



# Advanced Topics in Single Cell Omics

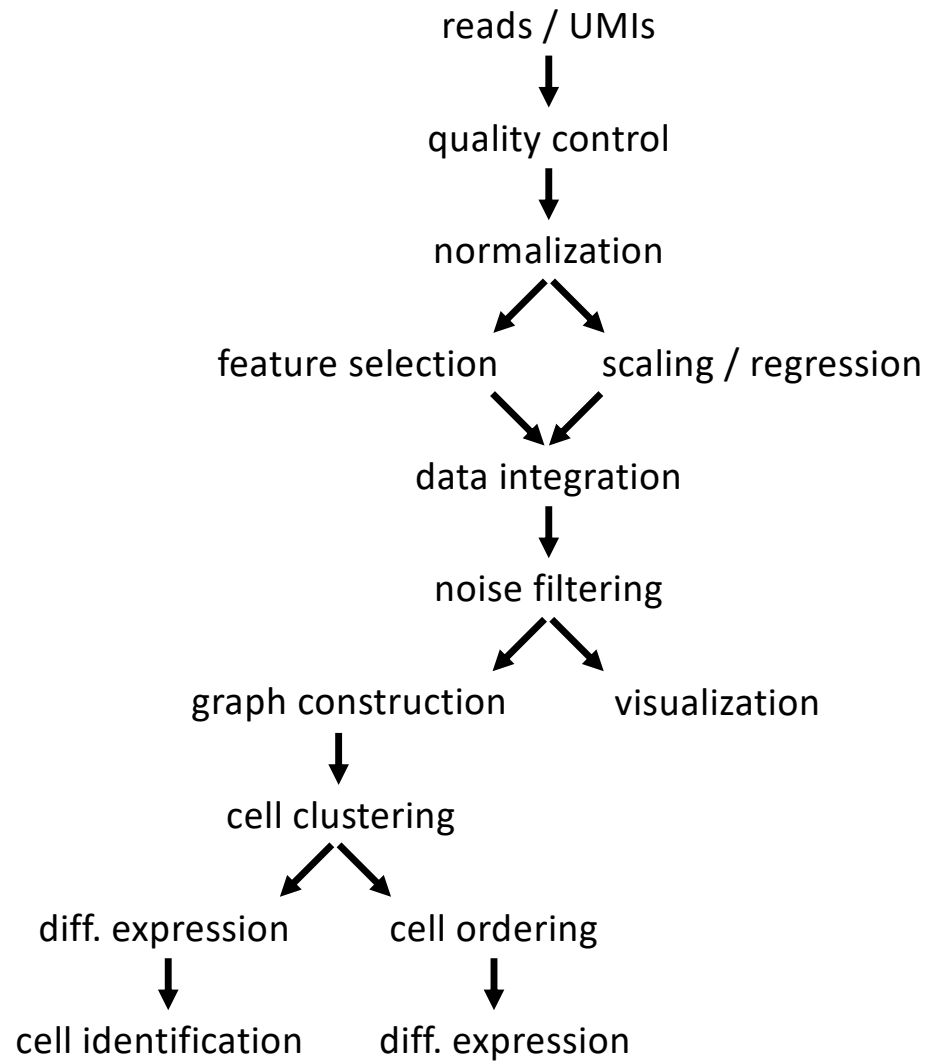
**Paulo Czarnewski**

Scientific Coordinator for the Human Developmental Cell Atlas (HDCA Sweden)

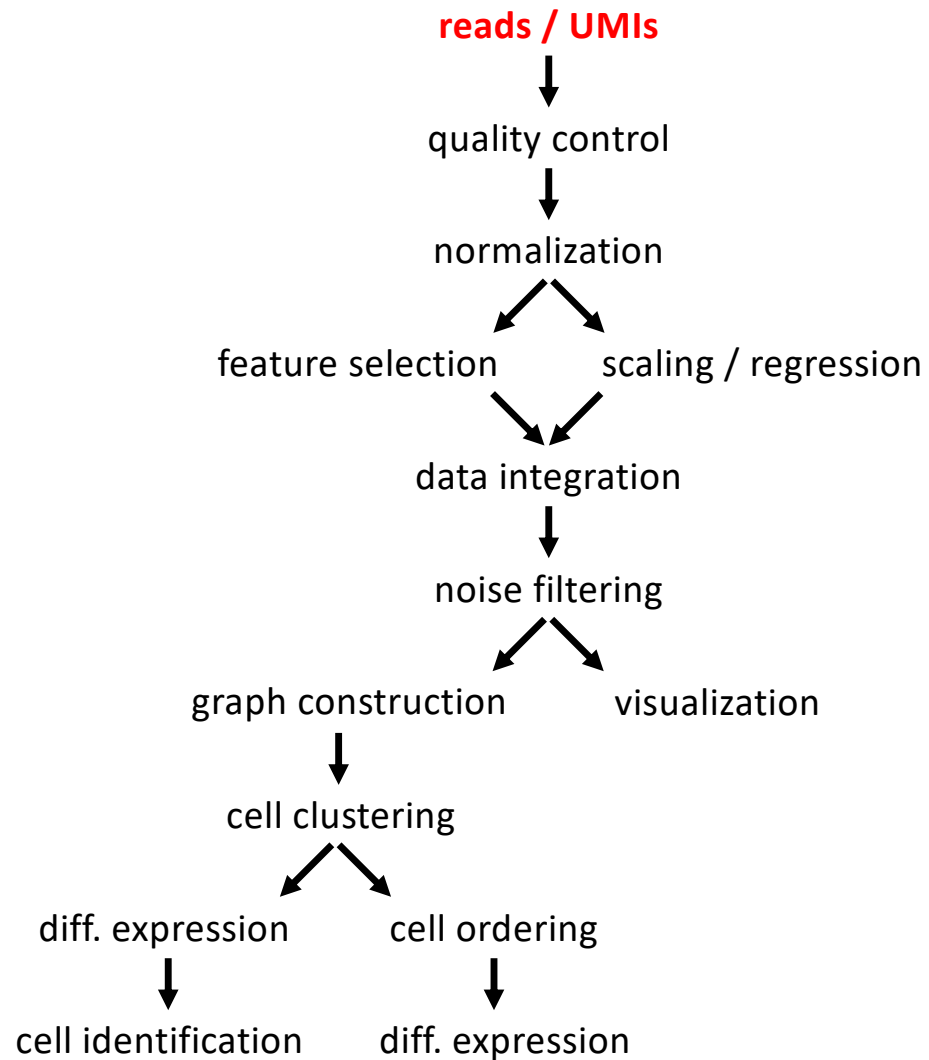
Senior Bioinformatician at the National Bioinformatics Infrastructure Sweden (NBIS)

2021-08-30

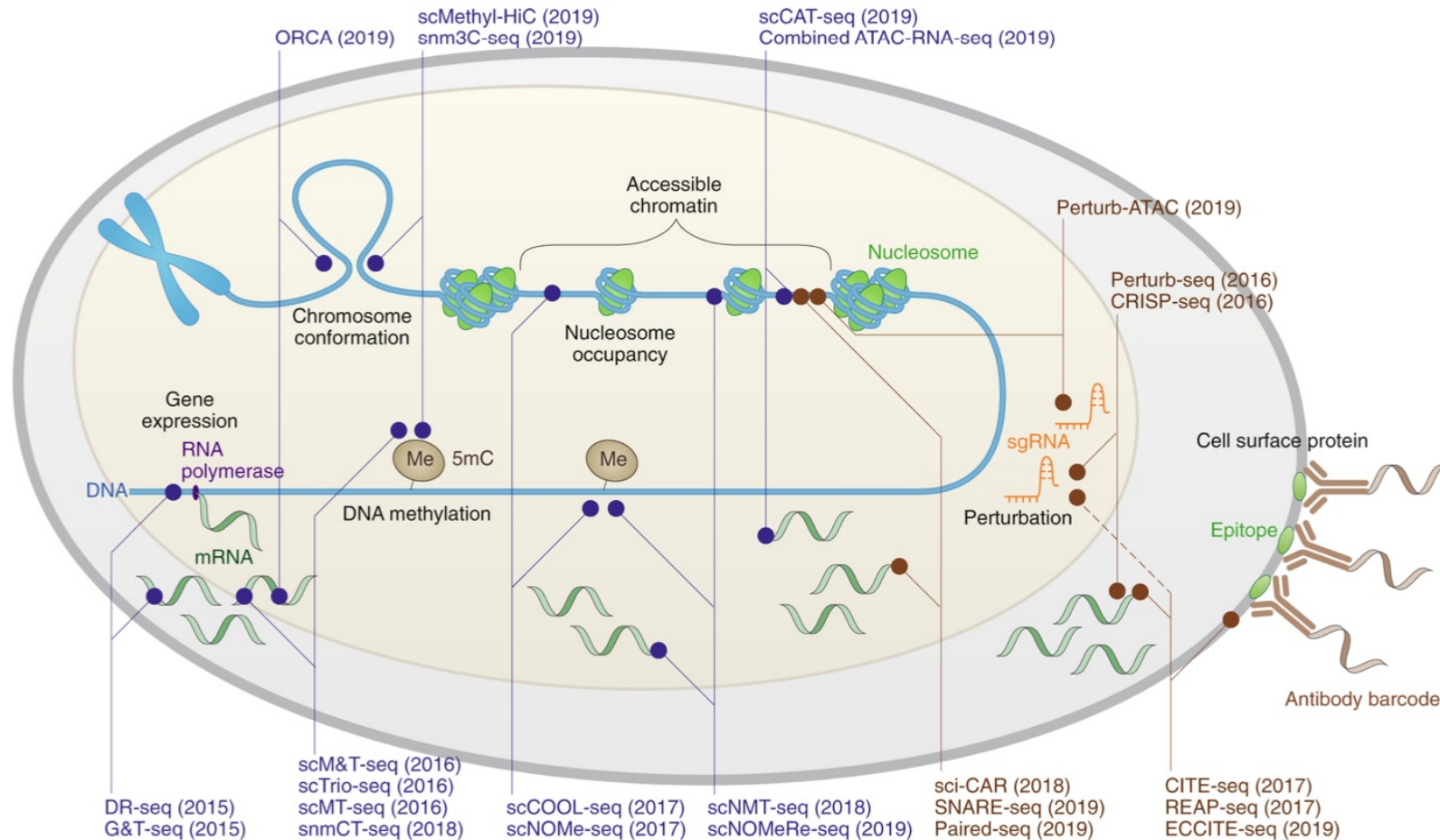
# scRNA-seq analysis workflow



# scRNA-seq analysis workflow



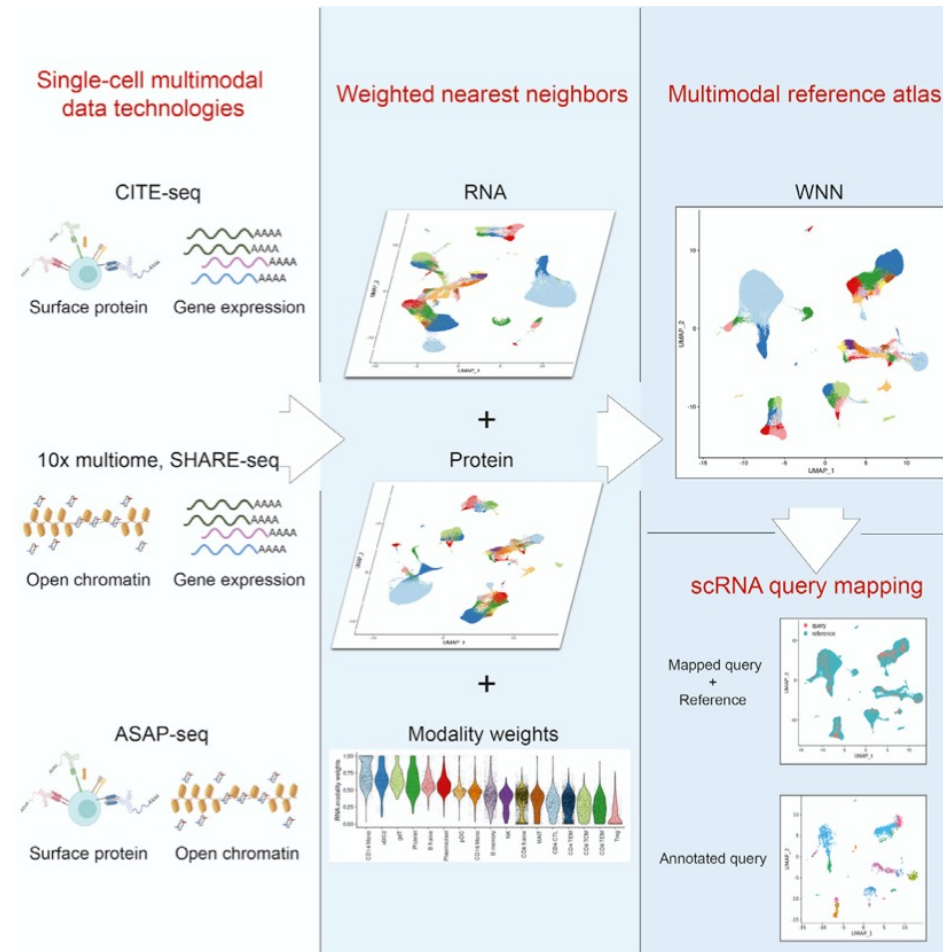
# scRNA-seq technologies



Zhu et al, Comment in Nature Methods, 2020

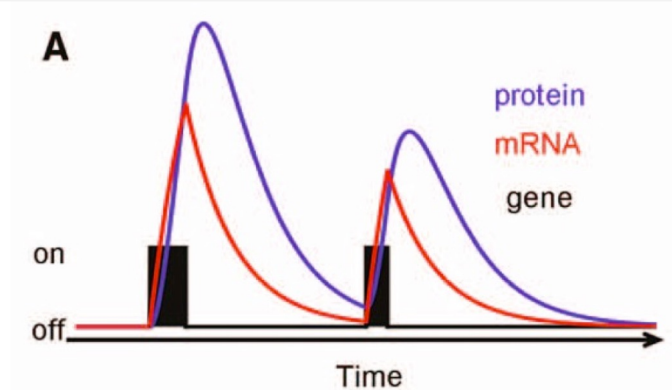


# scRNA-seq technologies

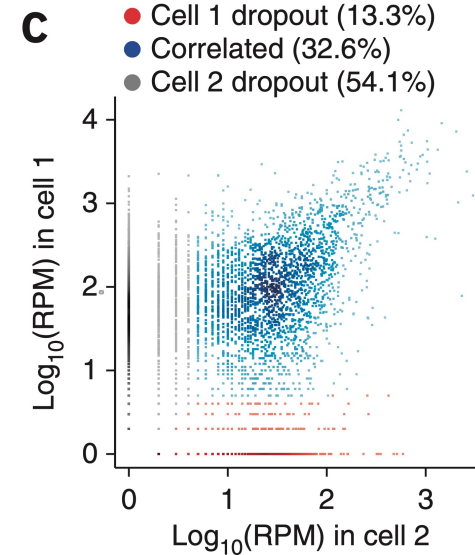


# scRNA-seq biases

- Amplification bias
- Drop-out rates
- Transcriptional bursting
- Background noise
- Bias due to cell-cycle, cell size and other factors
- Often clear batch effects
- Dissociation protocols may introduce transcriptional artifacts
- Ambient RNA

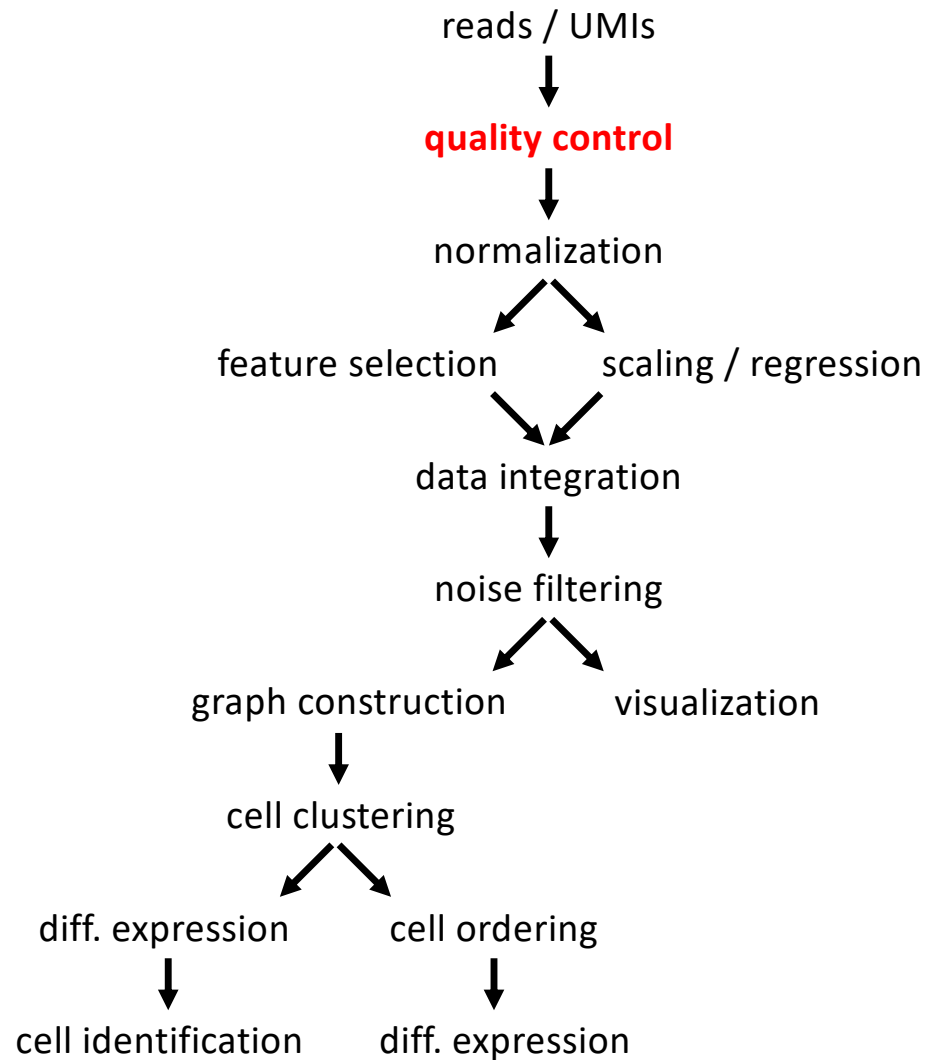


Suter et al. *Science* 2011



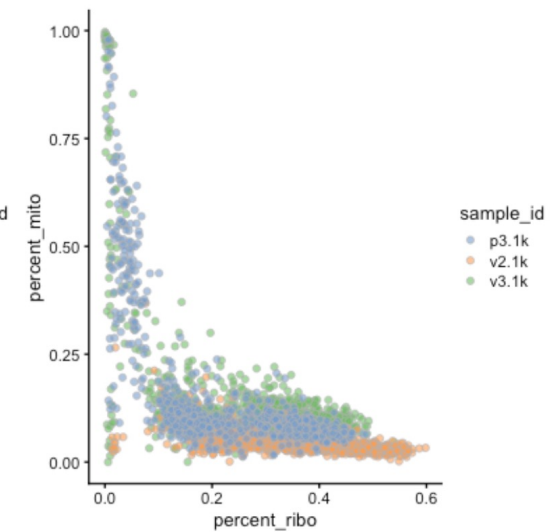
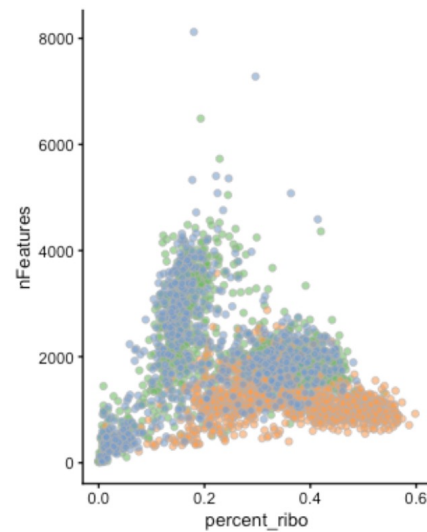
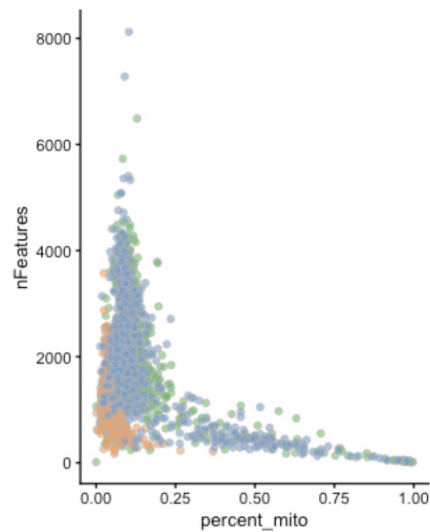
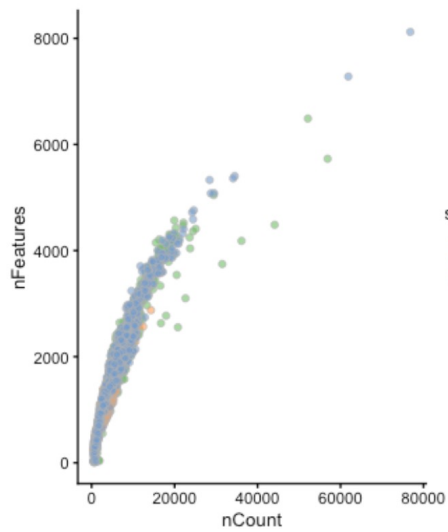
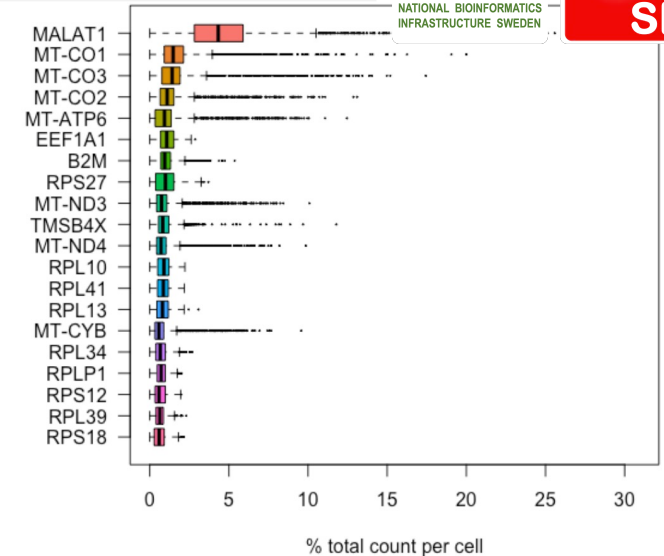
Karchenko et al. *Nature Methods* 2014

# scRNA-seq analysis workflow

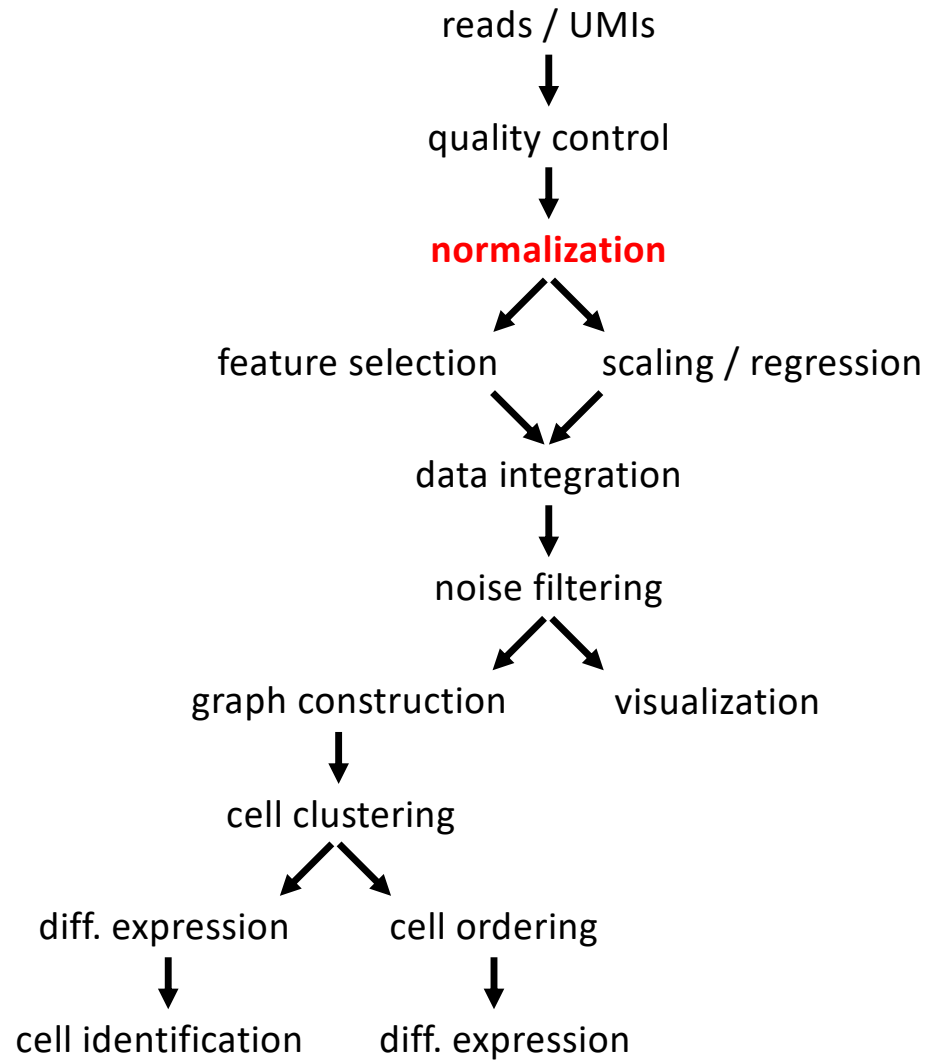


# scRNA-seq quality control

- Mapping statistics (% uniquely mapping)
- Cell cycle biases
- 3' bias – for full length methods like SS2
- mRNA-mapping read percentage
- Number of UMIs/read counts
- Number of detected genes
- Spike-in detection
- Mitochondrial percentage
- ribosomal percentage
- Protein-coding percentage



# scRNA-seq analysis workflow



# scRNA-seq normalization

## Count normalization (UMI and read counts)

for uneven sequencing depth

- CPM -  $\log[\text{CP}10\text{K}+1]$

## Gene length normalization (read counts)

for differences in gene detection due to gene length

- TPM (closer to UMI counts)
- FPKM

## Drop-out rate normalization (UMI and read counts)

for differences in RNA content / drop-out rates

- Deconvolution/Scran(Pooling-Across-Cells)
- SCnorm(Expression-DepthRelation)
- SCTransform
- Census
- Linnorm
- ZINB-WaVE
- ...

bulk

$$CPM = \log\left(\frac{\text{counts}}{\text{library}_{size}} \cdot 10^6 + 1\right)$$

single-cell

$$\log[TP10K + 1] = \log\left(\frac{\text{counts}}{\text{library}_{size}} \cdot 10^4 + 1\right)$$

Most common for UMI data / fast

$$FPKM = \log\left(\frac{\text{counts}}{\text{library}_{size} \cdot \text{transcript}_{length}} \cdot 10^4 + 1\right)$$

$$TPM = \log\left(\frac{\text{counts}}{\text{transcript}_{length}} \cdot \frac{10^4}{\sum \frac{\text{counts}}{\text{transcript}_{length}}} + 1\right)$$

# scRNA-seq analysis workflow



## Count normalization (UMI and read counts)

for uneven sequencing depth

- CPM -  $\log[CP10K+1]$

## Gene length normalization (read counts)

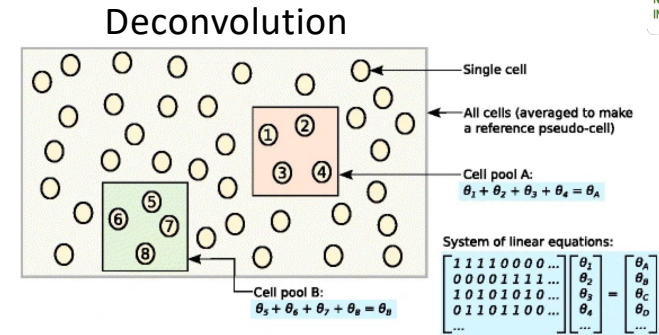
for differences in gene detection due to gene length

- TPM (closer to UMI counts)
- FPKM

## Drop-out rate normalization (UMI and read counts)

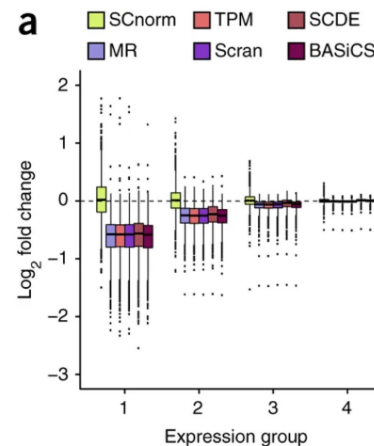
for differences in RNA content / drop-out rates

- Deconvolution/Scran (Pooling-Across-Cells)
- SCnorm (Expression-DepthRelation)
- SCTransform
- Census
- Linnorm
- ZINB-WaVE
- ...



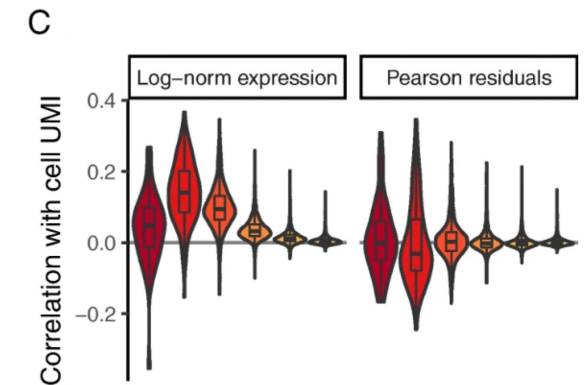
Lun et al. Genome Biol. 2016

## SCnorm



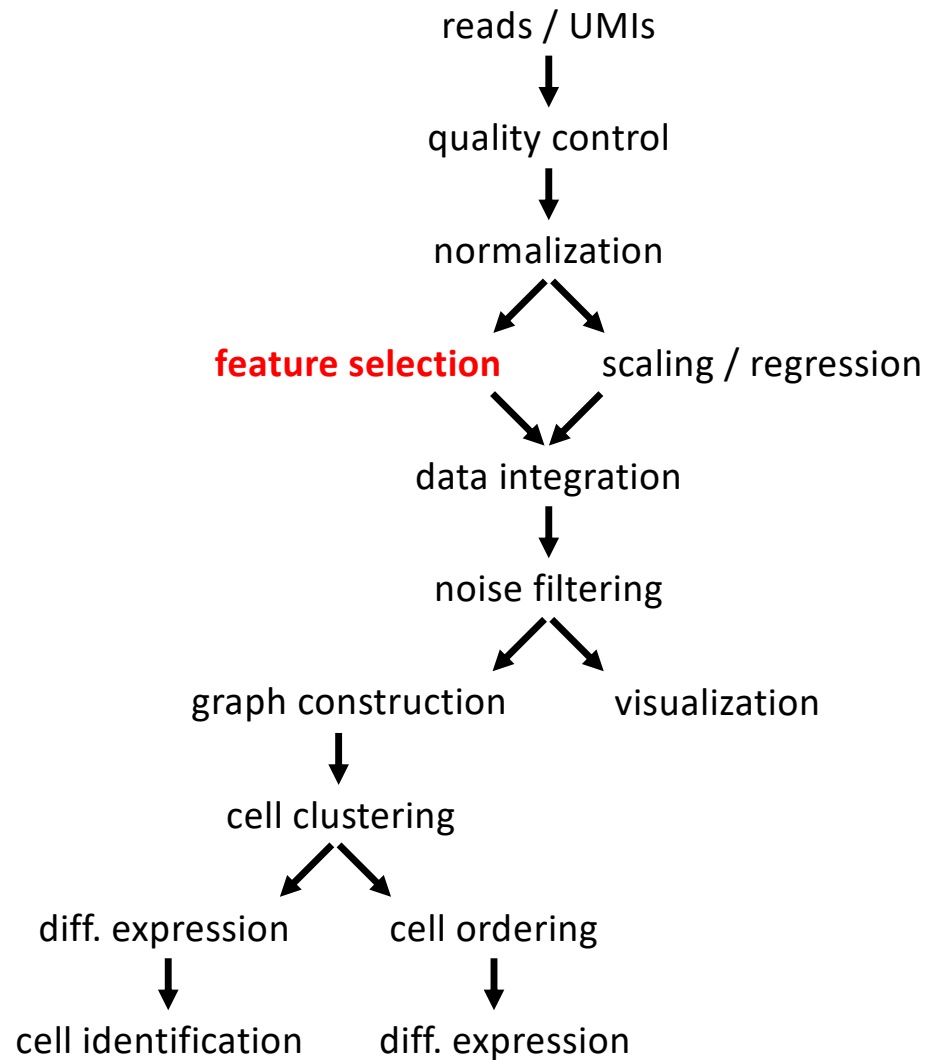
Bacher et al. Nature Methods 2017

## SCTransform



Hafmeister & Satija Genome Biology 2019

# scRNA-seq analysis workflow

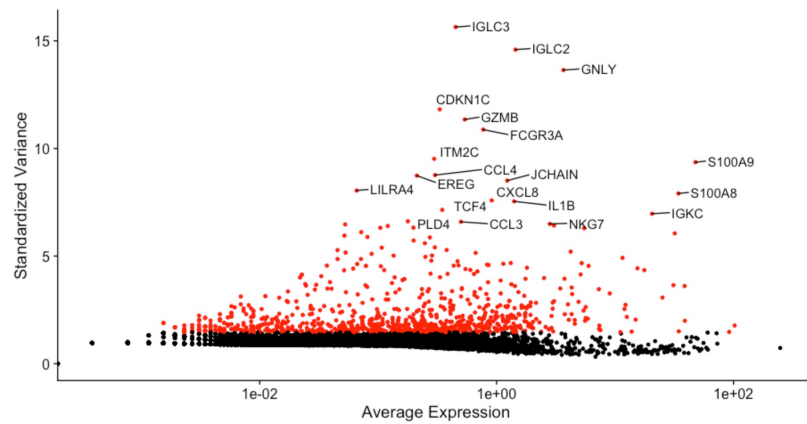




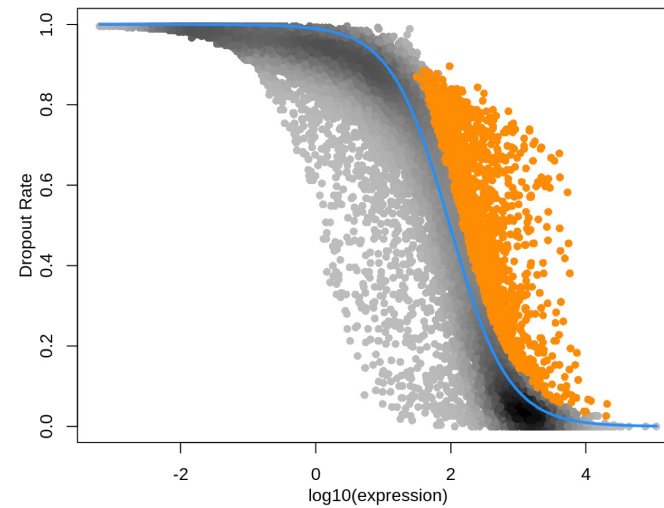
# scRNA-seq feature selection

Not all genes are important to define you cell types

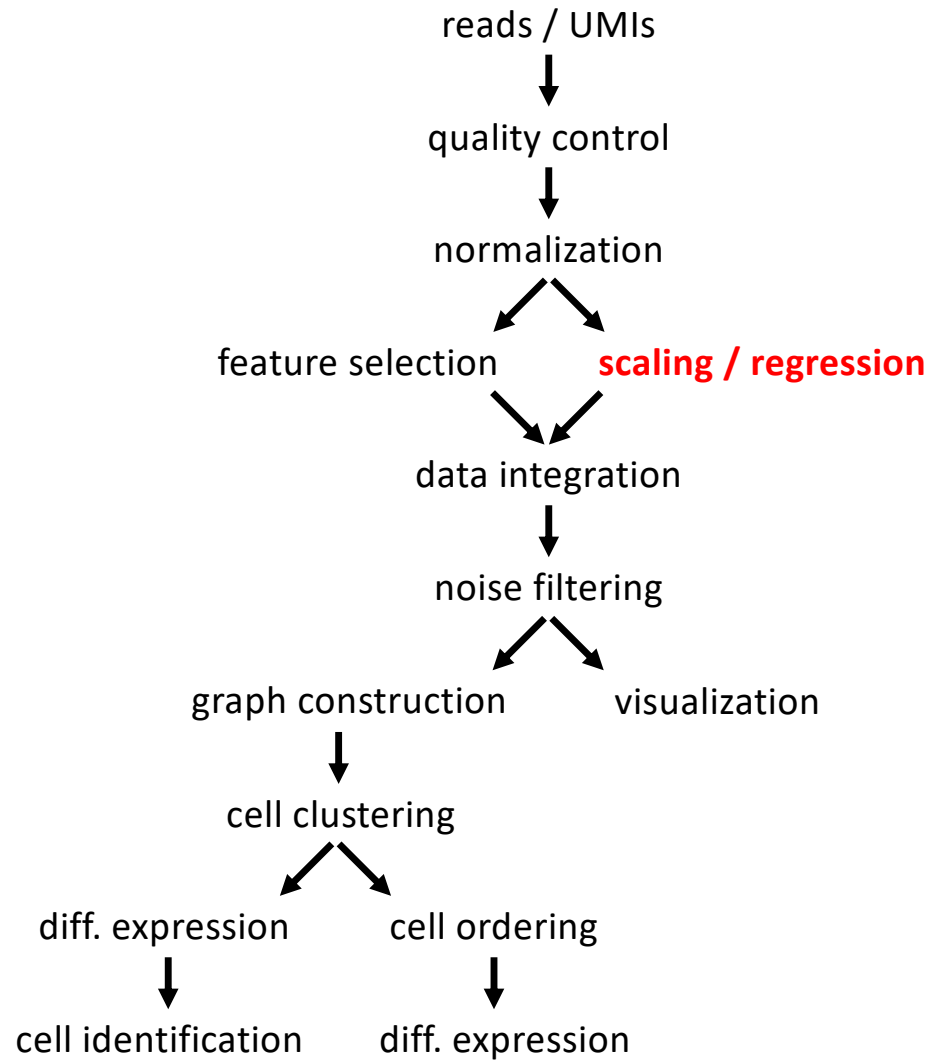
$$HVG = \frac{\text{variance}}{\log(\text{meanExpression})}$$



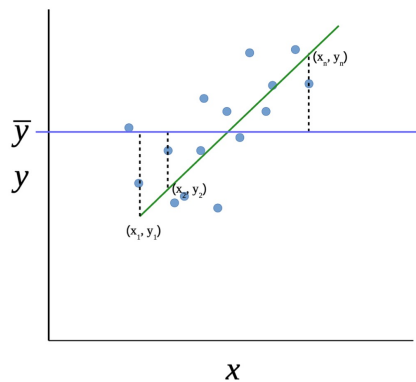
$$HVG = \frac{\log(\text{meanExpression})}{\text{dropout}_{rate}}$$



# scRNA-seq analysis workflow



# scRNA-seq scaling and regression of biases



Any source of variation that you do not expect to give separation of the cell types can be regressed out.

- Fit a line to the gene expression vs variable of interest
- Calculate residuals
- Remove variance explained by the variable of interest by taking the residuals.
- Linear / Negative Binomial / Poisson distributions



**fast**

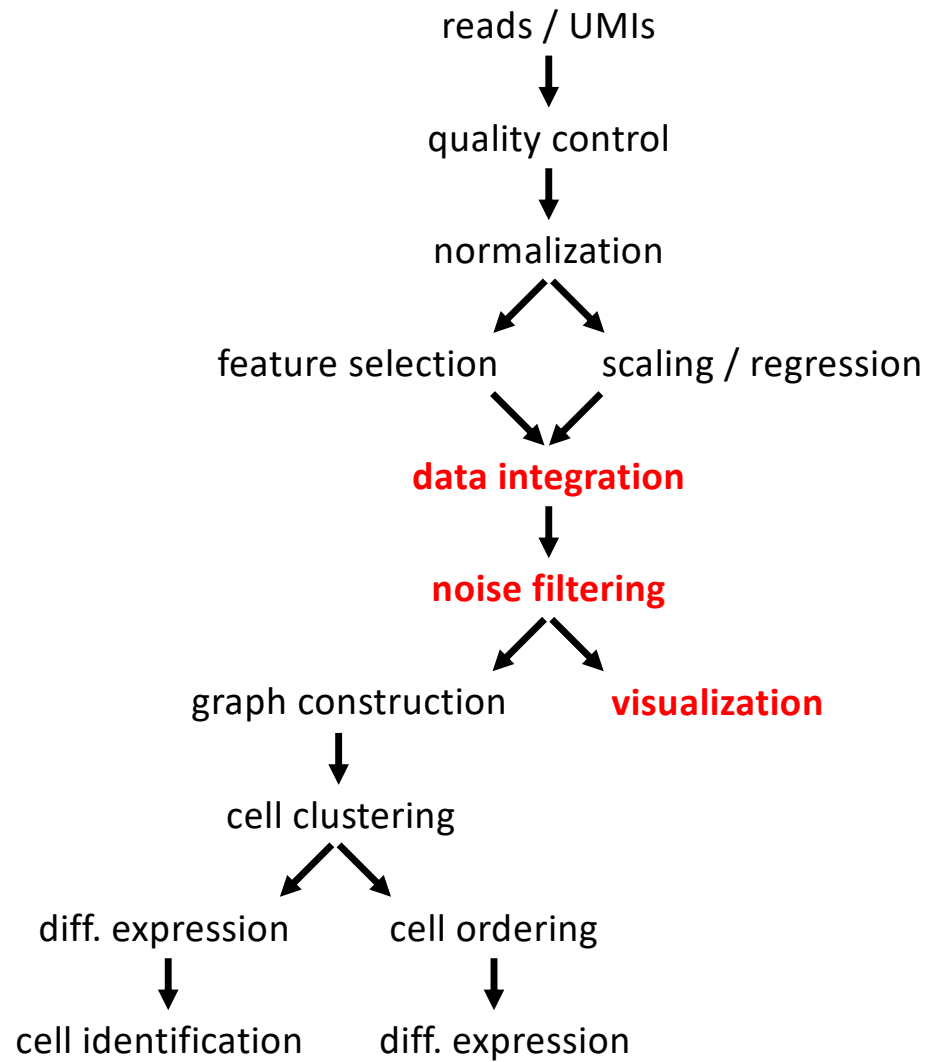
logNormalized counts follows a  
log-linear distribution



**slower (but more accurate)**

Regressing counts directly is better with count-based  
distributions

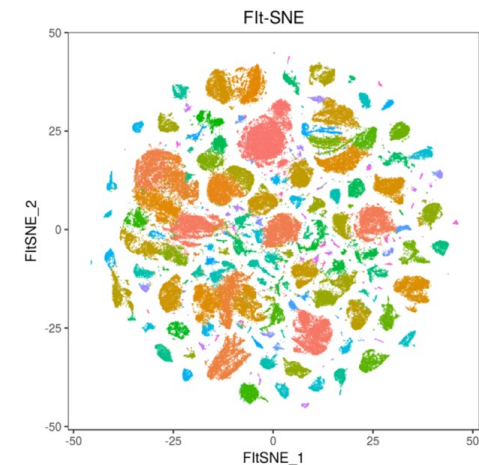
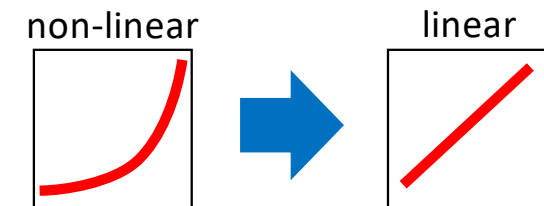
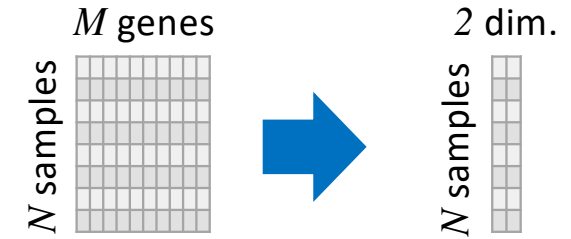
# scRNA-seq analysis workflow



# scRNA-seq dimensionality reduction

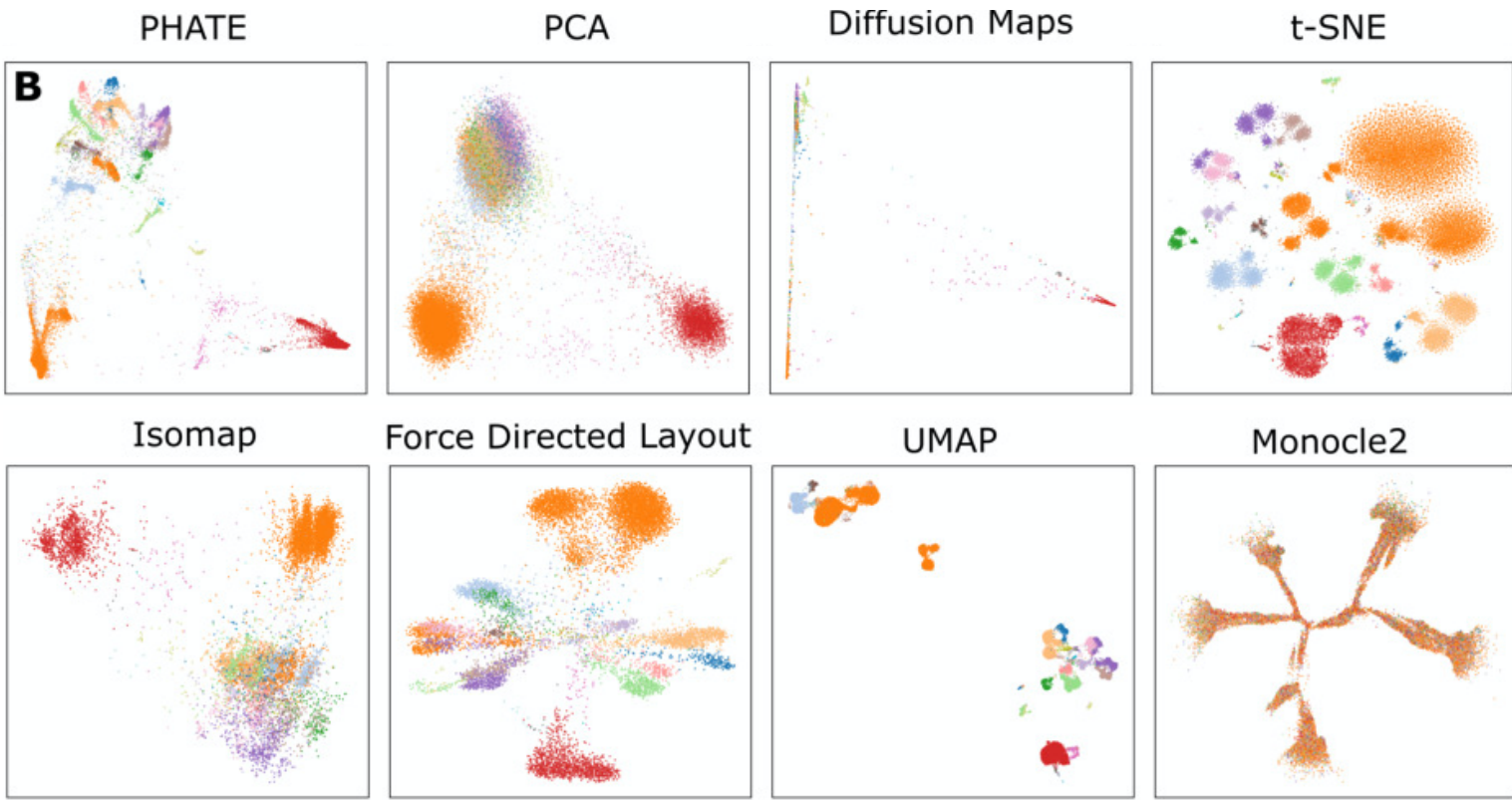
- Simplify complexity, so it becomes easier to work with.  
Reduce number of features (genes)  
In some: Transform non-linear relationships to linear
- “Remove” redundancies in the data
- Identify the most relevant information (find and filter noise)
- Reduce computational time for downstream procedures
- Facilitate clustering, since some algorithms struggle with too many dimensions
- Data visualization

... and more ...



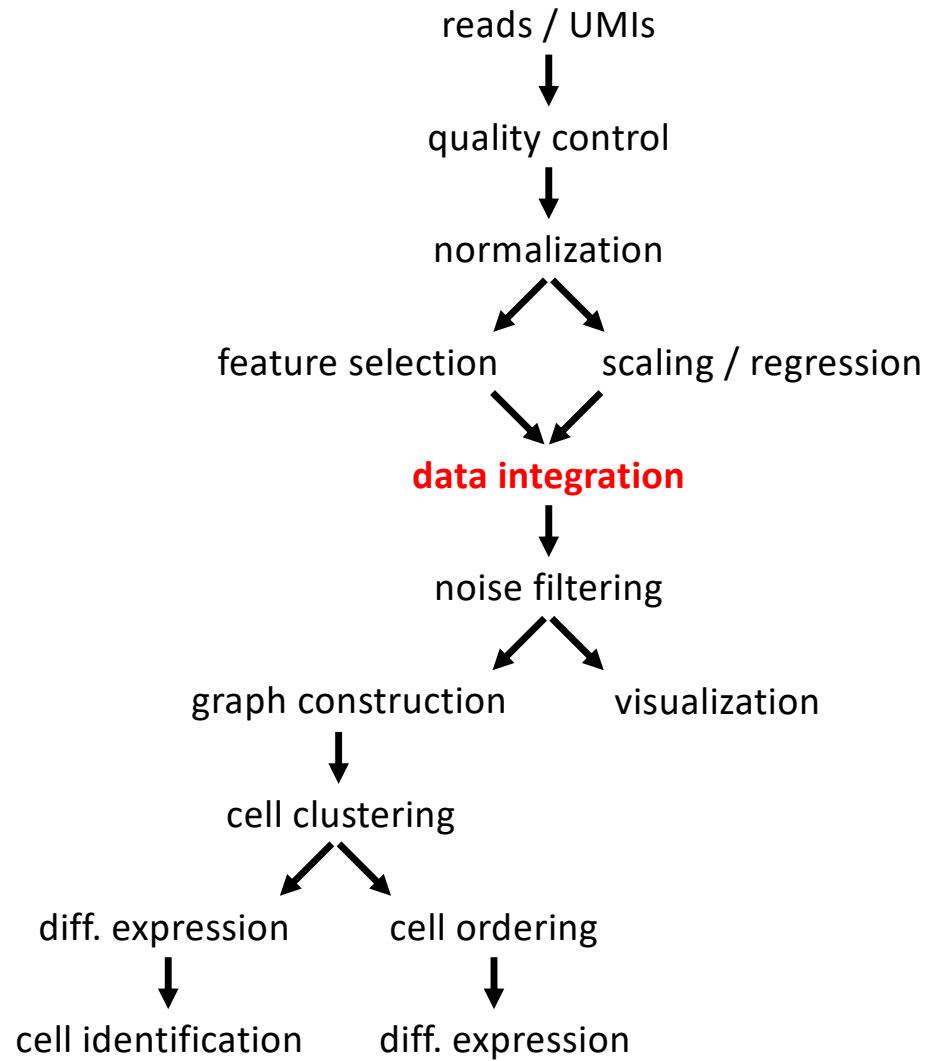
# Some dimensionality reduction methods

Shekhar et al. (2016)  
(n=6,174)



17	BC1A	BC5D
18	BC1B	BC6
19	BC2	BC7
21	BC3A	BC8/9_1
23	BC3B	BC8/9_2
24	BC4	Cone PR
25	BC5A	Muller Glia
26	BC5B	Rod BC
Amacrine_1	BC5C	Rod PR
Amacrine_2		

# scRNA-seq analysis workflow



# scRNA-seq data integration

We wish to obtain corrected data where the following goals are met:

## Goal:

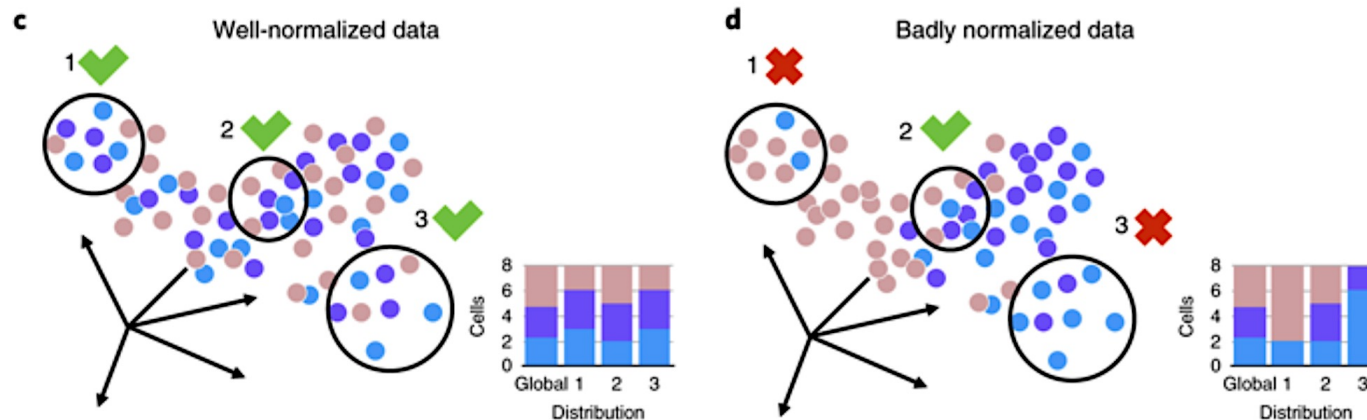
1. The batch-originating variance is erased
2. Meaningful heterogeneity is preserved
3. No artefactual variance is introduced

## What it practically means:

Similar cell types are intermixed across batches

We are not mixing distinct cell types (across or within batches)

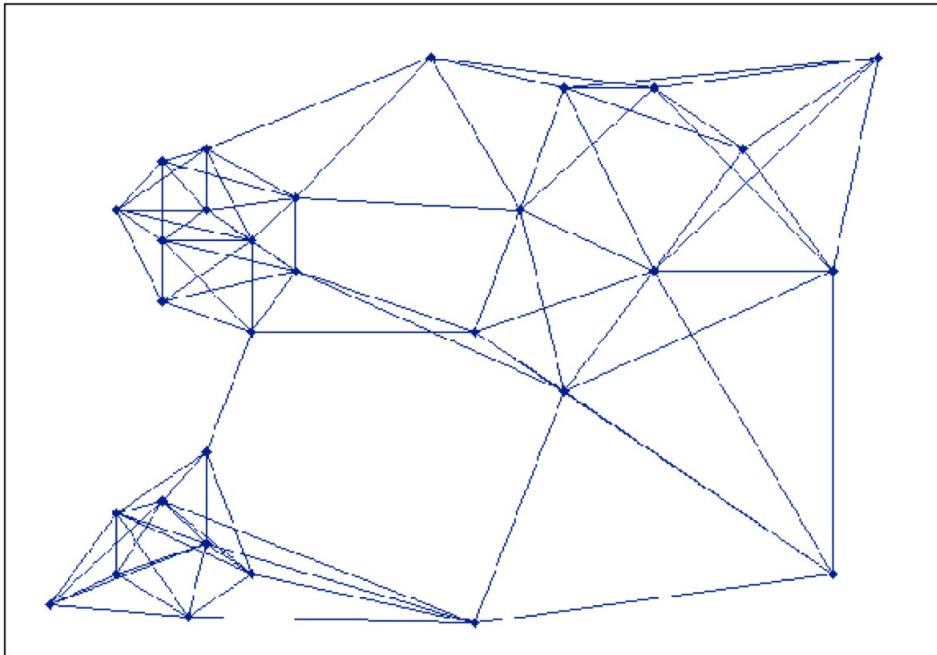
We do not separate similar cells within batches



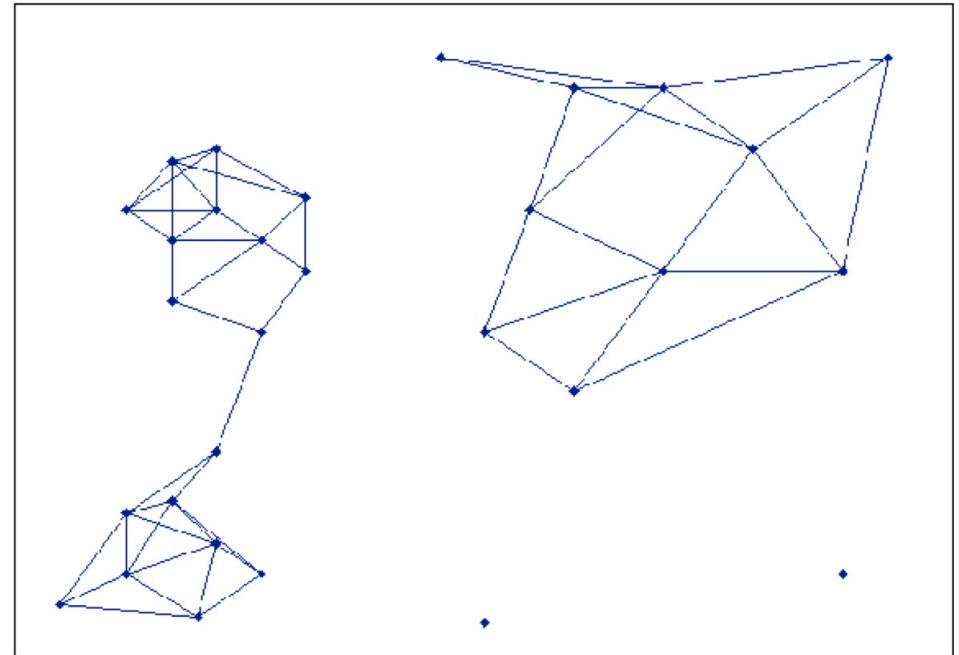


# scRNA-seq graph construction

The ***k*-Nearest Neighbor (*k*NN)** graph is a graph in which two vertices  $p$  and  $q$  are connected by an edge, if the distance between  $p$  and  $q$  is among the  $k$ -th smallest distances from  $p$  to other objects from  $P$ .



The **Shared Nearest Neighbor (SNN)** graph has weights that defines proximity, or similarity between two edges in terms of the number of neighbors (i.e., directly connected vertices) they have in common.

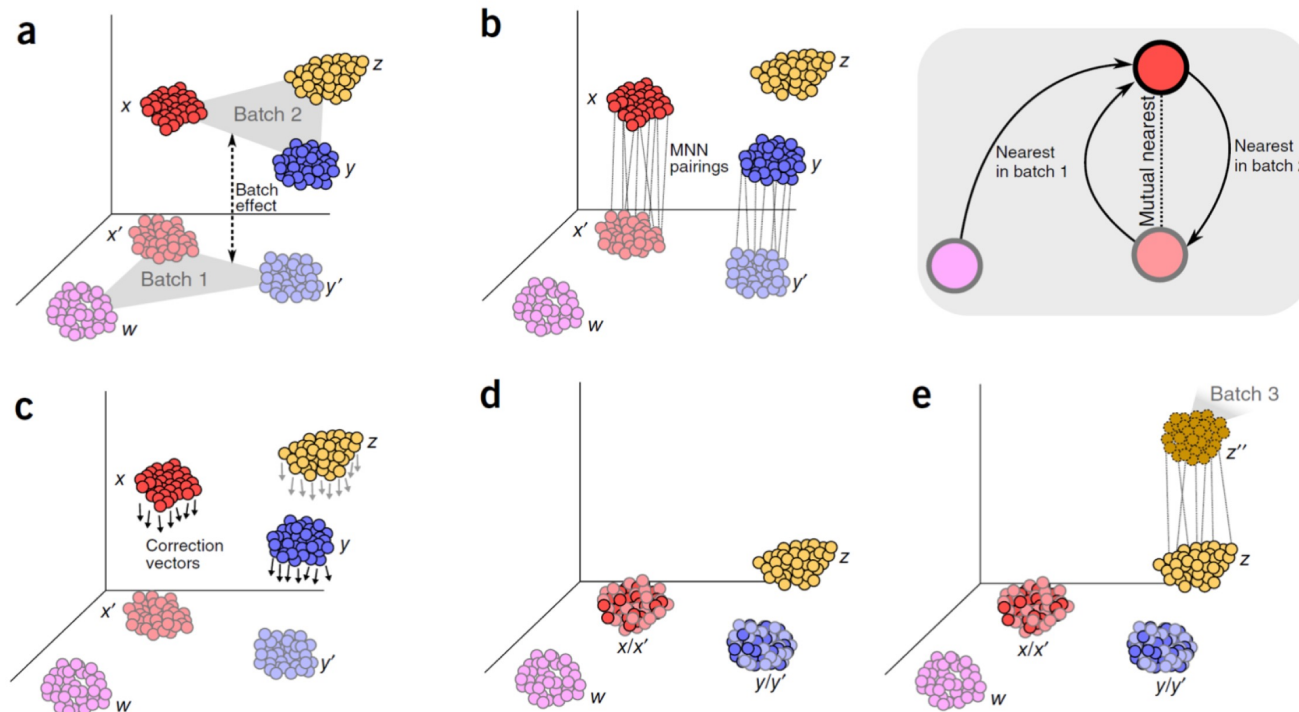


# scRNA-seq analysis workflow

Regression based bulk-RNAseq batch correction methods are slow and assume the batch is constant across cells

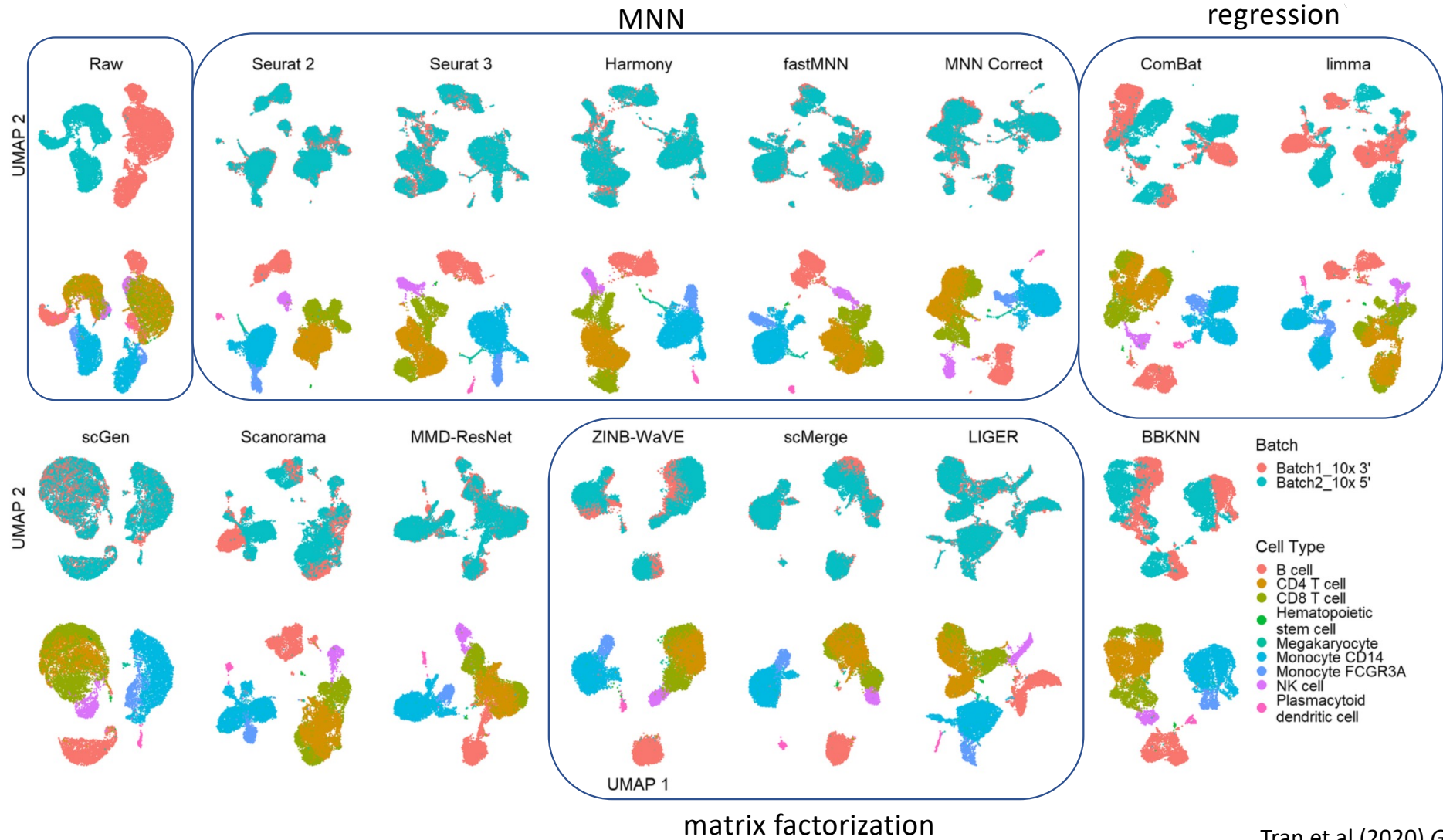
Modern data integration methods are based on the same principle:

- find MNN (mutual nearest neighbours) across datasets and correct each cell individually
- Done on a graph: much faster

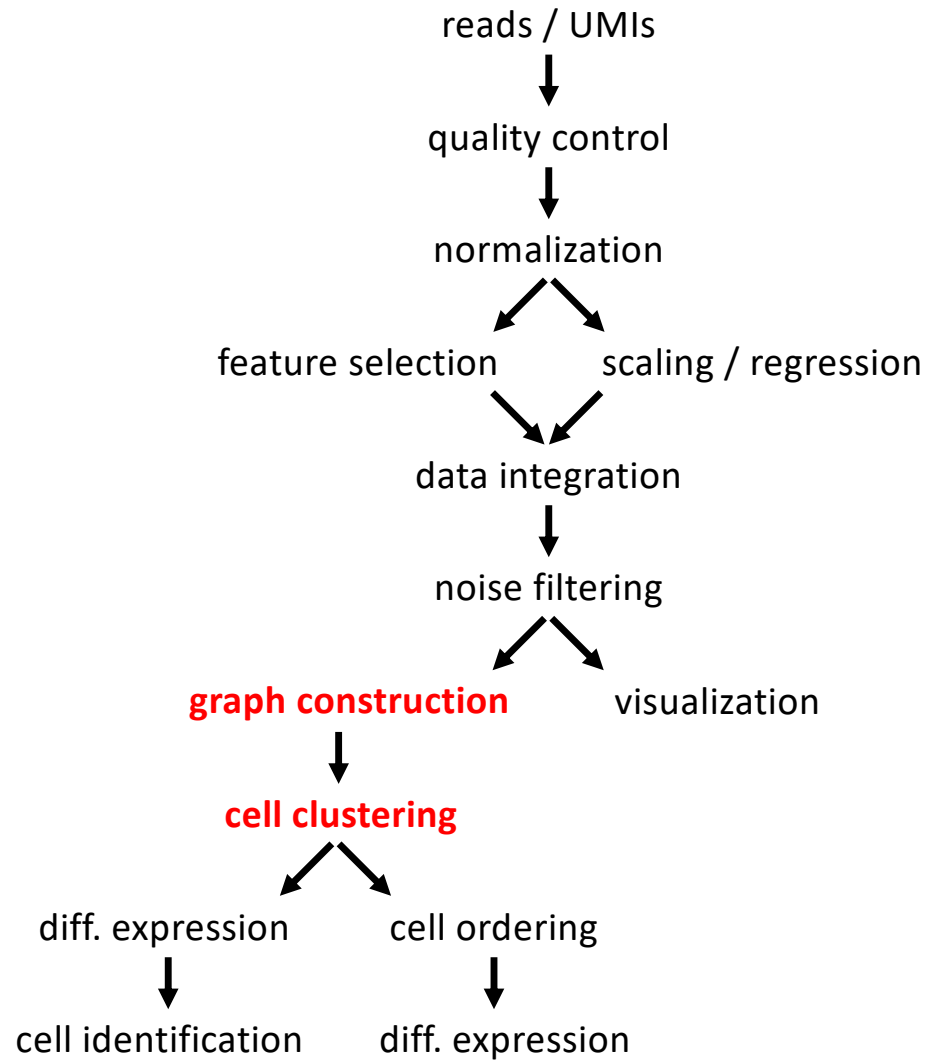


Haghverdi et al (2017) Nat Biotechnology

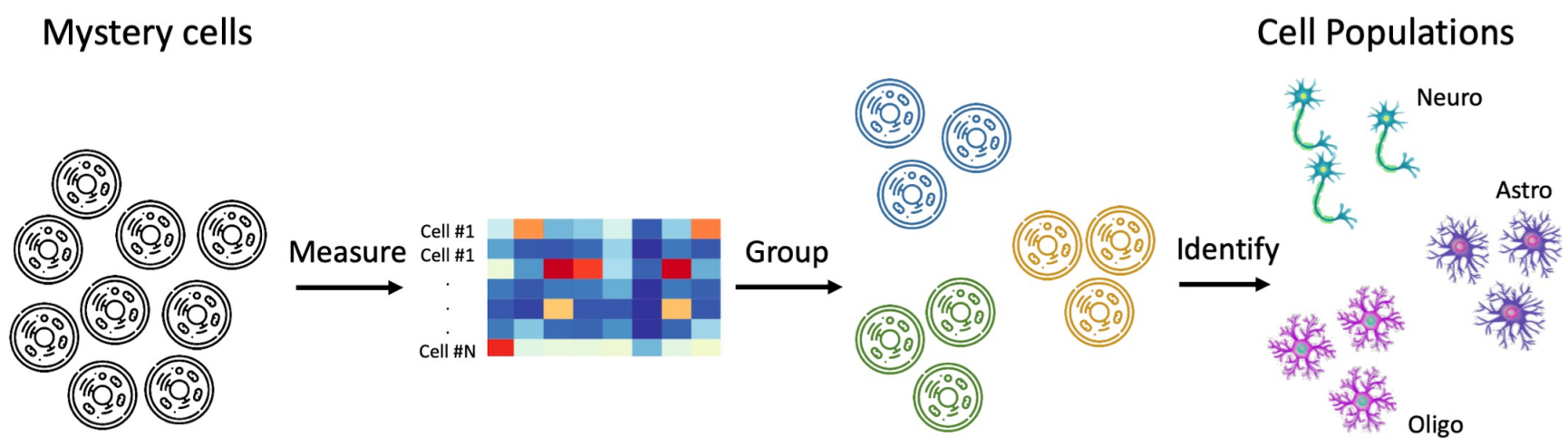
# scRNA-seq analysis workflow



# scRNA-seq analysis workflow

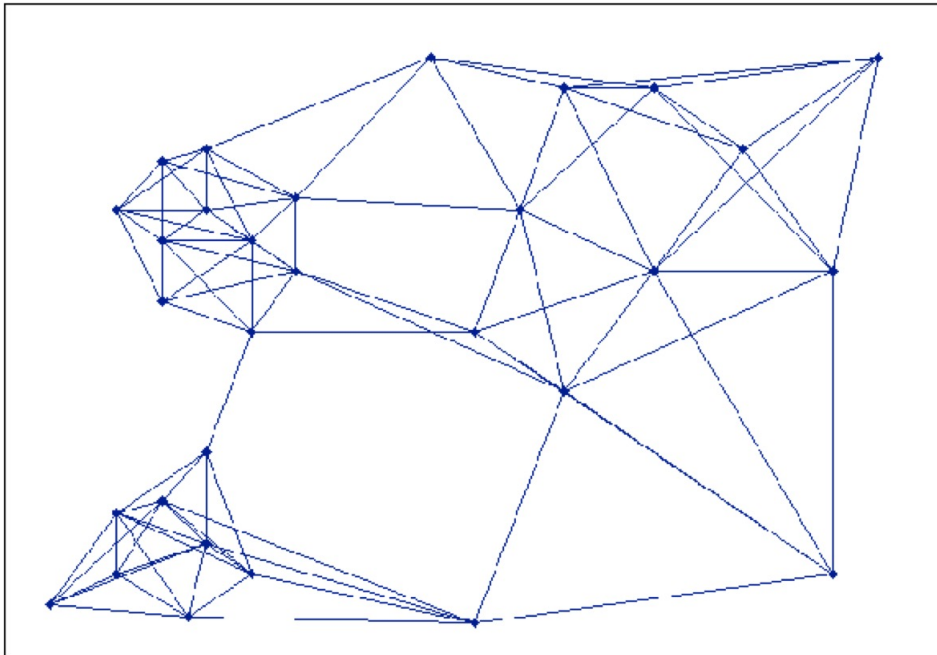


# scRNA-seq clustering

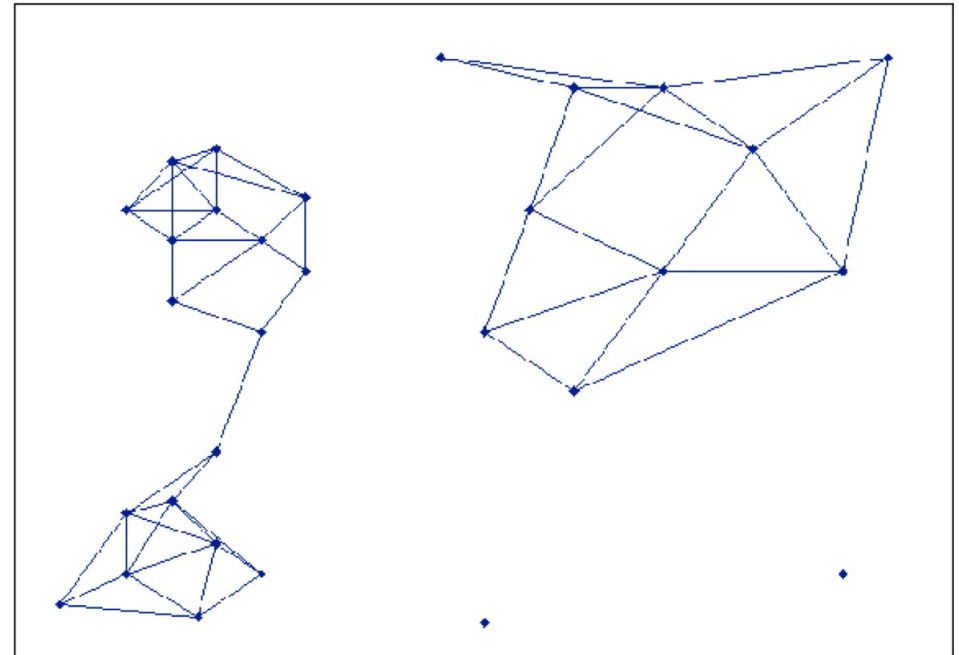


# scRNA-seq graph construction

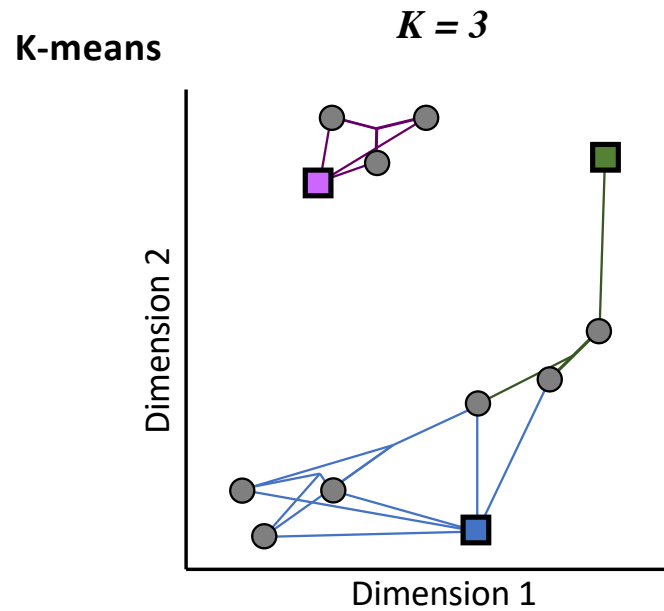
The ***k*-Nearest Neighbor (*k*NN)** graph is a graph in which two vertices  $p$  and  $q$  are connected by an edge, if the distance between  $p$  and  $q$  is among the  $k$ -th smallest distances from  $p$  to other objects from  $P$ .



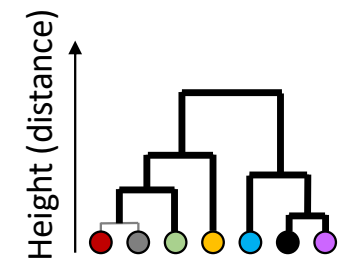
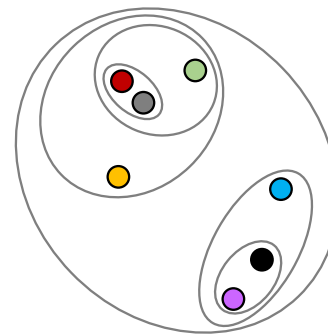
The **Shared Nearest Neighbor (SNN)** graph has weights that defines proximity, or similarity between two edges in terms of the number of neighbors (i.e., directly connected vertices) they have in common.



# scRNA-seq clustering



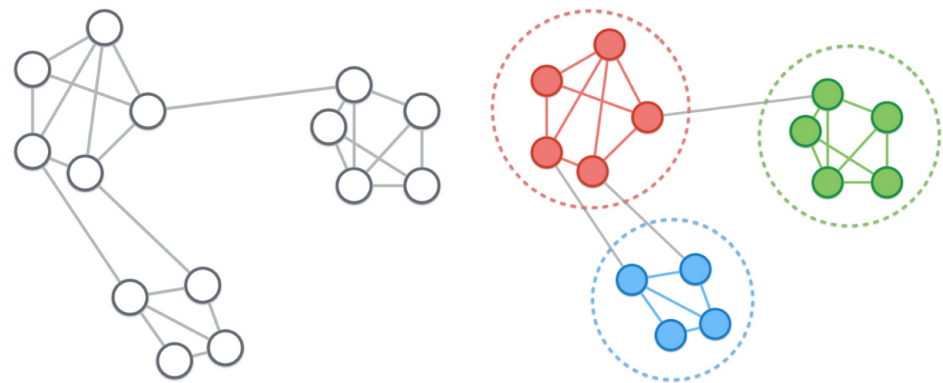
## Hierarchical Clustering



## GRAPH

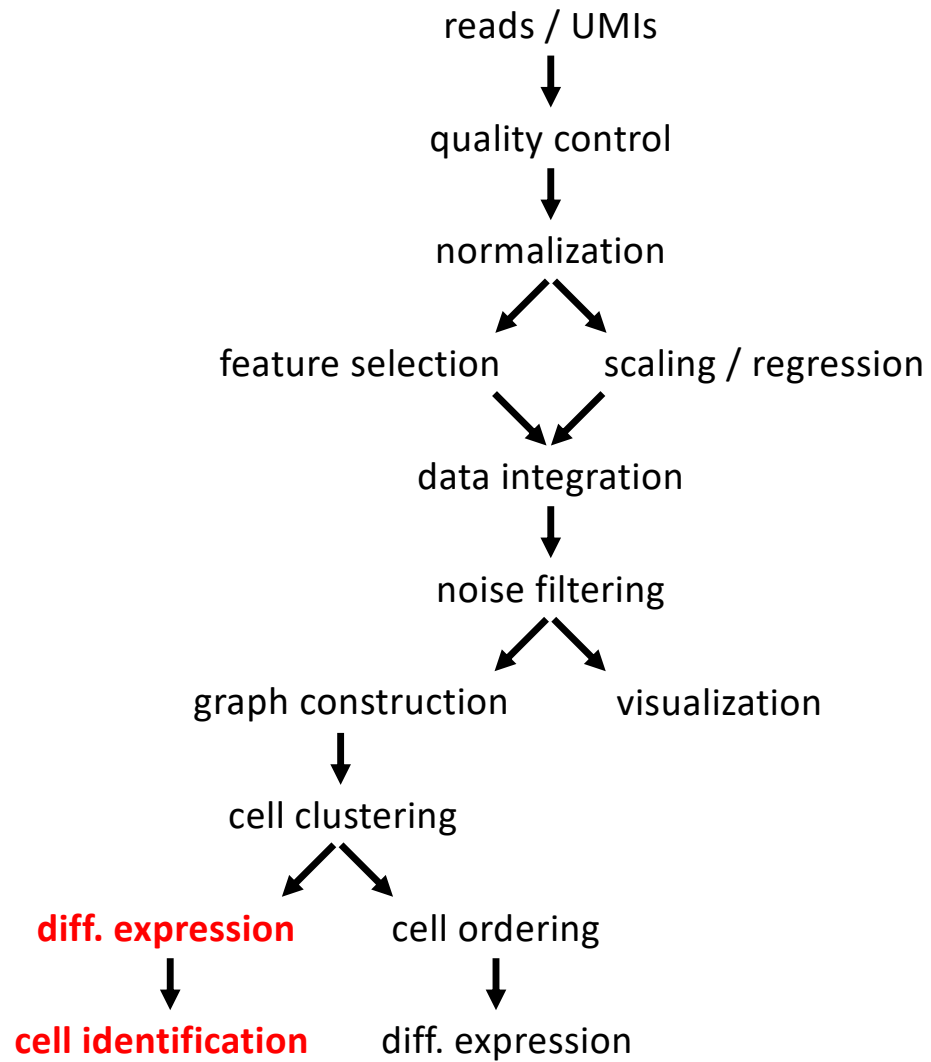
### Louvain / Leiden community detection

Communities, or clusters, are usually groups of vertices having higher probability of being connected to each other than to members of other groups.





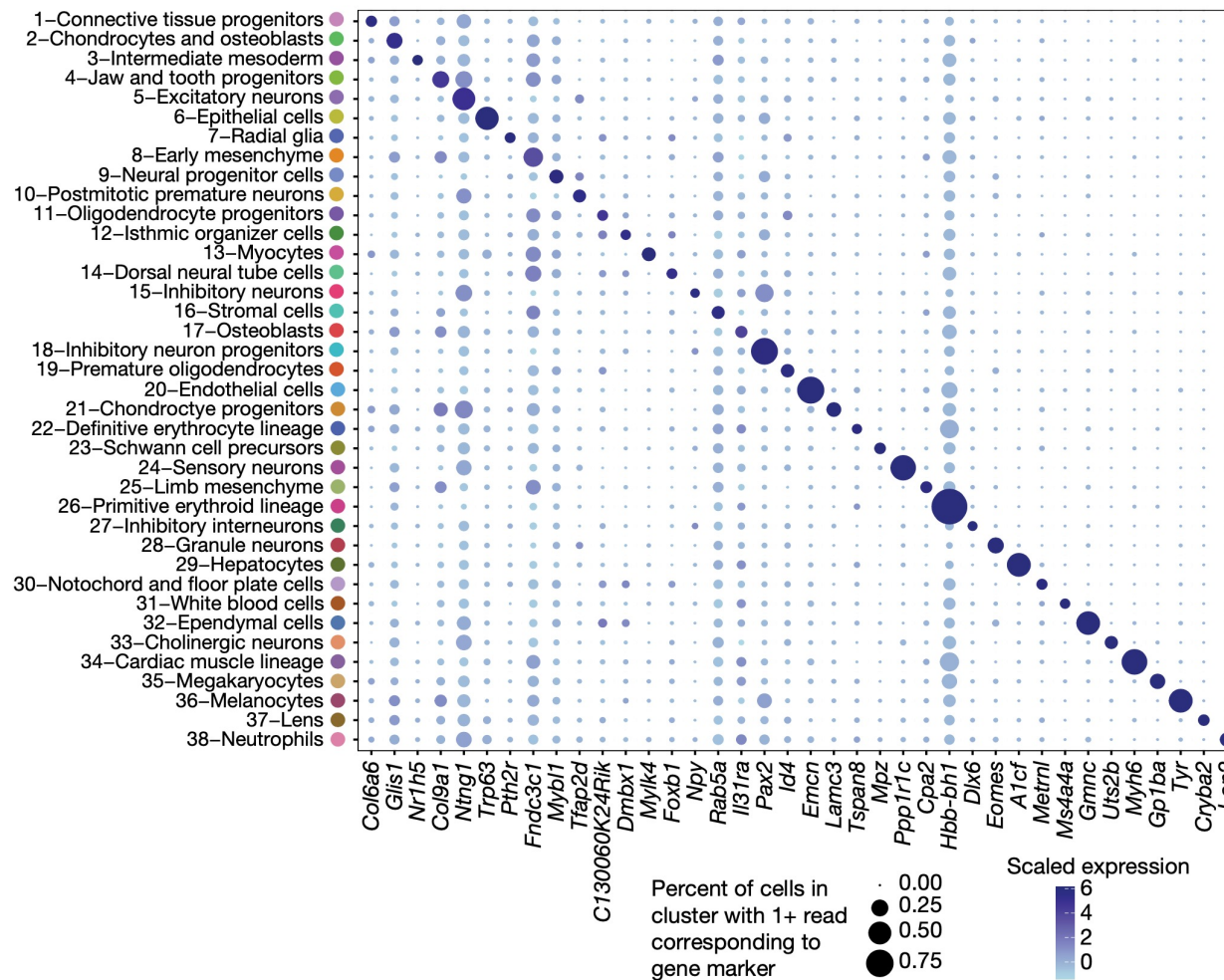
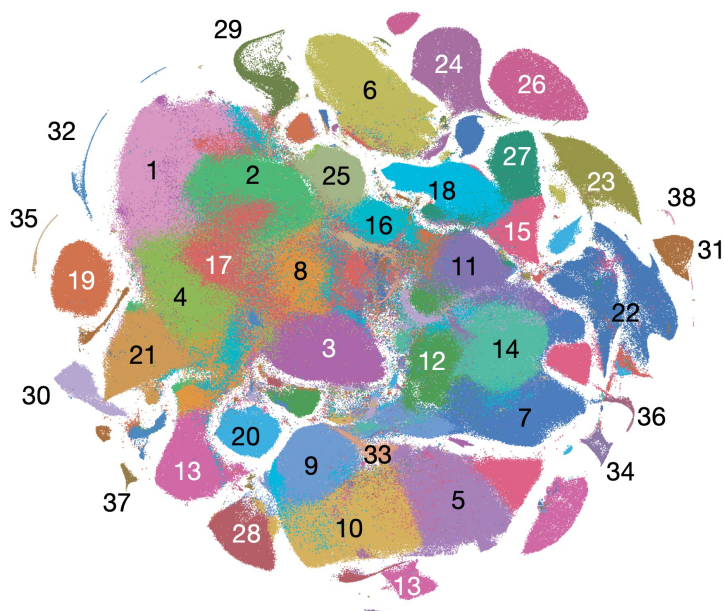
# scRNA-seq analysis workflow



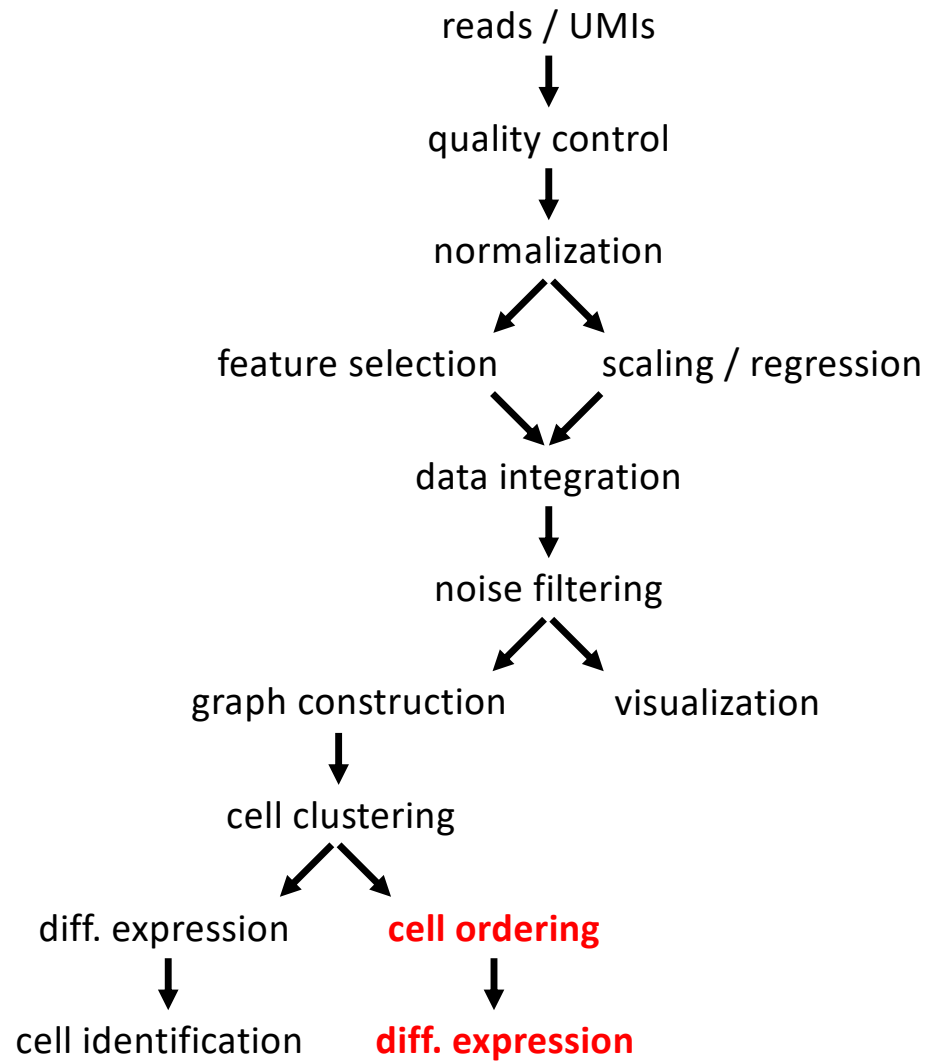




# scRNA-seq differential gene expression



# scRNA-seq analysis workflow



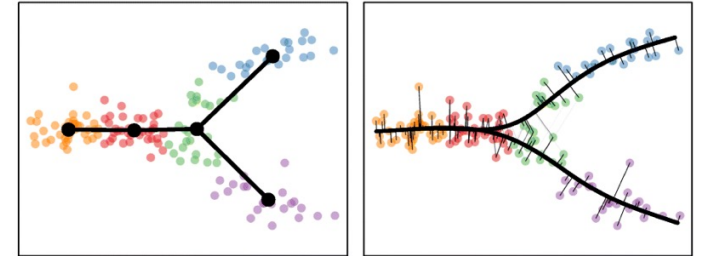
# scRNA-seq trajectory inference

Are you sure that you have a trajectory?

Do you have intermediate states?

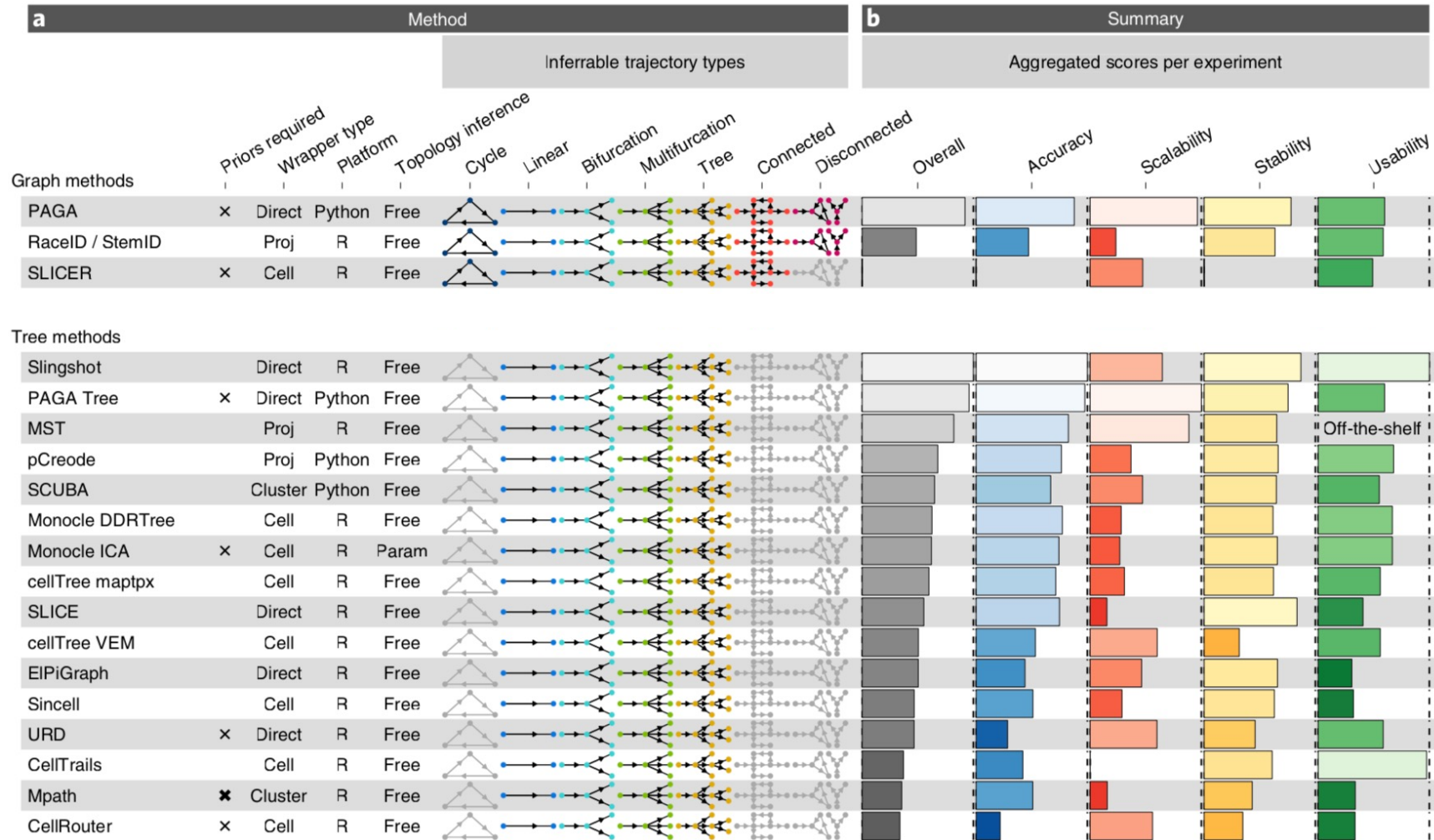
Do you believe that you have branching in your trajectory?

- ! Be aware, any dataset can be forced into a trajectory without any biological meaning!
- ! First make sure that gene set and dimensionality reduction captures what you expect.



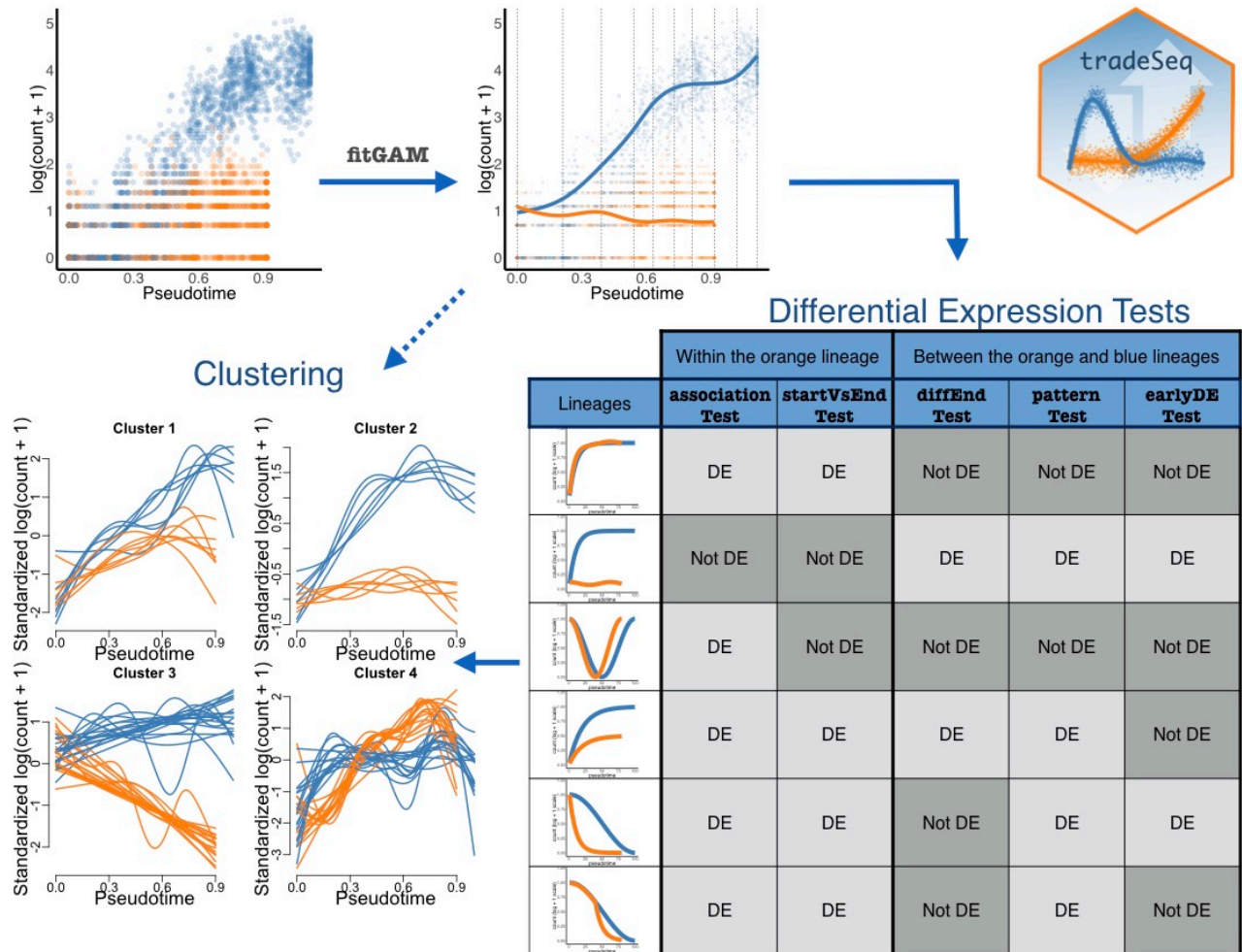
Street et al (2018) *BMC Genomics*

# scRNA-seq trajectory inference

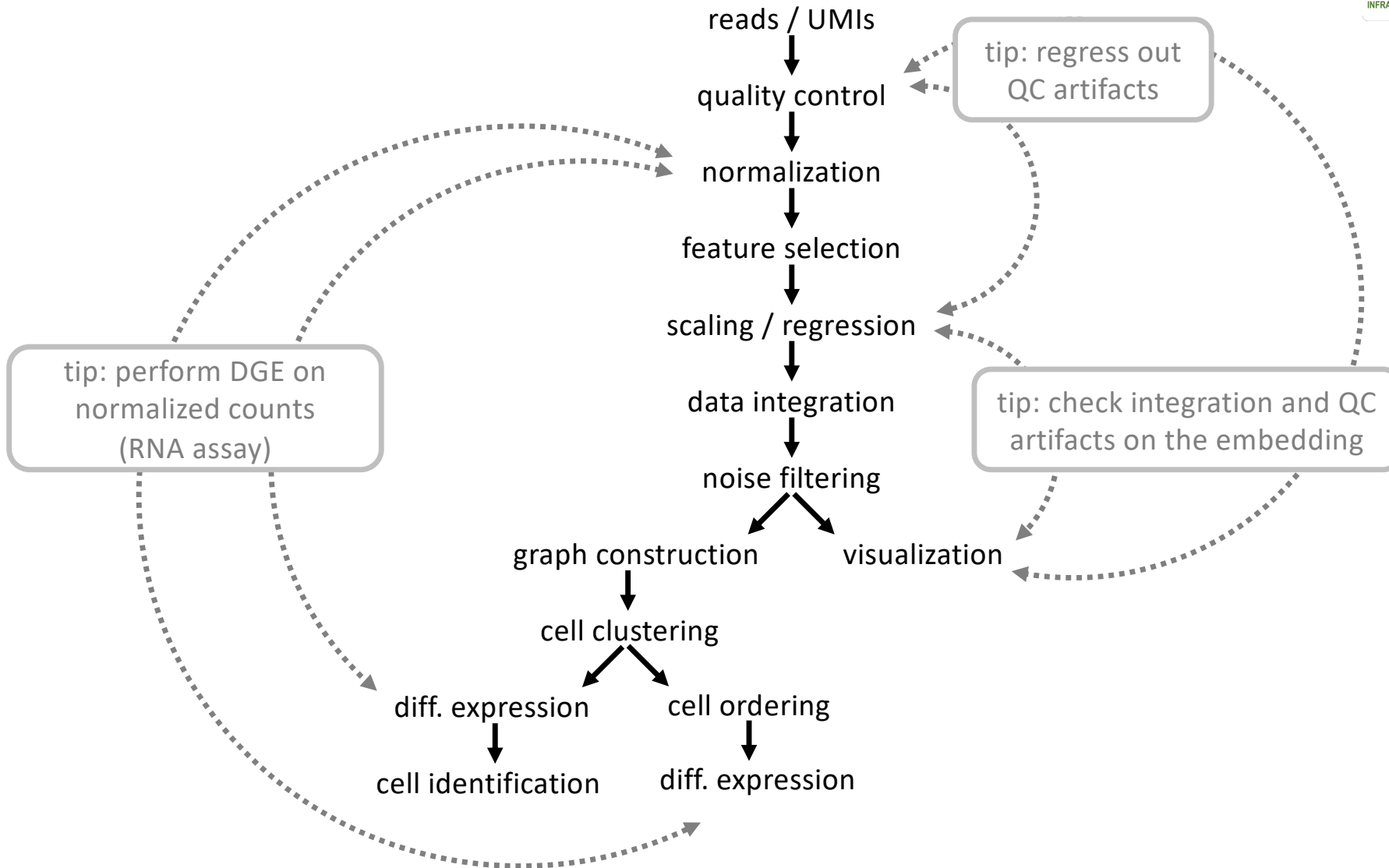




# scRNA-seq trajectory inference



# scRNA-seq analysis workflow



# scRNA-seq mini projects



Spatial  
transcriptomics



RNA  
velocity



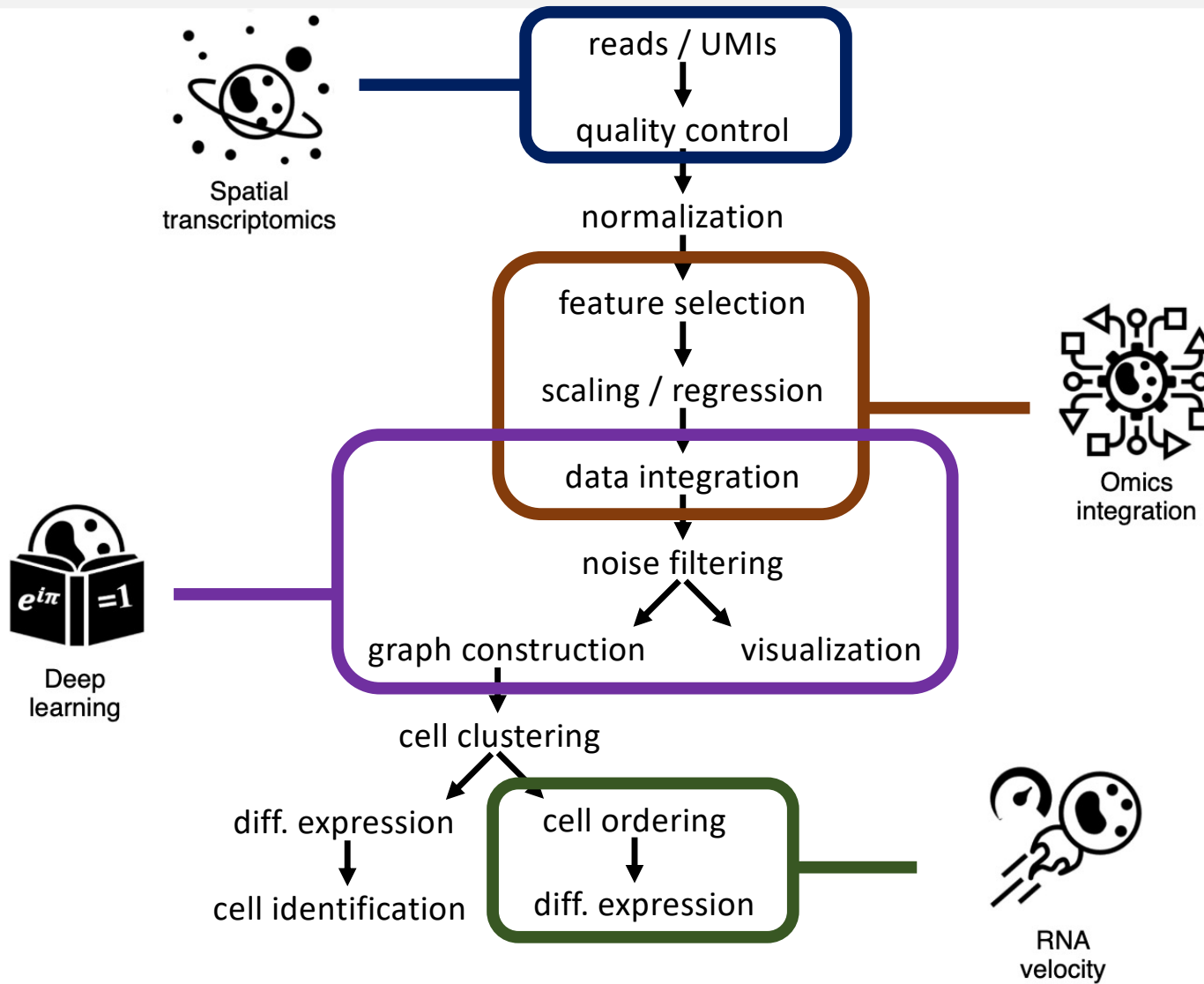
Omics  
integration



Deep  
learning



# scRNA-seq analysis workflow



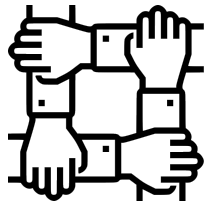
# Project-based learning (PBL)

Please read the material at:

[https://nbisweden.github.io/single-cell\\_sib\\_scilifelab\\_2021/projects.html](https://nbisweden.github.io/single-cell_sib_scilifelab_2021/projects.html)



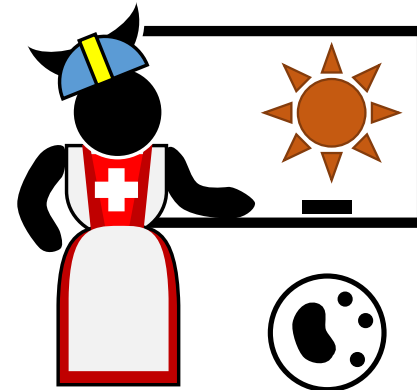
**Learning Strategy**



**Working in Groups**



**Tips for a good dynamic**



## Report.Rmd



### Load and merge datasets

- Consult the Glossary or additional sources for help
- Which file format do we have the data in?
- Describe in form of text the rationale for this step in your markdown report.



## Report.Rmd



### # Loading data

We first load the single cell RNA-seq dataset supplied from the ``.h5`` format in order to create a Seurat object.

```
```{r}
data <- Seurat::Read10X_h5( filename =
"data/colon_dataset.h5", use.names = T)
```
```



## Glossary



### Reading files

There are many formats available in which one can store single cell information, many of which cannot all be listed here. The most common formats are:

[...]

How to run it:

```
# From .csv .tsv .txt format
raw_matrix <- read.delim(
  file = "data/folder_sample1.csv",
  row.names = 1 )
```

```
# From .mtx format
sparse_matrix <- Seurat::Read10X(
  data.dir = "data/folder_sample1")
```

```
# From .h5 format
sparse_matrix <- Seurat::Read10X_h5(
  filename = "data/matrix_file.h5",
  use.names = T)
```

[...]



**Thank you!**

<https://czarnewski.github.io/czarnewski/index.html>