



Introduction to single cell RNA-seq

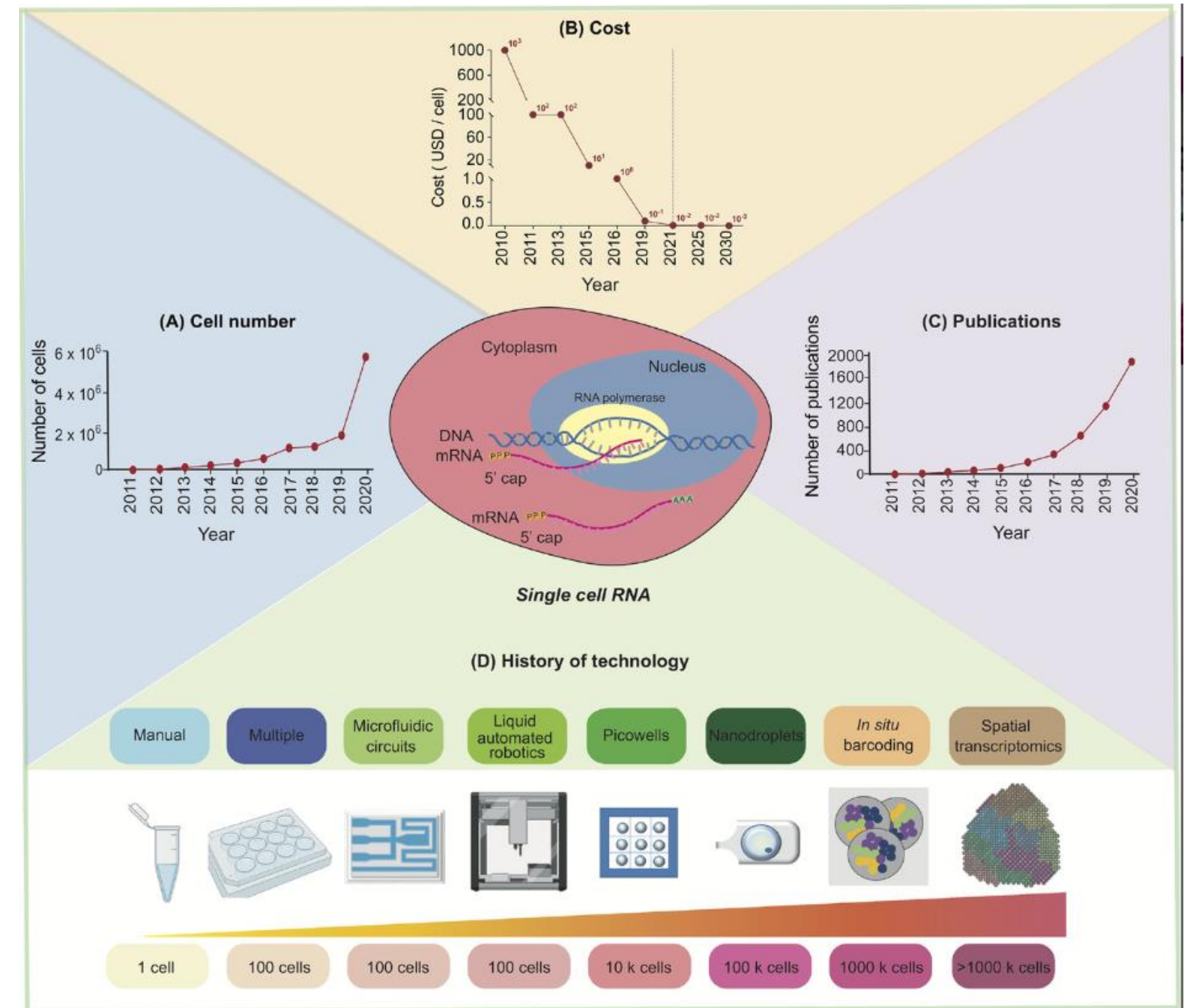
Xabier Bujanda Cundin
17/07/2023

Why single cell?

- **Cell:** the smallest unit that can live on its own and that makes up all living organisms and the tissues of the body
- Identification of cell types and state in a heterogeneous sample
- Assess differences between and within cell types

Challenges:

- Huge amount of data to handle
- Low read depth
- High variability between cells/samples
- Data processing, analysis, presentation and interpretation



scRNA-seq data can answer multiple biological questions

- Cell populations
- Expression differences cell clusters and conditions
- Time-series gene expression development

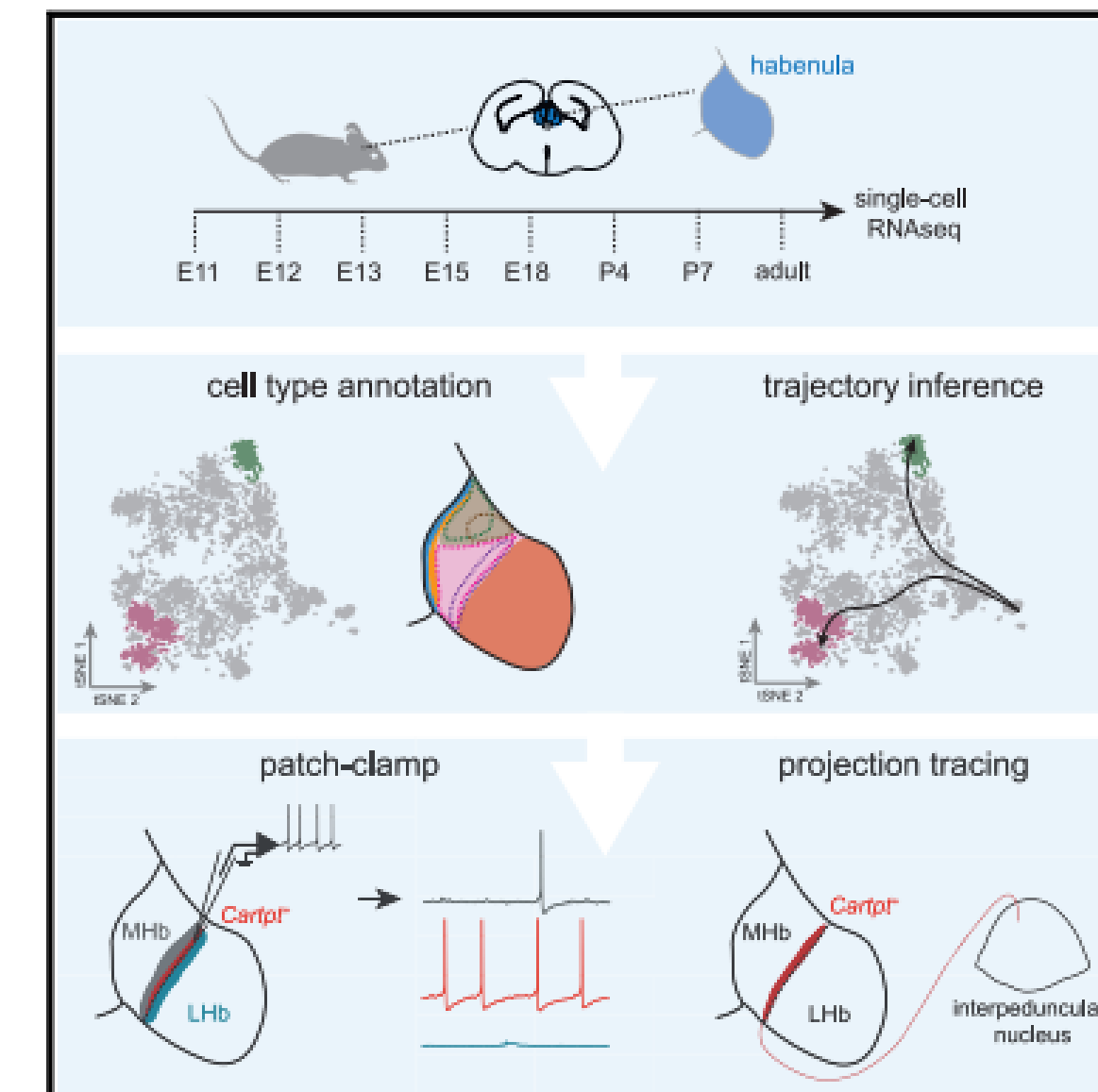
But...

- Low number of transcripts can be detected
- Some cell types may not be compatible for some scRNA-seq techniques
- Low expressed transcript may not be detected

Cell Reports

Molecular signatures and cellular diversity during mouse habenula development

Graphical abstract



Highlights

- scRNA-seq reveals cellular heterogeneity of the developing mouse habenula
- Trajectory inference reconstructs developmental paths of habenula cell types
- *Cartpt* marks physiologically distinct neurons that innervate the dorsal IPN
- Developing habenular cell types align to risk loci of human psychiatric disorders

Authors

Lieke L. van de Haar, Danai Riga, Juliska E. Boer, ..., Frank J. Meye, Onur Basak, R. Jeroen Pasterkamp

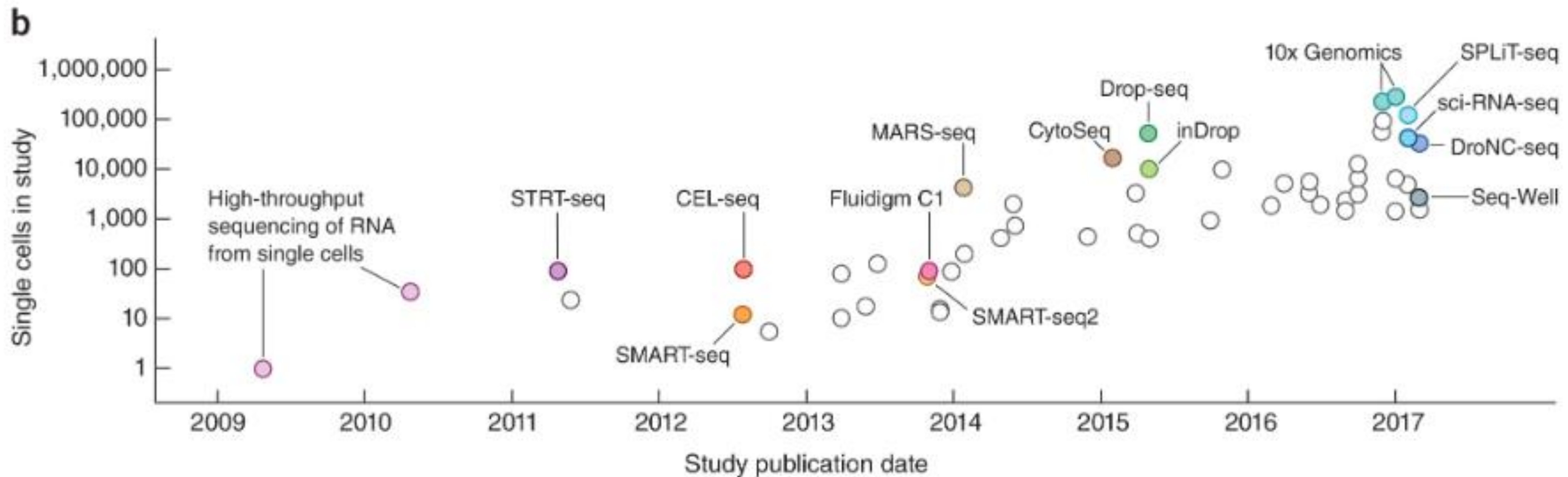
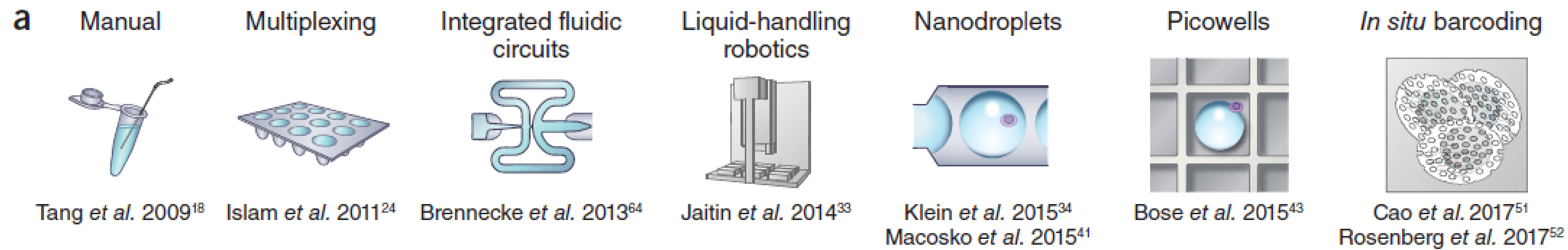
Correspondence

r.j.pasterkamp@umcutrecht.nl

In brief

van de Haar et al. characterize cell populations of the developing mouse habenula using single-cell RNA sequencing. They reconstruct the developmental trajectories of molecularly distinct neuronal subsets and establish physiological and connectivity properties of *Cartpt* neurons. Their analysis indicates a possible link between developing habenula neurons and human psychiatric disorders.

Single cell RNA-seq technologies



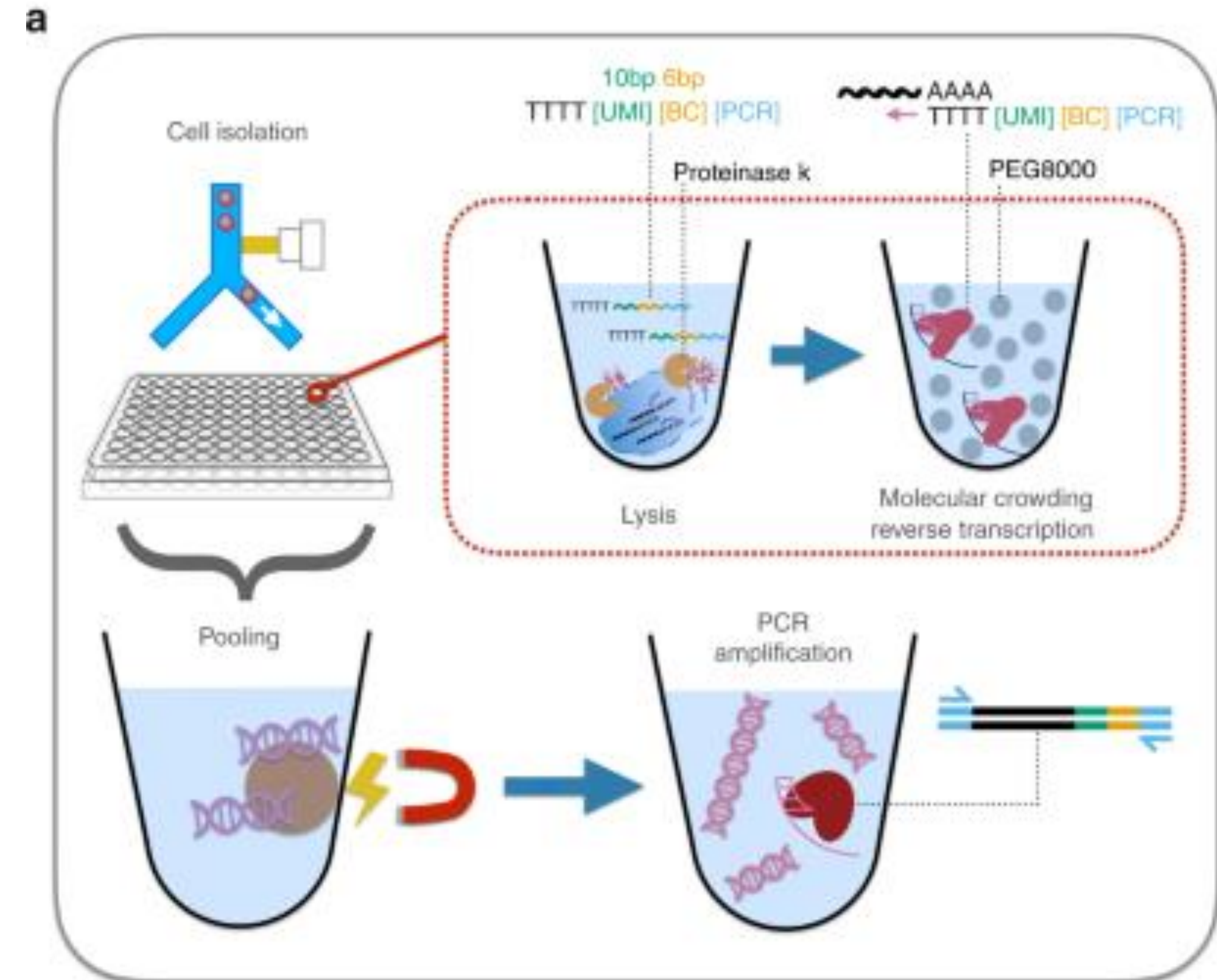
Single cell RNA-seq technologies

TABLE 1 Comparison of different single-cell RNA sequencing (scRNA-seq) technology and experimental protocols

Platforms	Isolation strategies	Tissue	Cell numbers	Targets	UMI	Amplification methods	Region	Published year
Smart-seq	FACS	Dissociated cell	Hundreds of cells	/	×	PCR	Full-length	2012
Smart-seq2	FACS	Dissociated cell	Hundreds of cells	/	×	PCR	Full-length	2013
Fluidigm C1	Micro-fluidic	Dissociated cell	Hundreds of cells	No poly(A) minus RNA detection	×	PCR	Full-length	2013
Drop-seq	Microdroplets	Dissociated cell	Large number of cell	No poly(A) minus RNA detection	√	PCR	3' end	2015
10x Genomics	Microdroplets	Dissociated cell	Large number of cells	No poly(A) minus RNA detection	√	PCR	3' end	2016
MATQ-seq	FACS	Dissociated cell	Hundreds of cells	No poly(A) minus RNA detection	√	PCR	Full-length	2017
Seq-Well	Micro-fluidic	Dissociated cell	Large number of cells	No poly(A) minus RNA detection	√	PCR	3' end	2017
CEL-seq	FACS	Dissociated cell	Hundreds of cells	No poly(A) minus RNA detection	√	IVT	3' end	2012
MARS-seq	FACS	Dissociated cell	Hundreds of cells	No poly(A) minus RNA detection	√	IVT	3' end	2014
inDrop-seq	Microdroplets	Dissociated cell	Large number of cell	No poly(A) minus RNA detection	√	IVT	3' end	2015
DNBelab C4	Microdroplets	Dissociated cell	Large number of cells	No poly(A) minus RNA detection	√	PCR	3' end	2019

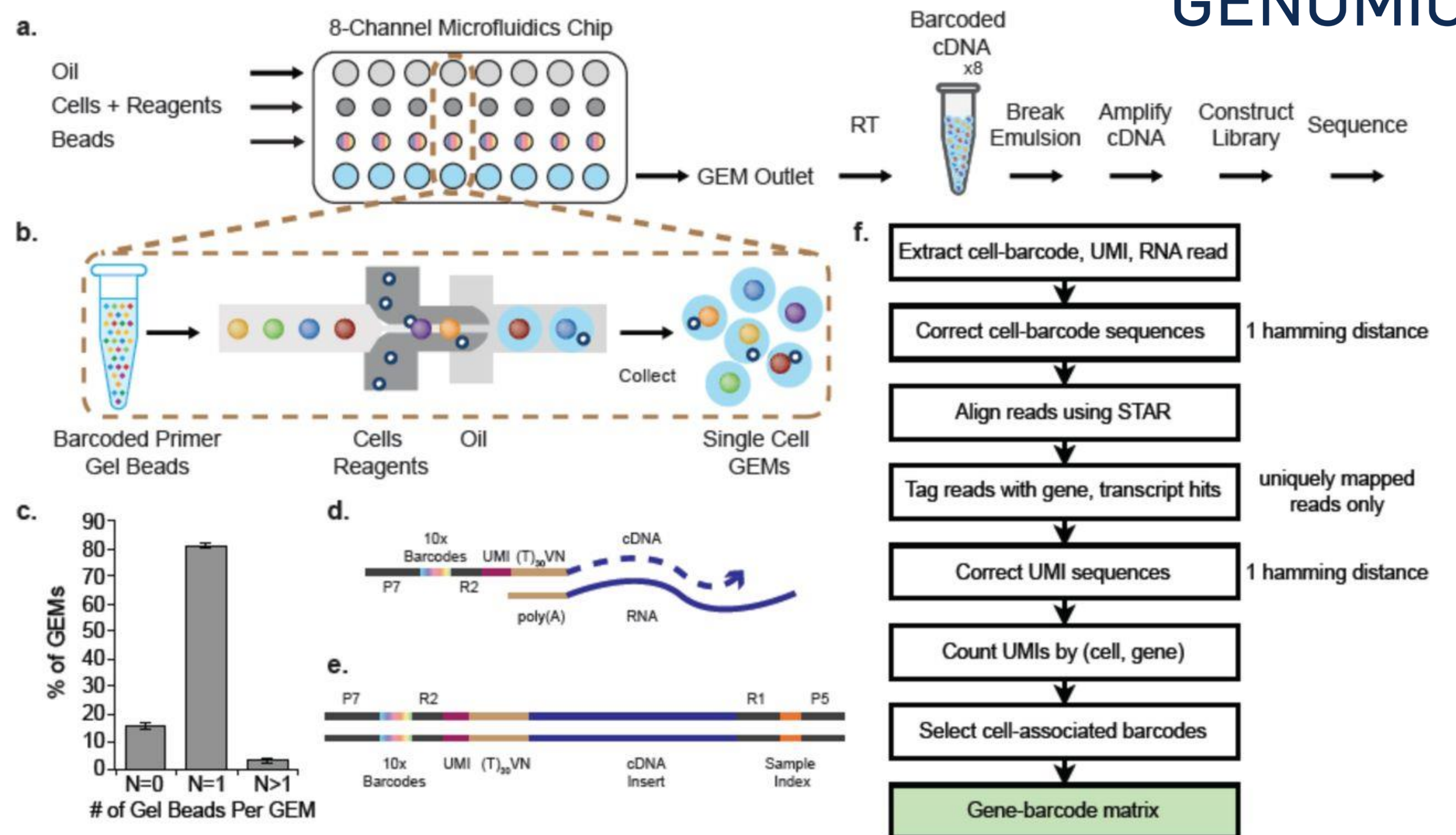
SMART-seq technology

- Cell separation through FACS
- Library preparation by cell
- Poor depth
- Whole gene sequencing (isoforms identification)



10X Genomics technology

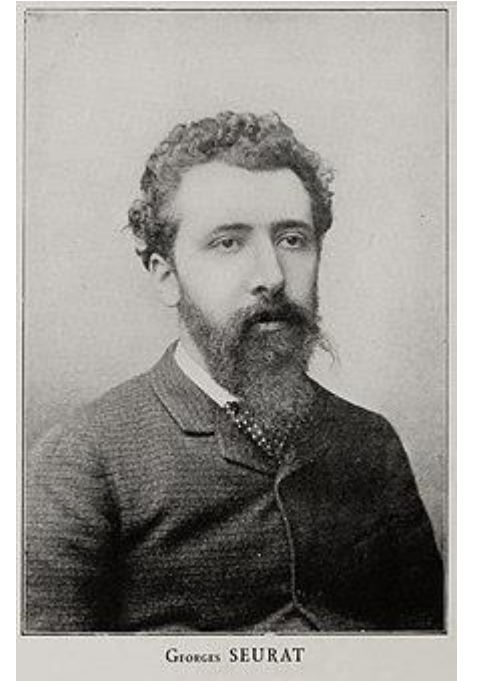
- High Sensitivity
- Lower technical noise
- High price
- Simpler data processing steps



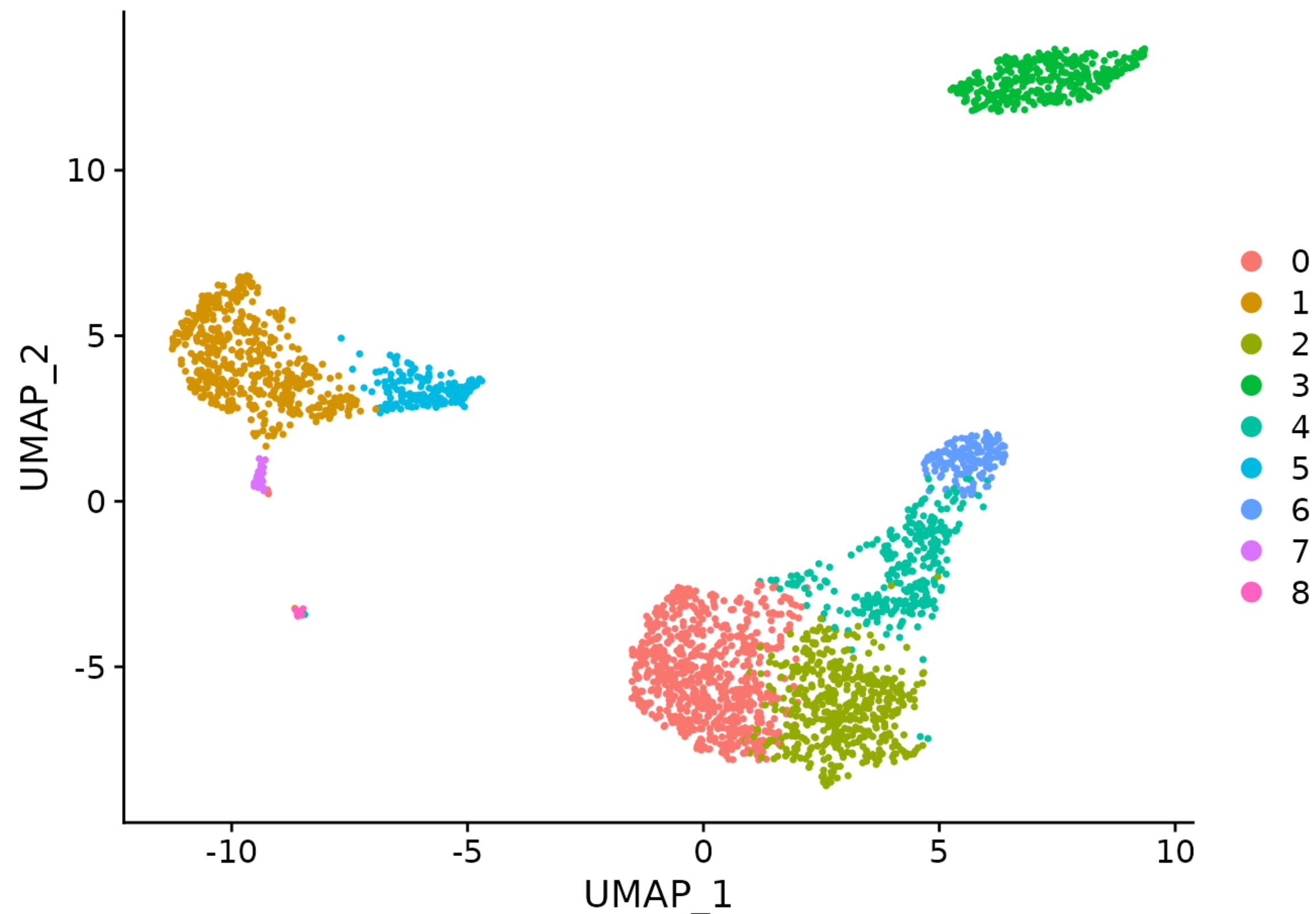
Platforms:

- Seurat (R)
- Scanpy (Python)
- Galaxy (webtool)
- Scater + scran (R, Bioconductor)
- Monocle3 and psupertime (R)





Georges SEURAT

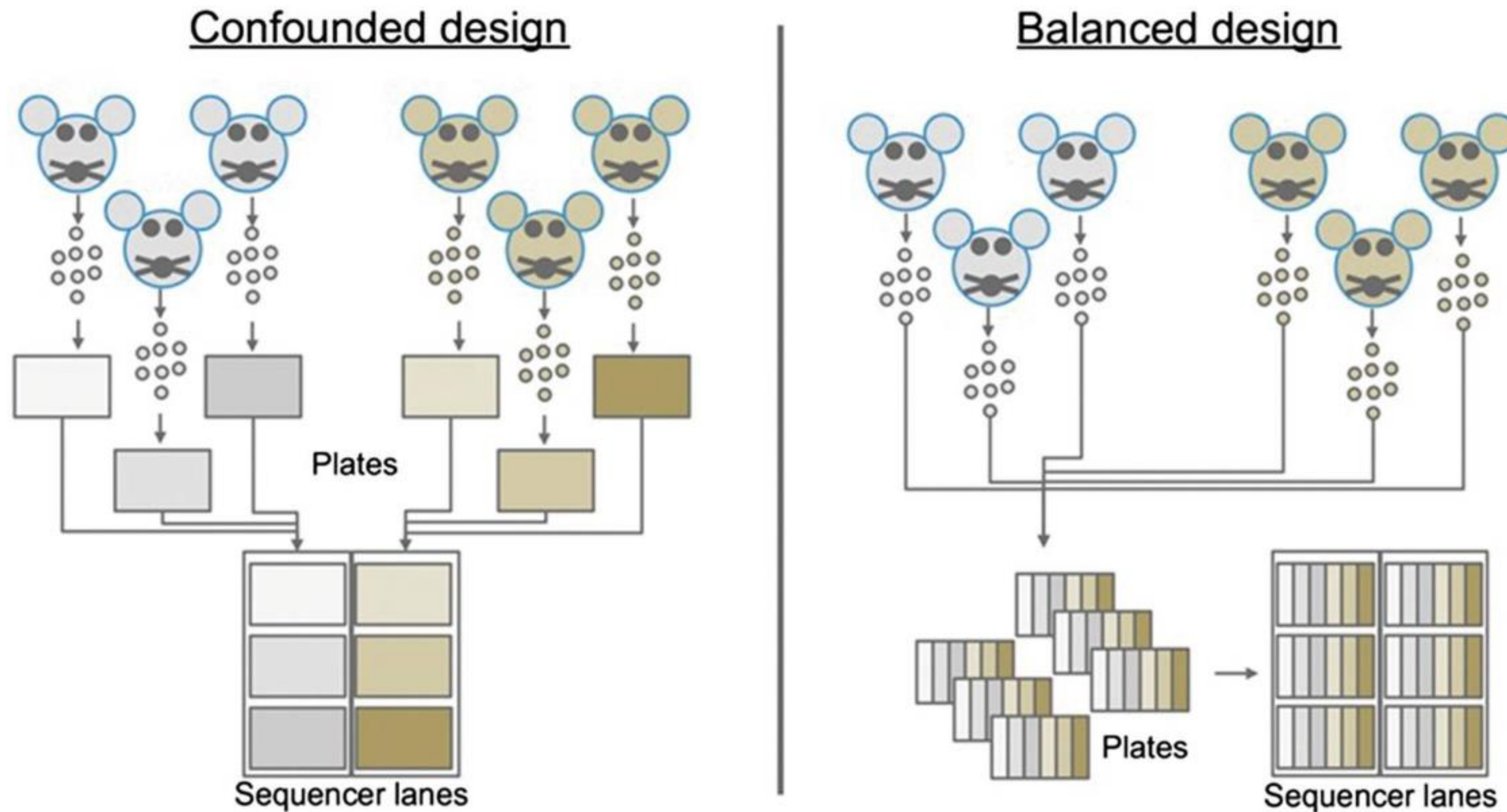


Why Seurat?

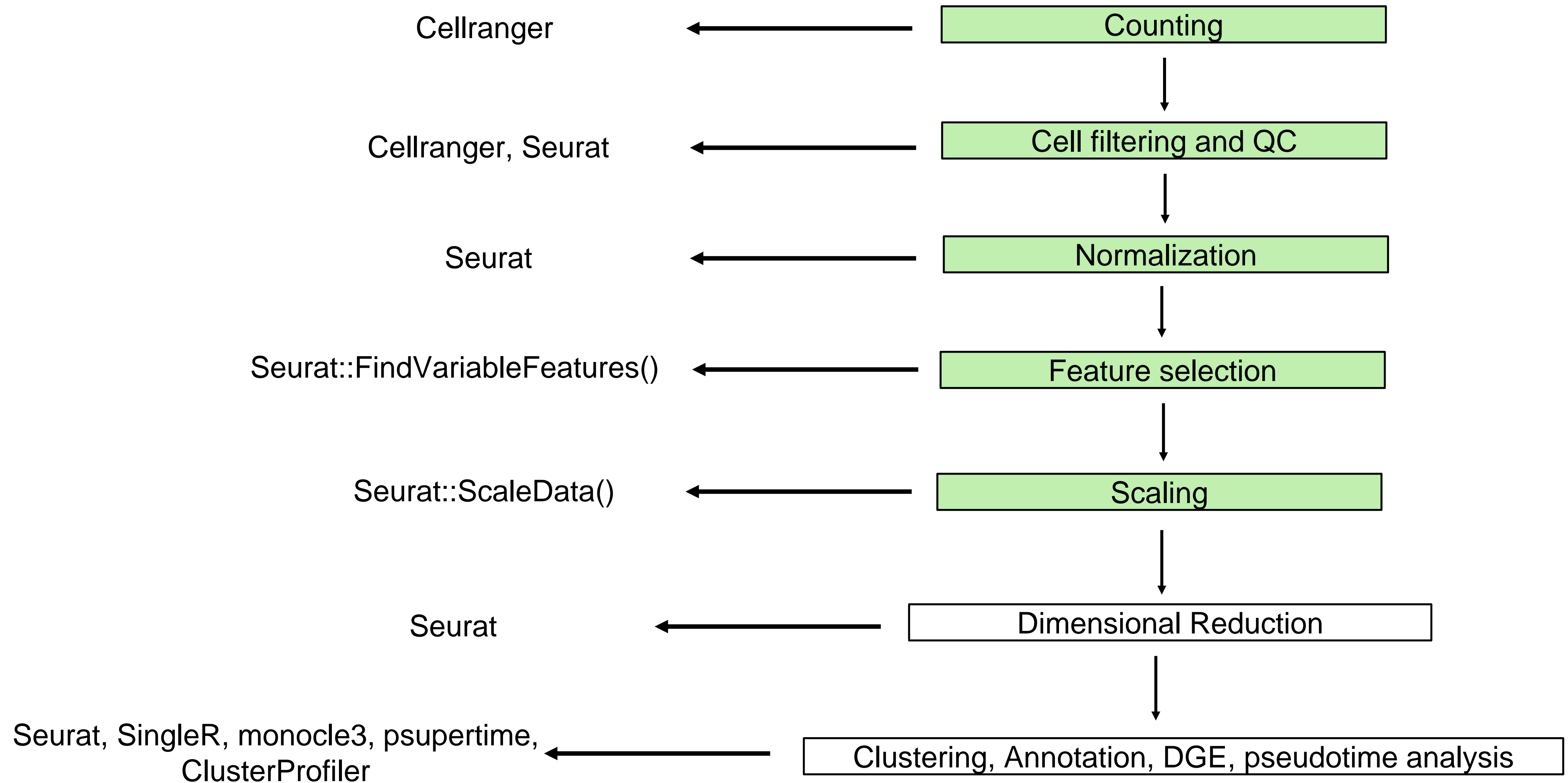
- Easy to use
- Very popular, supported and updated
- Deep documentation and good data model
- Several built-in functionalities
- Easy to extend with your own functionality
- Seurat objects extend to other packages

Experimental Design

- Replication and randomization
- Be aware of potential biases that can affect the results
- Record any factor for downstream correction



Data Analysis Overview



Counting

```
#Export path to make cellranger 7.1.0 available
export PATH=$PATH:/home/tigem/software/cellranger-7.1.0

# Define the path to the main directory containing the subfolders with input files
main_dir="/home/tigem/x.bujandac/Polishchuk_scrNAseq/MAN2-HRP_TC_sc"

# Loop through each subfolder
for subfolder in "$main_dir"/*; do

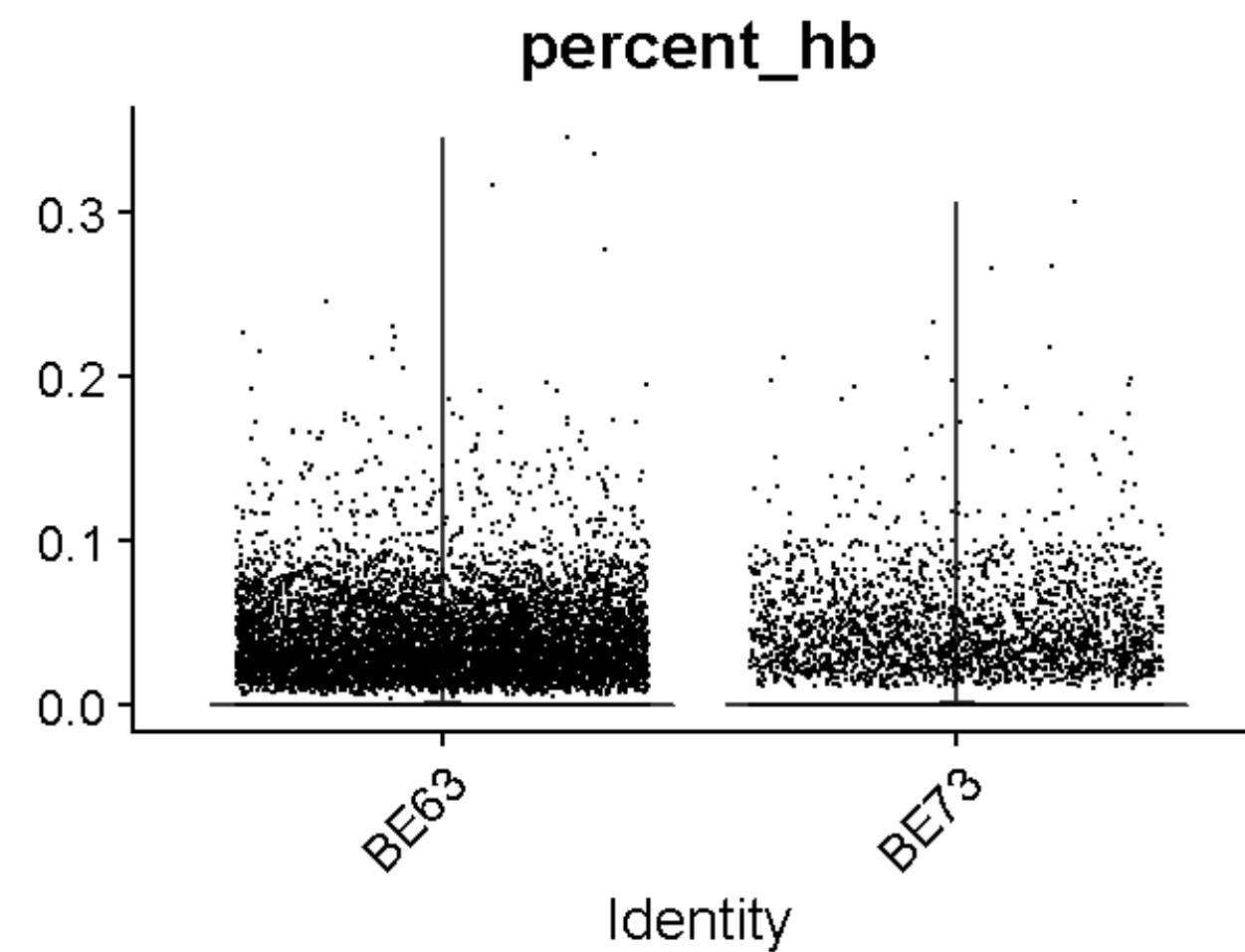
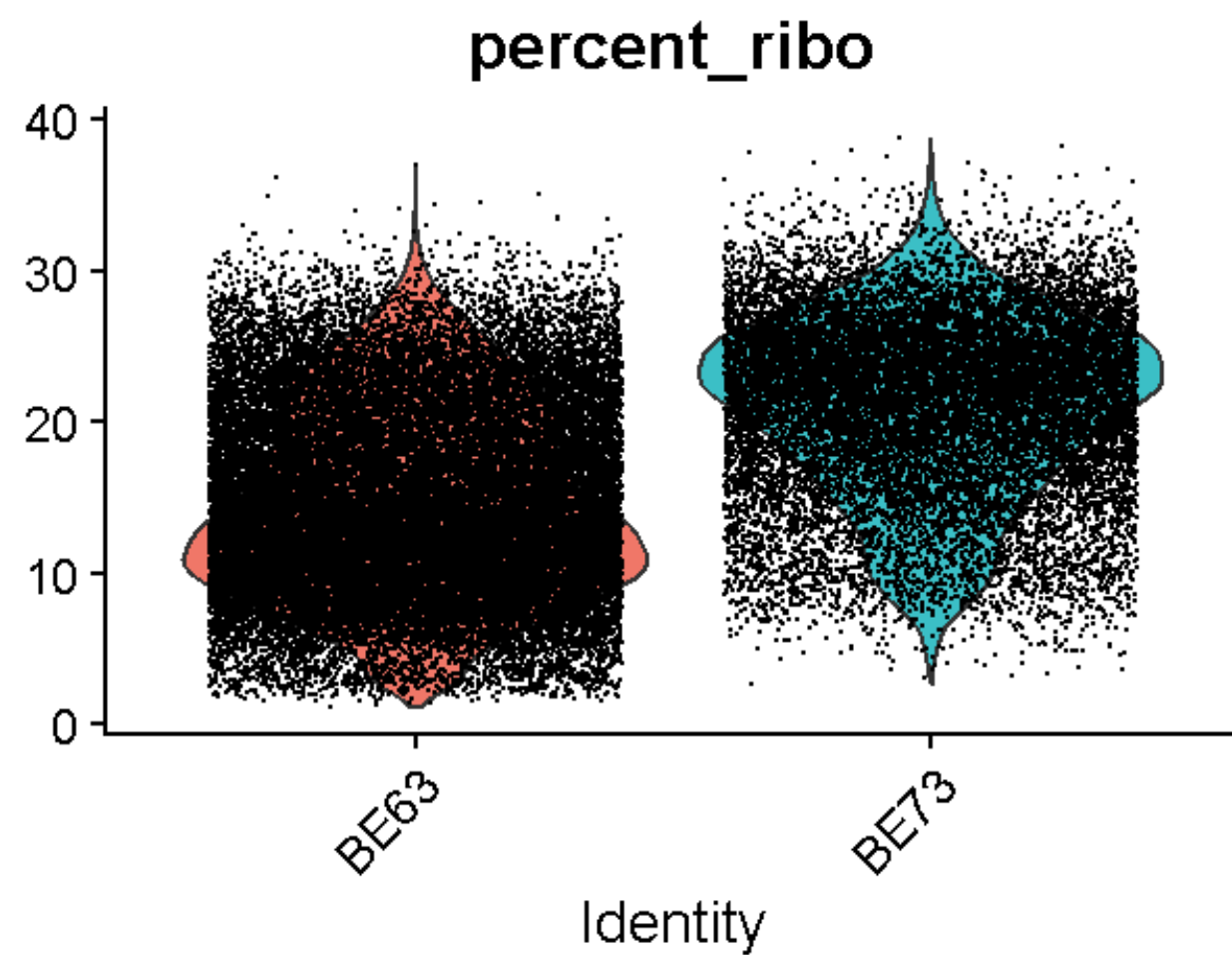
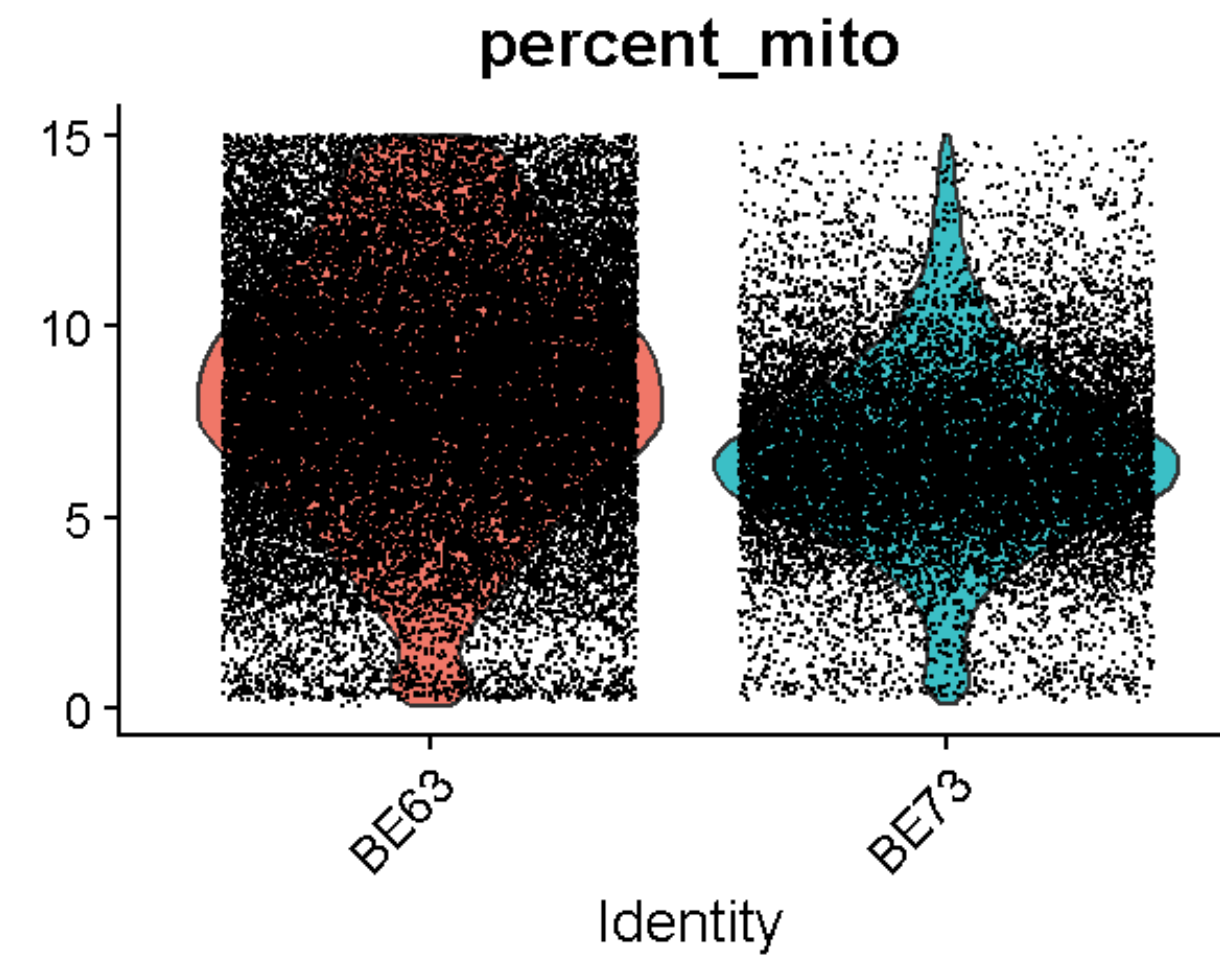
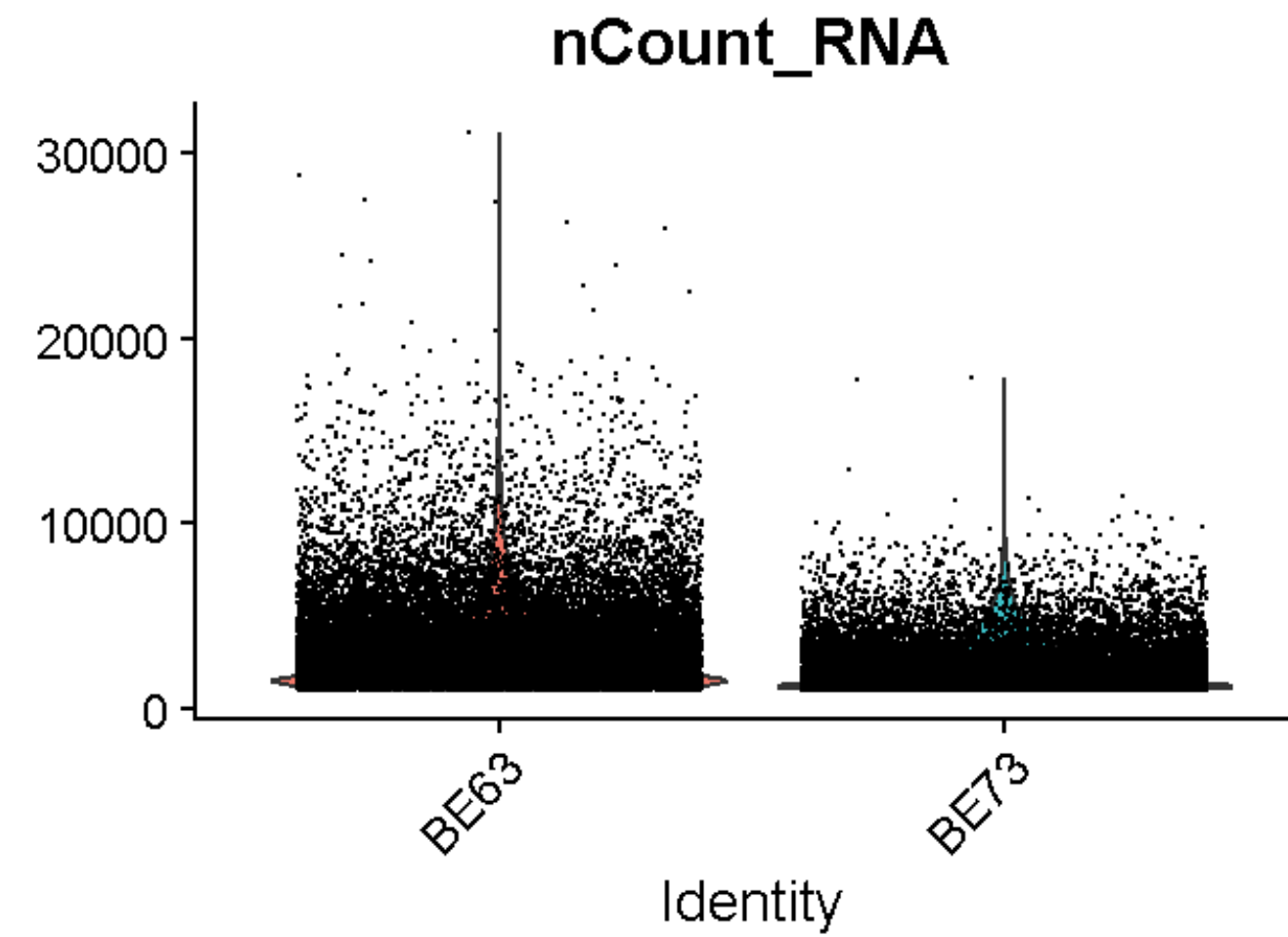
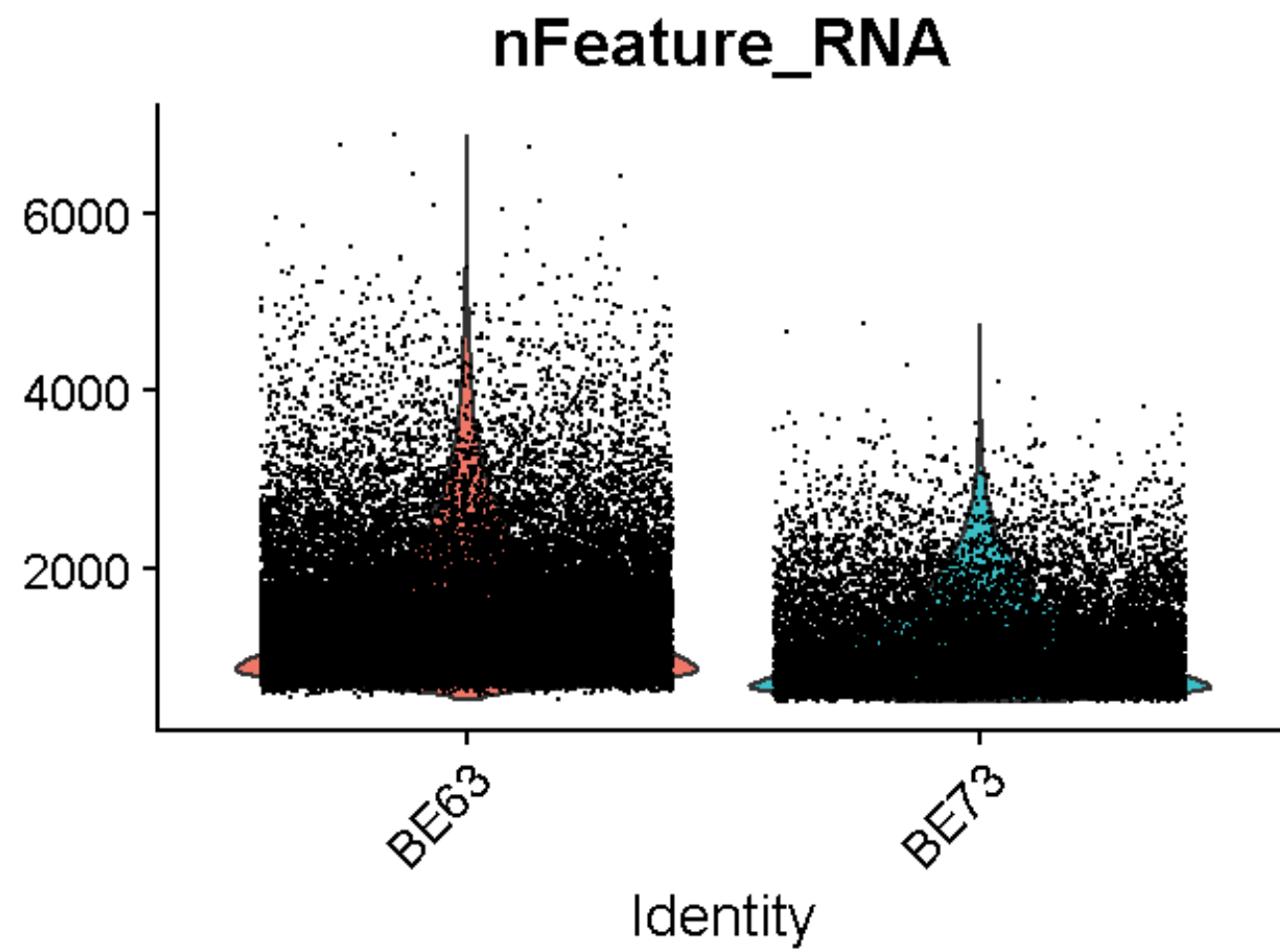
    echo "Processing files in $subfolder..."

    # Run cellranger count on the files in the subfolder
    cellranger count --id=$(basename "$subfolder") \
    --transcriptome=/home/tigem/x.bujandac/Polishchuk_scrNAseq/GRCh38 \
    --fastqs=$subfolder \
    --sample=$(basename "$subfolder") \
    --localcores 8 \
    --localmem 64

    echo "Finished processing files in $subfolder."
    echo
done
```

	Cell1	Cell2	...	CellN
Gene1	3	2	.	13
Gene2	2	3	.	1
Gene3	1	14	.	18
...
...
...
GeneM	25	0	.	0

Quality Control



- Differences among samples in:
- nFeatures and nCounts
 - Eventual apoptotic cells
 - Ceck for ribosomal genes and blood related genes



Quality Control

```

# Create Seurat Object
counts <- CreateSeuratObject(counts = counts, meta.data = metadata, project = "scRNAseqBE")
|
#===== Quality Control =====#

#calculate the % of mitochondrial genes for each cell
counts <- PercentageFeatureSet(counts, "^MT-", col.name = "percent_mito")

#calculate the % of ribosomal genes for each cell
counts <- PercentageFeatureSet(counts, "^RP[SL]", col.name = "percent_ribo")

# Percentage hemoglobin genes - includes all genes starting with HB except HBP.
counts <- PercentageFeatureSet(counts, "^HB[^P]", col.name = "percent_hb")

#Plot basic statistics
feats <- c("nFeature_RNA", "nCount_RNA", "percent_mito", "percent_ribo", "percent_hb")
VlnPlot(countsFiltered, group.by = "orig.ident", features = feats, pt.size = 0.1, ncol = 3) +
  NoLegend()

```

Data Normalization and Scaling

There are different ways to normalize data:

Logarithmic normalization: $x_i = \frac{\text{The read count of gene } X \text{ in cell } i}{\text{Total counts of cell } i} \times 10^4 \text{ (Eq.1)}$

$$f(x_i) = \ln(x_i + 1) \text{ (Eq.2)}$$

CPM, TPM, RPKM (old but gold)

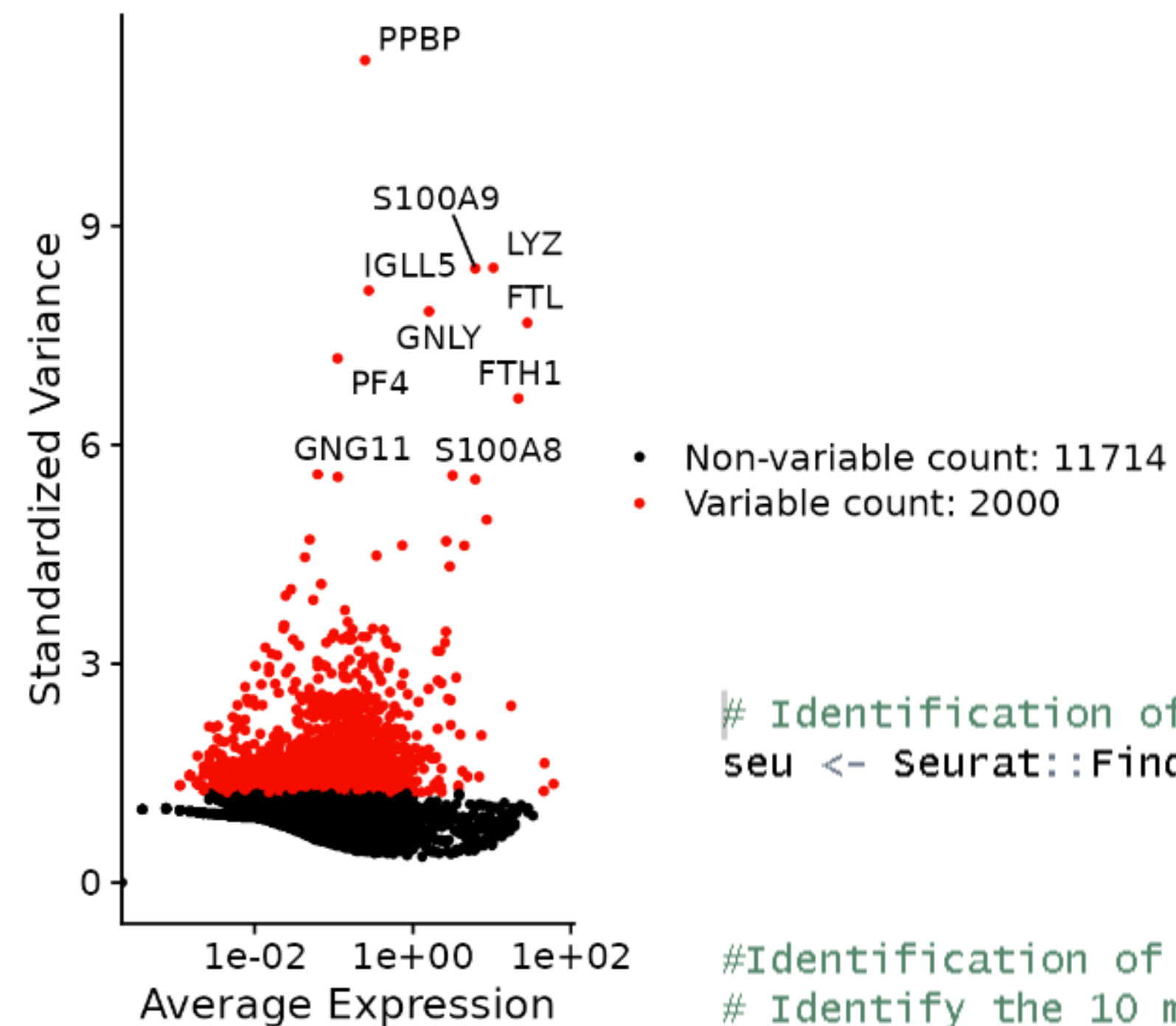
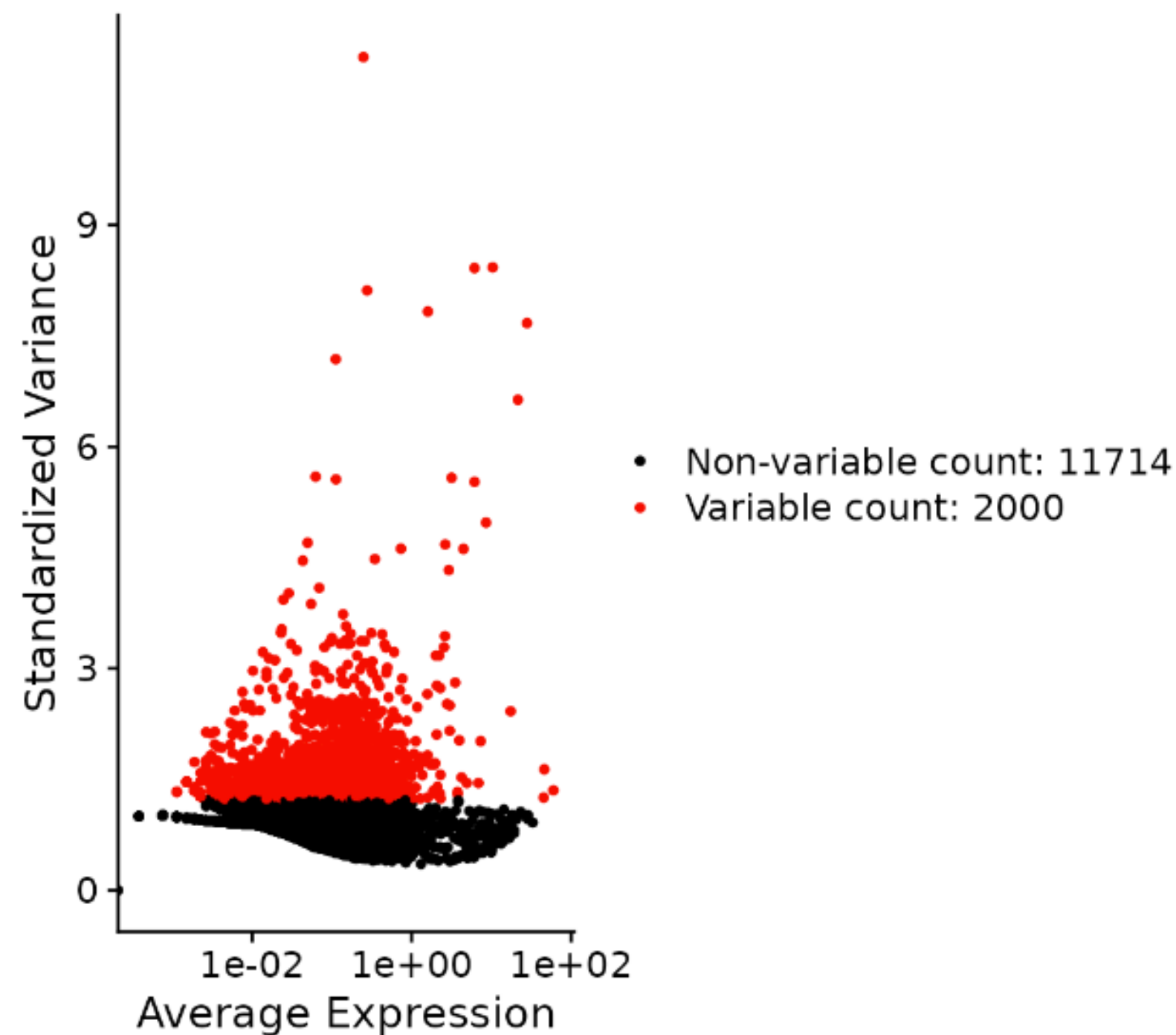
Work well under the assumption that the amount of RNA is the same in all cells and a uniform scaling factor is applicable for all genes

Sctransform: Normalization and variation control in one step



Provides an excellent tool to not only normalize data, but also stabilize variance and regress unwanted variation

Feature selection



- Identify several genes that exhibit high variability between cells
- Simple mathematical model to detect the most variable genes
- These most variable genes represent the features to prioritize for the downstream analysis

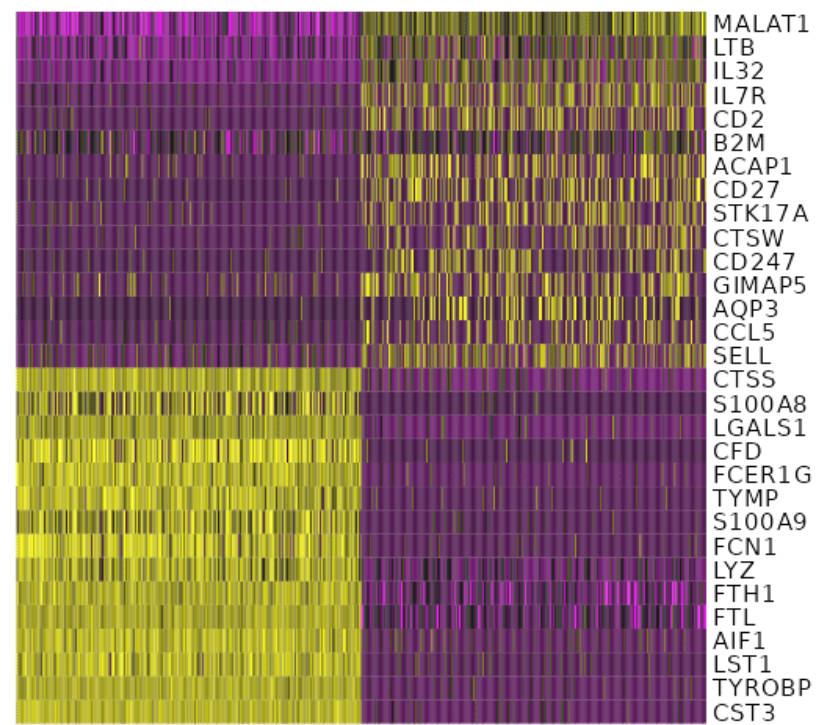
```
# Identification of highly variable features (feature selection)
seu <- Seurat::FindVariableFeatures(seu,
                                   selection.method = "vst",
                                   nfeatures = 2000)
```

```
# Identification of most expressed genes
# Identify the 10 most highly variable genes
top10 <- head(Seurat::VariableFeatures(seu), 10)
top10
```

```
#Plot
vf_plot <- Seurat::VariableFeaturePlot(seu)
Seurat::LabelPoints(plot = vf_plot,
                    points = top10, repel = TRUE)
```

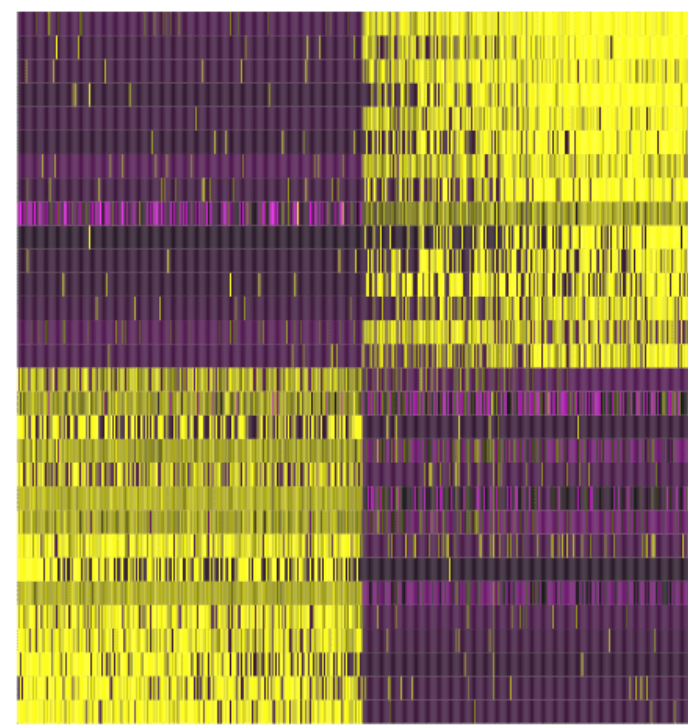

Linear Dimensional Reduction

PC_1



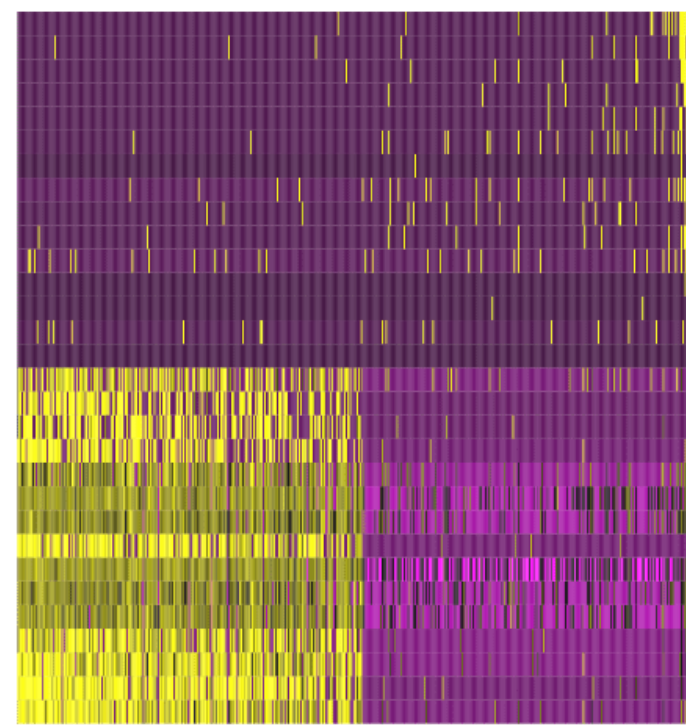
MALAT1
LTB
IL32
IL7R
CD2
B2M
ACAP1
CD27
STK17A
CTSW
CD247
GIMAP5
AQP3
CCL5
SELL
CTSS
S100A8
LGALS1
CFD
FCER1G
TYMP
S100A9
FCN1
LYZ
FTH1
FTL
AIF1
LST1
TYROBP
CST3

PC_2



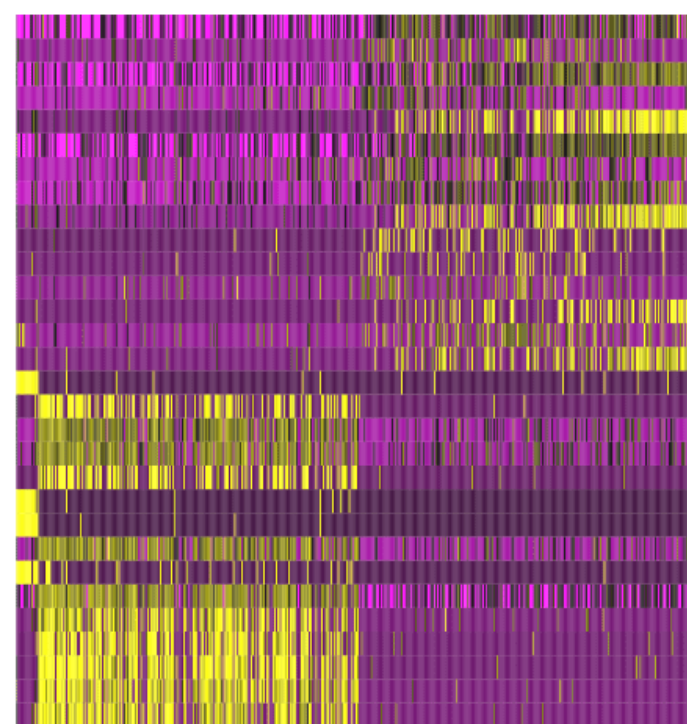
NKG7
PRF1
CST7
GZMB
GZMA
FGFBP2
CTSW
GNLY
B2M
SPON2
CCL4
GZMH
FCGR3A
CCL5
CD247
HLA-DRB5
CD37
HLA-DQA2
HLA-DPB1
HLA-DMA
CD74
HLA-DRB1
CD79B
LINC00926
HLA-DRA
HLA-DQB1
HLA-DQA1
TCL1A
MS4A1
CD79A

PC_3



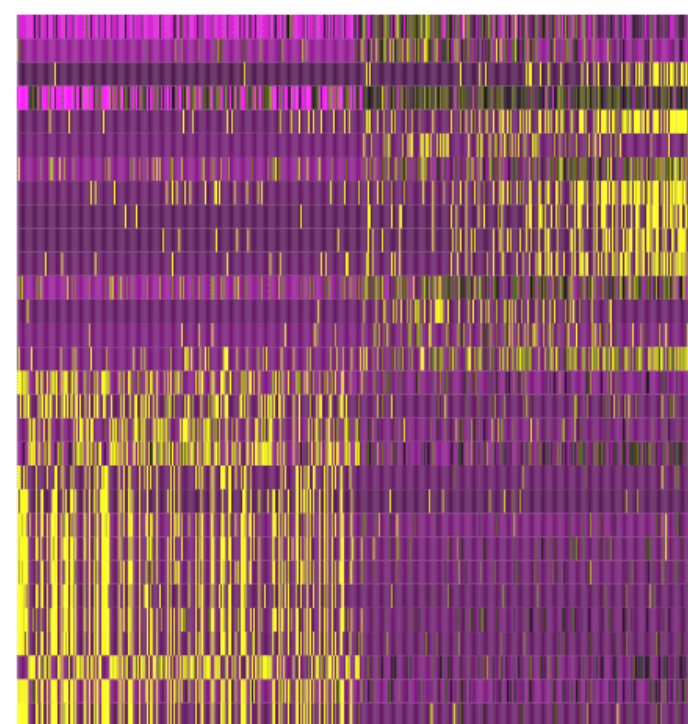
PPBP
PF4
SDPR
SPARC
GNG11
NRGN
GP9
RGS18
TUBB1
CLU
HIST1H2AC
APO01189.
ITGA2B
CD9
TMEM40
HLA-DMB
LINC00926
TCL1A
HLA-DQA2
HLA-DRB5
HLA-DRA
HLA-DRB1
MS4A1
CD74
HLA-DPA1
HLA-DPB1
HLA-DQB1
CD79B
CD79A
HLA-DQA1

PC_4



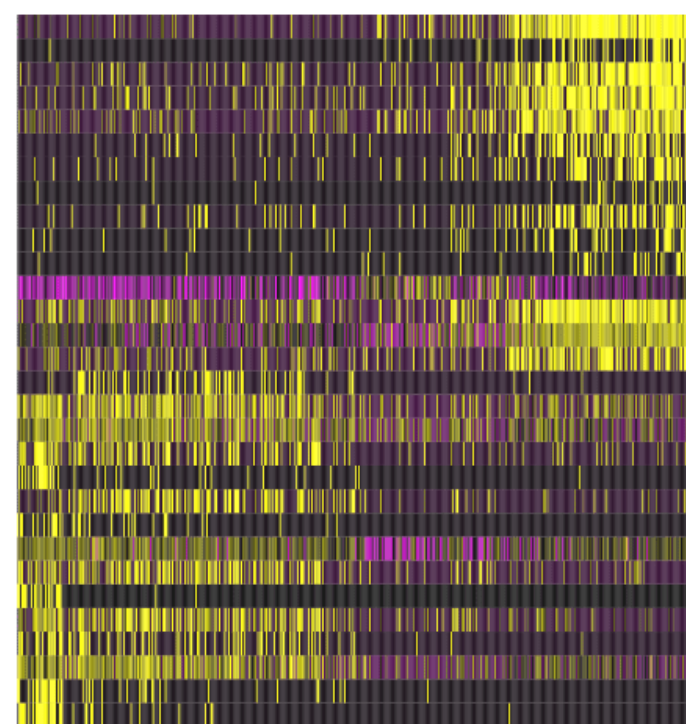
VIM
IL7R
S100A6
IL32
S100A8
S100A4
GIMAP7
S100A10
S100A9
MAL
AQP3
CD2
CD14
FYB
LGALS2
PPBP
HLA-DQA2
HLA-DPA1
HLA-DRB1
TCL1A
SDPR
PF4
HLA-DPB1
HIST1H2AC
CD74
HLA-DQB1
MS4A1
CD79A
CD79B
HLA-DQA1

PC_5



LTB
IL7R
CKB
VIM
MS4A7
AQP3
CYTIP
RP11-290F20.3
SIGLEC10
HMOX1
LILRB2
PTGES3
MAL
CD27
HN1
CTSW
CCL3
LGALS2
S100A9
GZMH
SPON2
GZMA
PRF1
CST7
CCL4
GNLY
FGFBP2
S100A8
NKG7
GZMB

PC_6



FCGR3A
CKB
MS4A7
RP11-290F:
RHOC
LILRA3
SIGLEC10
VMO1
HMOX1
CTD-2006K
PPM1N
MALAT1
IFITM3
IFITM2
LYN
FOLR3
GRN
GSTP1
ALDH2
CD1C
CD14
CACNA2D3
VIM
MS4A6A
SERPINF1
LGALS2
ID1
GPX1
CLEC10A
FCER1A

PCA is a statistical technique for reducing the dimensionality of the dataset by linearly transforming the data into a new coordinate system where most of the variation in the data can be described

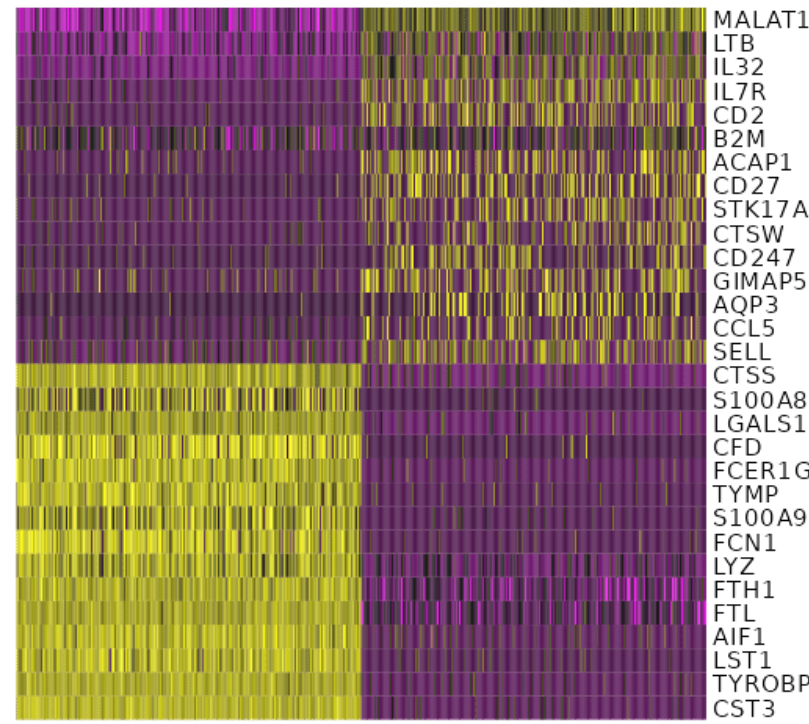
```
#===== Dimensional Reduction =====#
```

```
# PCA
```

```
countsFiltered <- RunPCA(countsFiltered, features = VariableFeatures(object = countsFiltered))
```

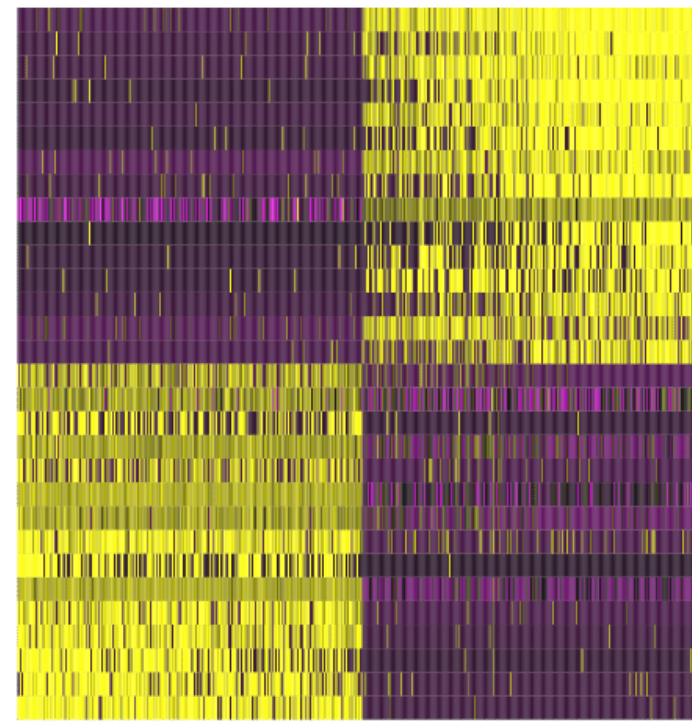
Linear Dimensional Reduction

PC_1



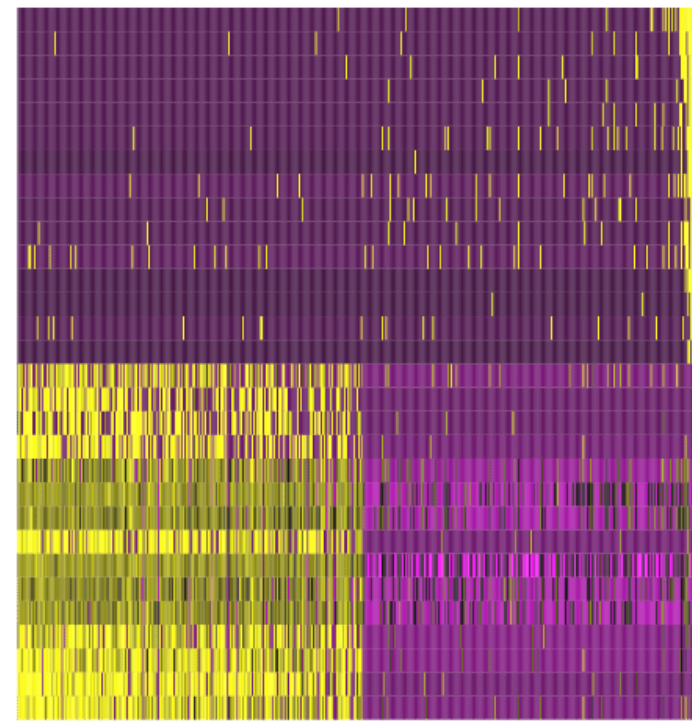
MALAT1
LTB
IL32
IL7R
CD2
B2M
ACAP1
CD27
STK17A
CTSW
CD247
GIMAP5
AQP3
CCL5
SELL
CTSS
S100A8
LGALS1
CFD
FCER1G
TYMP
S100A9
FCN1
LYZ
FTH1
FTL
AIF1
LST1
TYROBP
CST3

PC_2



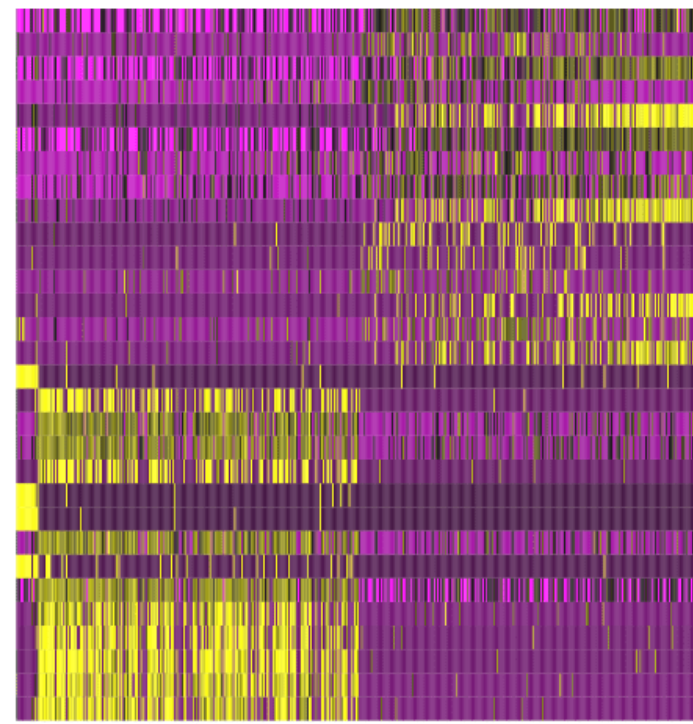
NKG7
PRF1
CST7
GZMB
GZMA
FGFBP2
CTSW
GNLY
B2M
SPON2
CCL4
GZMH
FCGR3A
CCL5
CD247
HLA-DRB5
CD37
HLA-DQA2
HLA-DPB1
HLA-DMA
CD74
HLA-DRB1
CD79B
LINC00926
HLA-DRA
HLA-DQB1
HLA-DQA1
TCL1A
MS4A1
CD79A

PC_3



PPBP
PF4
SDPR
SPARC
GNG11
NRGN
GP9
RGS18
TUBB1
CLU
HIST1H2AC
APO01189.
ITGA2B
CD9
TMEM40
HLA-DMB
LINC00926
TCL1A
HLA-DQA2
HLA-DRB5
HLA-DRA
HLA-DRB1
MS4A1
CD74
HLA-DPA1
HLA-DPB1
HLA-DQB1
CD79B
CD79A
HLA-DQA1

PC_4



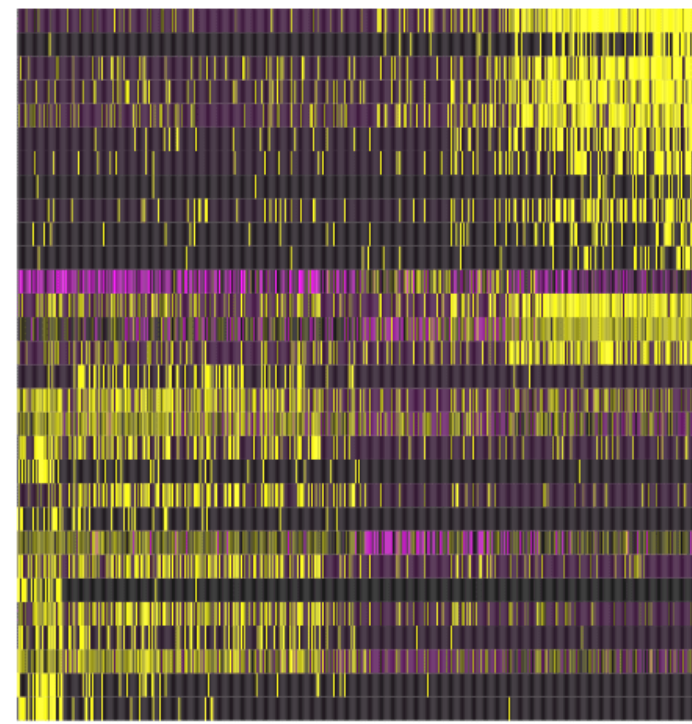
VIM
IL7R
S100A6
IL32
S100A8
S100A4
GIMAP7
S100A10
S100A9
MAL
AQP3
CD2
CD14
FYB
LGALS2
PPBP
HLA-DQA2
HLA-DPA1
HLA-DRB1
TCL1A
SDPR
PF4
HLA-DPB1
HIST1H2AC
CD74
HLA-DQB1
MS4A1
CD79A
CD79B
HLA-DQA1

PC_5

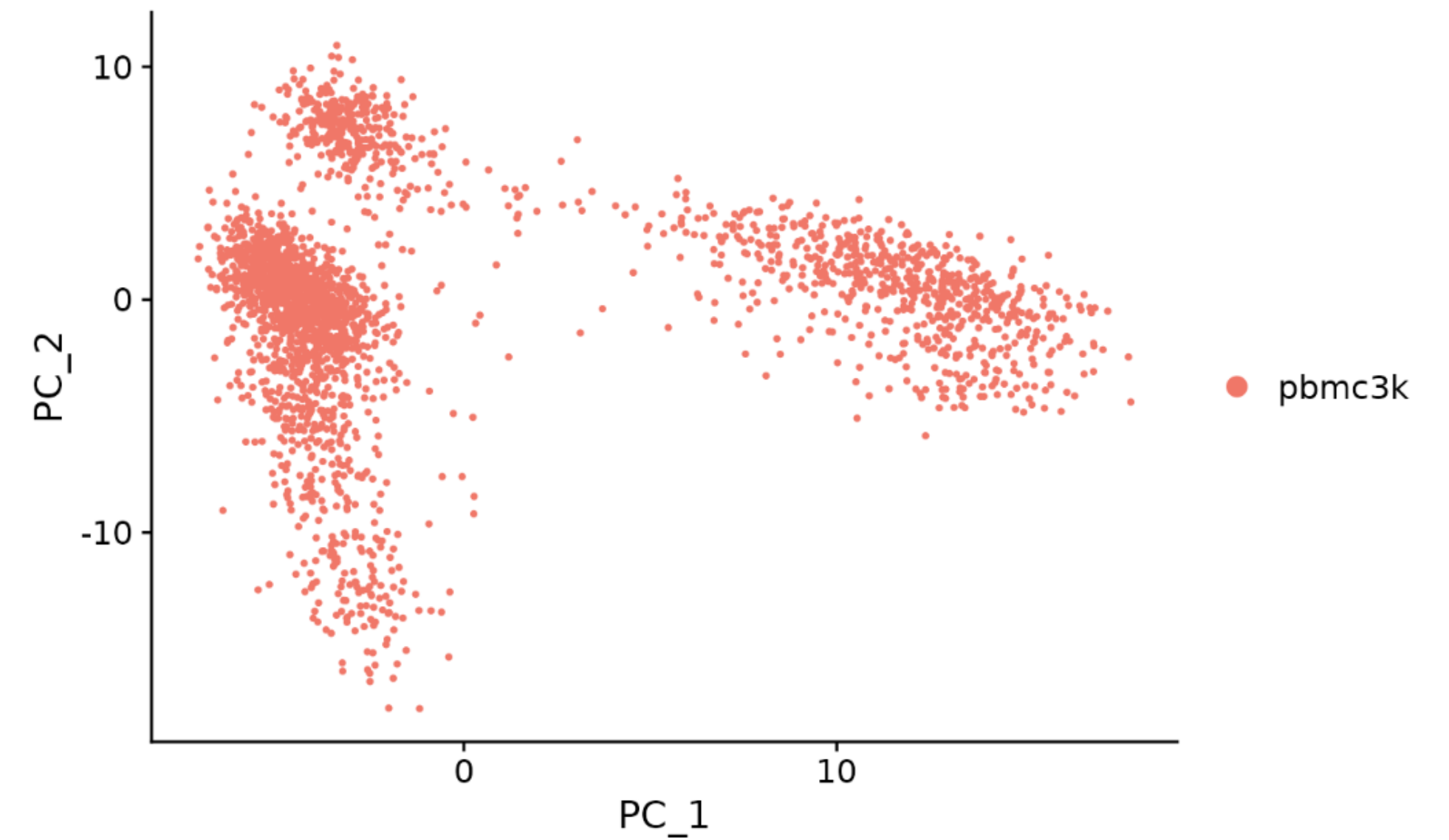
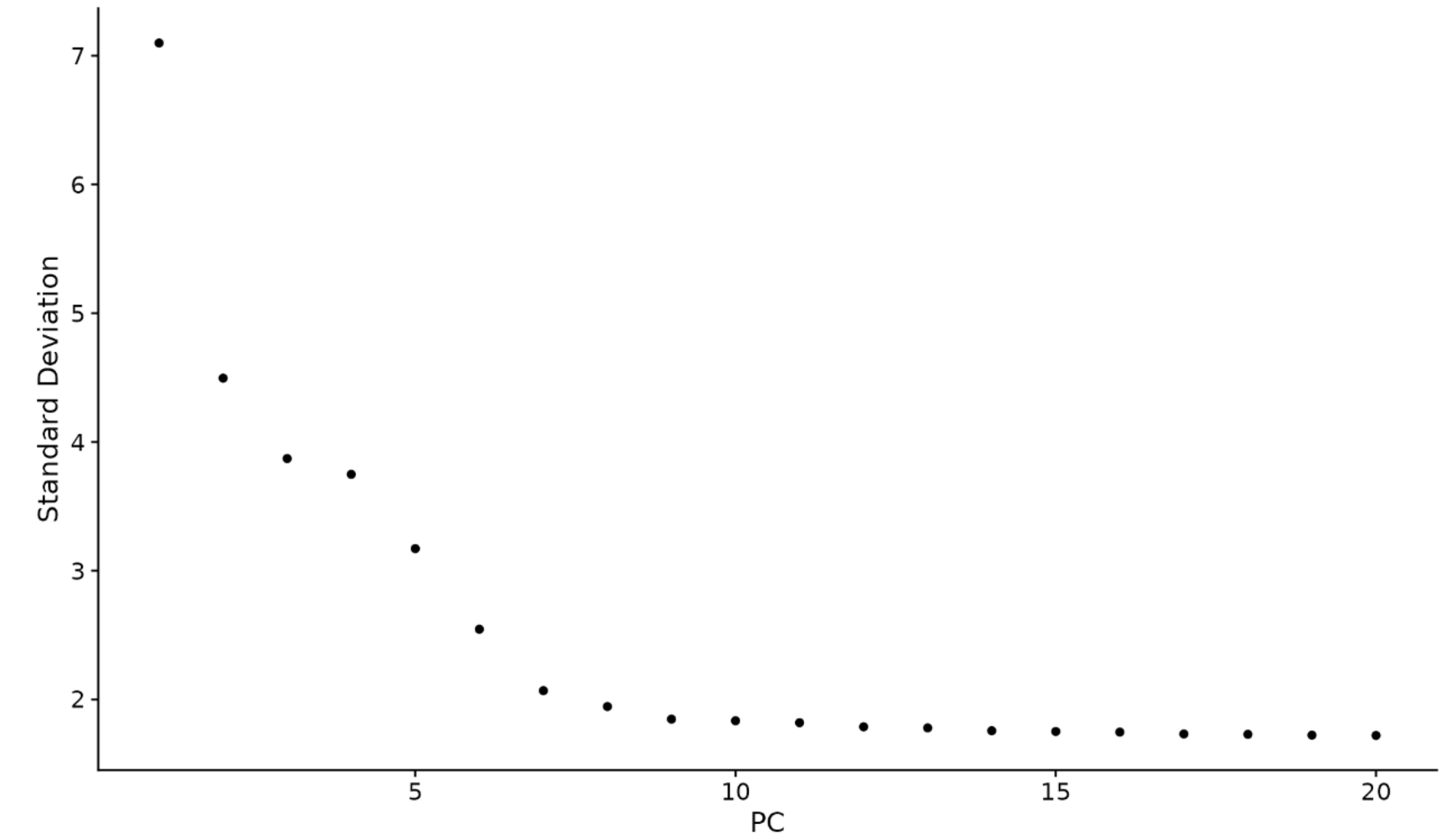


LTB
IL7R
CKB
VIM
MS4A7
AQP3
CYTIP
RP11-290F20.3
SIGLEC10
HMOX1
LILRB2
PTGES3
MAL
CD27
HN1
CTSW
CCL3
LGALS2
S100A9
GZMH
SPON2
GZMA
PRF1
CST7
CCL4
GNLY
FGFBP2
S100A8
NKG7
GZMB

PC_6



FCGR3A
CKB
MS4A7
RP11-290F:
RHOC
LILRA3
SIGLEC10
VMO1
HMOX1
CTD-2006K
PPM1N
MALAT1
IFITM3
IFITM2
LYN
FOLR3
GRN
GSTP1
ALDH2
CD1C
CD14
CACNA2D3
VIM
MS4A6A
SERPINF1
LGALS2
ID1
GPX1
CLEC10A
FCER1A

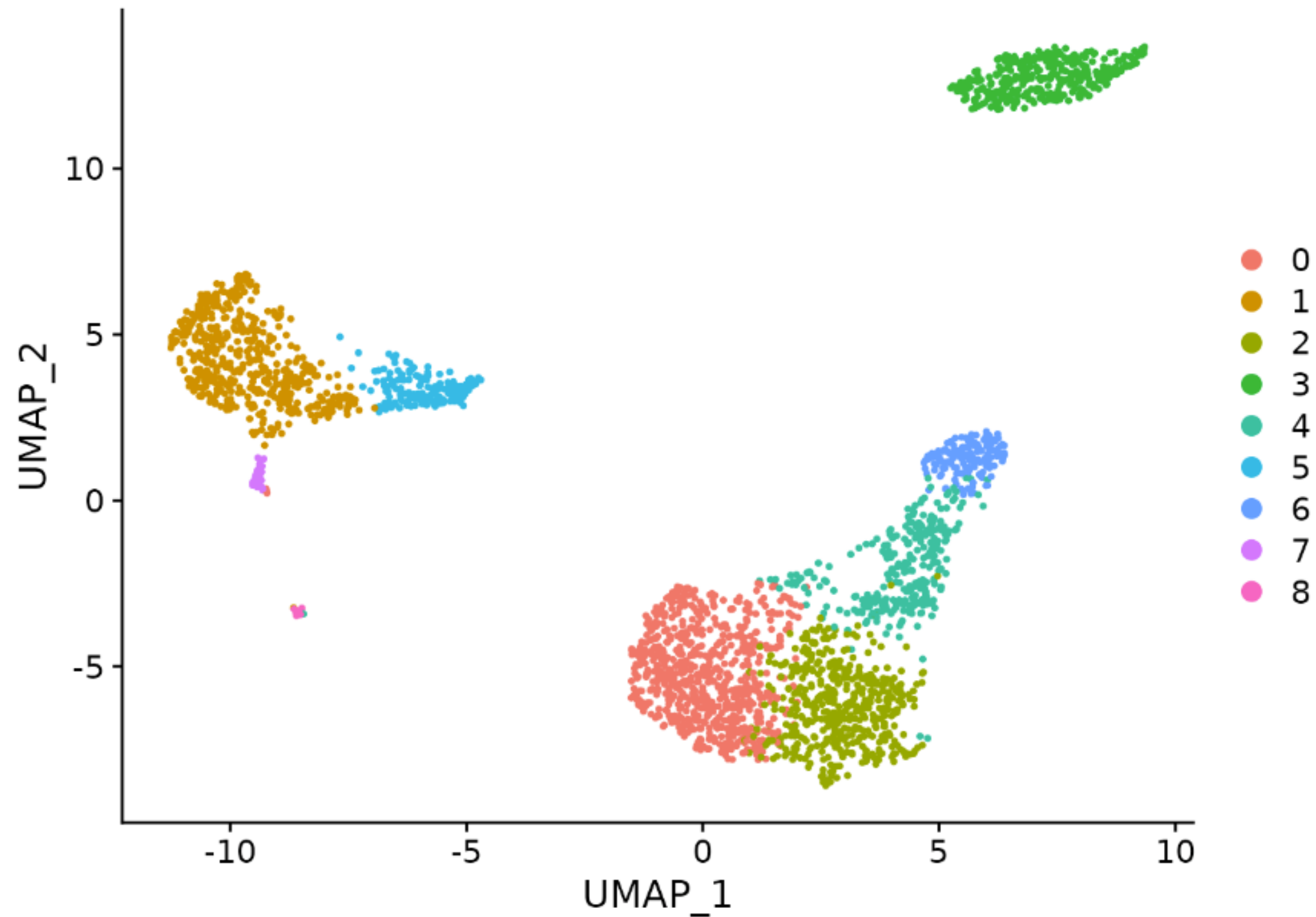


#===== Dimensional Reduction =====#

PCA

```
countsFiltered <- RunPCA(countsFiltered, features = VariableFeatures(object = countsFiltered))
```


Non - Linear Dimensional Reduction



```

##### Cluster the cells #####
countsFiltered <- FindNeighbors(countsFiltered, dims = 1:10)
countsFiltered <- FindClusters(countsFiltered, resolution = 0.5)
head(Idents(countsFiltered), 5)

##### Non-linear dimensional reduction|#####
#Run UMAP
countsFiltered <- RunUMAP(countsFiltered, dims = 1:6)

# individual clusters
DimPlot(countsFiltered, reduction = "umap", label = TRUE, repel = TRUE)

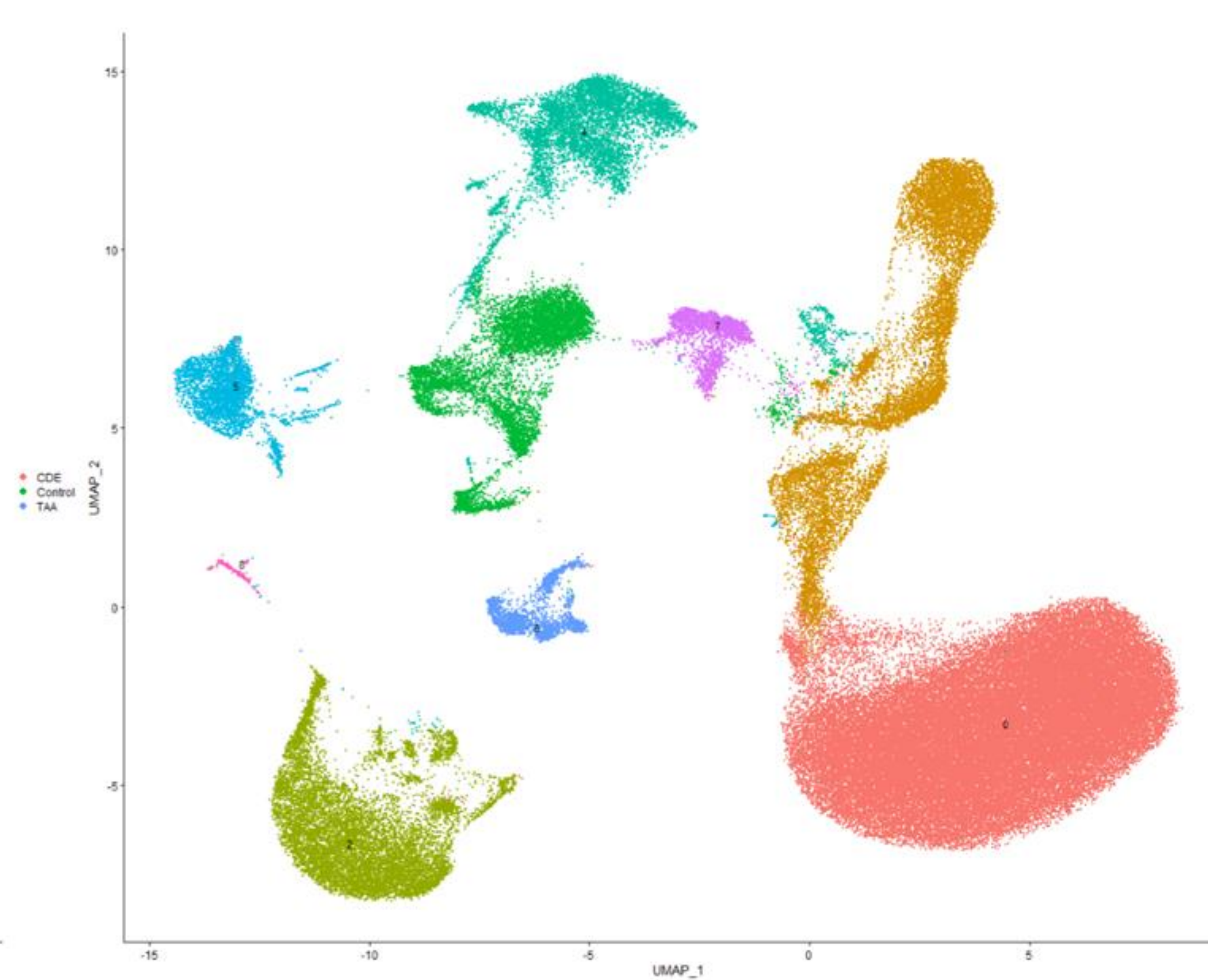
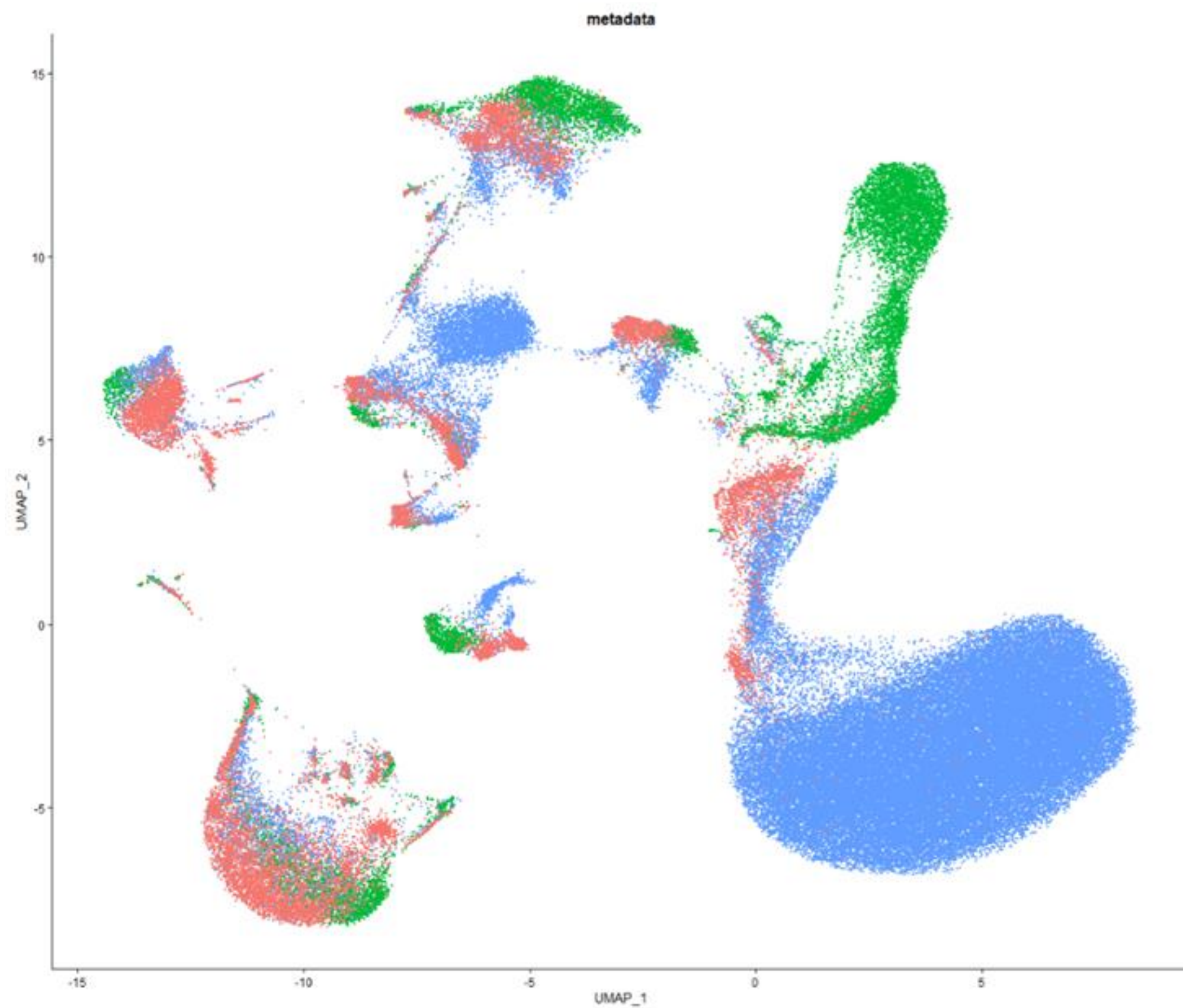
```

First we need to cluster cells!

Clustering groups a set of objects in a way that objects in the same group (**cluster**) are more similar to each other than to those in other groups

Non - Linear Dimensional Reduction

We can take advantage of these methods to identify potential batch effects



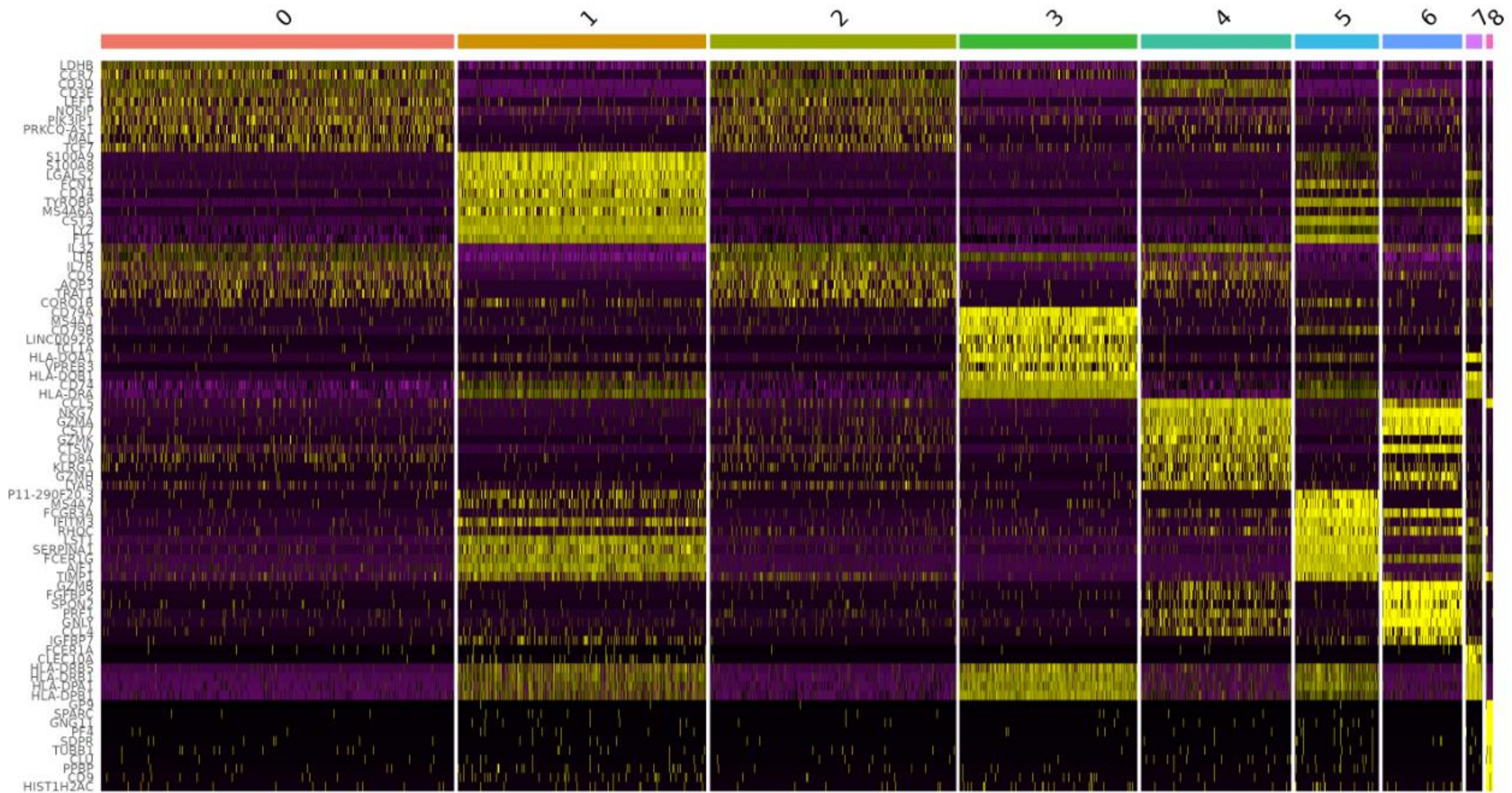
Clusters are driven by samples



Batch correction



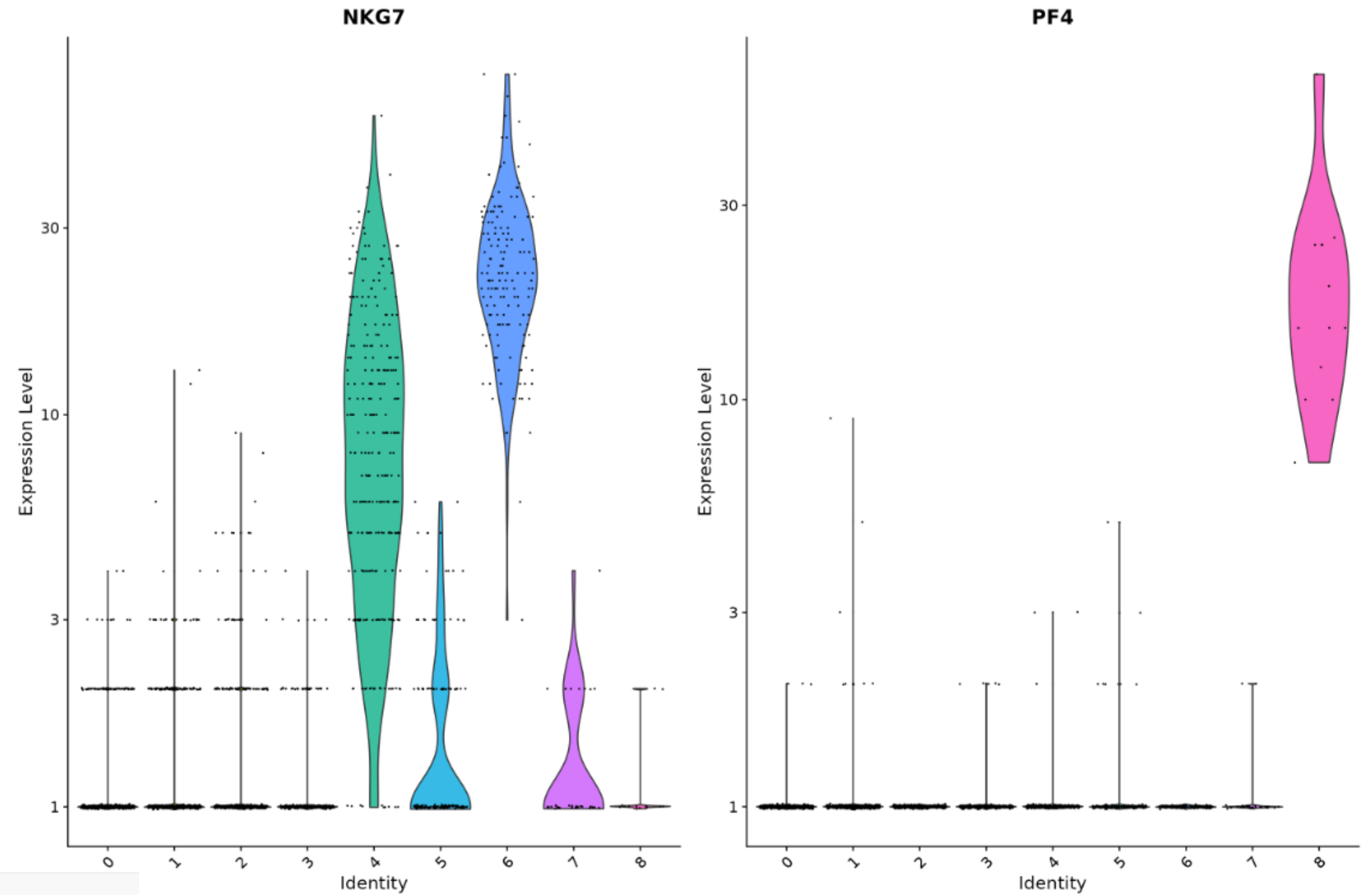
Finding Differentially Expressed Features



```
# find all markers of cluster 8
```

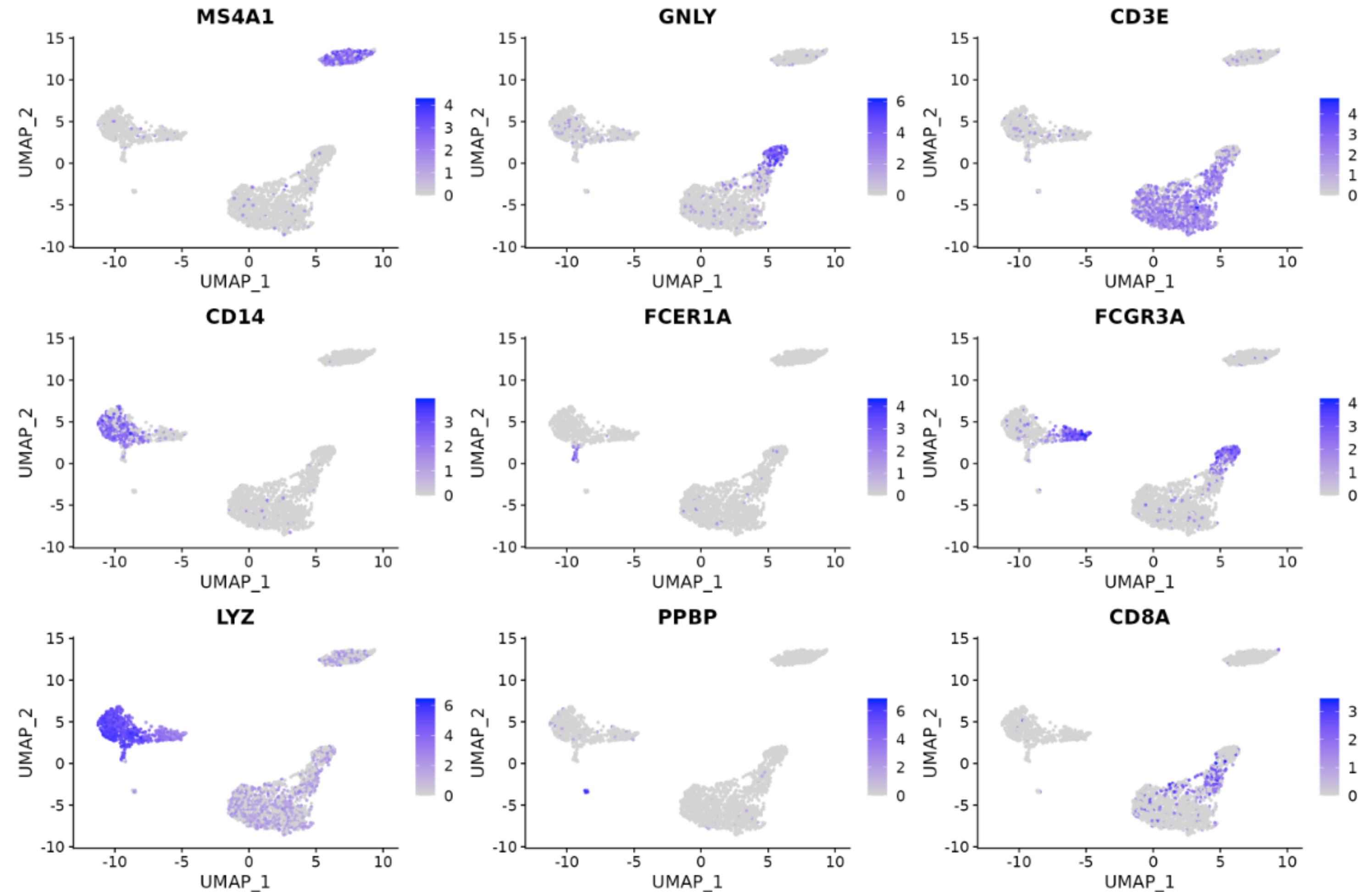
```
cluster8.markers <- FindMarkers(countsFiltered, ident.1 = 8, min.pct = 0.25)
```


Finding Differentially Expressed Features



```
# you can plot raw counts as well
VlnPlot(pbmc, features = c("NKG7", "PF4"), slot = "counts", log = TRUE)
```

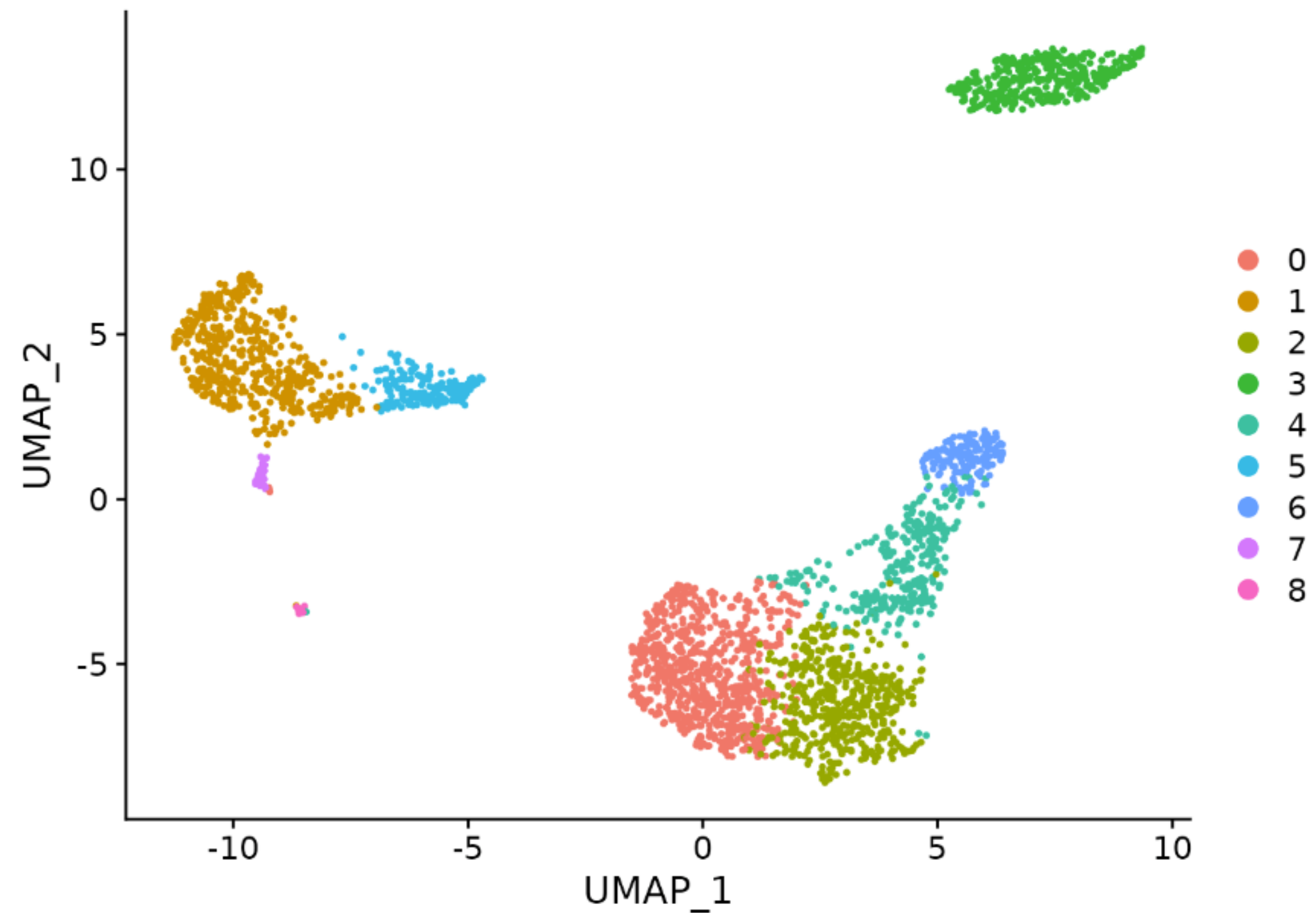

Finding Differentially Expressed Features



```
FeaturePlot(pbm, features = c("MS4A1", "GNLY", "CD3E", "CD14", "FCER1A", "FCGR3A", "LYZ", "PPBP", "CD8A"))
```

Adapted from https://satijalab.org/seurat/articles/pbmc3k_tutorial.html

Assigning Cell Identities to Clusters

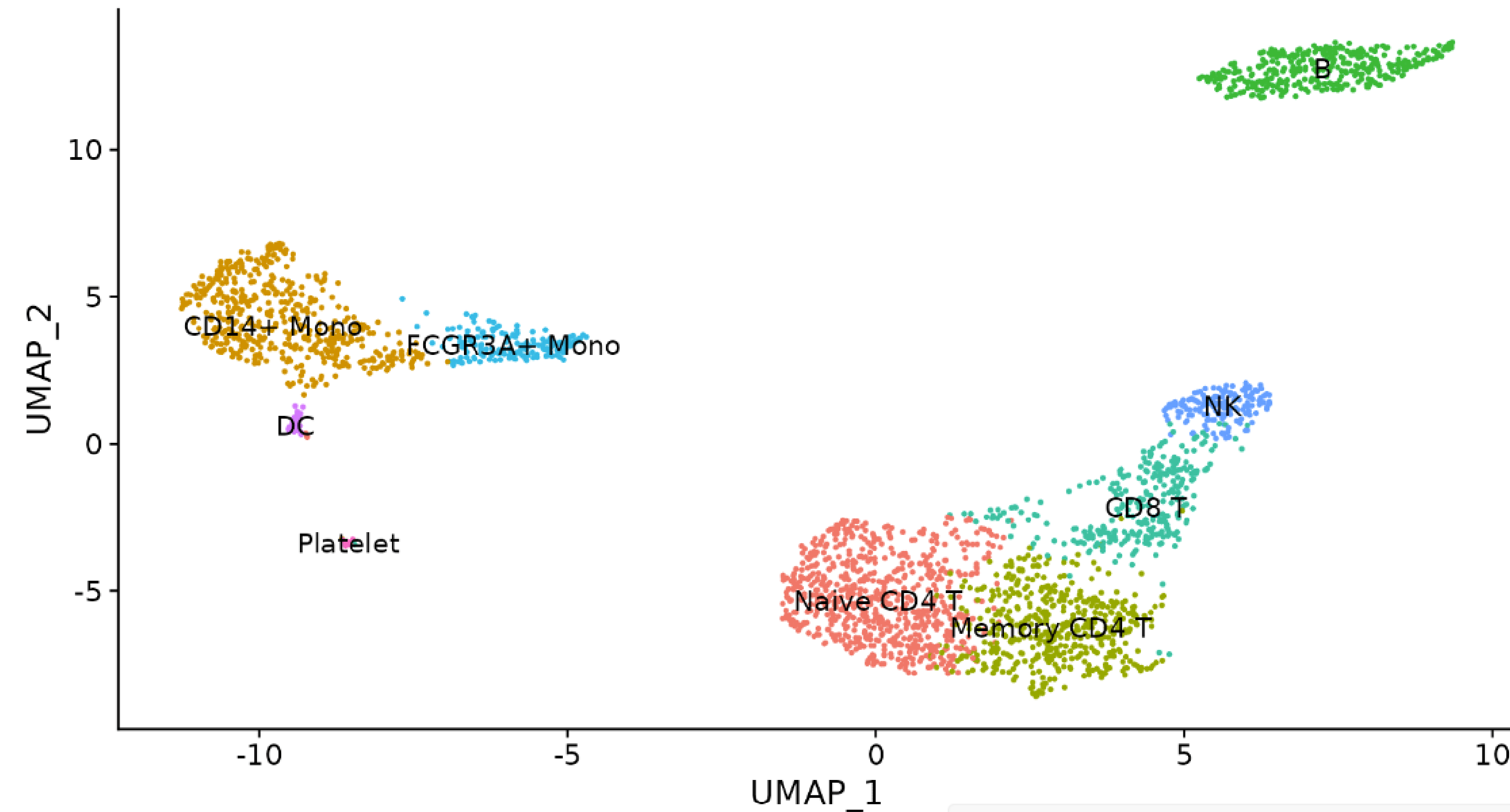


How to assign cell identities:

- Manually
- SingleR
- Using Machine Learning techniques

```
new.cluster.ids <- c("Naive CD4 T", "CD14+ Mono", "Memory CD4 T", "B", "CD8 T", "FCGR3A+ Mono",
  "NK", "DC", "Platelet")
names(new.cluster.ids) <- levels(pbmc)
pbmc <- RenameIdents(pbmc, new.cluster.ids)
DimPlot(pbmc, reduction = "umap", label = TRUE, pt.size = 0.5) + NoLegend()
```

Assigning Cell Identities to Clusters

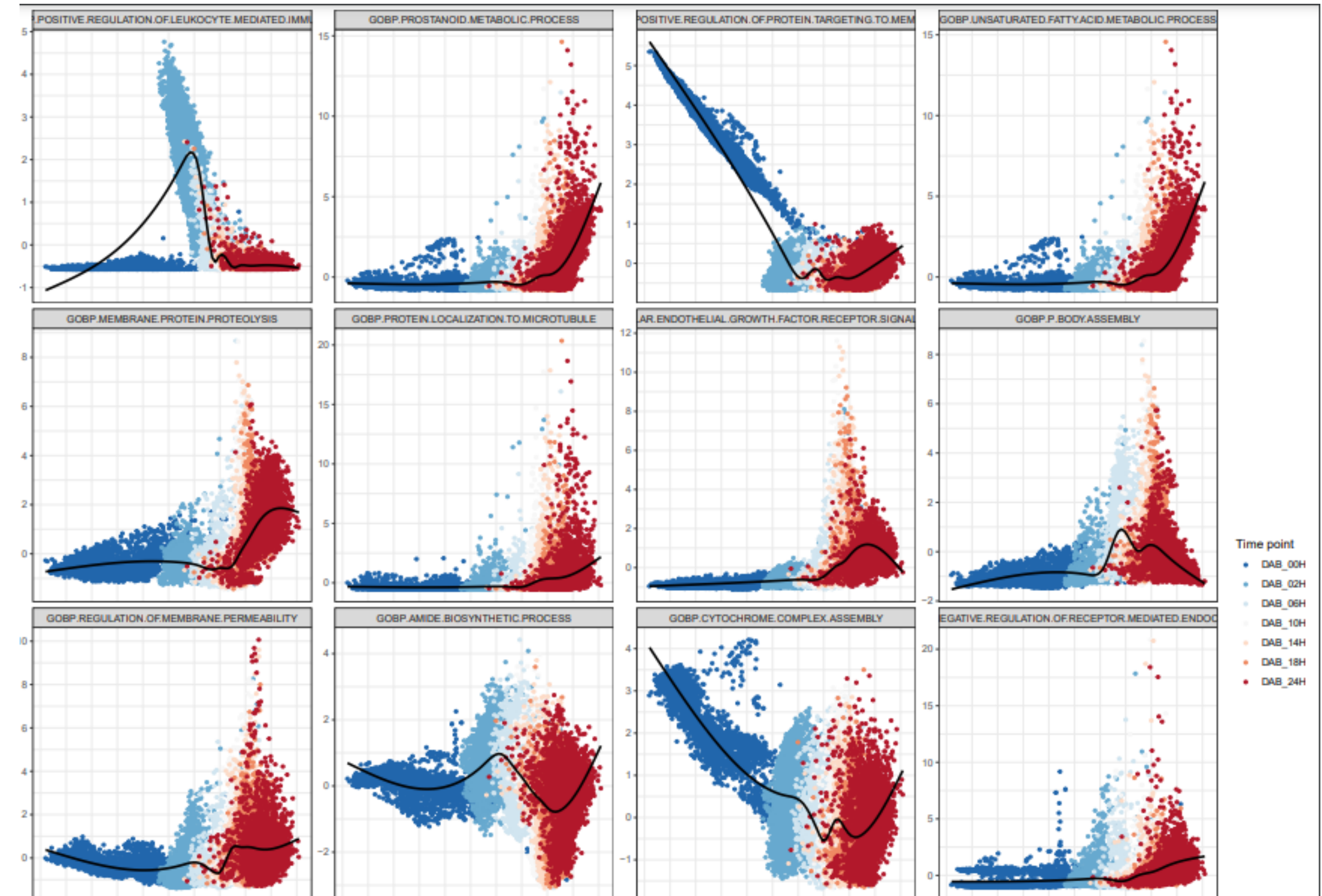
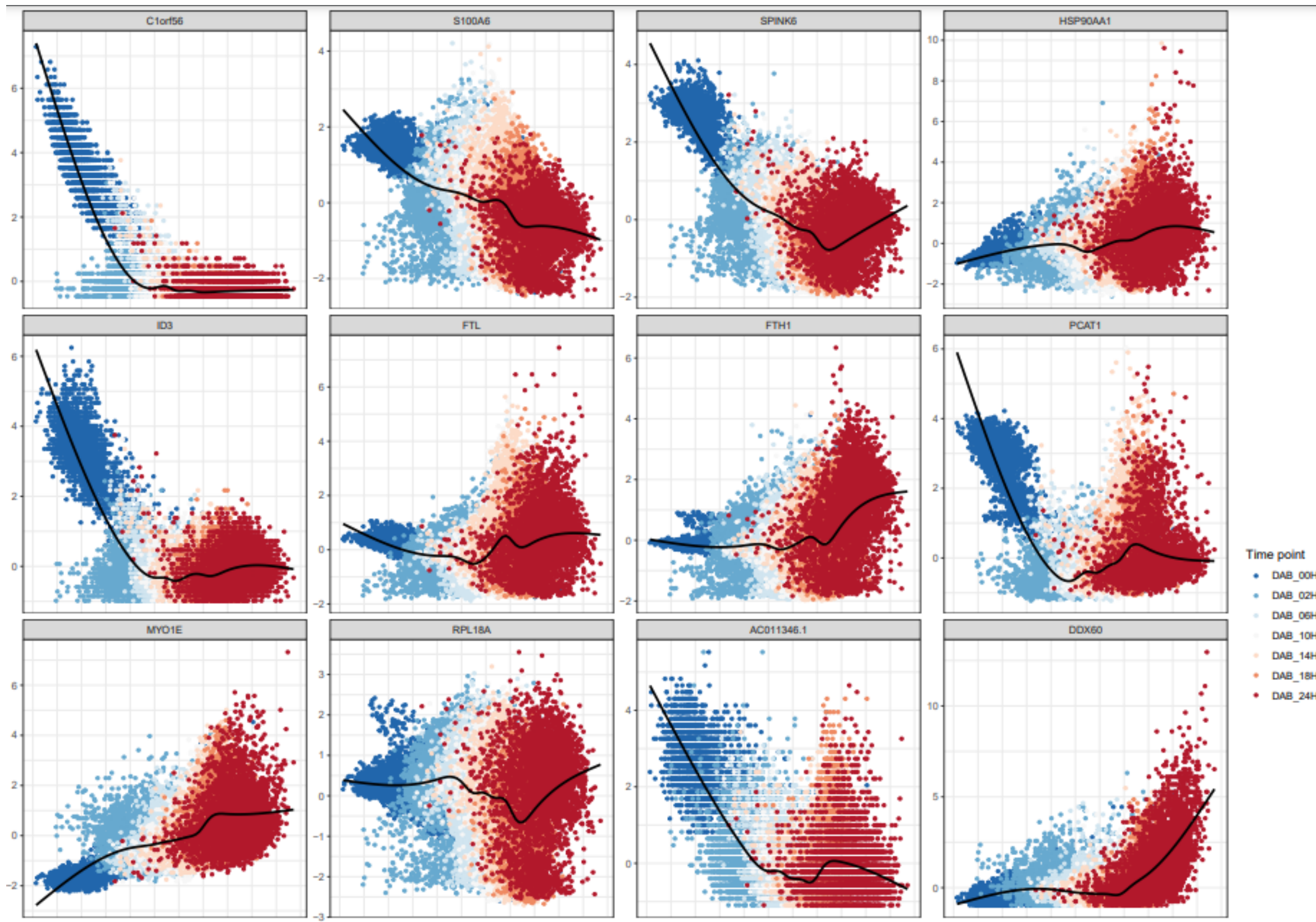


How to assign cell identities:

- Manually
- SingleR
- Using Machine Learning techniques

```
new.cluster.ids <- c("Naive CD4 T", "CD14+ Mono", "Memory CD4 T", "B", "CD8 T", "FCGR3A+ Mono",
                    "NK", "DC", "Platelet")
names(new.cluster.ids) <- levels(pbmc)
pbmc <- RenameIdents(pbmc, new.cluster.ids)
DimPlot(pbmc, reduction = "umap", label = TRUE, pt.size = 0.5) + NoLegend()
```


Pseudotime Analysis





Introduction to single cell RNA-seq

Xabier Bujanda Cundin
17/07/2023