

# **SINGLE-CELL analysis with 10x genomics**

## Introduction to scRNA-seq analysis

25.10.2023

Joana P. Bernardes

Institute for Clinical Molecular Biology-Kiel



**IKMB**

Institute of **Clinical**  
**Molecular Biology** Kiel

C COMPETENCE  
C CENTRE FOR  
G GENOMIC  
A ANALYSIS KIEL

A large orange circle is positioned on the left side of the slide, covering approximately one-third of the vertical space.

What do  
you know  
so far?

- Overview single-cell technologies
- Single-cell techniques in understanding disease
- Precision medicine: Present and future

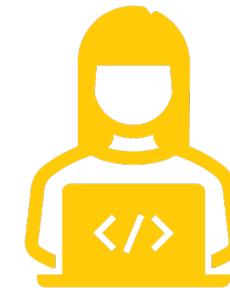


# Objectives for the next 3 days



## Theory: *Introduction to scRNA-seq analysis*

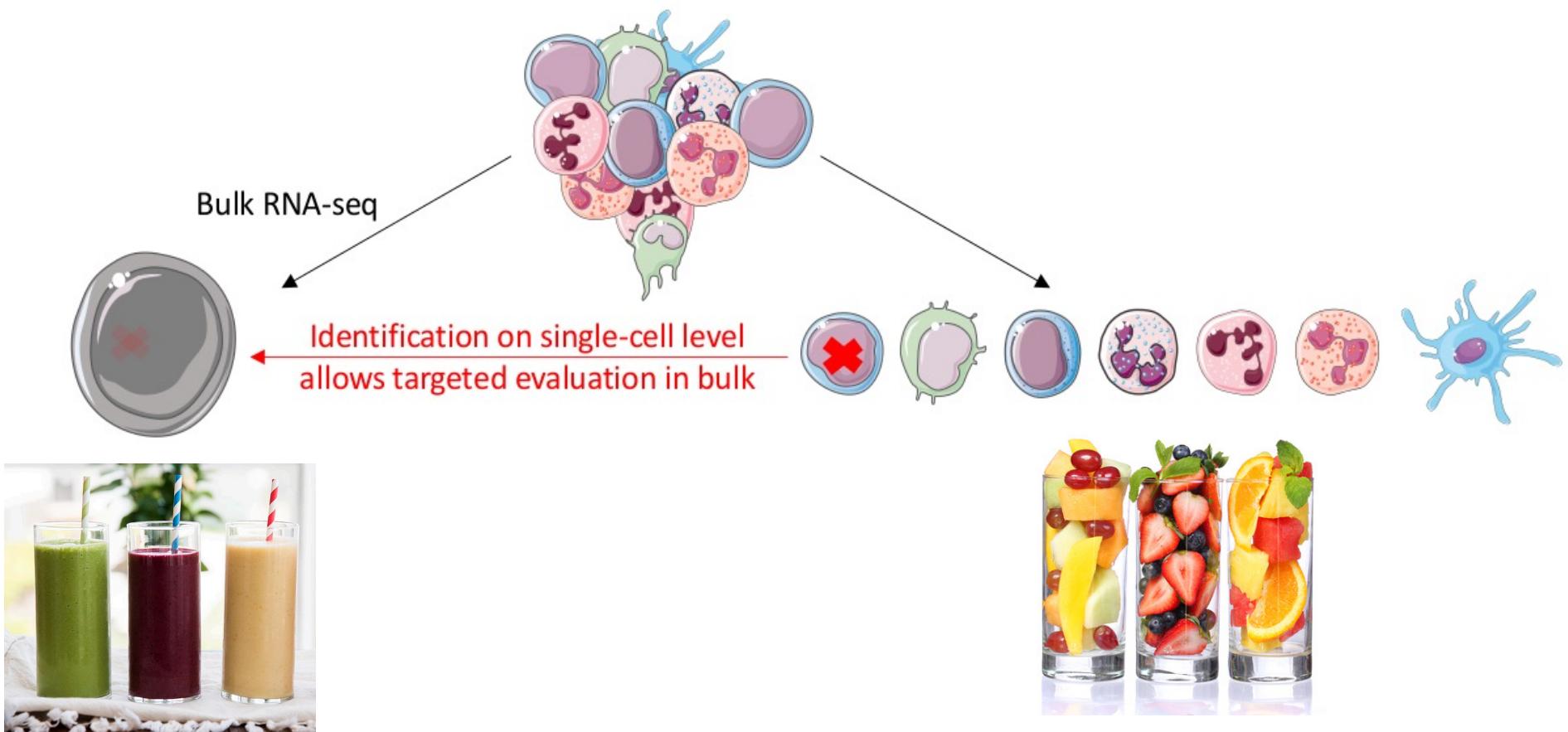
- Concepts of single cell RNA-seq technology
- Workflow of single cell analysis
- Limitations of the technology



## Workshop: *Practical Single cell transcriptome analysis 101*

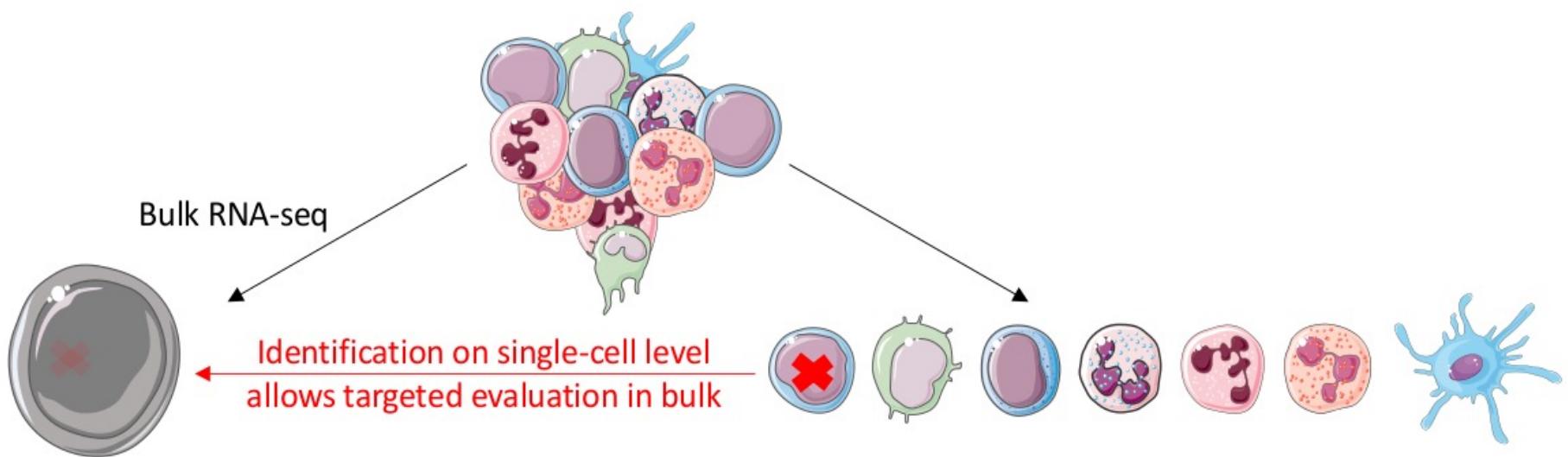
- Advances in medical life sciences using single cell technology
- COVID-19 pandemic as example

# In simple ‘words’:



Modified from Servier Medical Art by Jonas Schulte-Schrepping  
<http://365-smoothie-rezepte.de/obst-beeren-und-gemuese-sorgfaeltig-auswaehlen/>

# In simple ‘words’:



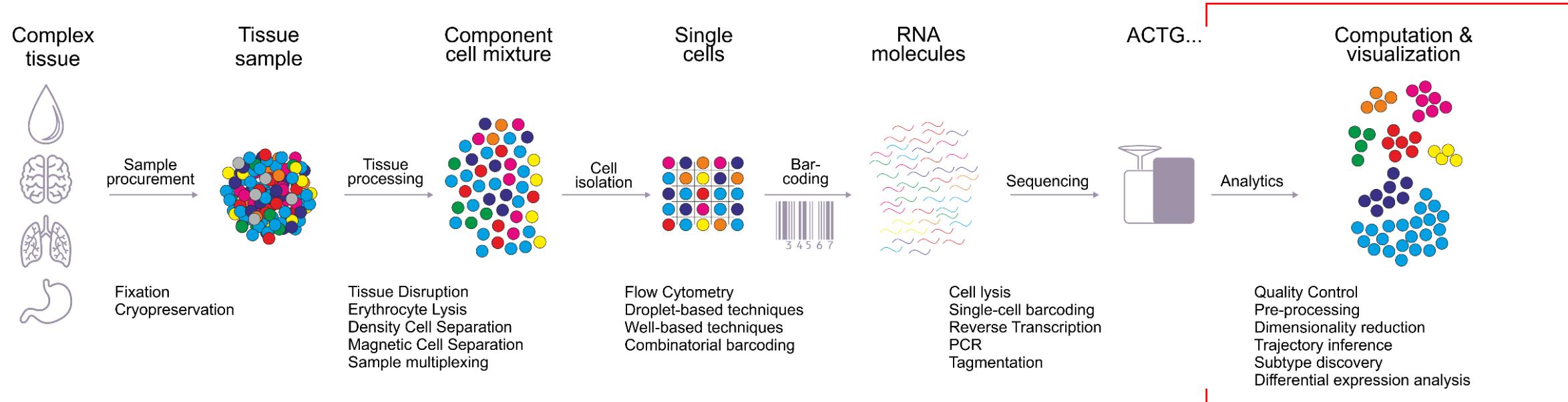
## e.g. *bulk RNA-seq*

- Average expression
- ‘Homogeneous’ cell population

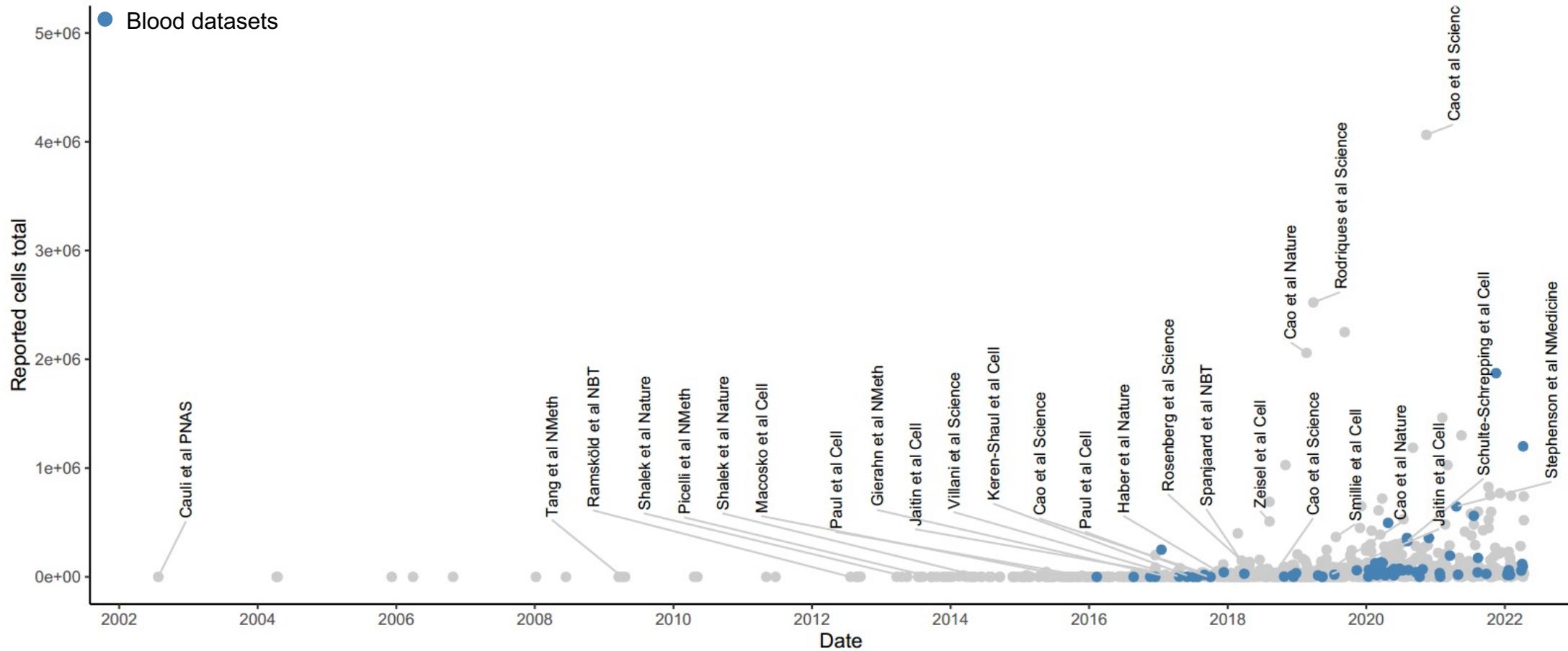
## e.g. *single cell RNA-seq*

- Cell specific expression
- Heterogenous cell population

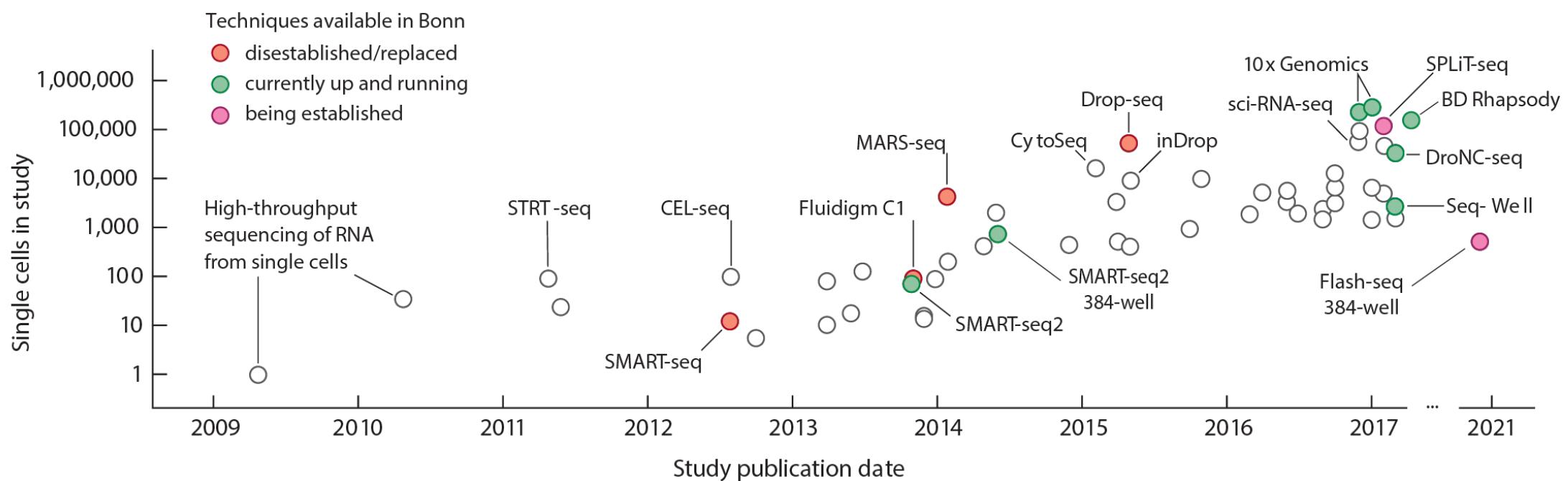
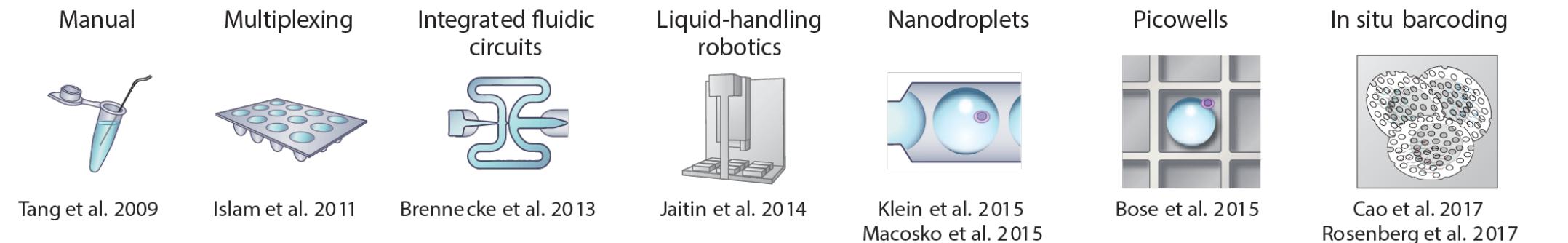
# Single-cell RNA-seq experiments are highly complex



# The evolution of single-cell RNA-seq

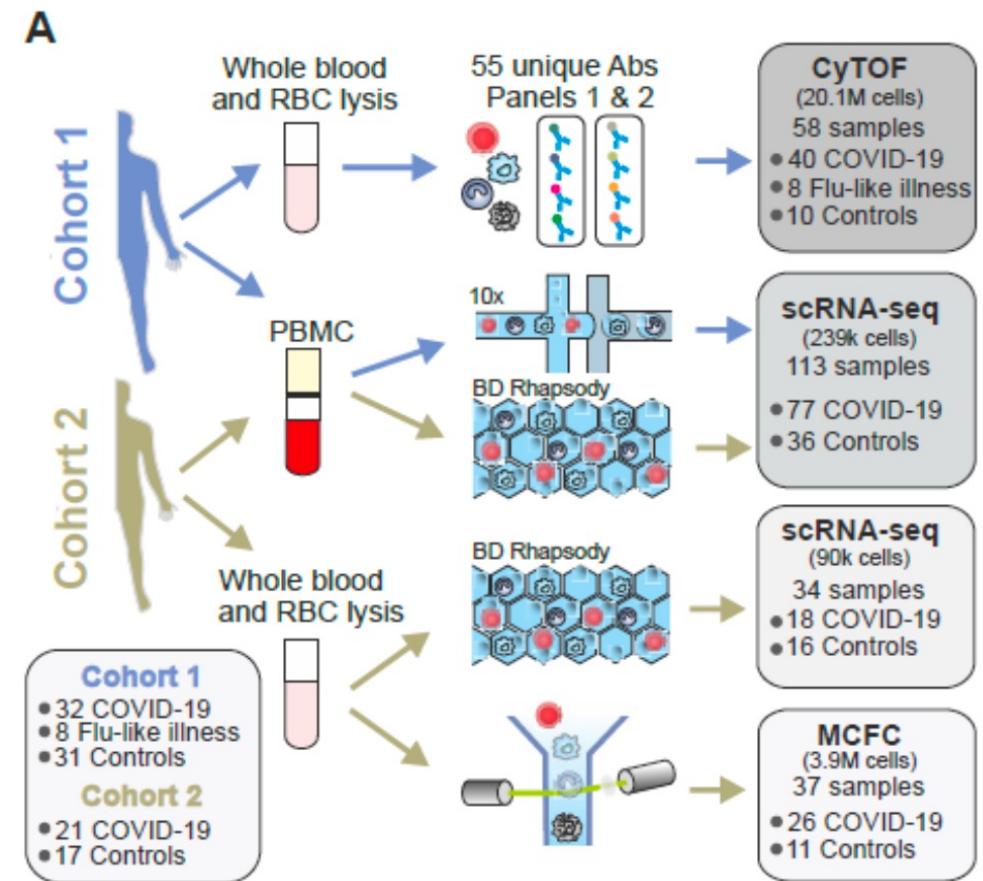


# Single-cell: what can we do?

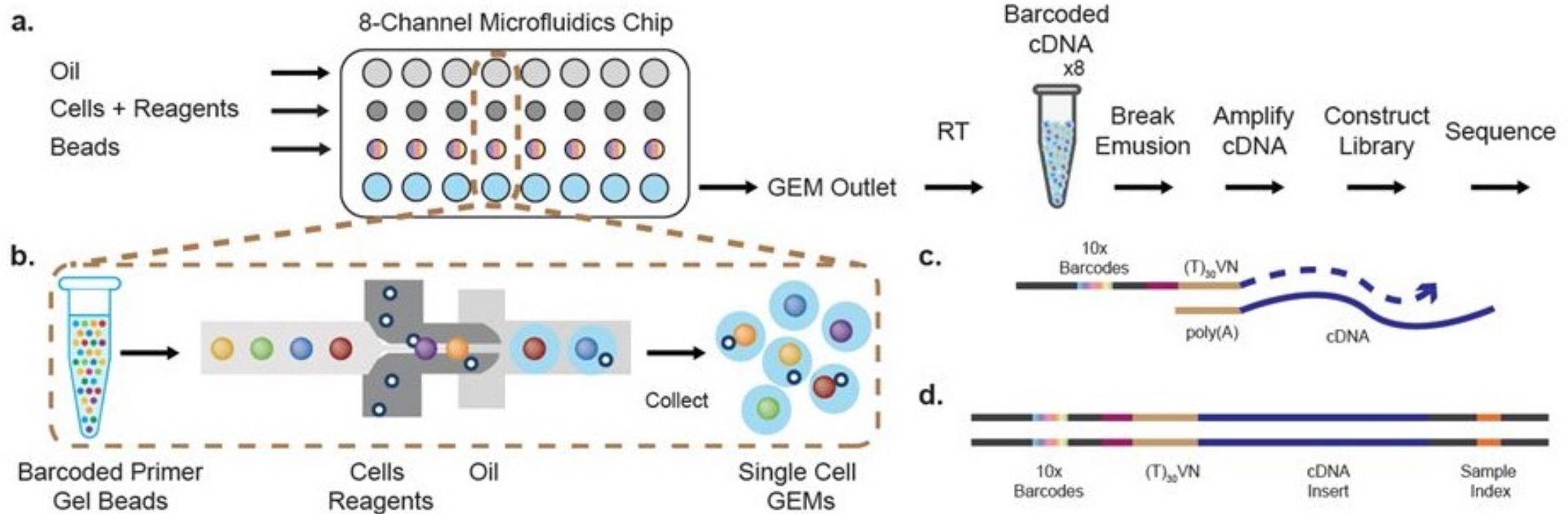


# Single cell technology available

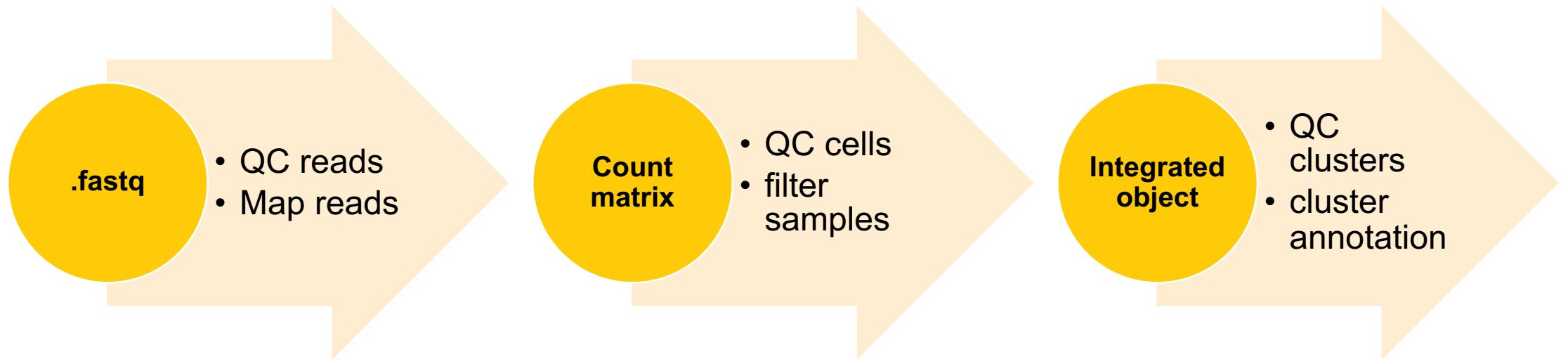
- BD Rhapsody
- **10x Genomics**
  - **Gene expression:** 3' and 5'
  - VDJ' (TCR and BCR)
  - CITE-seq
  - ATAC-seq
  - Multiome
  - **Spatial transcriptomics**



# 10x Genomics- single cell gene expression



# Processing workflow



# From reads to count matrix

- cell ranger

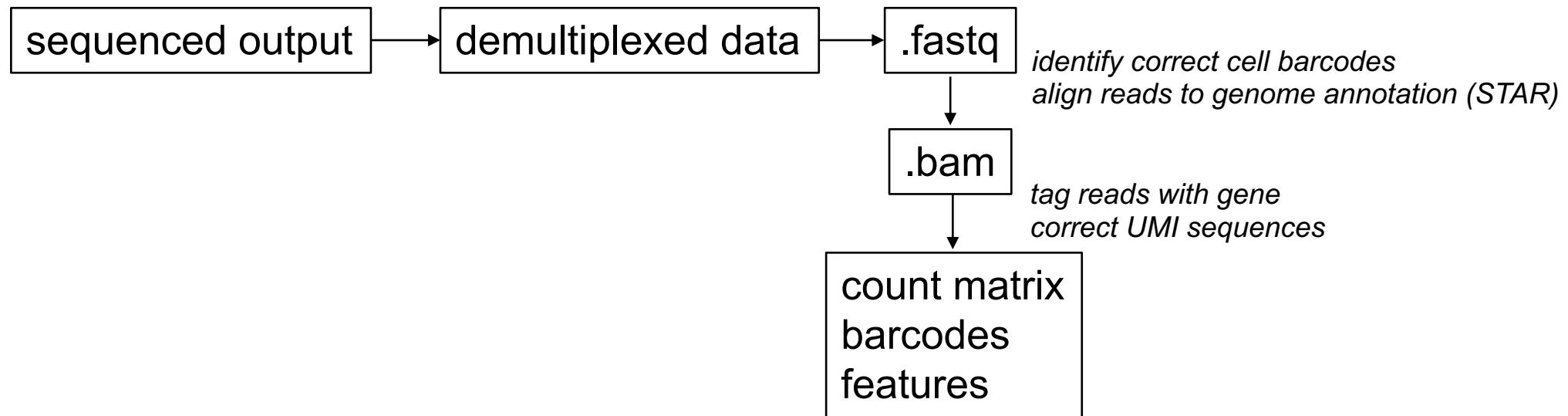


.fastq

- QC reads
- Map reads

# From reads to count matrix

- cell ranger

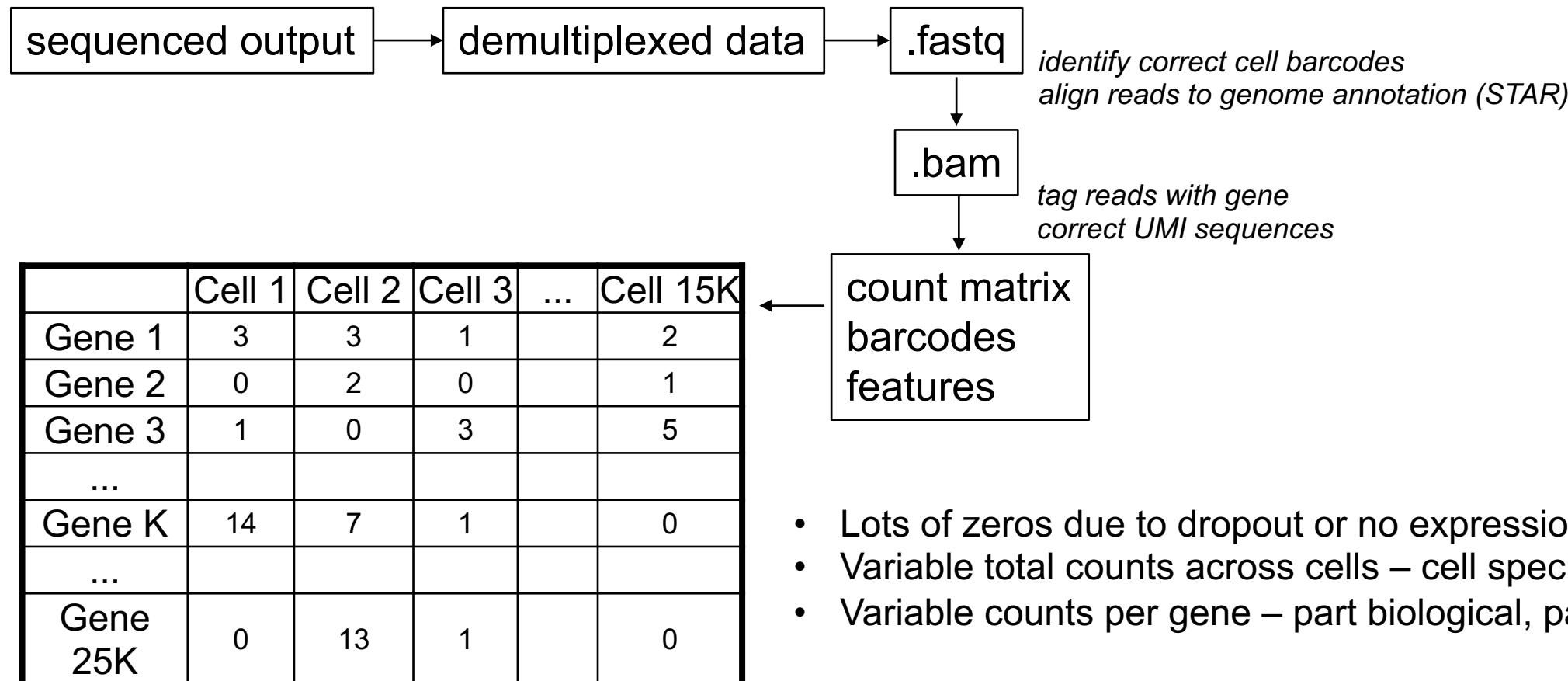


- QC reads
- Map reads

# From reads to count matrix

- 
- QC reads
  - Map reads

- cell ranger



# From matrix to sample objects

- Create object

*R packages*

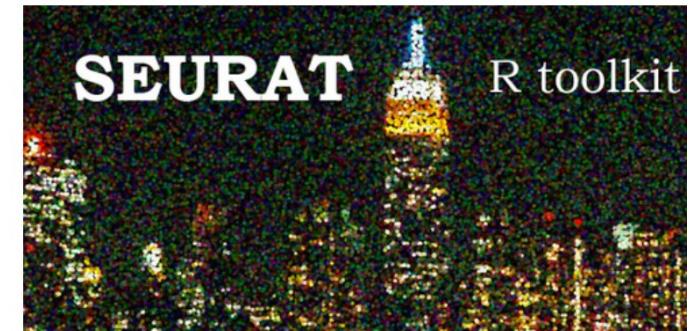
- **Seurat V4 > S4 seurat object**
- scater V1.18.3 > SingleCellExperiment
- Monocle3 > derived SingleCellExperiment

*Python packages*

- scanpy V1.7.0 > Anndata
- loom 3.0.0 > h5ad



- QC cells
- prepare samples



# The Seurat Object

- Each type of data has it's own slot
- Facilitates streamlined analysis, easy execution of functions

e.g. raw counts:

```
pbmc@assay$counts
```

```
GetAssayData(object = pbmc,  
             slot = 'counts')
```

```
> str(pbmc)  
Formal class 'Seurat' [package "Seurat"] with 12 slots  
  ..@ assays      :List of 1  
  ...$ RNA:Formal class 'Assay' [package "Seurat"] with 12 slots  
    ..@ counts      :Formal class 'dgCMatrix' [pri...  
    ...@ i          : int [1:2491554] 21 26 28 :  
    ...@ p          : int [1:1177] 0 2618 4423 !  
    ...@ Dim         : int [1:2] 15246 1176  
    ...@ Dimnames   :List of 2  
      ...$ : chr [1:15246] "AL627309.1" "AL627309.2"  
      ...$ : chr [1:1176] "AAACCCAAGGAGAGTA"  
    ...@ x          : num [1:2491554] 1 1 1 1 2  
    ...@ factors    : list()  
    ...@ data        :Formal class 'dgCMatrix' [pri...  
    ...@ i          : int [1:2491554] 21 26 28 :  
    ...@ p          : int [1:1177] 0 2618 4423 !  
    ...@ Dim         : int [1:2] 15246 1176  
    ...@ Dimnames   :List of 2  
      ...$ : chr [1:15246] "AL627309.1" "AL627309.2"  
      ...$ : chr [1:1176] "AAACCCAAGGAGAGTA"  
    ...@ x          : num [1:2491554] 1 1 1 1 2  
    ...@ factors    : list()  
    ...@ scale.data  : num[0 , 0 ]  
    ...@ key        : chr "rna_"  
    ...@ var.features: logi(0)  
    ...@ meta.features:'data.frame': 15246 obs.  
    ...@ misc       : NULL  
  ..@ meta.data   :'data.frame': 1176 obs. of  3 :  
  ...$ orig.ident : Factor w/ 1 level "10X_PBMC": 1 1  
  ...$ nCount_RNA : num [1:1176] 8286 5509 4280 2754 :  
  ...$ nFeature_RNA: int [1:1176] 2618 1805 1559 1225 :  
  ..@ active.assay: chr "RNA"  
  ..@ active.ident: Factor w/ 1 level "10X_PBMC": 1 1 1 :  
  ...- attr(*, "names")= chr [1:1176] "AAACCCAAGGAGAGTA"  
  ..@ graphs      : list()  
  ..@ neighbors   : list()  
  ..@ reductions  : list()  
  ..@ project.name: chr "10X_PBMC"  
  ..@ misc        : list()  
  ..@ version     :Classes 'package_version', 'numeric_version'  
  ...$ : int [1:3] 3 0 2  
  ..@ commands    : list()  
  ..@ tools       : list()  
  |
```

# From matrix to sample objects

- Quality control/filtering of cells:

- Number of genes [200; 5.000]

- remove ambient RNA
  - remove doublets

- % Mitochondria (<25%)

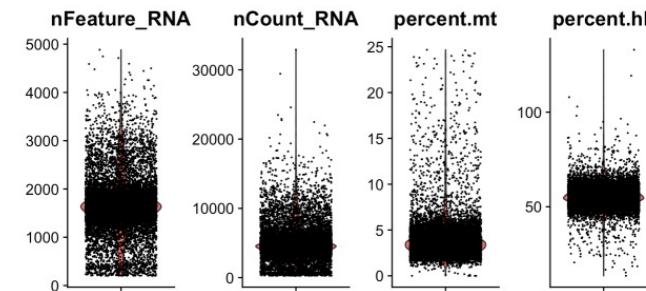
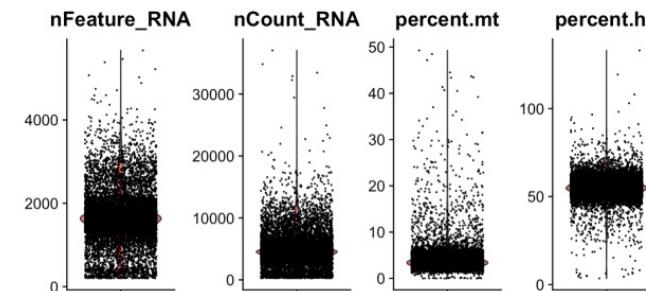
- remove disrupted cells

- % housekeeping genes (>30%)

- global quality measure



- QC cells
- prepare samples



# From matrix to sample objects

- Add the metadata- each *cell with information*

% mitochondria

% housekeeping genes

Sample code

Gender

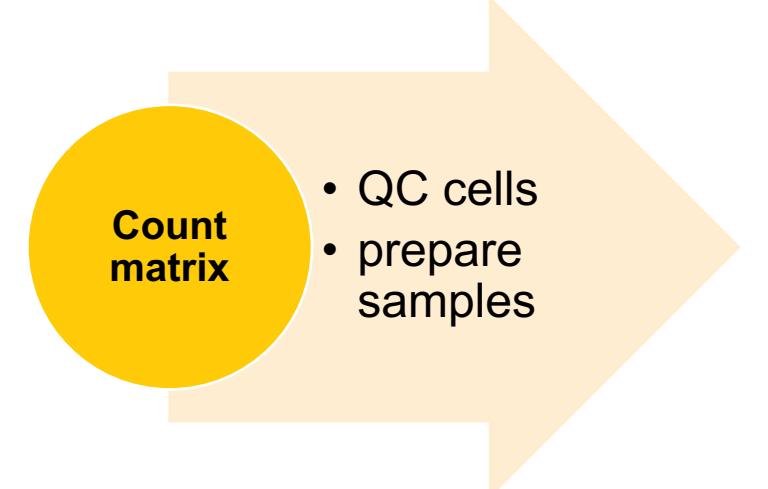
Treatment

Disease

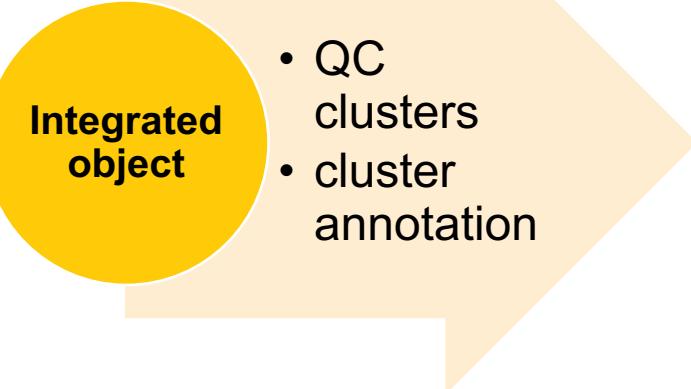
Timepoint

....

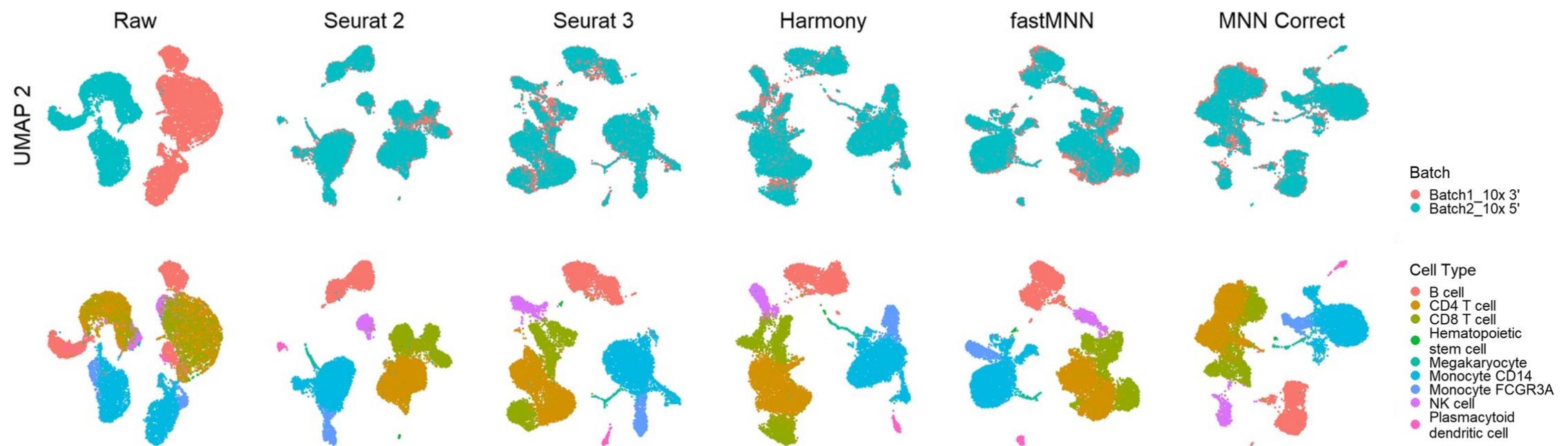
Cell annotation



# From samples to one object



- Sample integration
  - 1. Simple merge

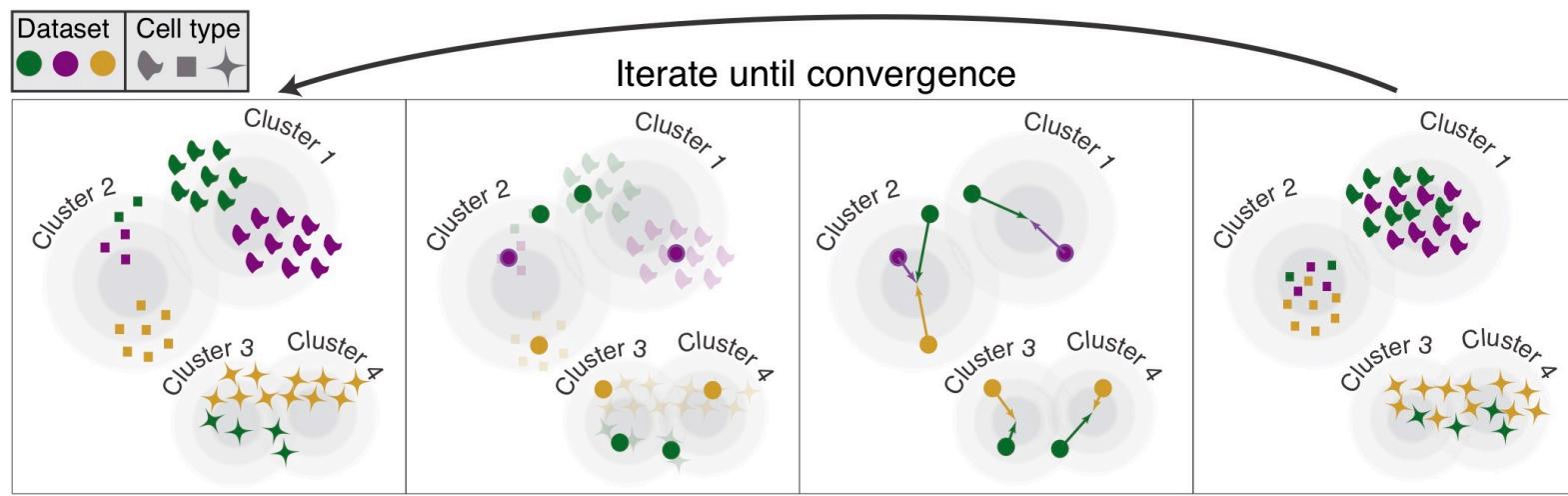


# From samples to one object



- QC clusters
- cluster annotation

- Sample integration
  - 1. Integration with correction
    - a. Harmony integration



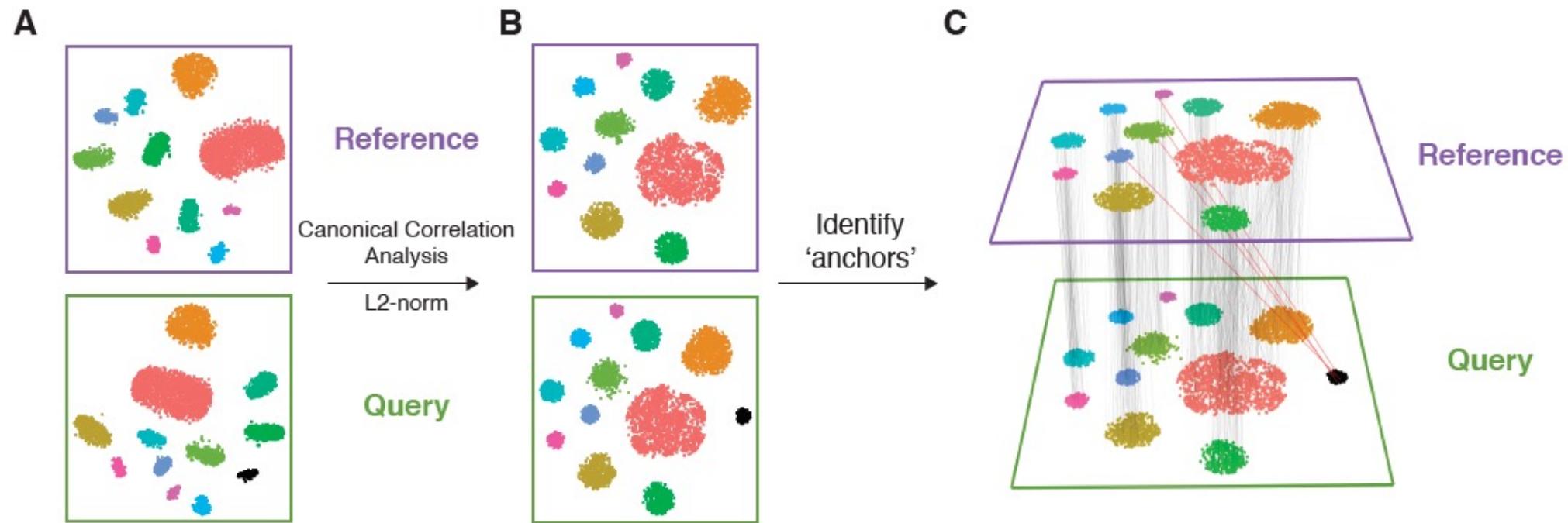
<https://portals.broadinstitute.org/harmony/articles/quickstart.html>

# From samples to one object



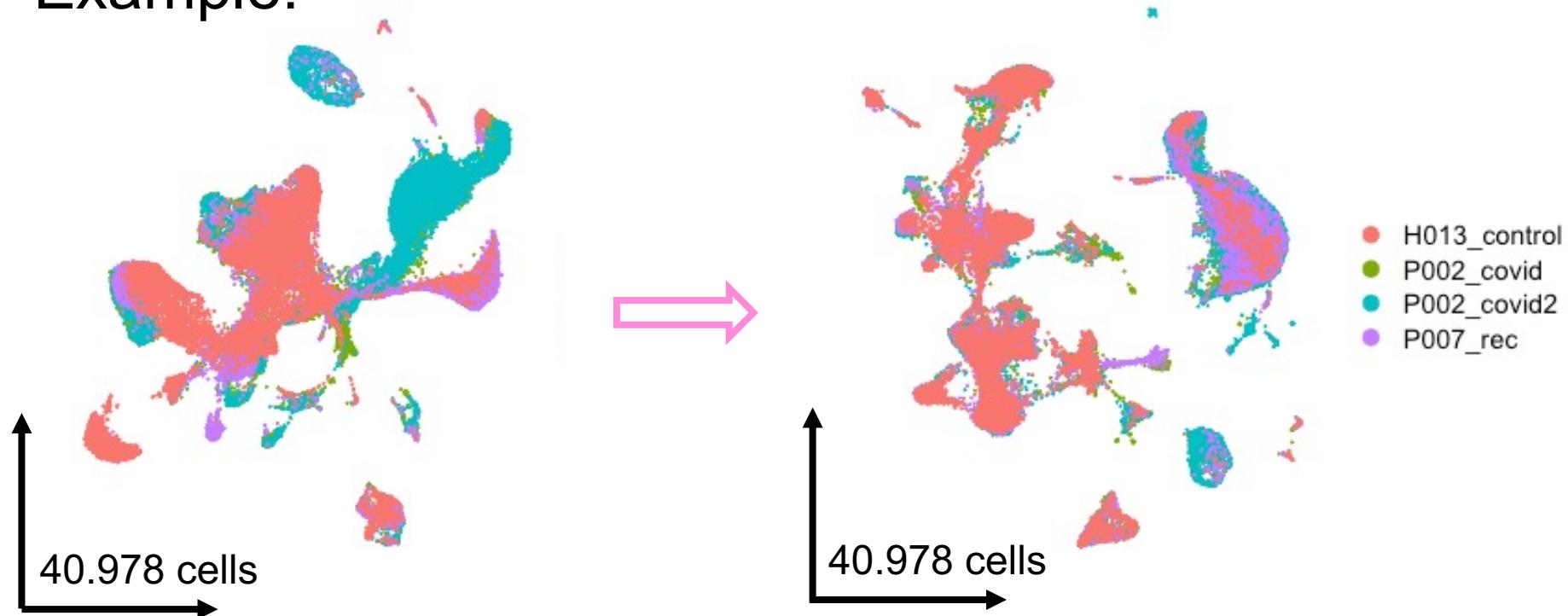
- QC clusters
- cluster annotation

## b. Seurat sample integration



# From samples to one object

- Example:

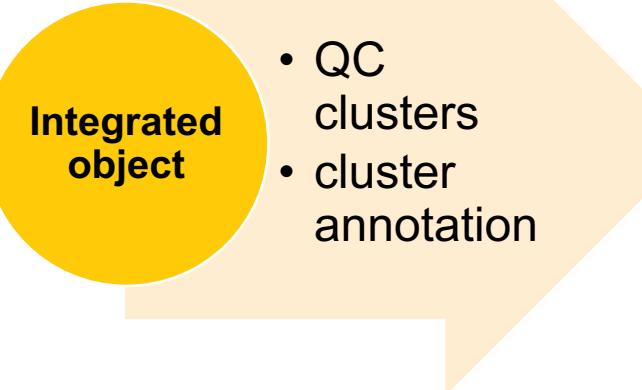


Integrated  
object

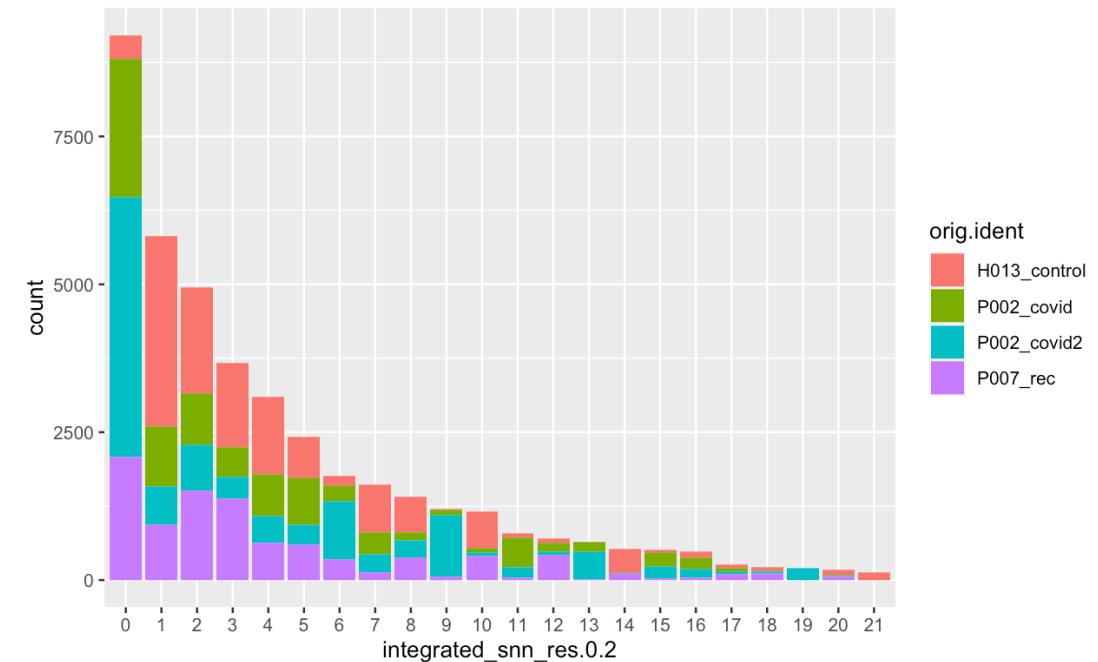
- QC clusters
- cluster annotation

! Goal is not to lose biological relevant features.

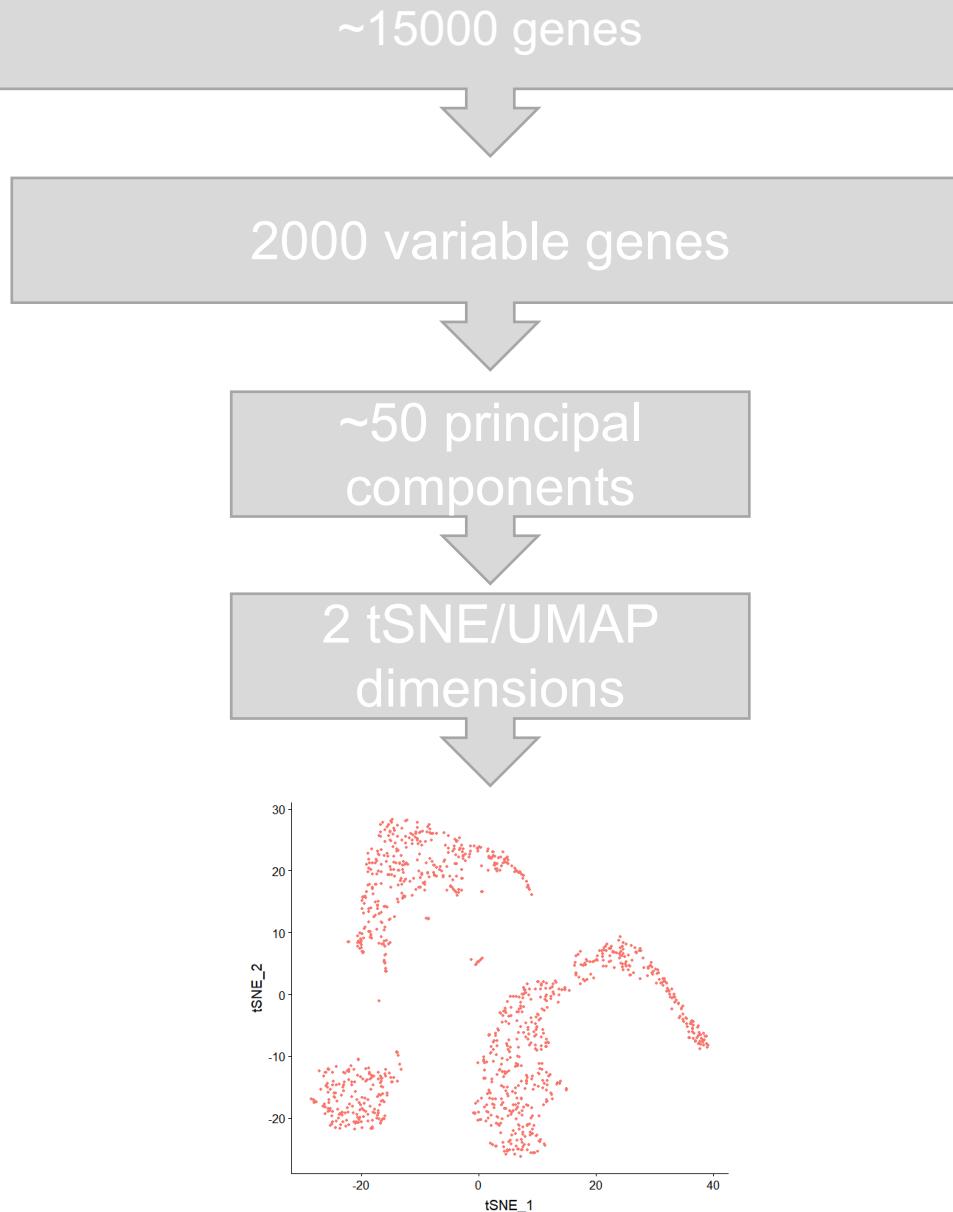
# From samples to one object



- Doublet removal:
  - Non-informative cells removal
  - actual doublets
  - cells sequenced too shallow
  - clusters with sample enrichment

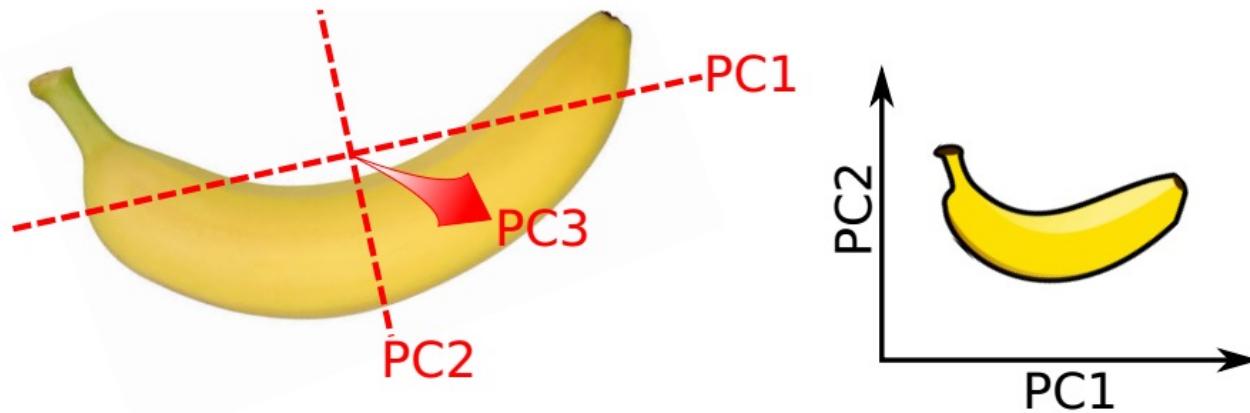


# Dimensionality reduction



# Dimensionality reduction with PCA

- Identifies axes of maximal variance in high-dimensional data
- Each subsequent principal component explains less variance



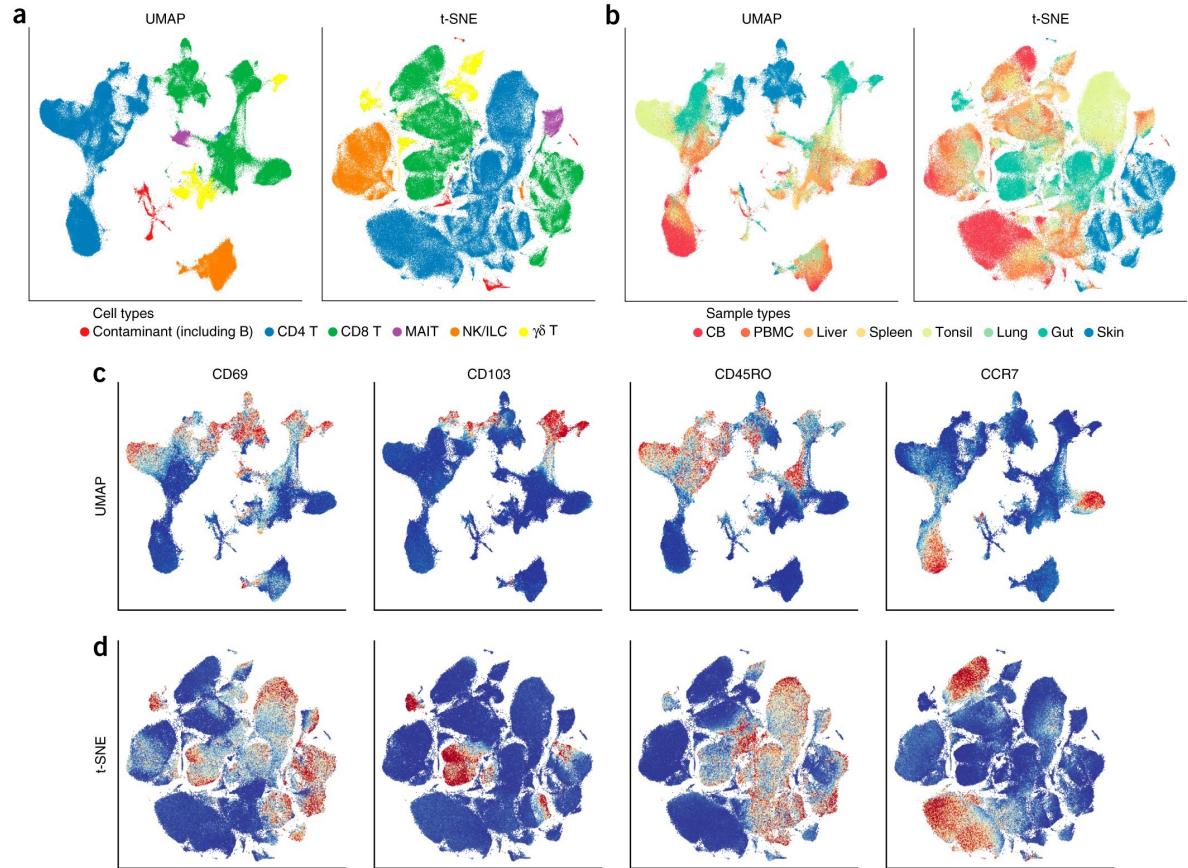
- Use the first few PCs as a “summary” of the relevant variance in the data
- Speed up downstream analyses by reducing dimensionality
  - Focus on biology, remove random noise in later PCs

<https://www.youtube.com/watch?v=FgakZw6K1QQ>

Slide adapted from Single-cell Analysis Workshop by Aaron Lun, CRUK Cambridge Institute

# Uniform Manifold Approximation and Projection for Dimension Reduction

UMAP by McInnes et al. is an incredibly powerful tool in the data scientist's arsenal, and offers a number of advantages over t-SNE. While both UMAP and t-SNE produce somewhat similar output, the increased speed, better preservation of global structure, and more understandable parameters make UMAP a more effective tool for visualizing high dimensional data. Finally, it's important to remember that no dimensionality reduction technique is perfect - by necessity, we're distorting the data to fit it into lower dimensions - and UMAP is no exception.



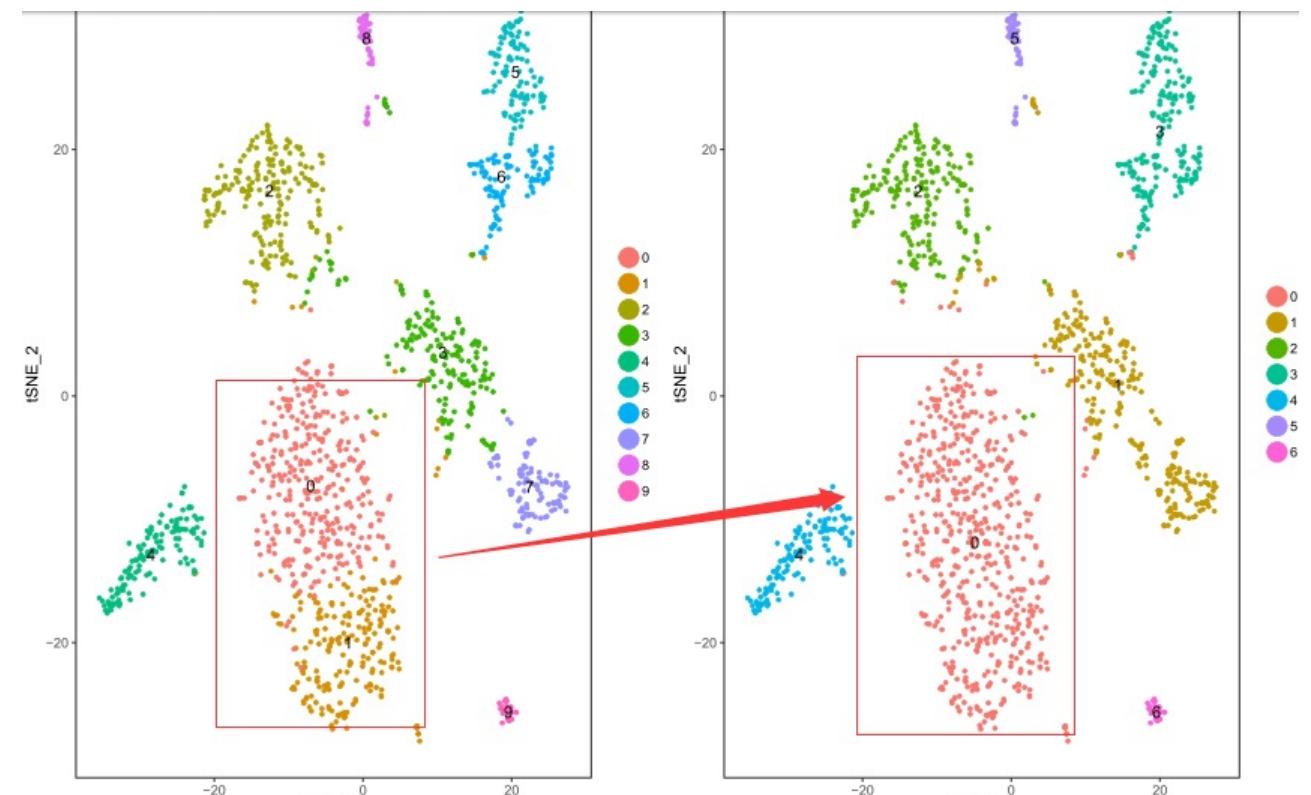
# How many clusters?

Cluster resolution is a parameter that controls the granularity of clustering in Seurat.  
Higher resolution > more clusters  
Lower resolution > fewer clusters

Keep in mind that the optimal cluster resolution may vary between datasets and experiments, and it's often a matter of experimentation and fine-tuning to strike the right balance between granularity and interpretability of the results.

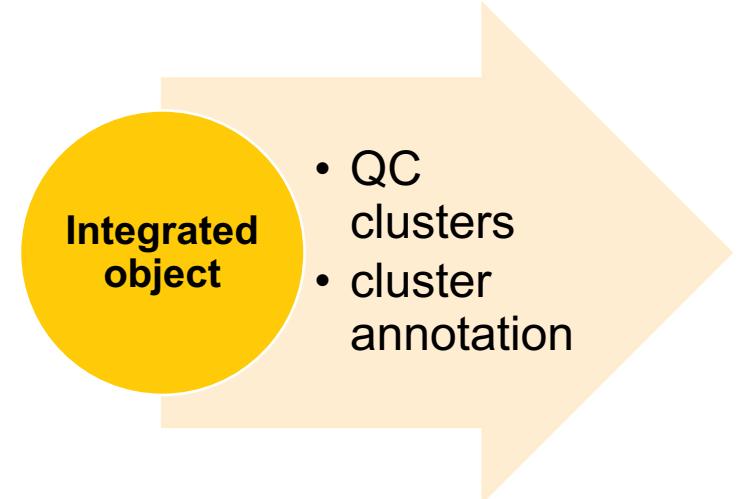


- QC clusters
- cluster annotation

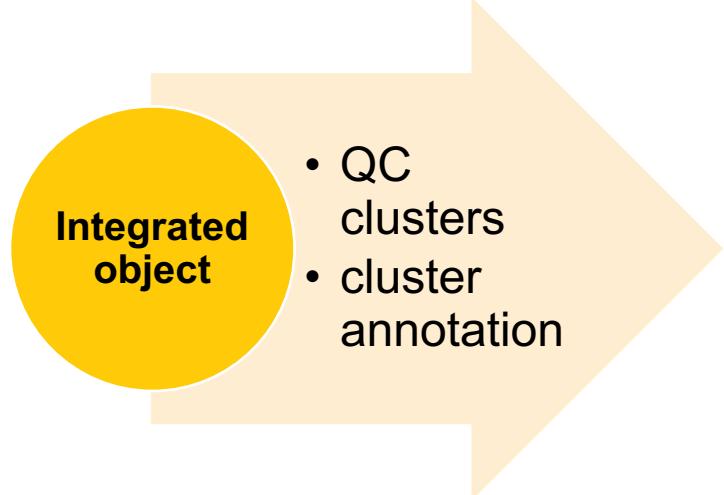


# From samples to one object

- Cluster annotation:
  1. Reference based method
  2. Package based method
  3. Literature or *in-house* biomarkers



# From samples to one object



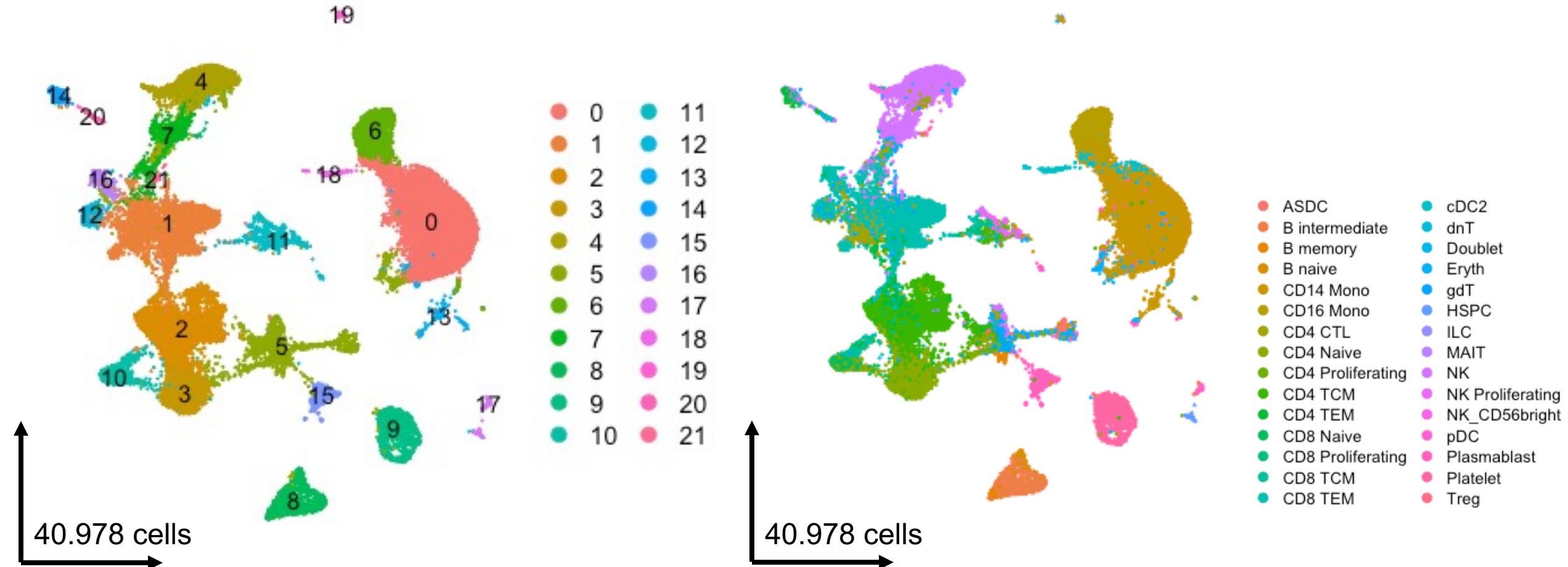
- Reference based method

- Choose a reference available in literature (e.g. spleen or PBMCs)
- Use the manuscript annotated cells as anchors/references for your own object/query.

Advantages:

Disadvantages:

# Reference based annotation (Seurat V4)



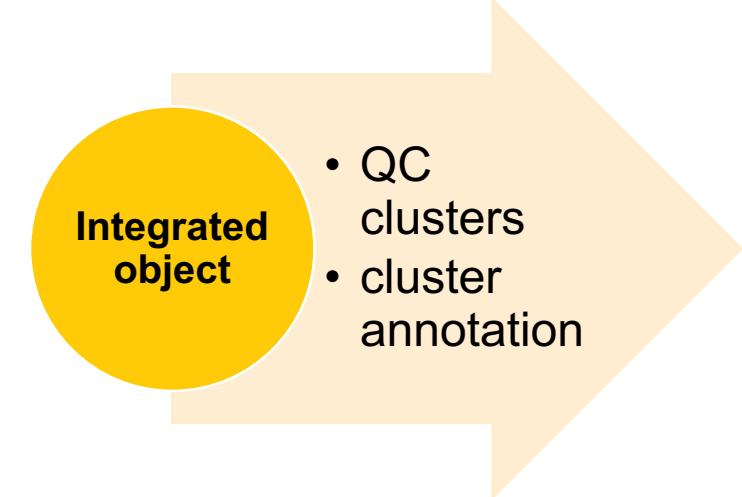
# From samples to one object

- Package based method

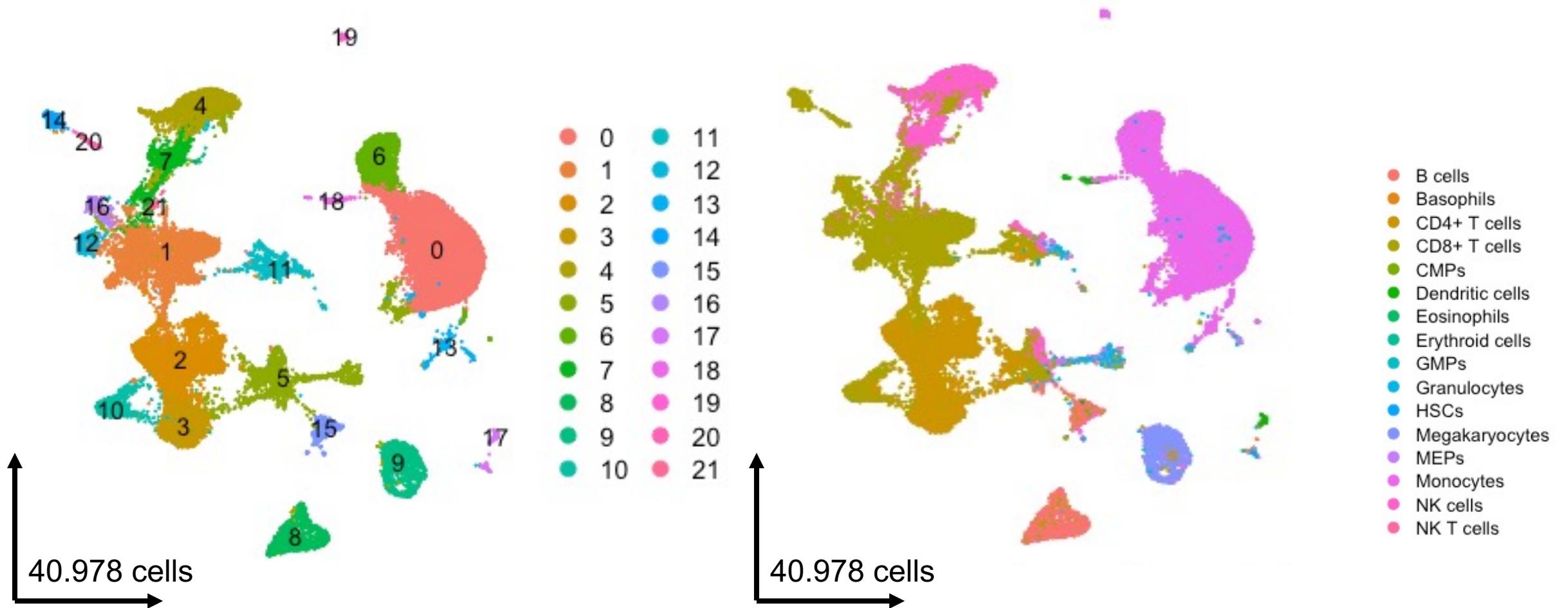
- Use packages available (SingleR)
- Use your cells expression and map to databases (single cell and flow cytometry)

Advantages:

Disadvantages:



## Package based annotation (SingleR- dmap)



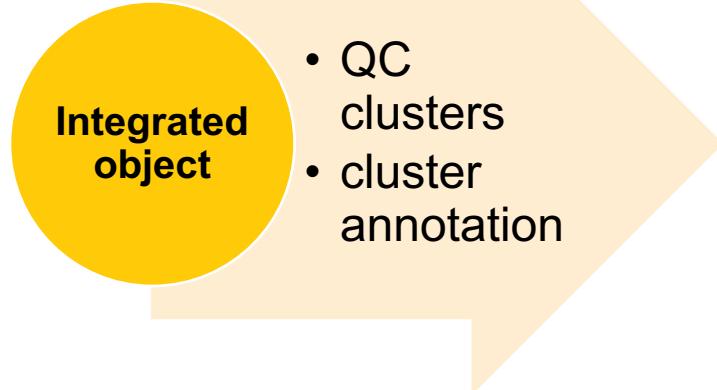
# From samples to one object

- Literature or *in-house* biomarkers

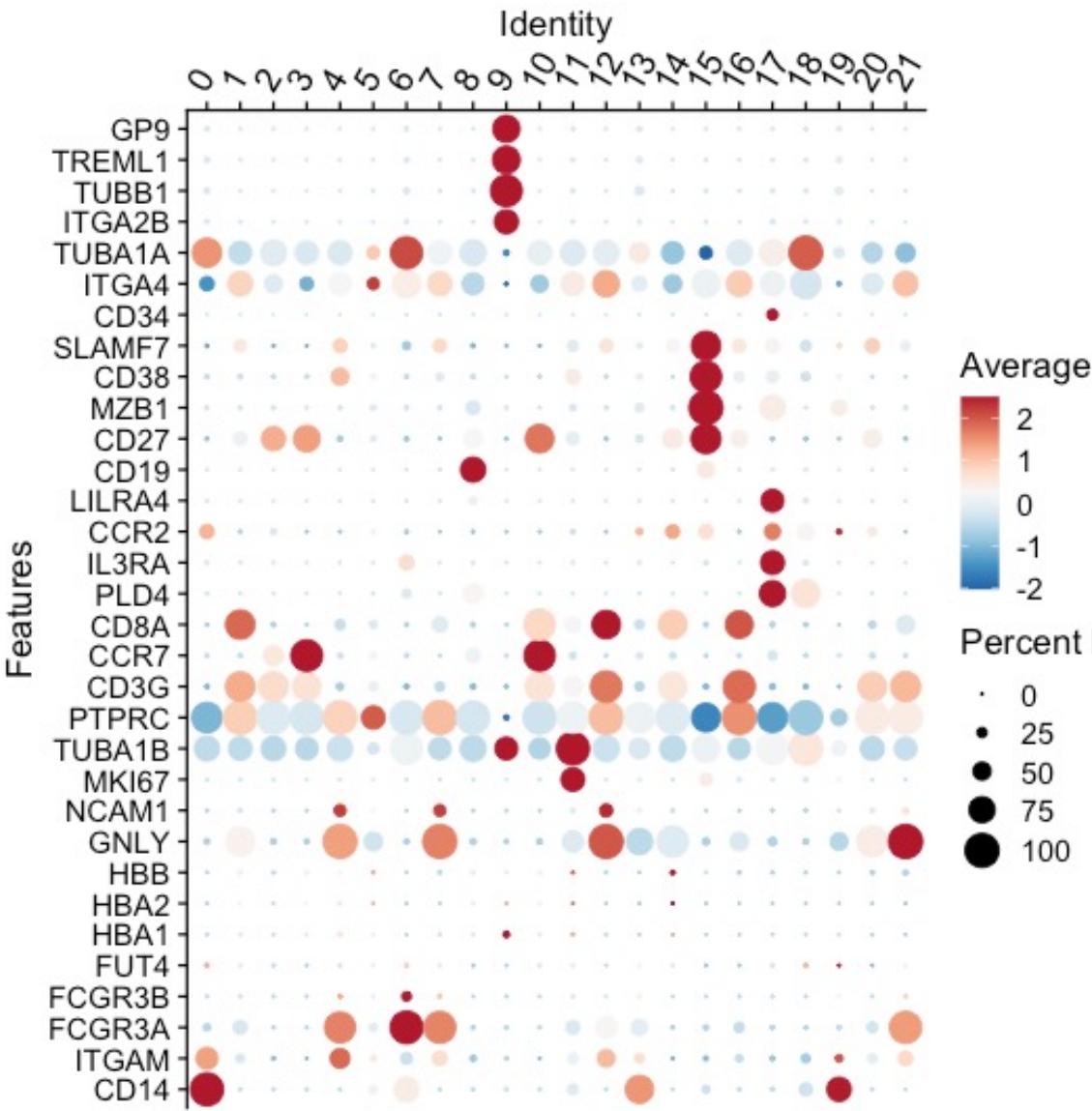
- Scientists have been using cell expression to identify cells for decades.
- Use (e.g. flow cytometry) markers to identify your cells based on cluster expression.

Advantages:

Disadvantages:



# Cluster based annotation (in-house/literature knowledge)



0: CD14 Monocytes

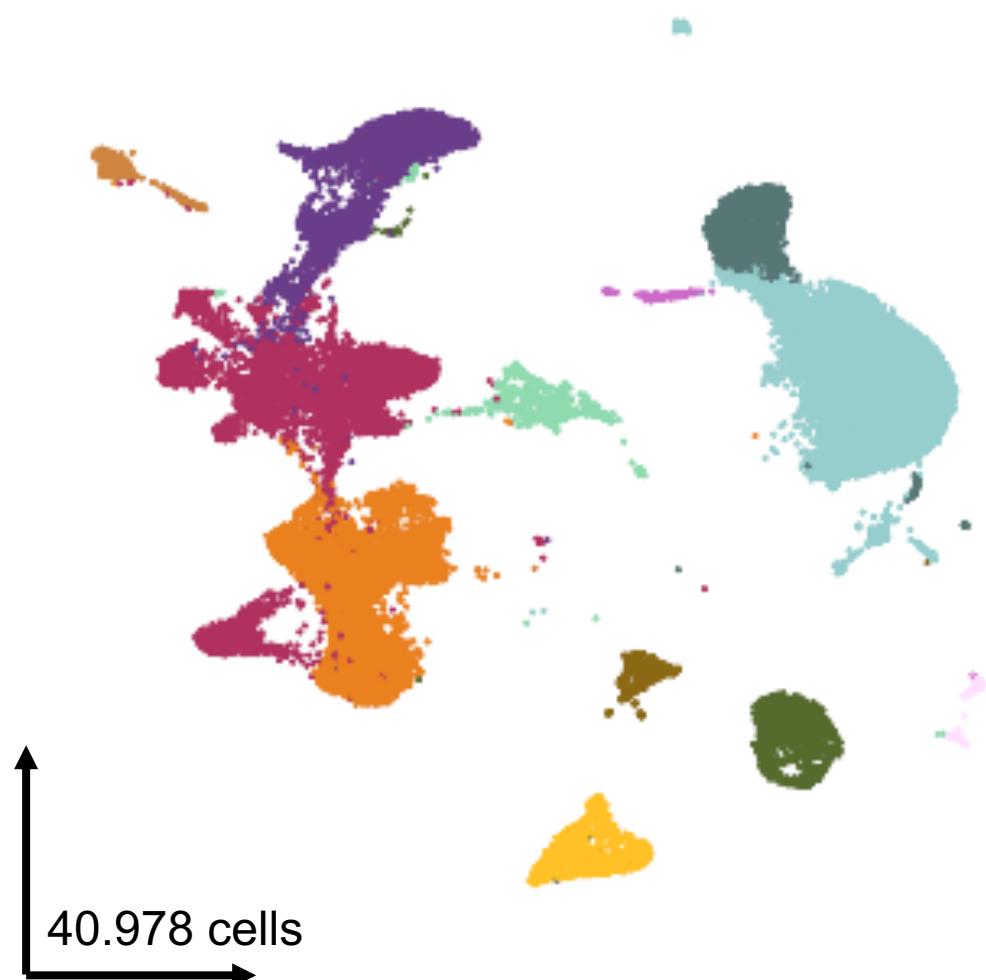
1: CD8 T cells

...

8: B cells

9: Megakaryocytes

# Your object



Integrated  
object

- QC clusters
- cluster annotation

# Recommendations

The object is big- depending on the cell number

- Work smart: parallelize operations
- Let your IT department/responsible person know

The focus should be on individual cell types/ compartments

- We normally do not want to ‘discover new cell types’
- Base the cell annotations on the experts in the field
- Histology books can be very useful

# Transcriptome analysis

- Finally we can start dissecting the biological signals:

cell proportion analysis	Differential gene expression	Coexpression analysis	pseudotime analysis
<ul style="list-style-type: none"><li>• scCODA</li><li>• ANOVA</li><li>• Mixed linear models</li></ul>	<ul style="list-style-type: none"><li>• Wilcoxon</li><li>• MAST</li></ul>	<ul style="list-style-type: none"><li>• hdWGCNA</li></ul>	<ul style="list-style-type: none"><li>• monocle3</li><li>• RNA velocity</li></ul>

**e.g. treated vs. non-treated or throughout timepoints**  
(covid-19 paper)

## cell proportion analysis

- scCODA
- ANOVA
- Mixed linear models

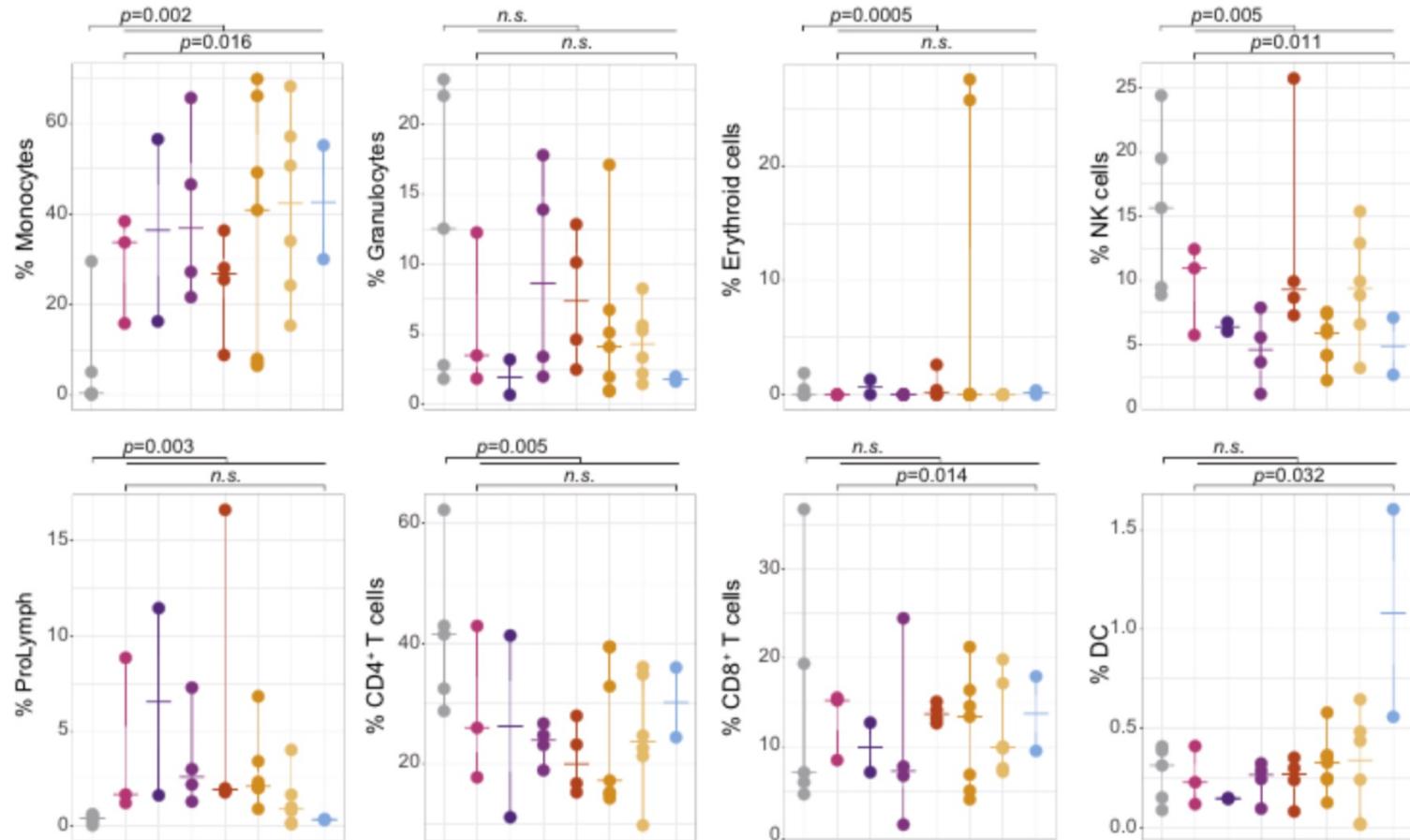
# Cell proportion analysis

- Identifies significant changes in cell proportions based on treatment or other parameters.

Disadvantages:

# Cell proportion change through time: Mixed linear model

## Cell proportion change between healthy and covid19: ANOVA

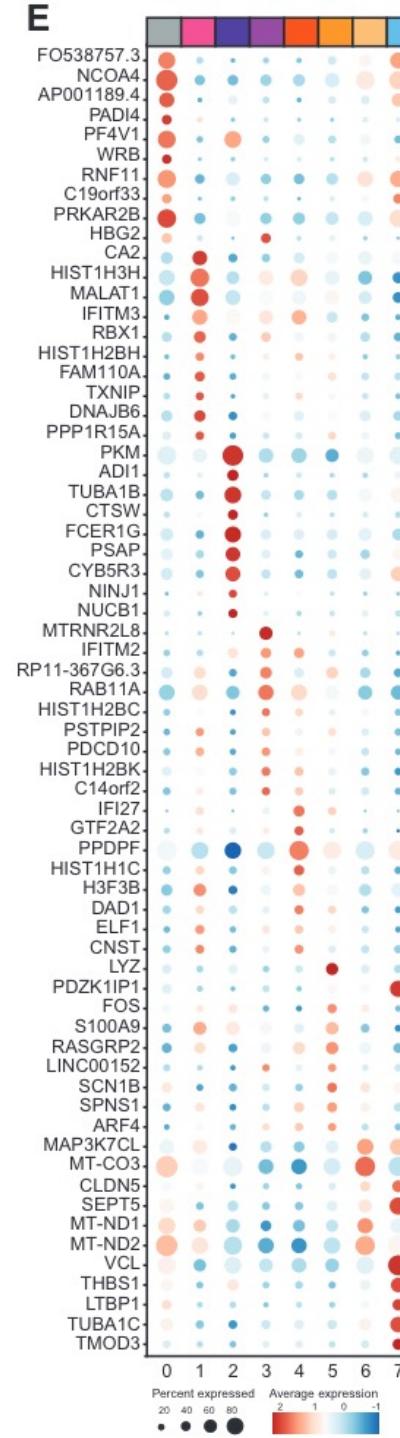
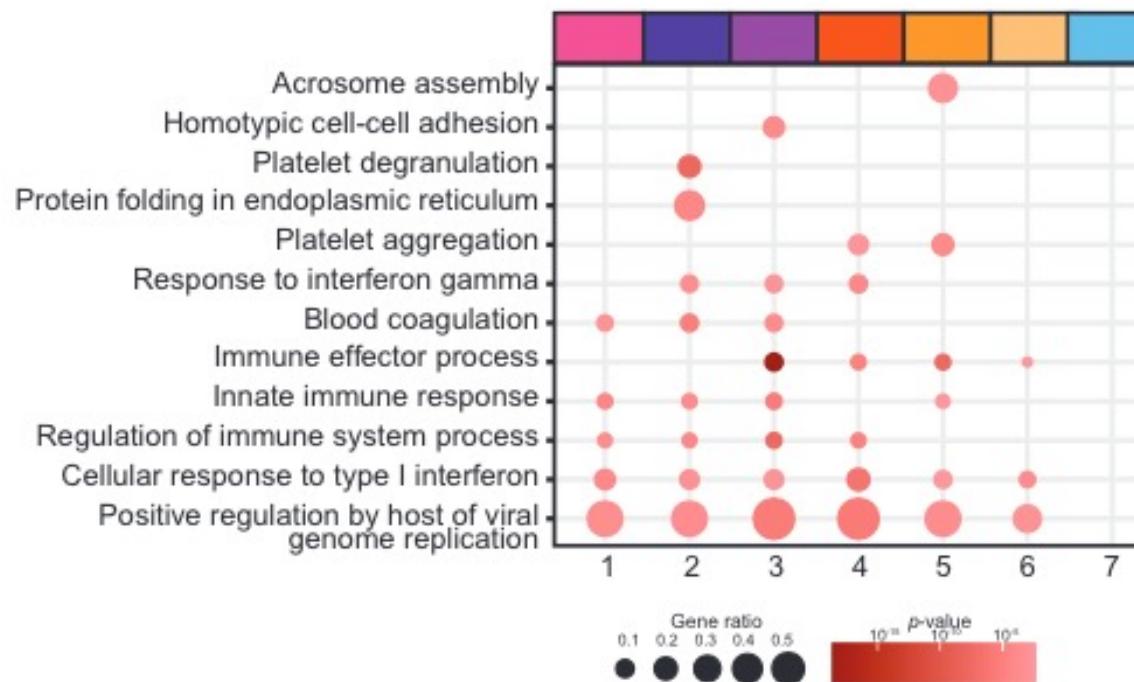


# Differential expression analysis

- Identifies expression changes between treatments, for individual genes within a cell type. Differential expression can be based on:
  - Signature genes: compare one against all other groups/clusters.
  - Pairwise comparisons: compare one group (e.g. treatment) to another group.

Disadvantages:

Differentially expressed genes can also be used to inform GO terms enriched throughout timepoints



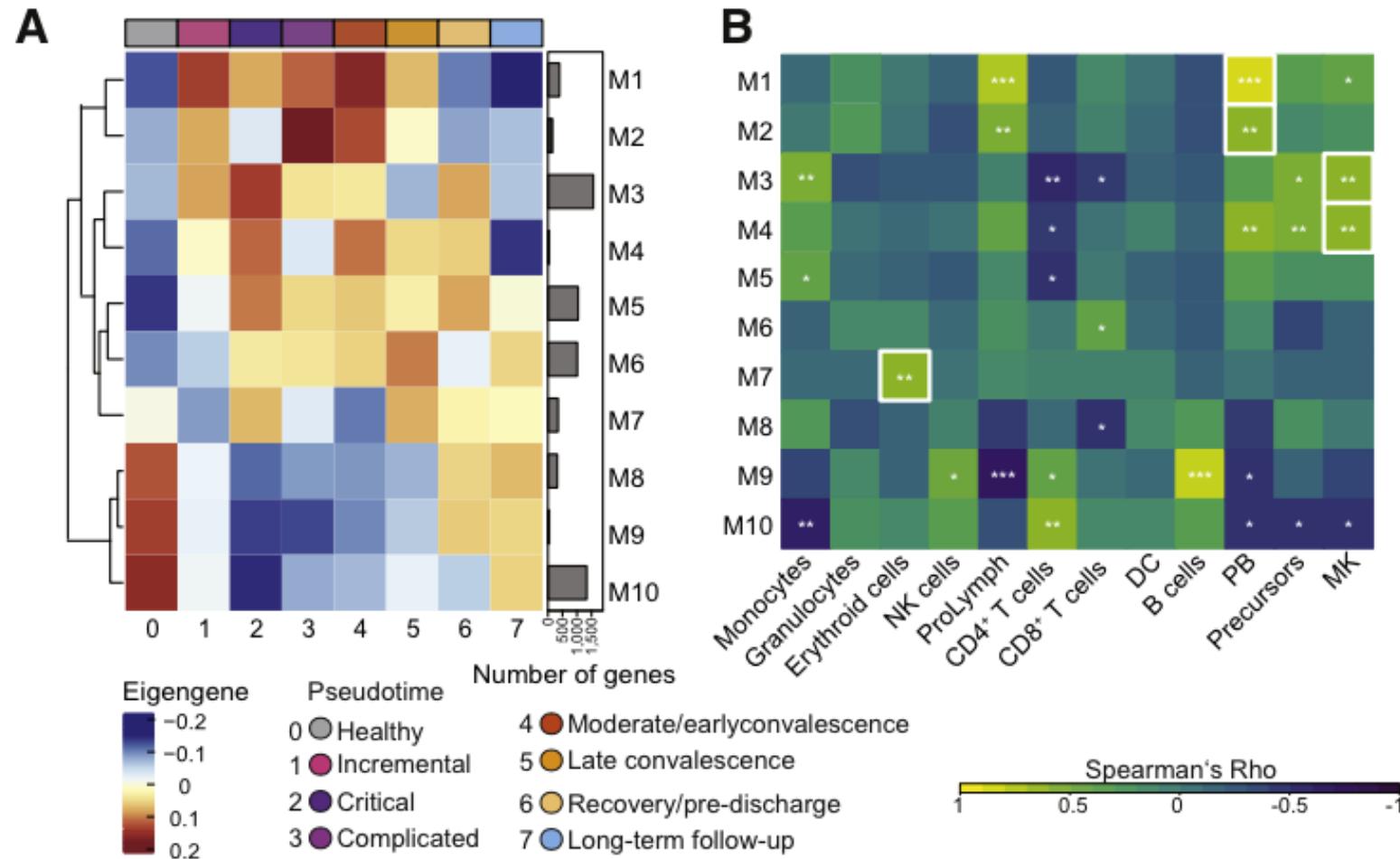
## Coexpression analysis

# Coexpression analysis

- Specially useful for longitudinal studies. Groups genes into modules that change similarly throughout time for a particular cell-type.
- Based on differentially expressed genes
- Relate modules with clinical parameters

Disadvantages:

# Coexpression modules can be calculated in bulk RNA-seq and correlated with cell-type proportions



## pseudotime analysis

- monocle3
- RNA velocity

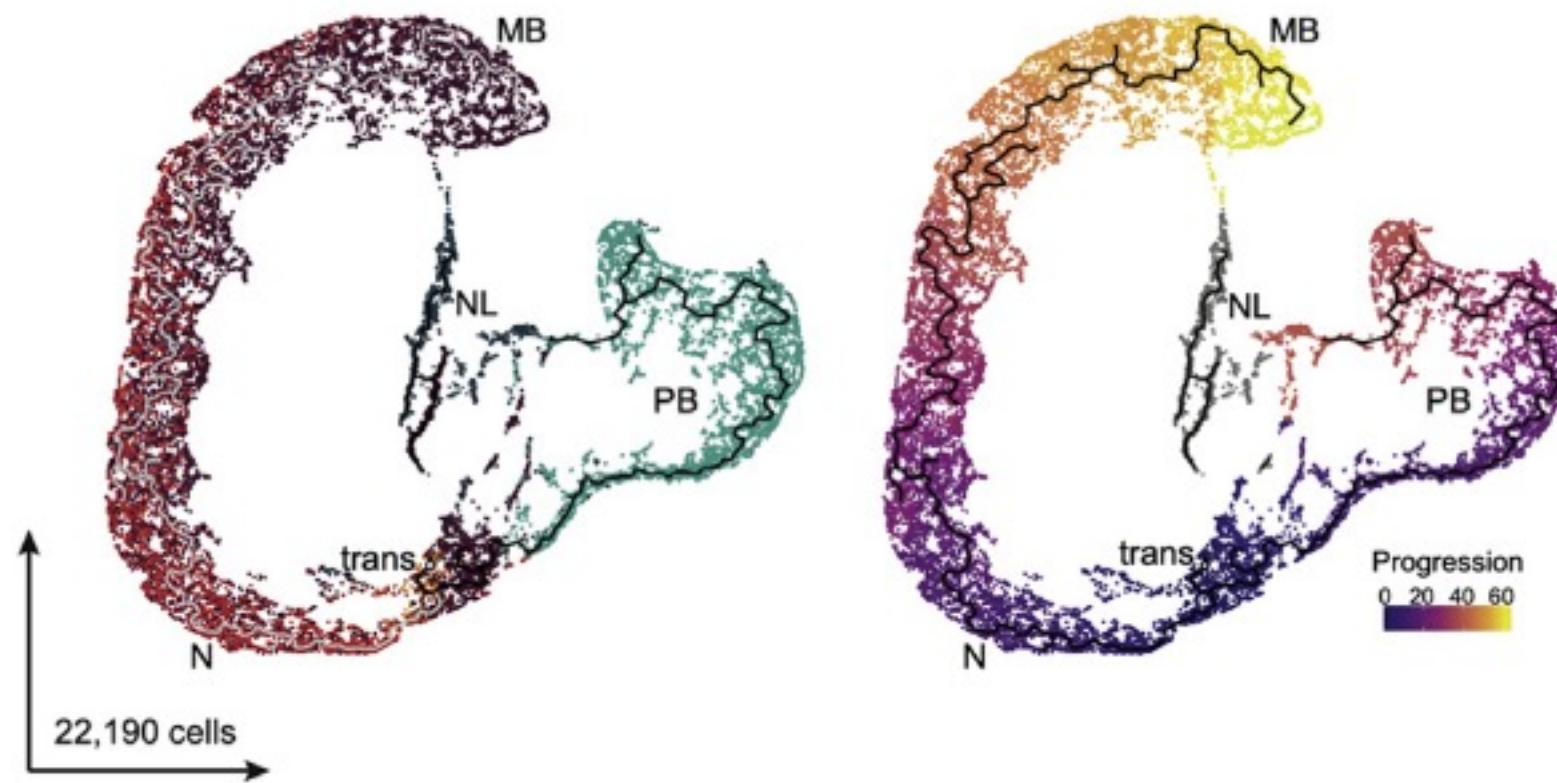
# Pseudotime or trajectory analysis

- Allocate cells to lineages and then orders them based on pseudotimes within these lineages.

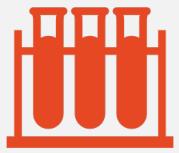
- Rooted (Monocle3)
- Unrooted: scVelo

Disadvantages:

## Cell-type lineage can resolve differentiation



# What we learned



Preparation of samples into the single cell object.

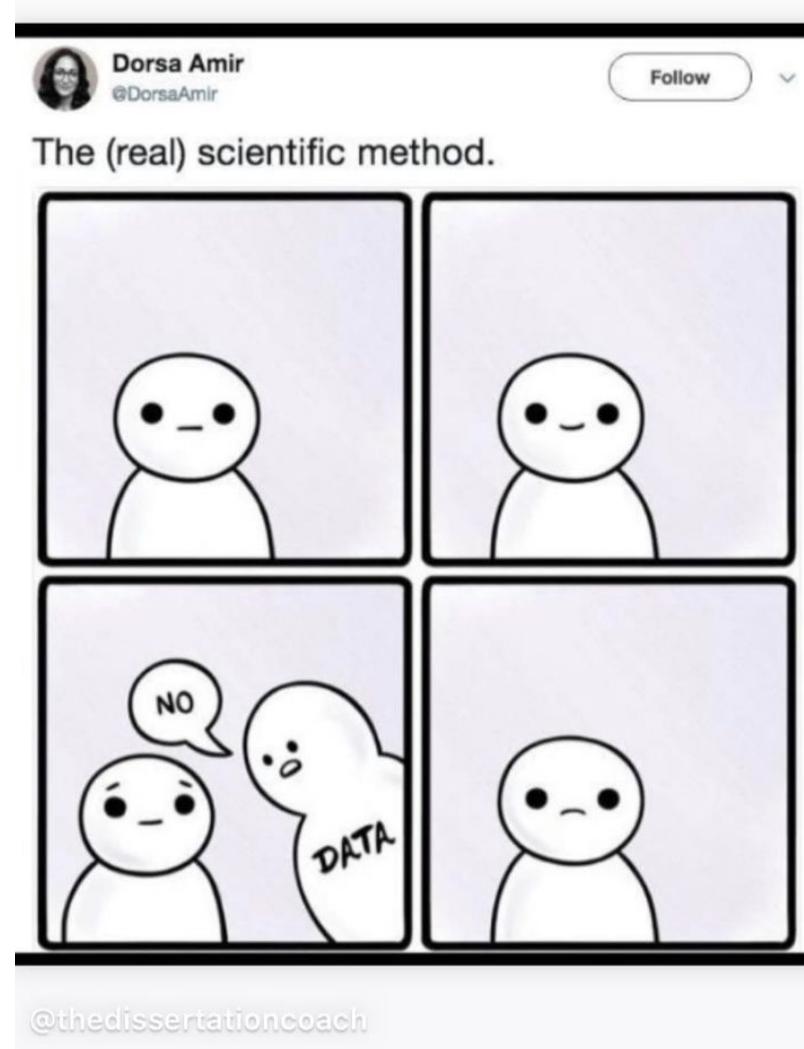


How to annotate the single cell object.



How to analyze the cell types annotated.

# Questions?



# Where can I get public data?

- [www.fastgenomics.org](http://www.fastgenomics.org)
- <https://www.ncbi.nlm.nih.gov/geo/>
- <https://support.10xgenomics.com/single-cell-gene-expression/datasets/>
- [https://portals.broadinstitute.org/single\\_cell](https://portals.broadinstitute.org/single_cell)
- <https://preview.data.humancellatlas.org/>
- <http://jinglebells.bgu.ac.il/>
- ....



Reducing barriers and accelerating single-cell research

