

Start Setting up for the Workshop

Workshop webpage: <https://go.tufts.edu/discworkshop>

Log on with Tufts Credentials to On Demand on Tufts Cluster <https://ondemand.pax.tufts.edu/>

Click on `Interactive Apps > RStudio Pax` and you will see a form to fill out to request compute resources to use RStudio on the Tufts HPC cluster. We will fill out the form with the following entries:

- `Number of hours : 5`
- `Number of cores : 1`
- `Amount of memory : 16GB`
- `R version : 4.0.0`
- `Reservation for class, training, workshop : Bioinformatics Workshops`

Single-cell RNA-Seq Bioinformatics

Day 2: Cell-type classification and Differential Expression

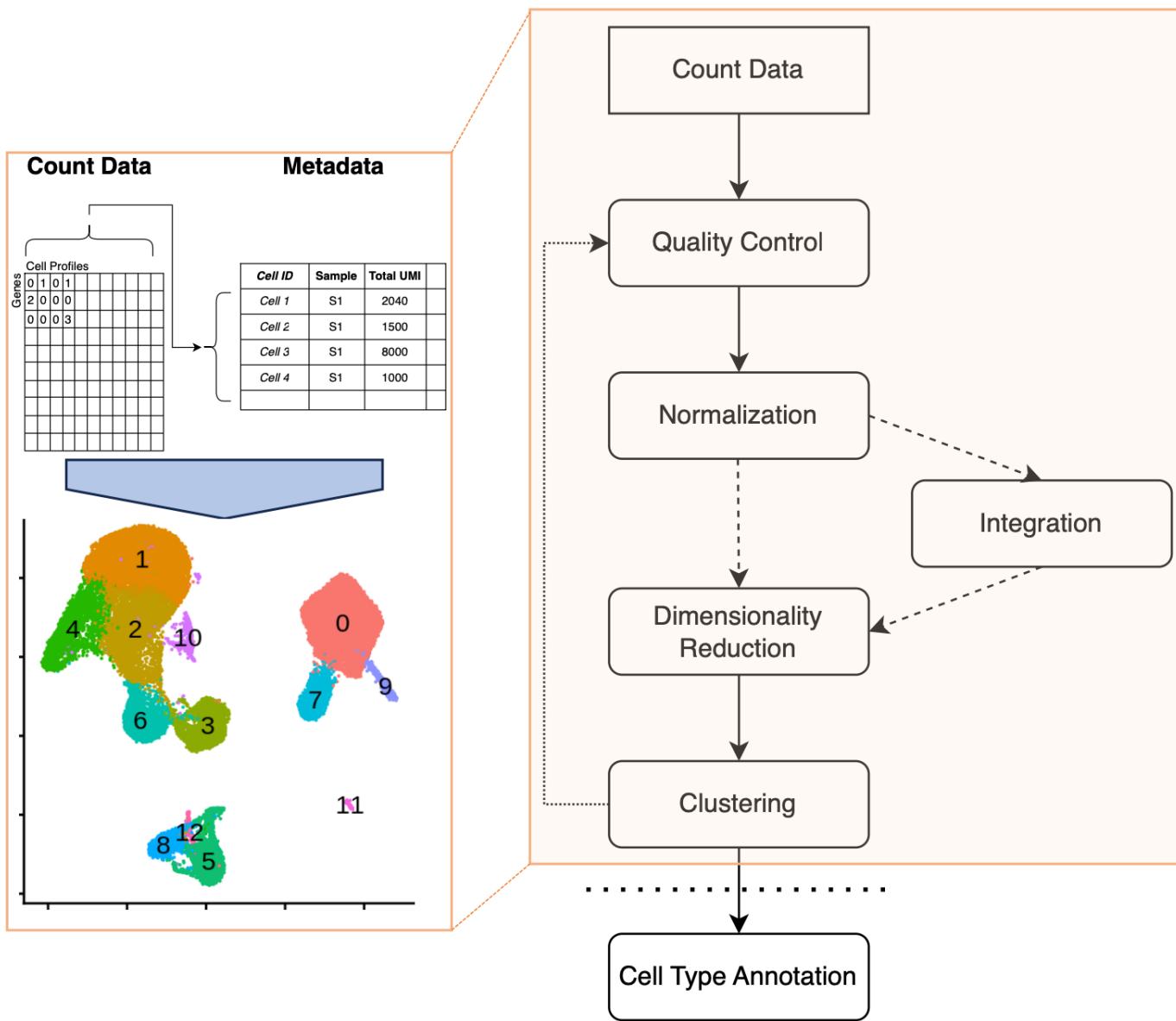
Data Intensive Studies Center (DISC) Single-cell series Oct-Nov 2023

Rebecca Batorsky, Data Scientist, Rebecca.Batorsky@tufts.edu

Eric Reed, Data Scientist, Eric.Reed@tufts.edu

Case Study

Workshop part I recap



scRNAseq Data

Peripheral Blood Mononuclear Cells (PBMCs)

Includes:

- B cells
- T cells
- Natural killer cells
- Monocytes
- Macrophages

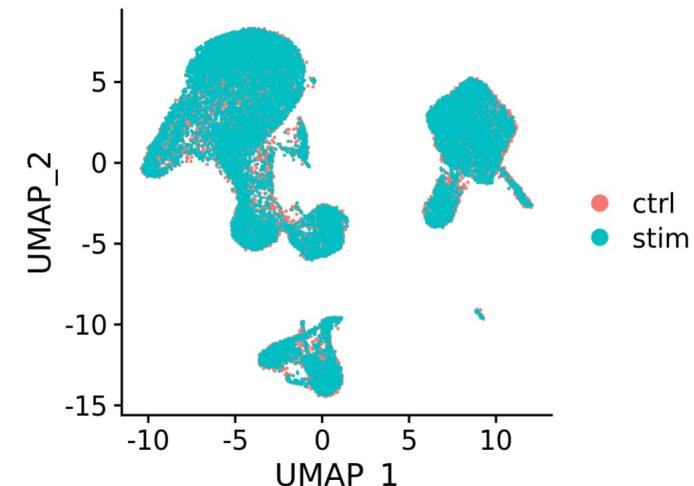
Doesn't include:

- Neutrophils
- Platelets
- Red blood cells

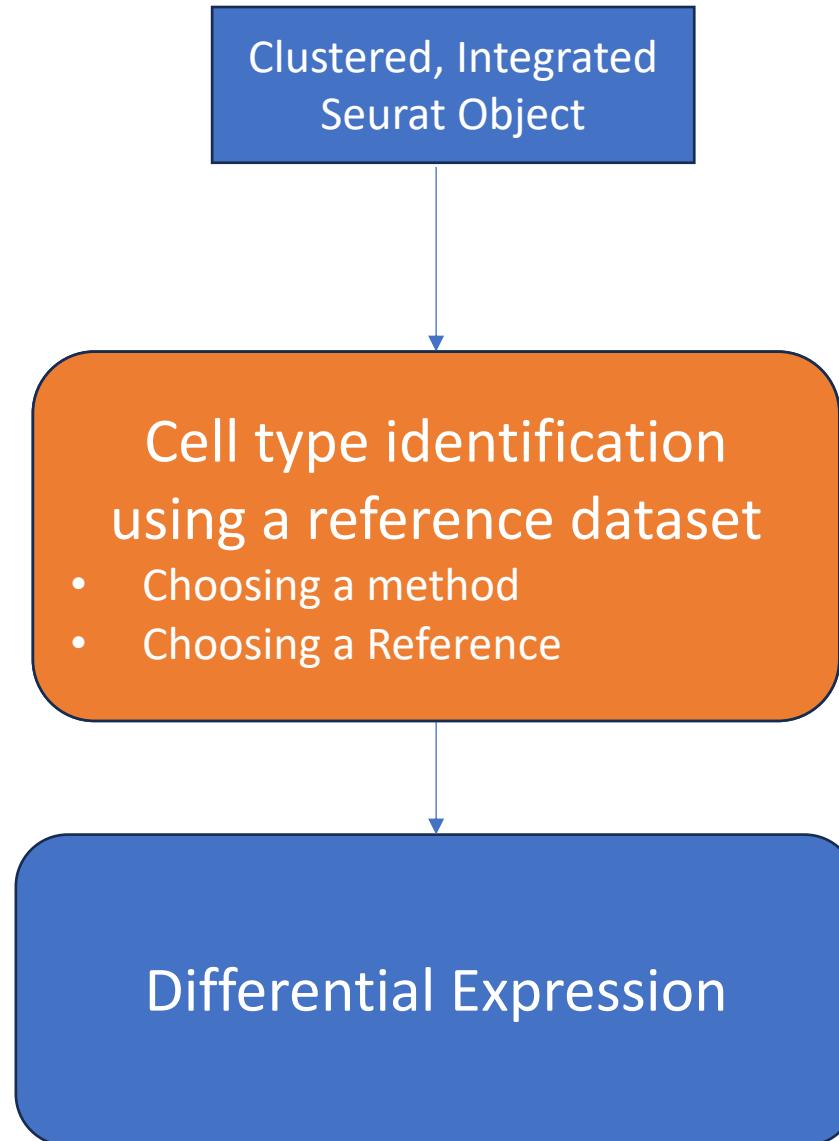
Two samples

Control: "ctrl"

Interferon beta-treated (stimulated): "stim"



Workshop day 2



Our Seurat Object

@metadata

| | Sample | nCount_RNA | nFeature_RNA | ... | integrated_snn_res.0.4 |
|--------|--------|------------|--------------|-----|------------------------|
| Cell-1 | ctrl | 2344 | 874 | | 1 |
| Cell-2 | ctrl | 3125 | 896 | | 2 |
| Cell-3 | stim | 2578 | 725 | | 3 |

@reductions

...

@assays

RNA

counts

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | 0 | 3 | 0 |
| gene2 | 1 | 1 | 0 |

counts

- SCT norm for sequencing depth

SCT

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | 0 | 2 | 0 |
| gene2 | 1 | 0 | 0 |

Integrated

data

- Norm
- log1p

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | 0 | 2.02 | 0 |
| gene2 | 2.42 | 1.15 | 0 |

data

- log1p

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | 0 | 1.09 | 0 |
| gene2 | 0.69 | 0 | 0 |

data

- integrated

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | -0.55 | 6.29 | 0.35 |
| gene2 | 1.7 | 0.01 | -0.36 |

- Scale.data
- SCTModel

- Scale.data
- Var.features

Adding Cell-type labels

Query dataset

| | Sample | nCount_RNA | nFeature_RNA | ... | integrated_snn_res.0.4 | Cell-type |
|--------|--------|------------|--------------|-----|------------------------|-----------|
| Cell-1 | ctrl | 2344 | 874 | | 1 | ? |
| Cell-2 | ctrl | 3125 | 896 | | 2 | ? |
| Cell-3 | stim | 2578 | 725 | | 3 | ? |

query@assays\$RNA@data

| | Cell-1 | Cell-2 | Cell-3 |
|-------|--------|--------|--------|
| gene1 | 0 | 2.02 | 0 |
| gene2 | 2.42 | 1.15 | 0 |

Reference dataset

Reference Metadata

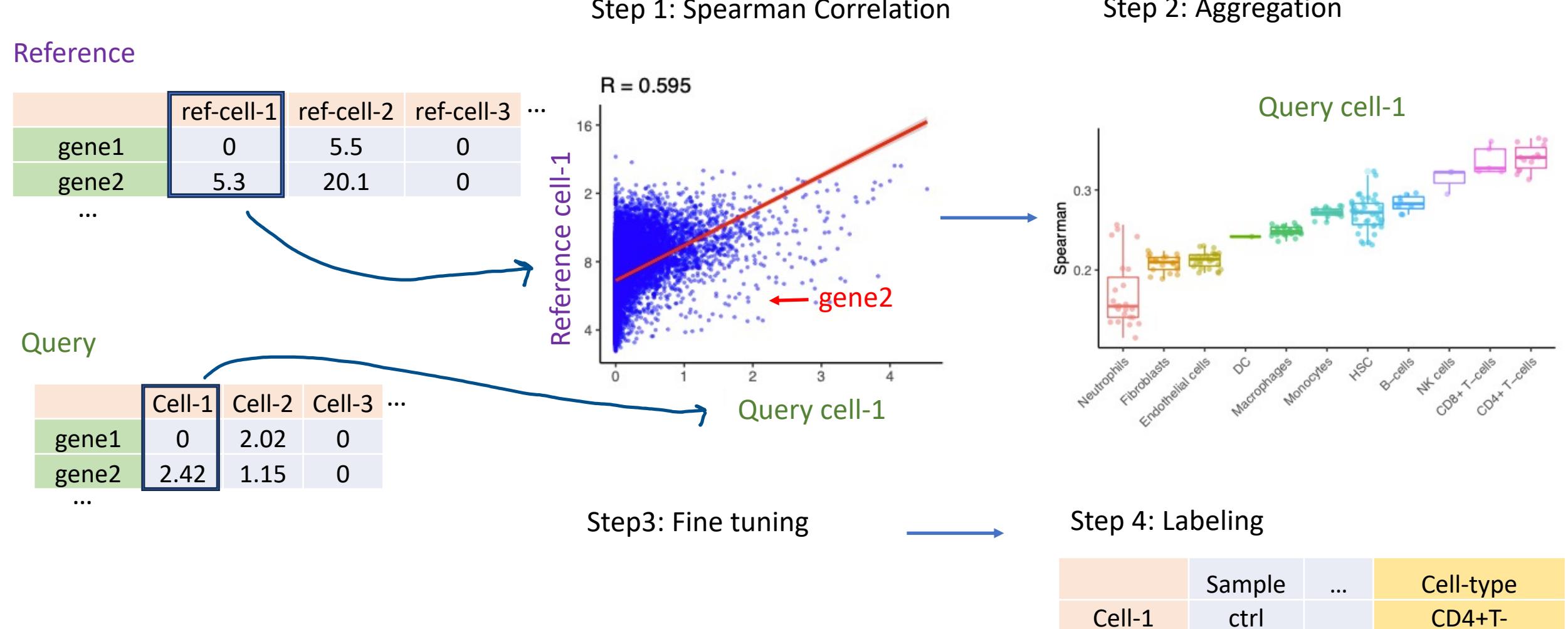
| | ... | Cell-type |
|------------|-----|------------|
| ref-cell-1 | | T-cell |
| ref-cell-2 | | B-cell |
| ref-cell-3 | | Macrophage |

Normalized data

| | ref-cell-1 | ref-cell-2 | ref-cell-3 |
|-------|------------|------------|------------|
| gene1 | 0 | 5.5 | 0 |
| gene2 | 5.3 | 20.1 | 0 |

SingleR Correlation Method

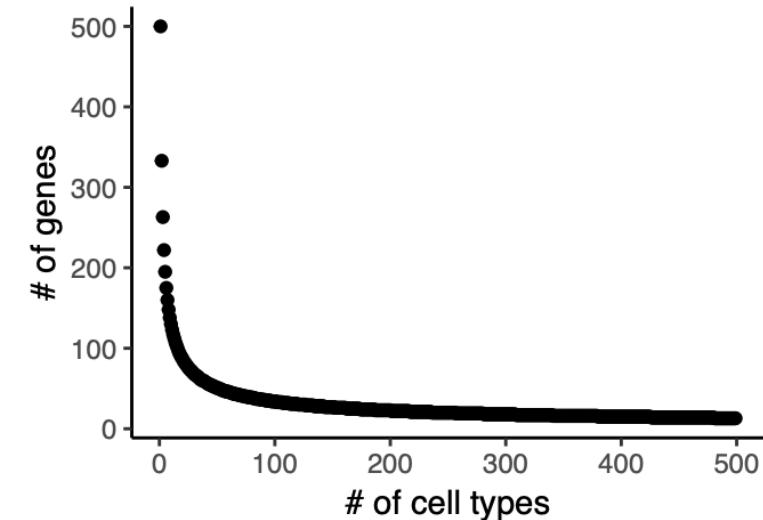
Using a reference database of cell-type gene expression, calculate correlation of gene expression for each cell



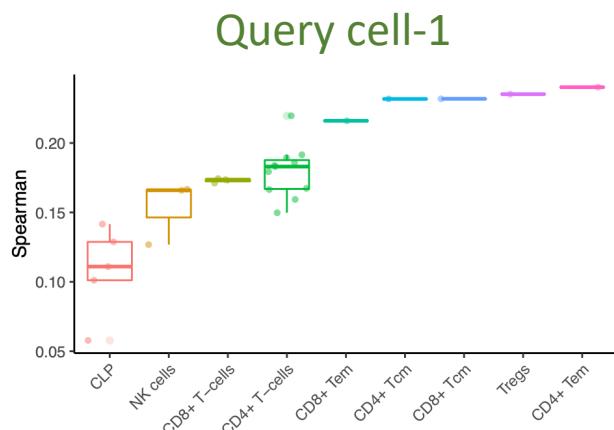
SingleR method

Step 3: Fine tuning

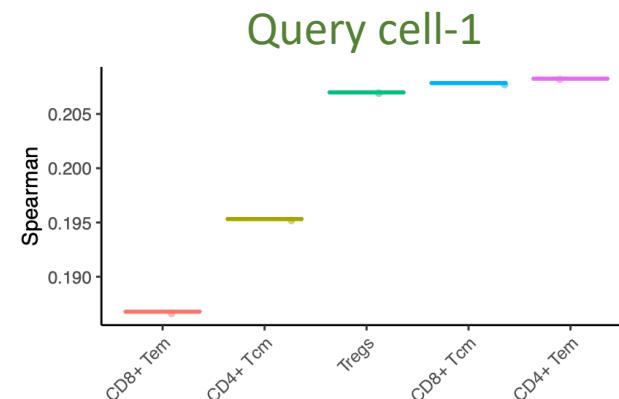
- Assignment is performed in multiple iterations.
- At each iterations, only the top N variable genes between cell-types are considered
- After each iteration, only best scoring cell types are kept



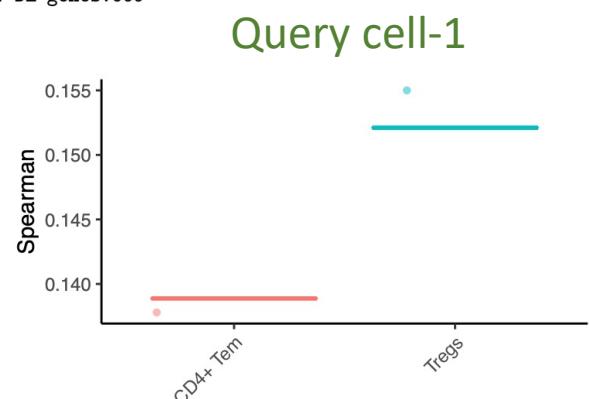
```
## [1] "Number of DE genes:1819"
```



```
## [1] "Number of DE genes:1287"
```



```
## [1] "Number of DE genes:666"
```



SingleR with celldex database of pure cell types

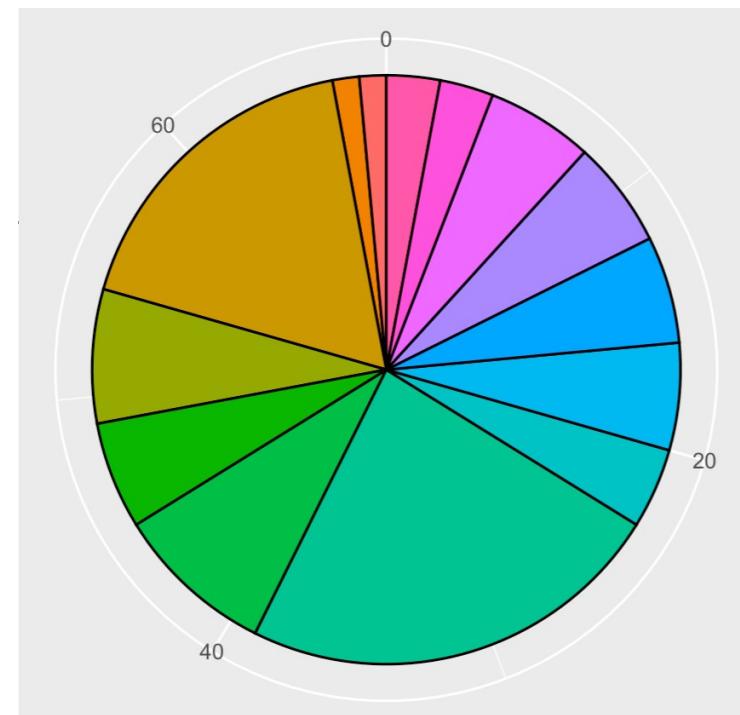
Celldex R library is a general purpose reference, giving convenient access to expression profiles for many cell types

| Some examples | Description |
|---------------------------------------|--|
| Human primary cell atlas (HPCA) | Microarray datasets derived from human primary cells |
| Blueprint/ENCODE | Bulk RNA-seq data for pure stroma and immune cells |
| Immunological Genome Project (ImmGen) | Microarray profiles of pure mouse immune cells |

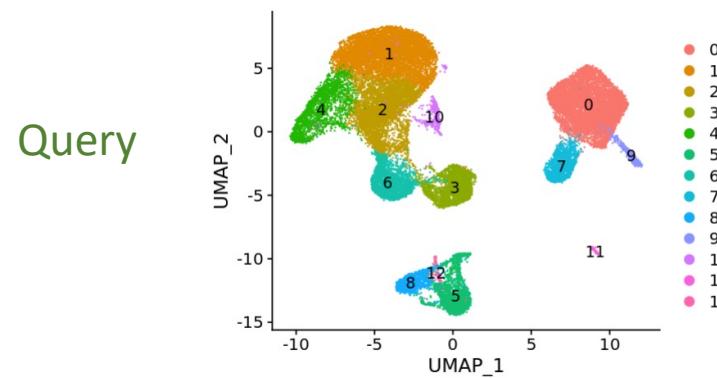
Labels are given at two resolutions, e.g. for HPCA **label.main.** = T cells

label.fine =

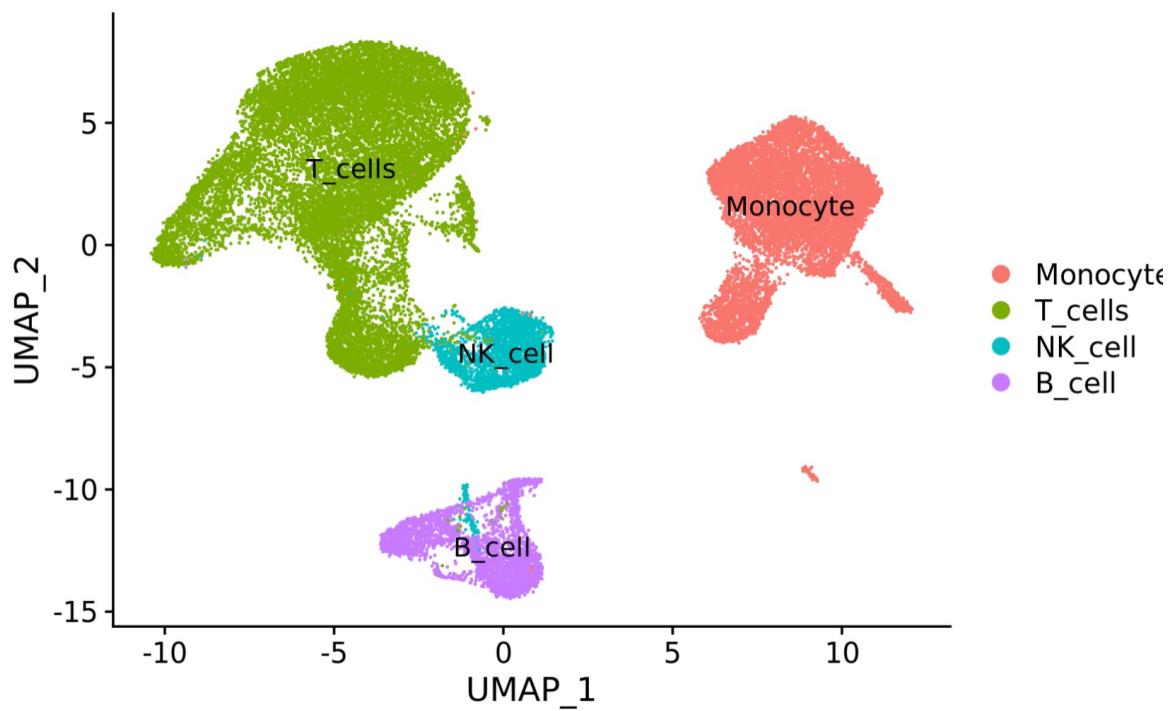
- █ T_cell:CCR10-CLA+1,25(OH)2_vit_D3/IL-12
- █ T_cell:CCR10+CLA+1,25(OH)2_vit_D3/IL-12
- █ T_cell:CD4+
- █ T_cell:CD4+_central_memory
- █ T_cell:CD4+_effector_memory
- █ T_cell:CD4+_Naive
- █ T_cell:CD8+
- █ T_cell:CD8+_Central_memory
- █ T_cell:CD8+_effector_memory
- █ T_cell:CD8+_effector_memory_RA
- █ T_cell:CD8+_naive
- █ T_cell:effector
- █ T_cell:gamma-delta
- █ T_cell:Treg:Naive



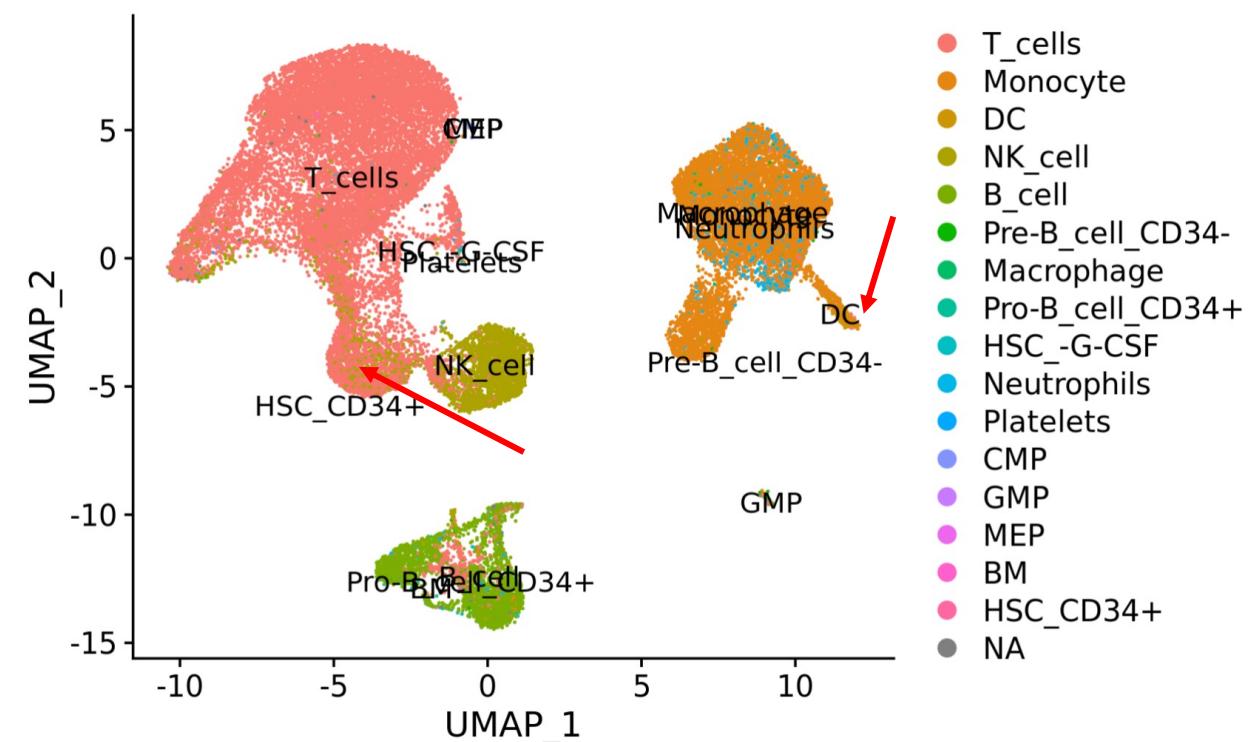
SingleR Results



Cluster level, HPCA, label.main



Cell level, HPCA, label.main



Seurat Integration mapping method



Products Research

Support > Single Cell Gene Expression > Datasets

3k PBMCs from a Healthy Donor

Single Cell Gene Expression Dataset by Cell Ranger 1.1.0

Peripheral blood mononuclear cells (PBMCs) from a healthy donor (same donor as pbmc6k).

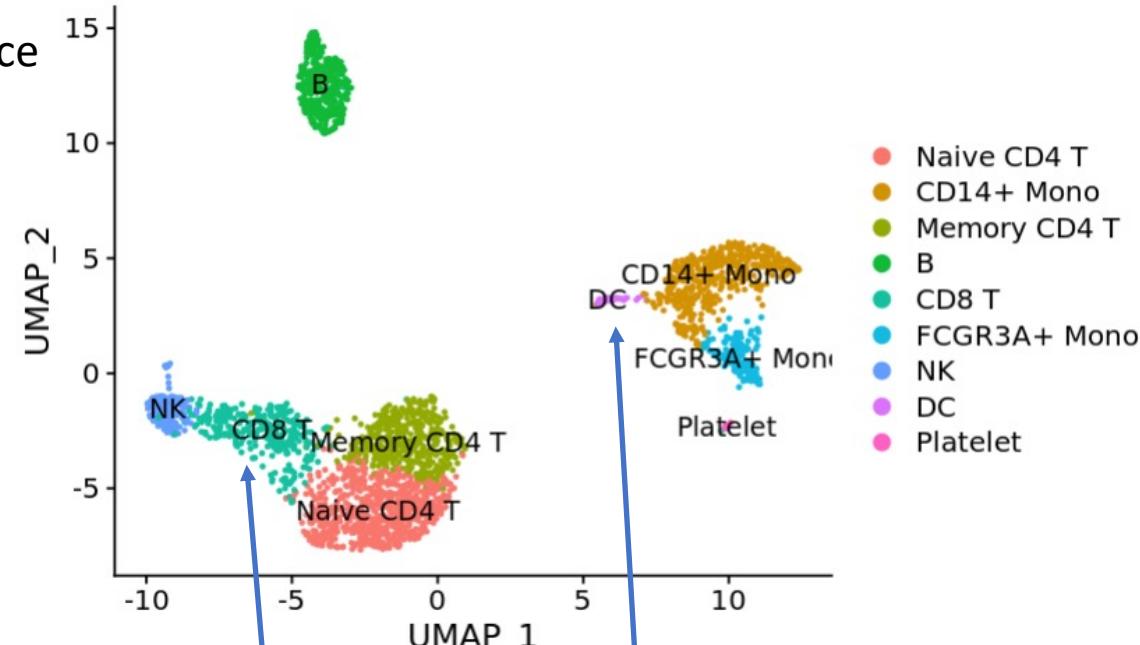
PBMCs are primary cells with relatively small amounts of RNA (~1pg RNA/cell).

- 2,700 cells detected
- Sequenced on Illumina NextSeq 500 with ~69,000 reads per cell
- 98bp read1 (transcript), 8bp I5 sample barcode, 14bp I7 GemCode barcode and 10bp read2 (UMI)
- Analysis run with --cells=3000

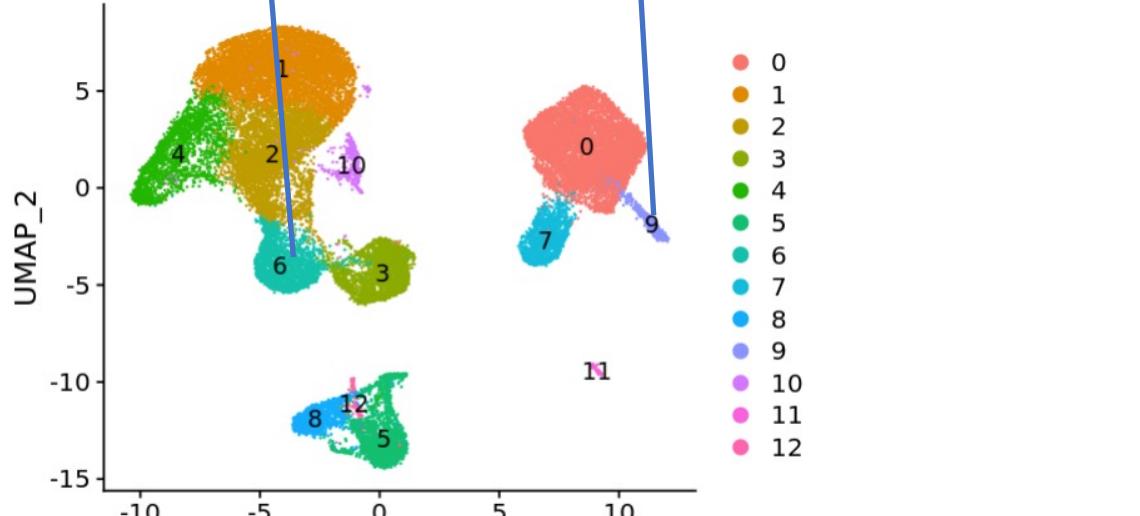
Published on May 26, 2016

This dataset is licensed under the Creative Commons Attribution license.

Reference

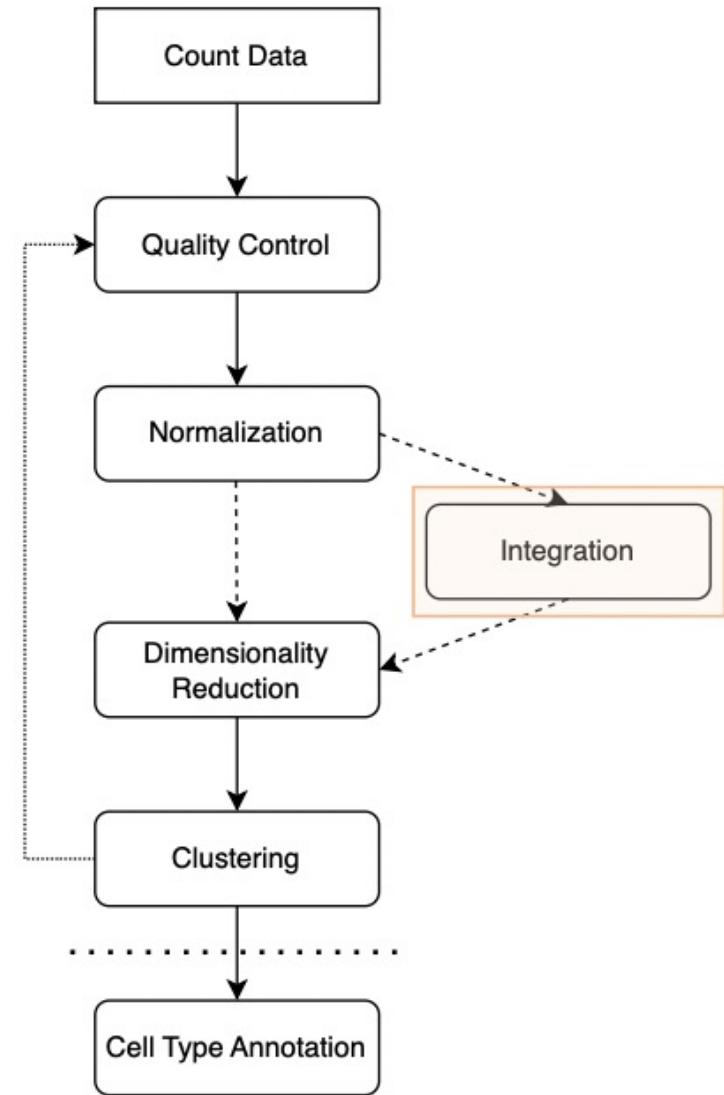
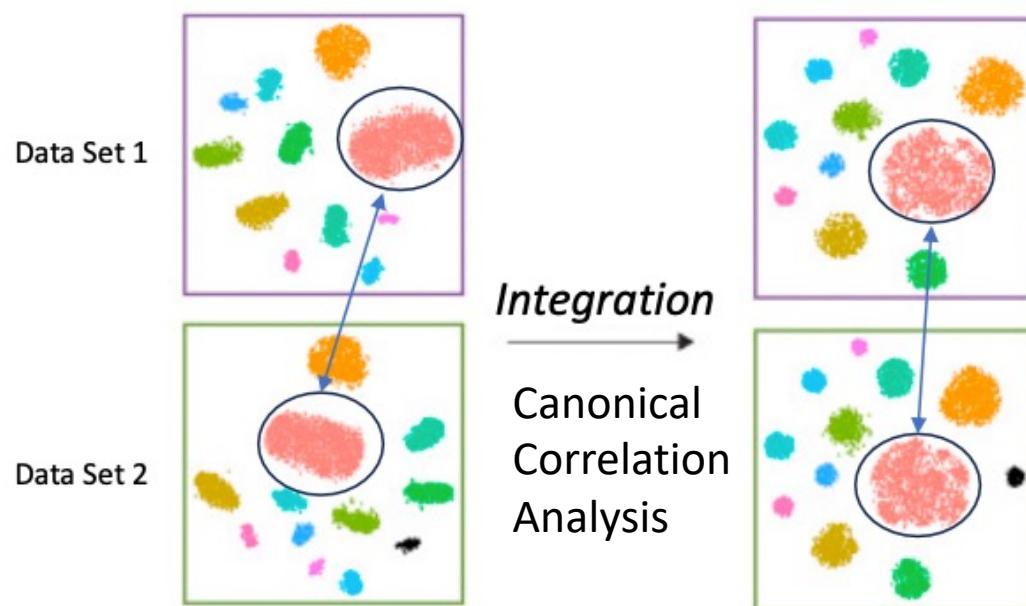


Query

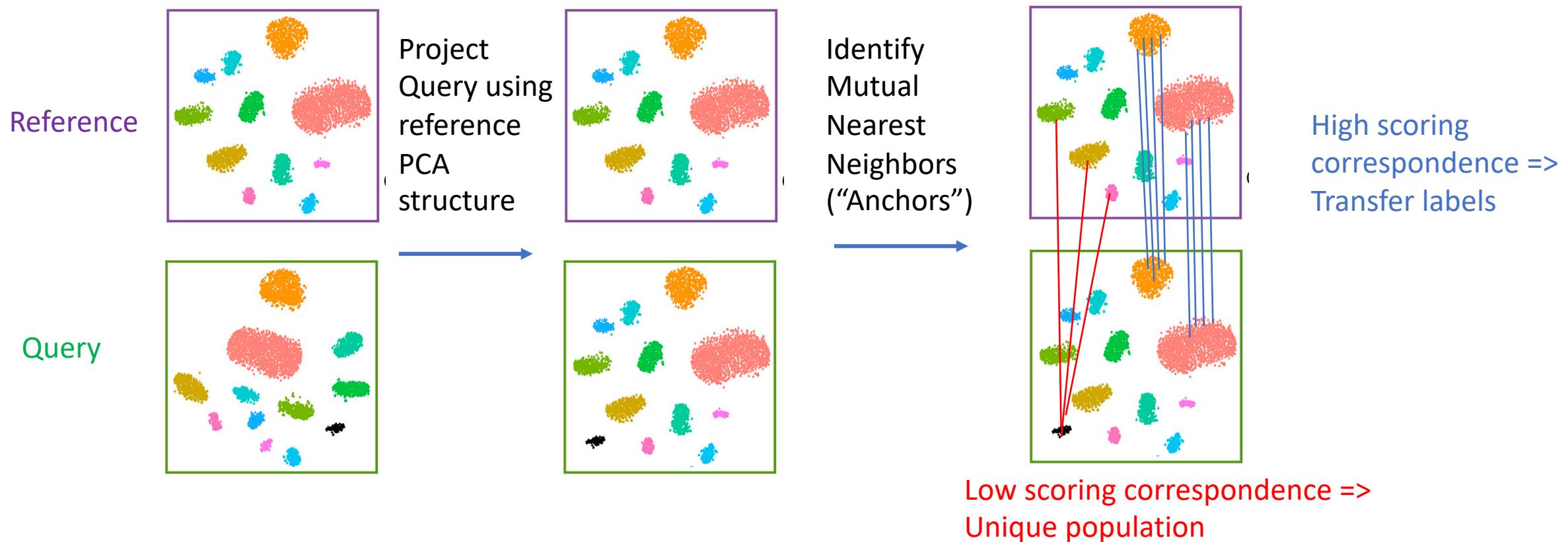


Integration (recap from workshop 1)

Transform our data such that the profiles of the same cell types are more similar across data sets

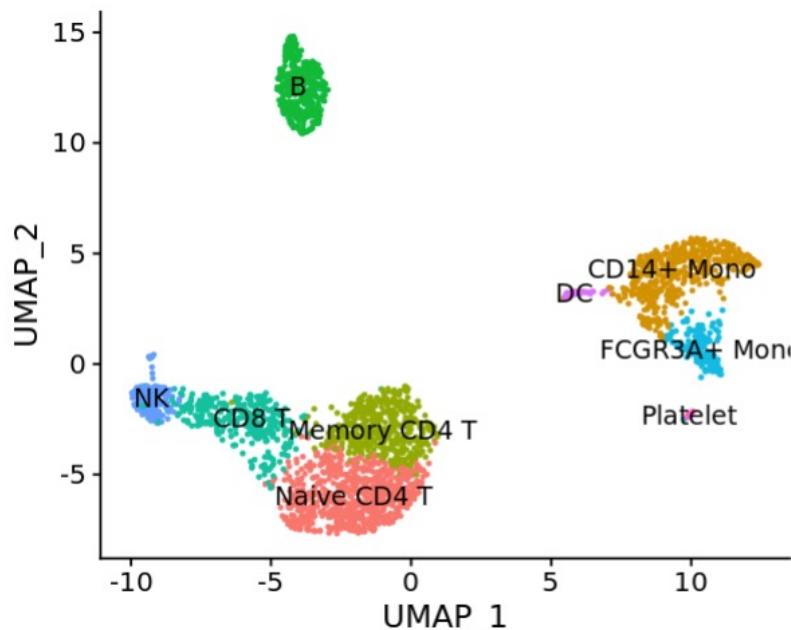


Seurat Integration mapping method



Where is the reference PCA?

Reference PBMC



@metadata

| | | |
|------------|-----|------------|
| | ... | Cell-type |
| ref-cell-1 | | T-cell |
| ref-cell-2 | | B-cell |
| ref-cell-3 | | Macrophage |

@assays

| RNA, SCT | | ref-cell-1 | ref-cell-2 | ref-cell-3 |
|----------|-------|------------|------------|------------|
| | | gene1 | 0 | 1.09 |
| gene1 | gene2 | 5.04 | 0 | 0 |
| gene2 | 1 | 0 | 0 | 0 |

@reductions

PCA

Feature loadings

| | PC_1 | PC_2 | PC_3 | ... |
|--------|---------------|--------------|---------------|-----|
| PPBP | -1.172699e-02 | 1.523301e-02 | -1.475722e-01 | |
| S100A9 | -1.186967e-01 | 2.094446e-02 | -2.730169e-02 | ... |
| IGLL5 | 8.678533e-03 | 5.451734e-02 | 5.480080e-02 | |
| LYZ | -1.195951e-01 | 1.637529e-02 | -1.491376e-02 | |

UMAP

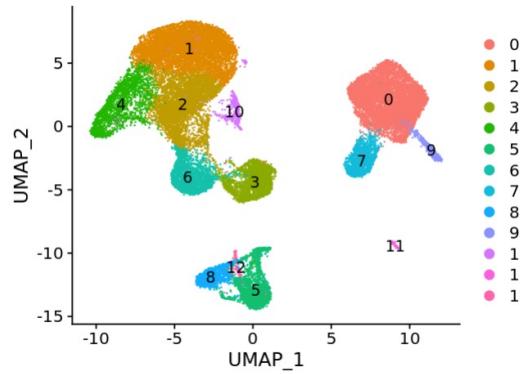
Cell embeddings

| | PC_1 | PC_2 | PC_3 | ... |
|------------------|------------|--------------|--------------|-----|
| AAACATACAACCAC-1 | 4.60604661 | -0.603719506 | -0.605242902 | ... |
| AAACATTGAGCTAC-1 | 0.16708090 | 4.544217123 | 6.451886683 | |

| | UMAP_1 | UMAP_2 |
|------------------|--------------|--------------|
| AAACATACAACCAC-1 | -3.198993068 | -4.242983556 |
| AAACATTGAGCTAC-1 | -4.163467985 | 10.733455443 |

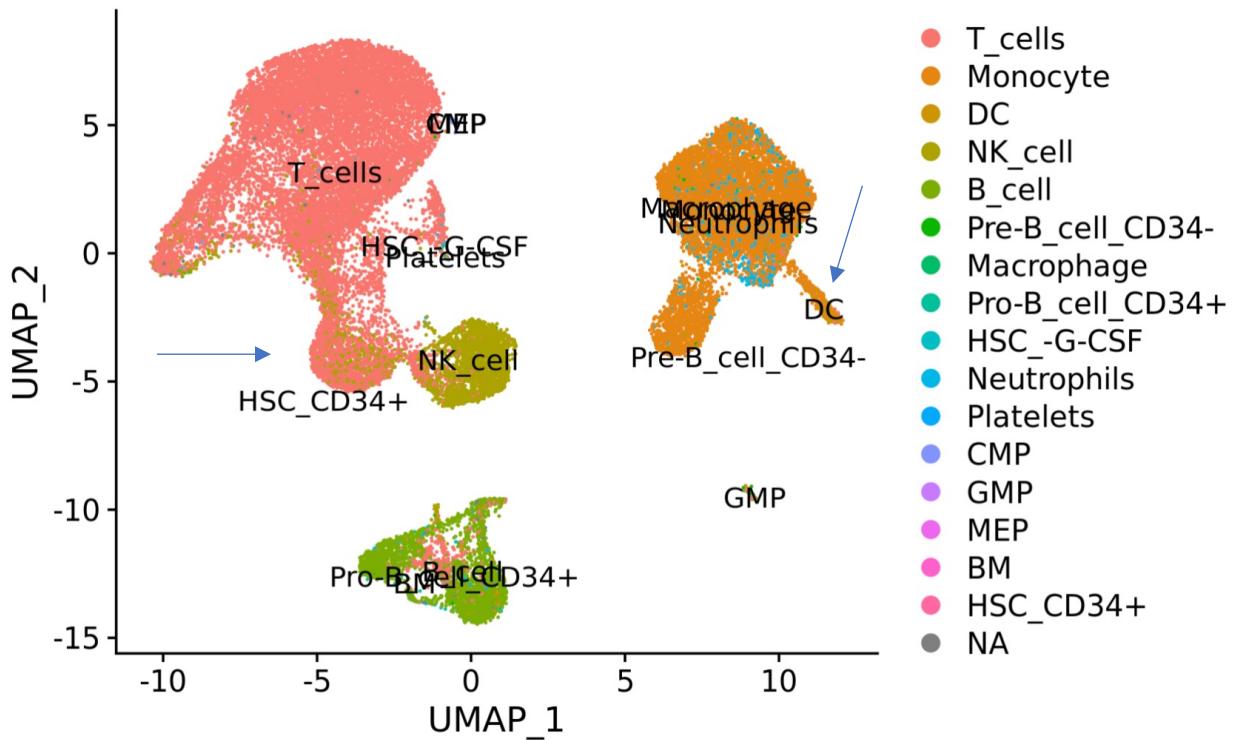
...

Query



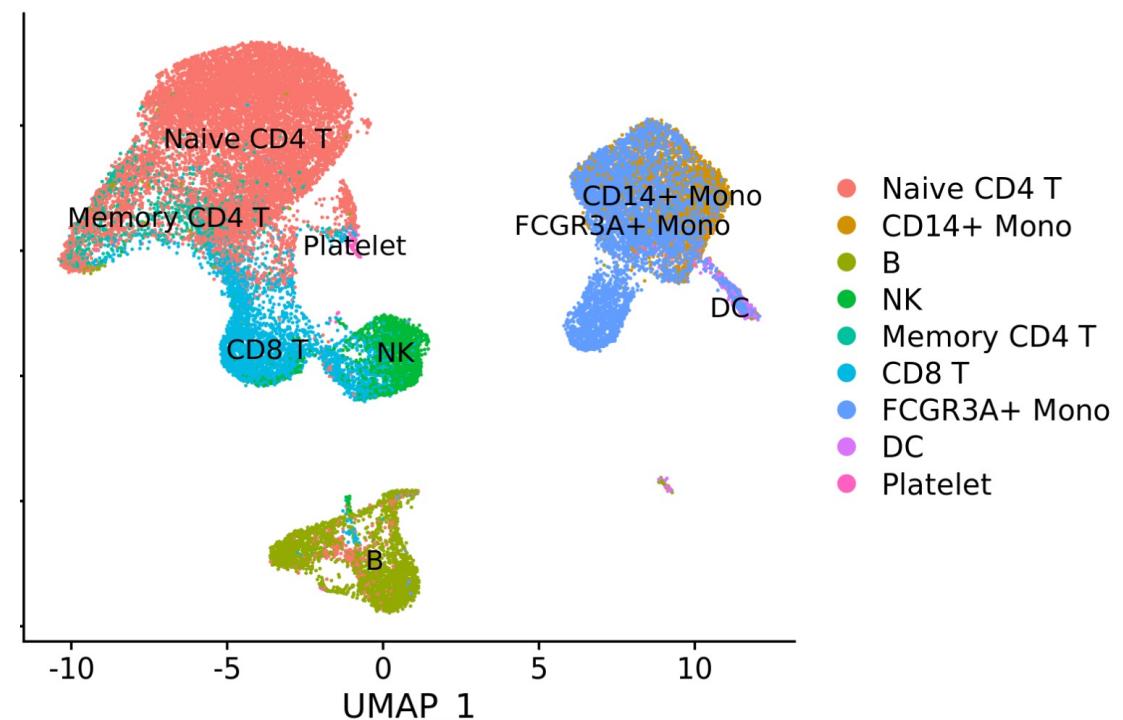
SingleR Results

Cell level, HPCA, label.main



Seurat Results

Human PBMC



Using single-cell datasets as a reference

Typically contains cells present in a single organism, tissue, condition and should be chosen to match your experiment as closely as possible.

10x Genomics reference database

The screenshot shows the 10x Genomics website interface. At the top, there's a dark header with the 10x Genomics logo, 'Products', and 'Research' links. Below the header, a breadcrumb navigation shows 'Support > Single Cell Gene Expression > Datasets'. The main content area features a title '3k PBMCs from a Healthy Donor' and a subtitle 'Single Cell Gene Expression Dataset by Cell Ranger 1.1.0'. It includes a brief description of PBMCs and a bulleted list of analysis details. At the bottom, it says 'Published on May 26, 2016' and 'This dataset is licensed under the Creative Commons Attribution license.'

Broad Single Cell Portal

The screenshot shows the Broad Institute Single Cell Portal homepage. The header includes the portal logo, a search bar, and links for 'Help & resources', 'Create a study', and 'Sign in'. A central circular graphic displays 'Featuring 605 studies 37,053,555 cells'. Below the header, there are search filters for 'Search studies' and 'Search genes', and a 'Search by text' input field. The main content area lists several study entries, such as 'Intertumoral lineage diversity and immunosuppressive transcriptional programs in well-differentiated gastroenteropancreatic neuroendocrine tumors' and 'Leptin receptor neurons in the dorsomedial hypothalamus input to the circadian feeding network'. Each entry includes a thumbnail, study ID, cell count, and a brief description.

Literature search

<https://tabula-sapiens-portal.ds.czbiohub.org/>

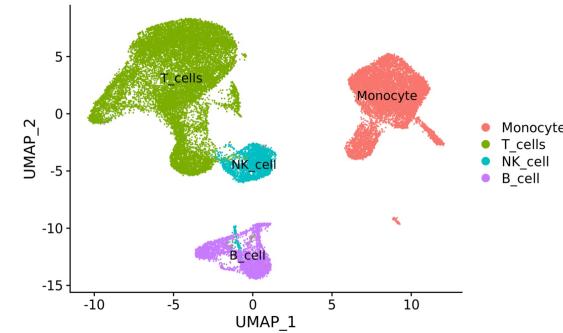
<https://chanzuckerberg.github.io/cellxgene-census/>

Pro: can match well

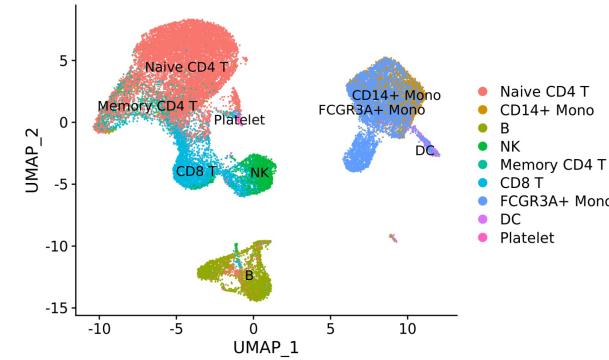
Con: may not be available, and takes a lot of pre-processing

Example workflow

1. Rough cell type mapping using label.main in Celldex.



2. Fine cell type mapping using a well-matched single-cell dataset.

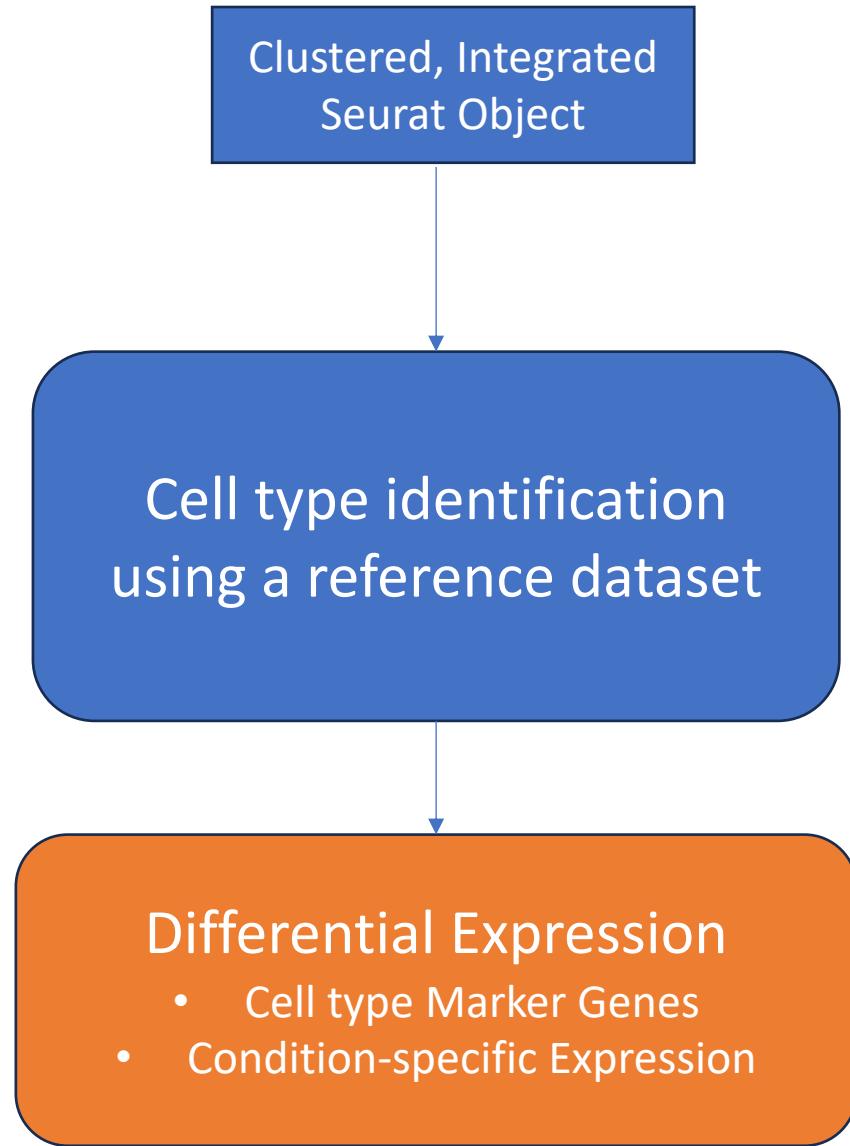


3. Known and calculated marker genes (next section)

What if there are no good references for your cell type?

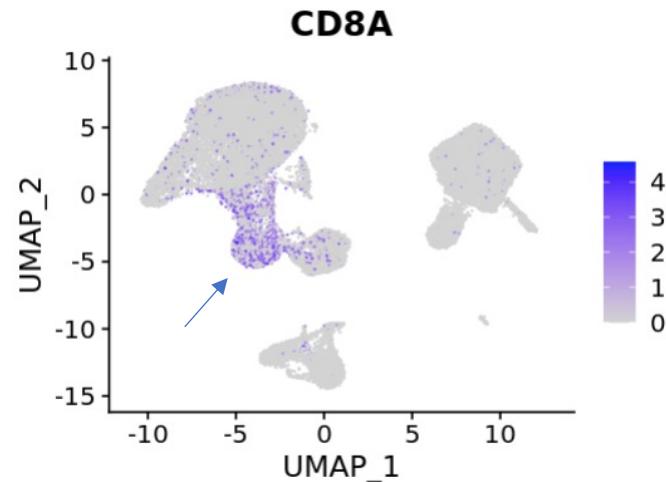
Questions?

Workshop day 2



Differential Expression

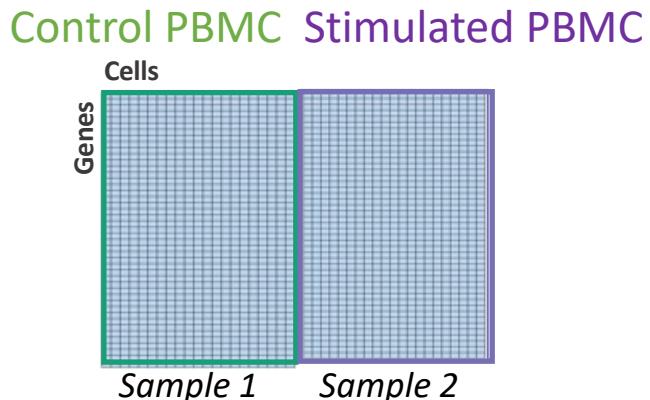
- 1) Genes that are overexpressed in one cell-type compared to all other cell-types



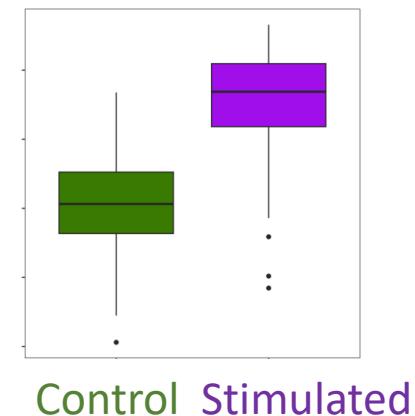
“Marker Genes”

- can be computed within a single sample

- 2) Gene that are statistically differently expressed between phenotypes or conditions, usually within a cell-type



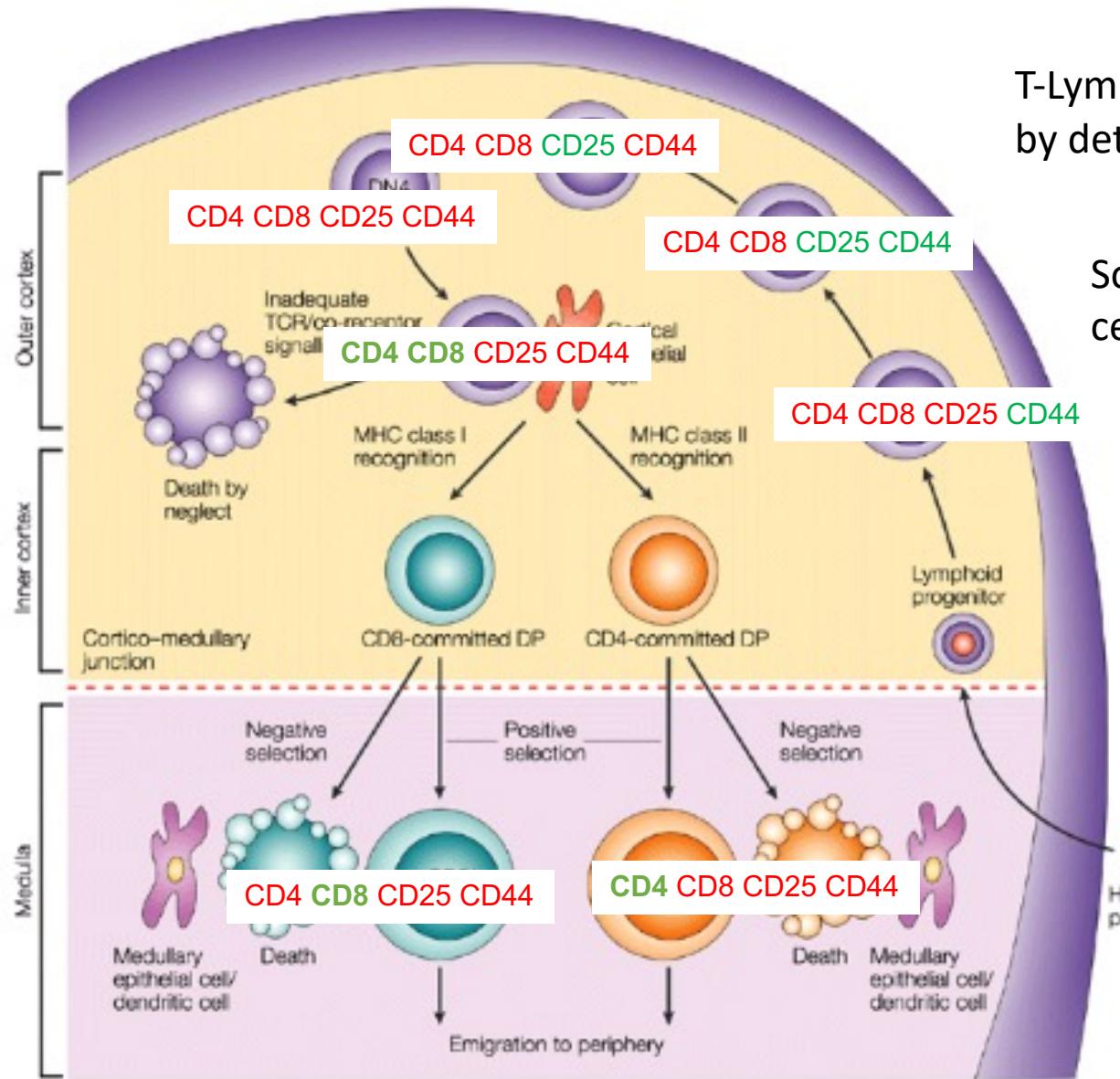
Gene1 Expression



“Differentially Expressed Genes”

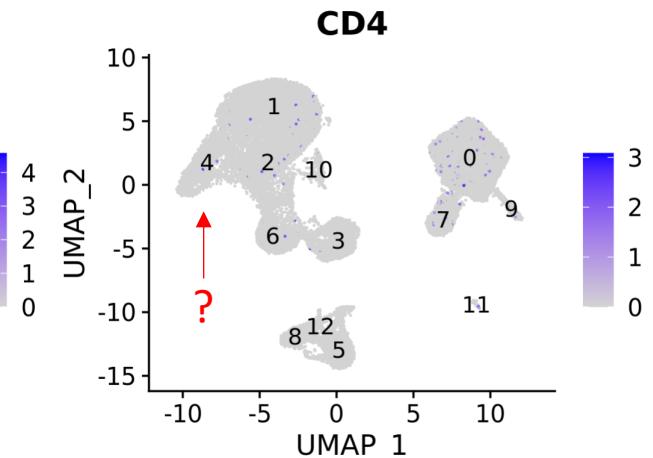
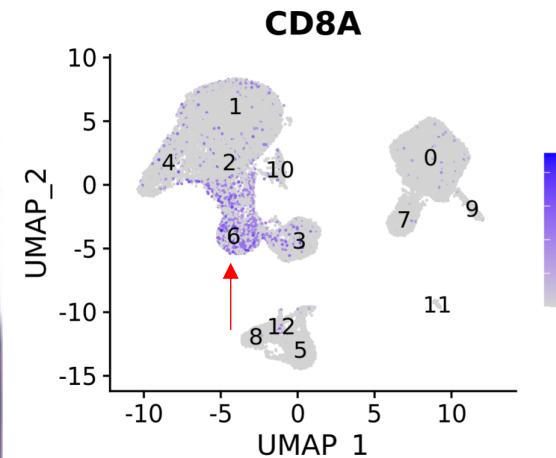
- typically requires biological replicates and modeling of inter-replicate variability

Cell-type Marker Genes



T-Lymphocyte development in the Thymus was historically defined by detection of 4 cell surface proteins **CD4 CD8 CD25 CD44**.

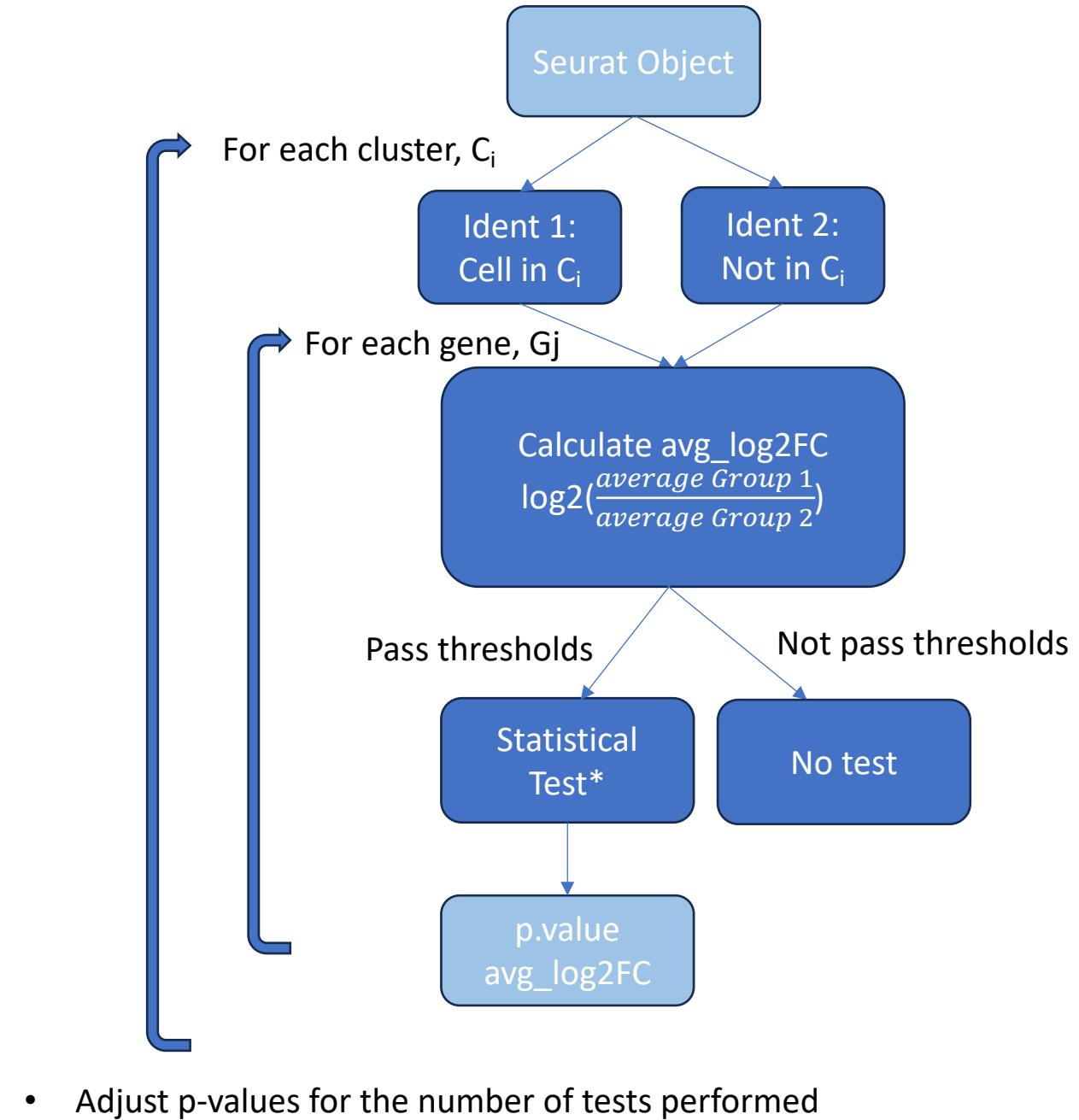
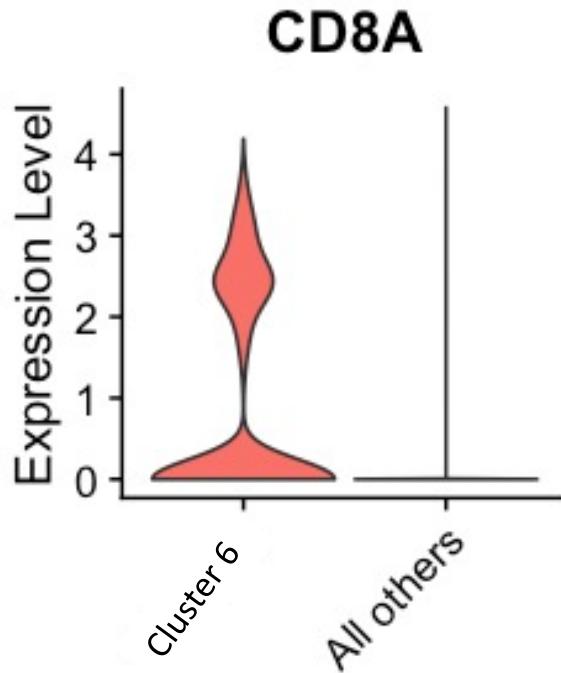
Some of these markers must be redefined for scRNAseq because cell surface protein quantity != gene transcript quantity



Find All Markers

Genes that are more highly expressed by cells in one cluster compared to all other cells

```
FindAllMarkers(seurat_object,  
    test.use = "wilcox",  
    logfc.threshold = 0.25,  
    min.pct = 0.1,  
    only.pos = TRUE,  
    ...)
```



*more on this soon

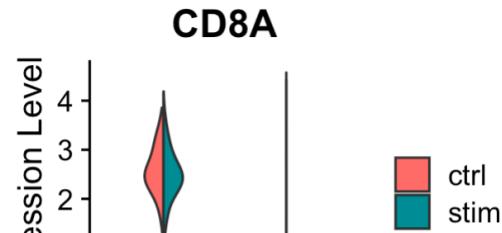
Find Conserved Markers

- Find Genes that are markers in two groups

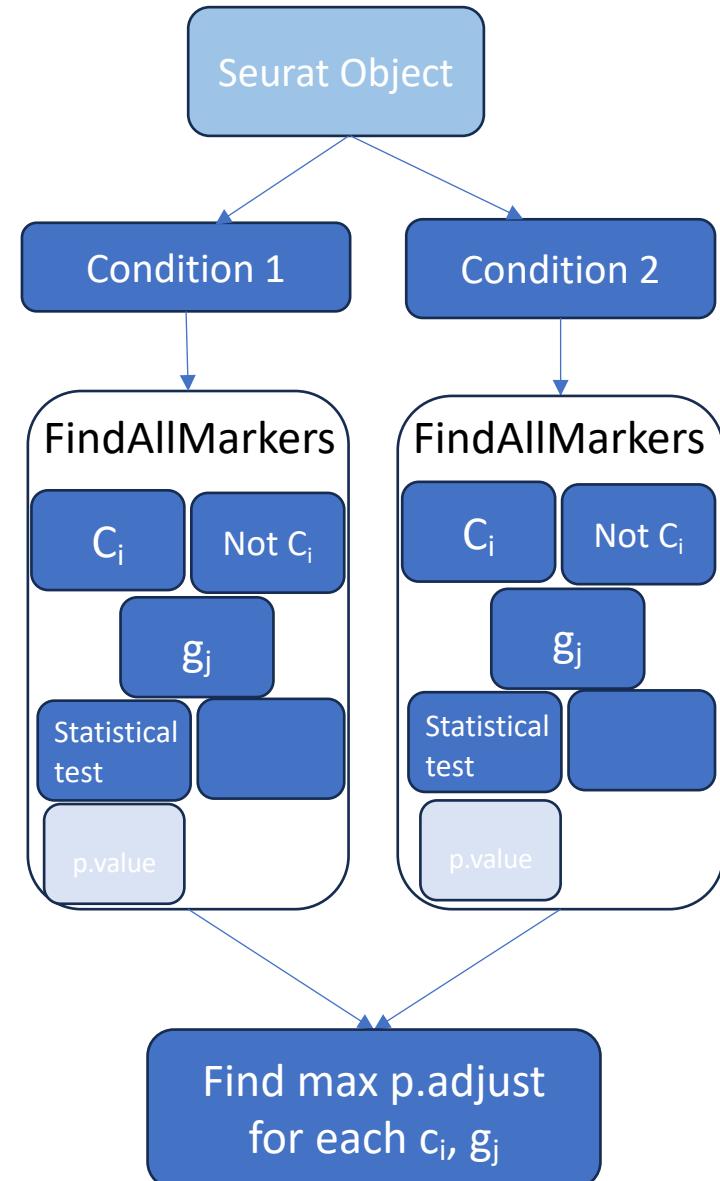
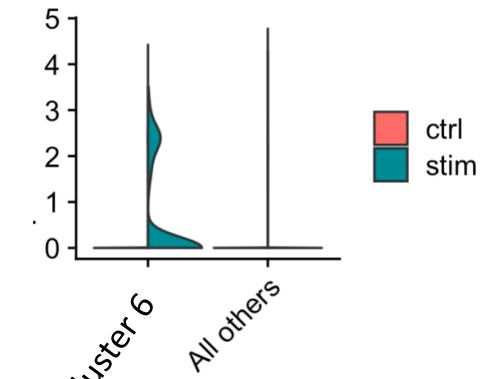
```
FindConservedMarkers(seurat_object,  
    ident.1 = 6,  
    grouping.var = "sample",  
    ...)
```

- Split cells by grouping.var

Cluster 6 marker
in **both conditions**



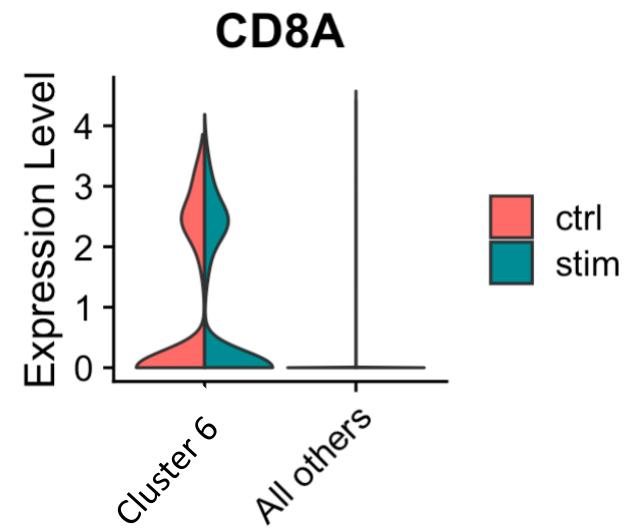
Cluster 6 marker
in **stim only**



Markers: Which statistical test?

```
FindConservedMarkers(seurat_object,  
  ident.1 = 6,  
  grouping.var = "sample",  
  test.use = "wilcox",  
  ...)
```

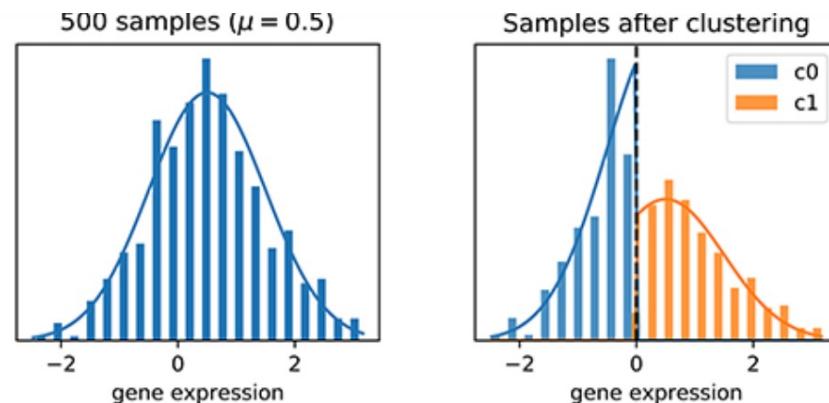
There are 9 options, we'll discuss two.



Interpretation of p-values

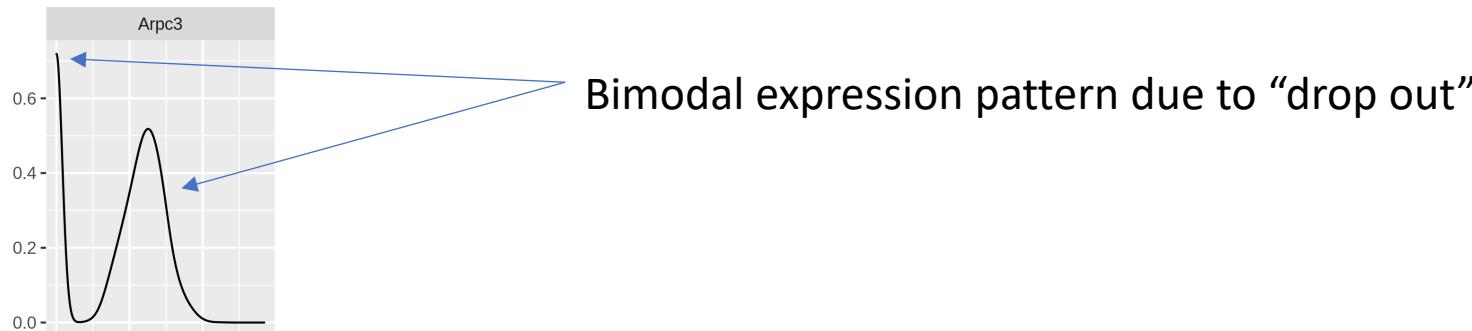


Most statistical tests are designed for hypotheses to be defined before experiments are carried out. In finding cluster marker genes, testing is performed **on the same data that is used to extract the clusters in the first place**. Thus, there will be some genes that are differentially expressed between clusters by construction.



Markers: Which statistical test?

- Wilcoxon Rank Sum test - default
 - Assumes independence of samples (cells)
 - no assumption of statistical distribution
- MAST (Model-based Analysis of Single-cell Transcriptomics):
 - Assumes independence of samples (cells)
 - Assumes counts are zero-inflated, *usually* true for single-cell data (though see ref)



- Allows adjusting for additional covariates such as batch, sex, sample, through the `latent.vars` parameter

Markers: Results

```
cluster_6_markers = FindConservedMarkers(seurat_object,
                                         ident.1 = 6,
                                         grouping.var = "sample",
                                         only.pos = TRUE)
```

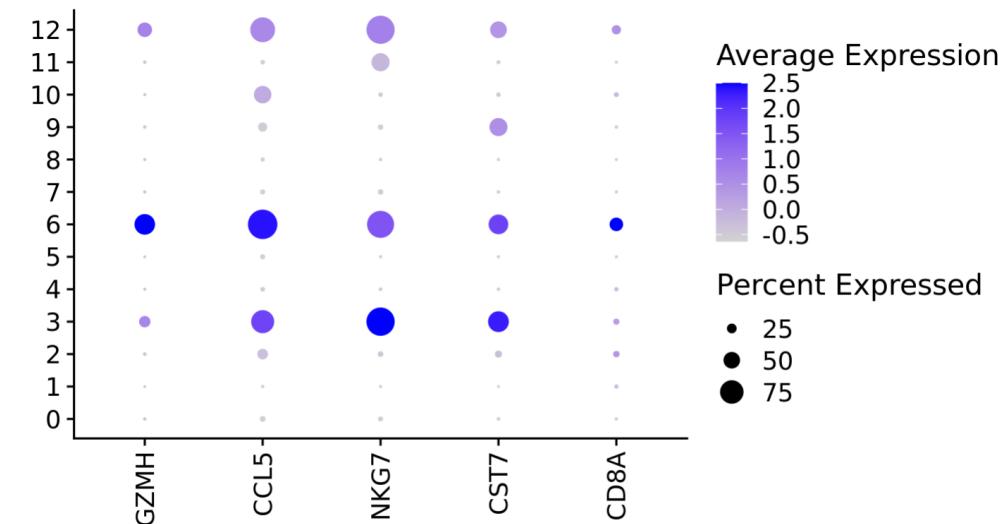
| | stim_p_val | stim_avg_log2FC | stim_pct.1 | stim_pct.2 | stim_p_val_adj | ctrl_p_val | ctrl_avg_log2FC | ctrl_pct.1 | ctrl_pct.2 | ctrl_p_val_adj | max_pval | minimump_p_val |
|-------------|--------------|-----------------|------------|------------|----------------|---------------|-----------------|------------|------------|----------------|---------------|----------------|
| GNLY | 0.000000e+00 | 1.2250346 | 0.666 | 0.127 | 0.000000e+00 | 0.000000e+00 | 1.1372563 | 0.603 | 0.111 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 |
| CD8A | 0.000000e+00 | 2.0559922 | 0.393 | 0.055 | 0.000000e+00 | 0.000000e+00 | 2.1068908 | 0.387 | 0.058 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 |
| PRF1 | 0.000000e+00 | 1.6820287 | 0.578 | 0.110 | 0.000000e+00 | 1.916110e-165 | 1.3894355 | 0.260 | 0.046 | 2.695008e-161 | 1.916110e-165 | 0.000000e+00 |

Multiple p-values!

- `stim_p_val`, `ctrl_p_val` - the raw p-value for each group
- `stim_p_val_adj`, `ctrl_p_val_adj` - the FDR adjusted p-value for each group
- `max_pval` - largest p-value of the two groups
- `minimump_p_val` - combined p-value, using package `metap`

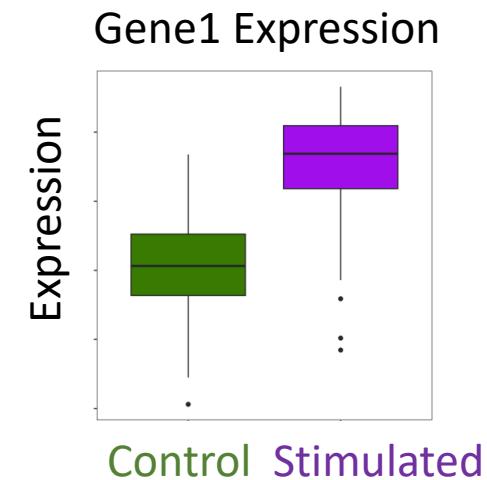
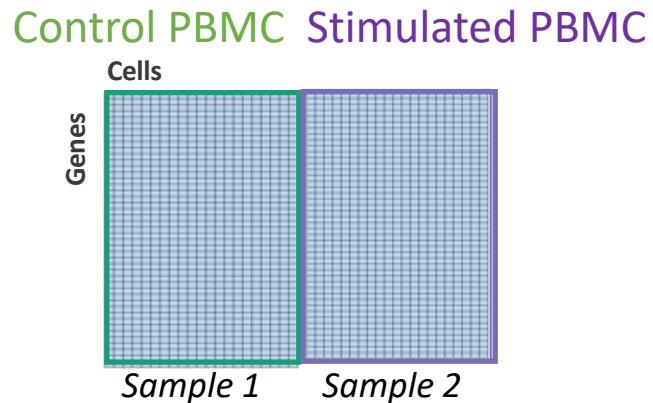
Fold changes

- `stim_avg_log2FC` and `ctrl_avg_log2FC`: log fold-change of the average expression between the two cluster 6 and all other clusters. Positive values indicate that the feature is more highly expressed in the cluster 6.



Differential Expression in scRNAseq

1) Genes that are overexpressed in one cell-type compared to all other cell-types

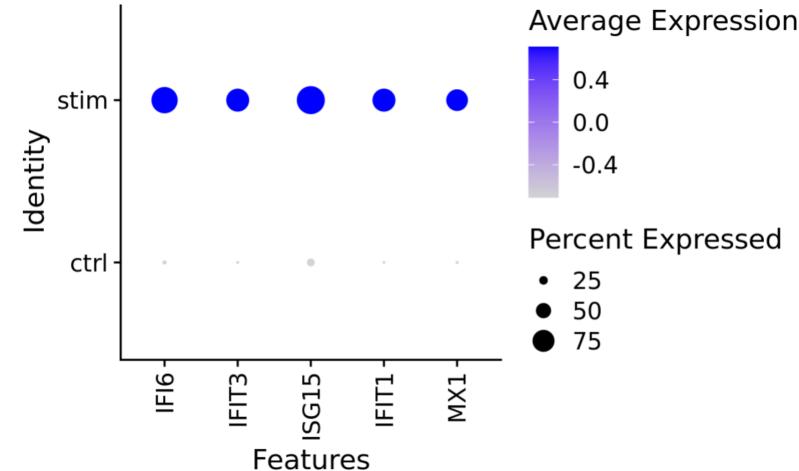


"Differentially Expressed Genes"
- typically requires biological replicates
and modeling of inter-replicate
variability

Differential Expression: Which statistical test?

With no replicates, naïve method: Seurat FindMarkers with test.use = "wilcox" or "MAST"

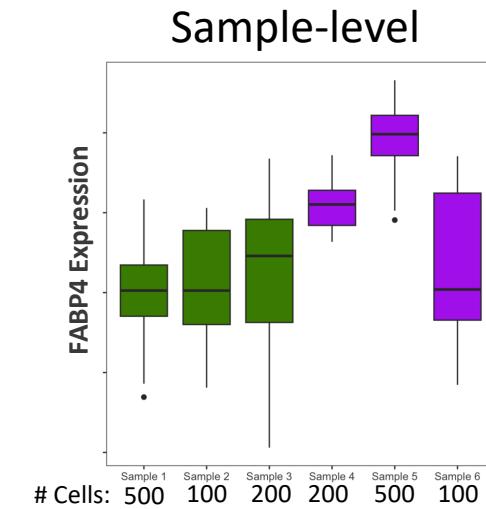
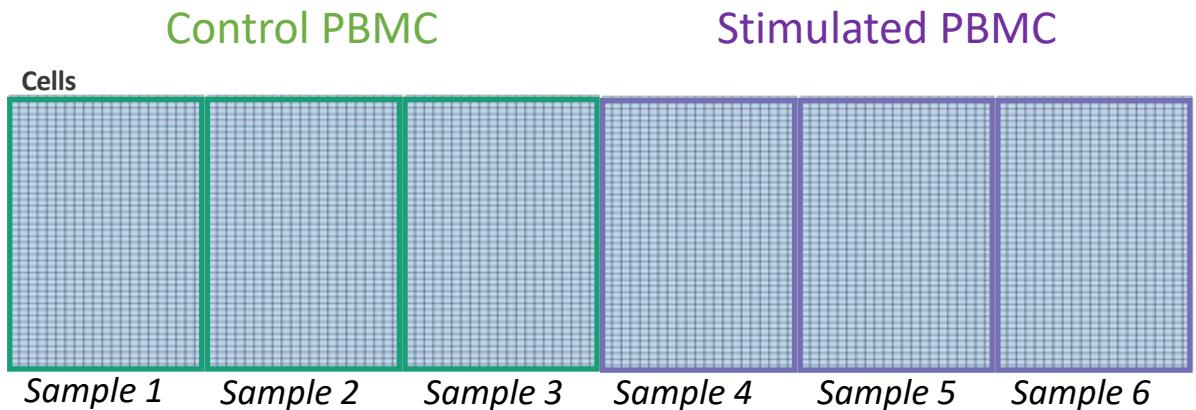
```
deg = FindMarkers(cluster_6,  
                  ident.1 = "stim",  
                  ident.2 = "ctrl",  
                  only.pos = FALSE)
```



Pseudoreplication bias

"Cells from the same individual share common genetic and environmental backgrounds and are not statistically independent; therefore, they are **subsamples or pseudoreplicates**. Thus, single-cell data have a hierarchical structure that many current single-cell methods do not address, leading to biased inference, highly inflated type 1 error rates, and reduced robustness and reproducibility."

A better study design



(Thanks Eric Reed, for the image)

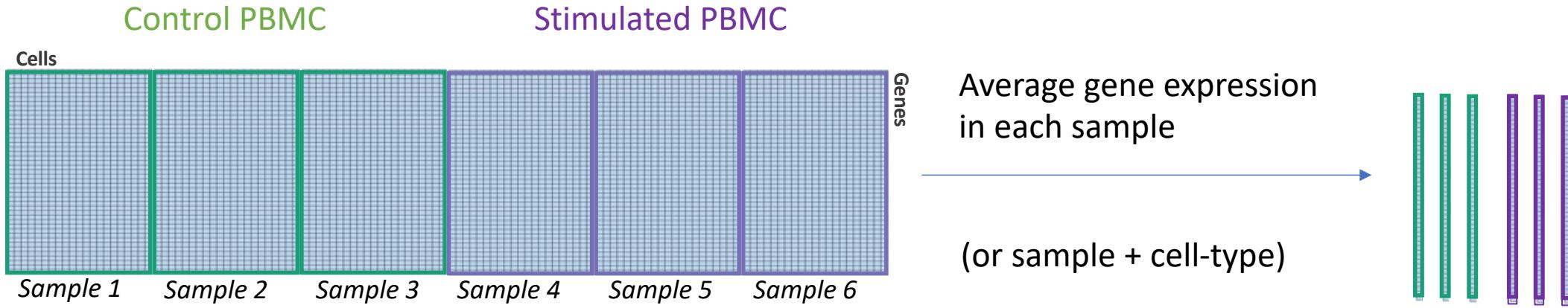
With replicates, there are options to model inter-individual variability:

- **Pseudobulk:** Transform data into bulk-RNAseq-like data by aggregating gene counts within each sample
- **Latent Variable:** Test whether the difference in gene expression between the groups can be explained by the difference in one or multiple latent variables (e.g. batch)
- **Mixed-model:** modeled as random effect



Helpful to talk to statistician to choose

Pseudobulk approach



Bulk-RNAseq differential expression tools can then be used, e.g. linear regression with covariates. For each gene, condition i:

Expression
Values

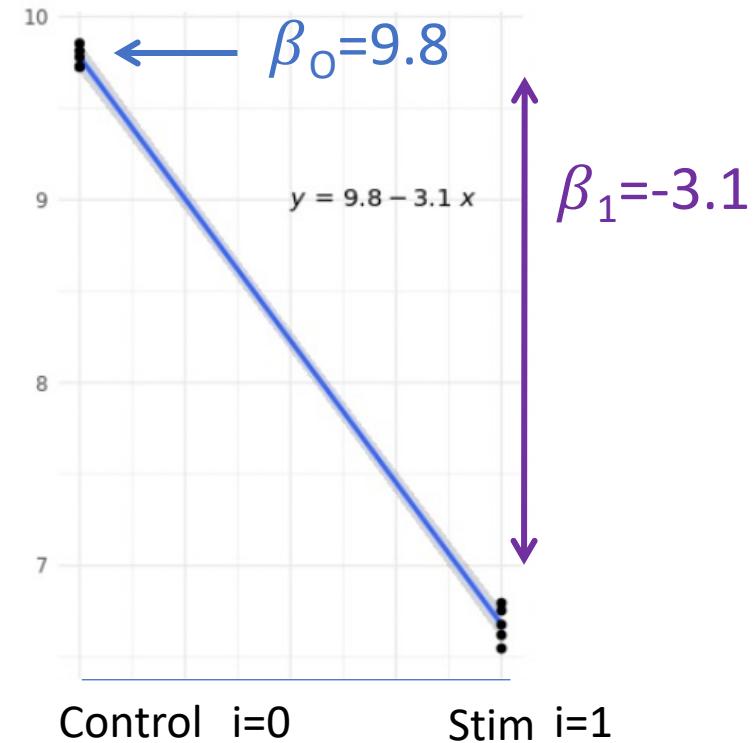
$$y \sim N(\mu_i, \sigma^2)$$

$$\mu_i = \beta_0 + \beta_1 X_i$$

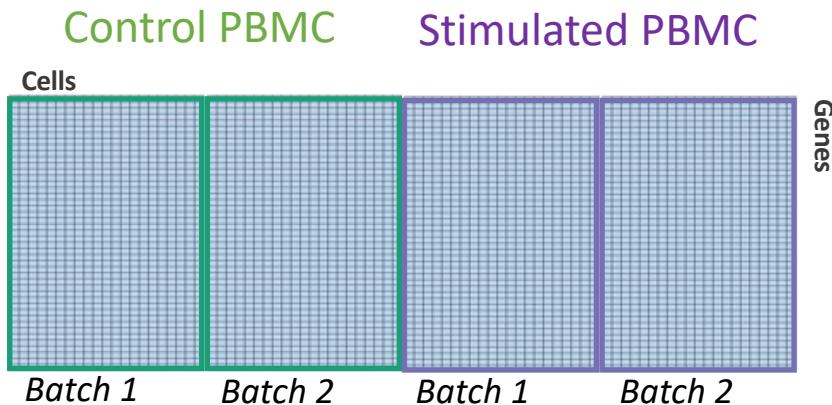
Intercept

Slope:
difference
between
conditions

Condition:
0-control
1-stim



Latent Variables for batch effects



In this case, our samples were grouped into batches.
e.g. subject tested in both conditions.

Can difference in gene expression between the conditions be explained by batch?

For each gene, the expression in cell c:

$$\text{Expression Values} \quad y_c \sim N(\mu_{ij}, \sigma^2)$$

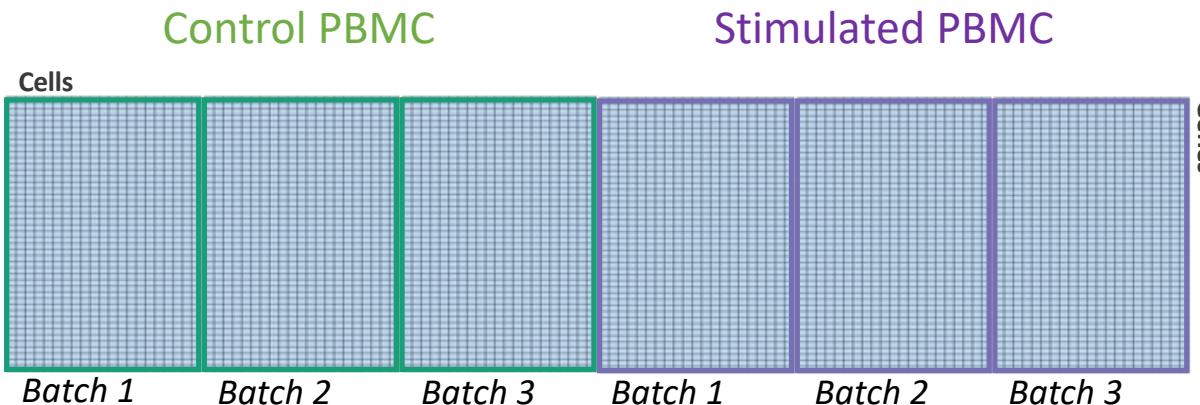
$$\mu_{ij} = \beta_o + \beta_1 X_i + \beta_2 X_j$$

| Intercept | Condition effect | Condition: 0-control | Batch effect | Batch 0 - Batch 1 | Batch 1 - Batch 2 |
|-----------|------------------|----------------------|--------------|-------------------|-------------------|
| | | | | | |

Seurat can do this:

```
deg = FindMarkers(cluster_6,  
                  ident.1 = "stim",  
                  ident.2 = "ctrl",  
                  test.use = "MAST",  
                  latent.vars = "batch"  
                  only.pos = FALSE)
```

Linear mixed-effect models are flexible



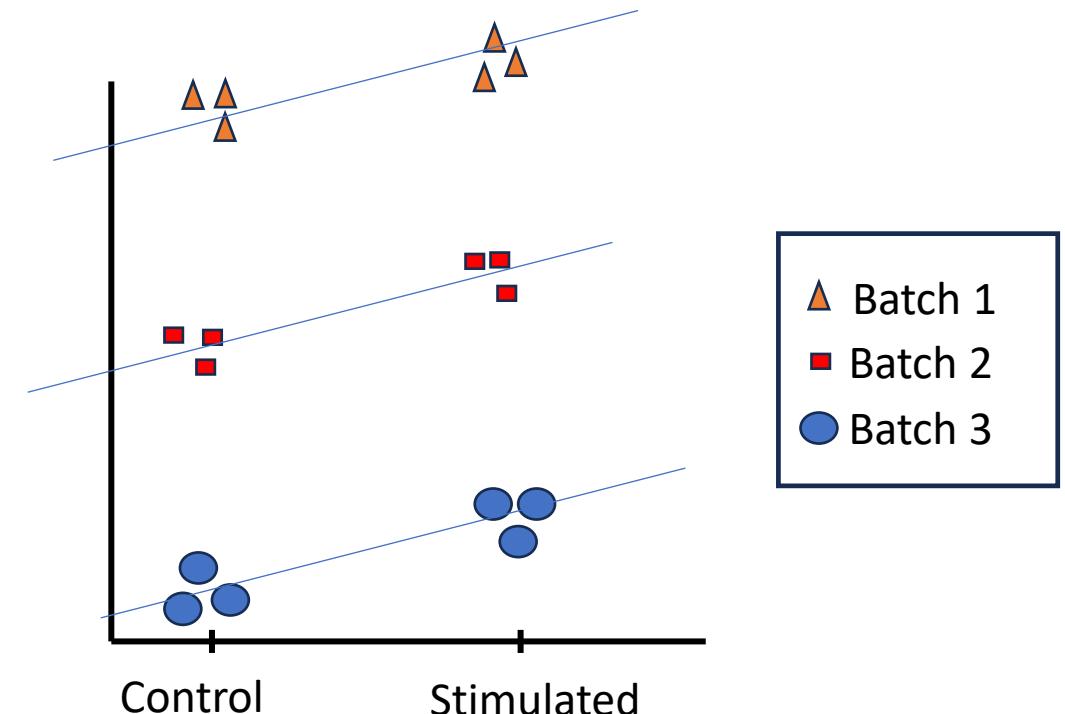
$$y_{ij} \sim N(\mu_{ij}, \sigma^2),$$
$$\gamma_j \sim N(0, \varphi^2)$$

$$\mu_{ij} = \beta_1 X_i + \gamma_j W_j$$

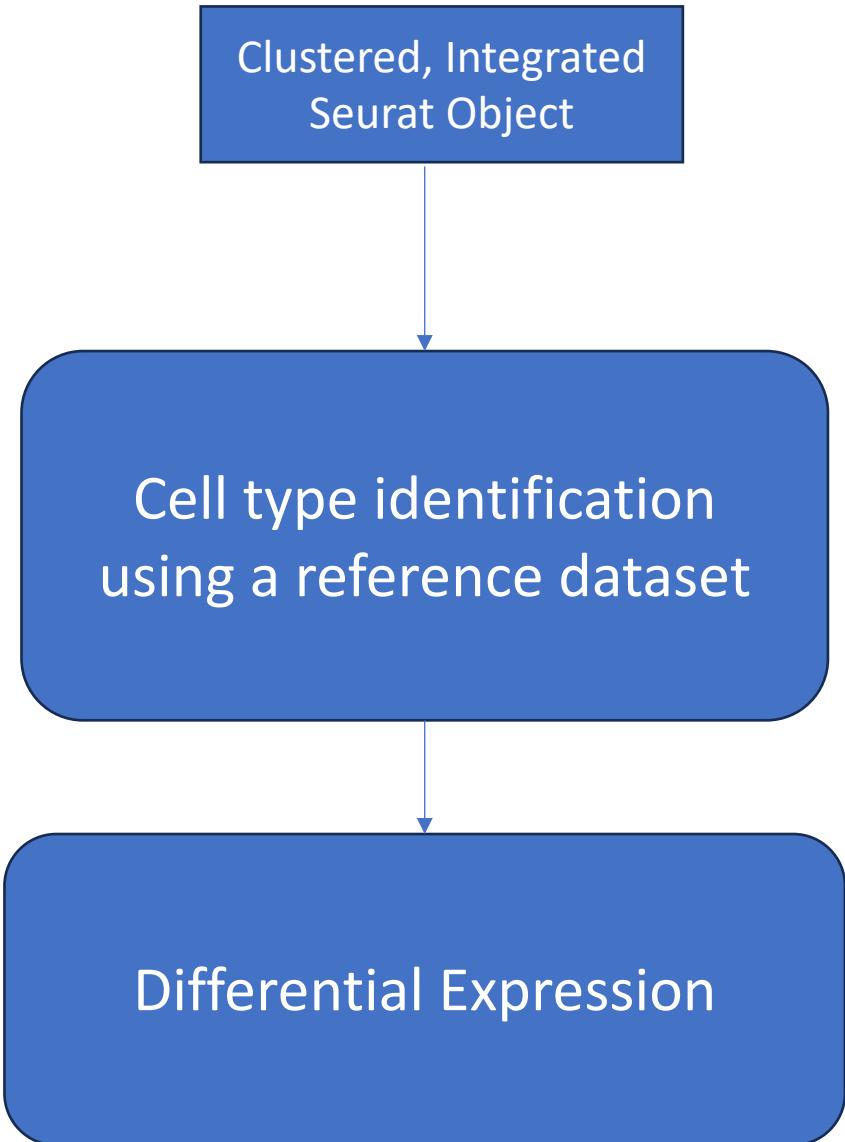
Random coefficient for batch j

Design matrix for random effects of batch j, condition i

Mixed-effect models allow us to model correlations within multiple batches and between conditions.



Wrap Up



- Choosing a method and reference
 - SingleR correlation with bulk database
 - Integration mapping with matching scRNAseq
- Cell type Marker Genes
- Condition-specific Expression



References for statistical methods

- Zhang et al Cell. Syst. 2020 [Valid post-clustering differential analysis for single-cell RNA-Seq.](#)
- Zimmerman 2021 [A practical solution to pseudoreplication bias in single-cell studies](#)
- Juntilla et al Briefings in Bioinformatics 2022 [Benchmarking methods for detecting differential states between conditions from multi-subject scRNAseq data.](#)
- Das 2021 [Statistical methods for analysis of single-cell RNA-sequencing data](#)

Bayesian mixed effects regression example:

Karagiannis et al 2023 [Multi-modal profiling of peripheral blood cells across the human lifespan reveals distinct immune cell signatures of aging and longevity](#)

[UCLA Stats: Mixed-effects model tutorial](#)