135817.

Miyamoto, D. T., Y. Zheng, B. S. Wittner, R. J. Lee, H. Zhu, K. T. Broderick, R. Desai, D. B. Fox, B. W. Brannigan, J. Trautwein, K. S. Arora, N. Desai, D. M. Dahl, L. V. Sequist, M. R. Smith, R. Kapur, C. L. Wu, T. Shioda, S. Ramaswamy, D. T. Ting, M. Toner, S. Maheswaran, and D. A. Haber. 2015. 'RNA-Seq of single prostate CTCs implicates noncanonical Wnt signaling in antiandrogen resistance', *Science*, 349: 1351-6.

Patel, A. P., I. Tirosh, J. J. Trombetta, A. K. Shalek, S. M. Gillespie, H. Wakimoto, D. P. Cahill, B. V. Nahed, W. T. Curry, R. L. Martuza, D. N. Louis, O. Rozenblatt-Rosen, M. L. Suvà, A. Regev, and B. E. Bernstein. 2014. 'Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma', *Science*, 344: 1396-401.

Reyfman, P. A., J. M. Walter, N. Joshi, K. R. Anekalla, A. C. McQuattie-Pimentel, S. Chiu, R. Fernandez, M. Akbarpour, C. I. Chen, Z. Ren, R. Verma, H. Abdala-Valencia, K. Nam, M. Chi, S. Han, F. J. Gonzalez-Gonzalez, S. Soberanes, S. Watanabe, K. J. N. Williams, A. S. Flozak, T. T. Nicholson, V. K. Morgan, D. R. Winter, M. Hinchcliff, C. L. Hrusch, R. D. Guzy, C. A. Bonham, A. I. Sperling, R. Bag, R. B. Hamanaka, G. M. Mutlu, A. V. Yeldandi, S. A. Marshall, A. Shilatifard, L. A. N. Amaral, H. Perlman, J. I. Sznajder, A. C. Argento, C. T. Gillespie, J. Dematte, M. Jain, B. D. Singer, K. M. Ridge, A. P. Lam, A. Bharat, S. M. Bhorade, C. J. Gottardi, G. R. S. Budinger, and A. V. Misharin. 2019. 'Single-Cell Transcriptomic Analysis of Human Lung Provides Insights into the Pathobiology of Pulmonary Fibrosis', *Am J Respir Crit Care Med*, 199: 1517-36.

Rizvi, A. H., P. G. Camara, E. K. Kandror, T. J. Roberts, I. Schieren, T. Maniatis, and R. Rabadan. 2017. 'Single-cell topological RNA-seq analysis reveals insights into cellular differentiation and development', *Nat Biotechnol*, 35: 551-60.

Saigusa, Ryosuke, and Klaus Ley. 2020. 'CITE-Seq Hits Vascular Medicine', *Clinical Chemistry*, 66: 751-53.

Satija, R., J. A. Farrell, D. Gennert, A. F. Schier, and A. Regev. 2015. 'Spatial reconstruction of single-cell gene expression data', *Nat Biotechnol*, 33: 495-502.

Team, R Core. 2014. "R: a language and environment for statistical computing. Version 3.1. 2 [computer program]. R Foundation for Statistical Computing, Vienna, Austria." In.

Treutlein, B., D. G. Brownfield, A. R. Wu, N. F. Neff, G. L. Mantalas, F. H. Espinoza, T. J. Desai, M. A. Krasnow, and S. R. Quake. 2014. 'Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq', *Nature*, 509: 371-5.

Villani, A. C., R. Satija, G. Reynolds, S. Sarkizova, K. Shekhar, J. Fletcher, M. Griesbeck, A. Butler, S. Zheng, S. Lazo, L. Jardine, D. Dixon, E. Stephenson, E. Nilsson, I. Grundberg, D. McDonald, A. Filby, W. Li, P. L. De Jager, O. Rozenblatt-Rosen, A. A. Lane, M. Haniffa, A. Regev, and N. Hacohen. 2017. 'Single-cell RNA-seq reveals new types of human blood dendritic cells, monocytes, and progenitors', *Science*, 356.

Wickham, H. 'ggplot2: elegant graphics for data analysis Springer; New York; 2009', *URL http://had. co. nz/ggplot2/book.[Google Scholar]*.

Witten, Ian H, and Eibe Frank. 2002. 'Data mining: practical machine learning tools and techniques with Java implementations', *Acm Sigmod Record*, 31: 76-77.

Zappia, L., B. Phipson, and A. Oshlack. 2017. 'Splatter: simulation of single-cell RNA sequencing data', *Genome Biol*, 18: 174.