# scRNA-seq : visualization

Bastien Job, Gustave Roussy, Villejuif

Nathalie Lehmann, Institut Pasteur, Paris

Audrey Onfroy, Institut Mondor, Créteil
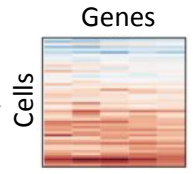
École de bioinformatique AVIESAN-IFB-INSERM 2024
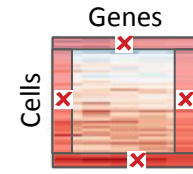
# scRNA-Seq pipeline overview



We want a <u>visual</u> <u>summary</u> of thousands cells' gene expression.

# How do we get to data visualization and clustering ?

Raw matrix

QC and filtering

Normalization

Clustering

# How do we get to data visualization and clustering ?



Raw matrix

QC and filtering

Normalization

We need to **summarize gene expression** to a few dimensions (20-100) to effectively separate the populations

Clustering

# Why an intermediary step is necessary ?



## scRNA-Seq data are sparse

> 70 % of the expression matrix is 0 : **not very informative**



http://cmdlinetips.com/wp-content/uploads/2018/03/Sparse_Matrix.png

```
prop(expr_mat == 0)
```

## Data are noisy

Some genes are more informative than some other.
There is **biological / technical noise** in gene expression.

## Computational time and ressources



We will summarize genes expression in few dimensions, before building the 2D projection.  5

# The right way to get to data visualization and clustering

# Our analyses goals



Dimension reduction → Clustering → Annotation

Cell type 1
Cell type 2
Cell type 3

# Challenges



❖ How to reduce the number of dimensions ?
❖ How many dimensions ?

❖ How to identify cell populations ?

❖ How ?

Genes

Dimensions

Cells

Cells

Cells

Dim. 2

Dim. 1

normalized matrix

reduced space

cells visualization

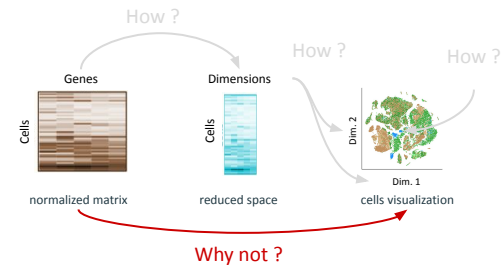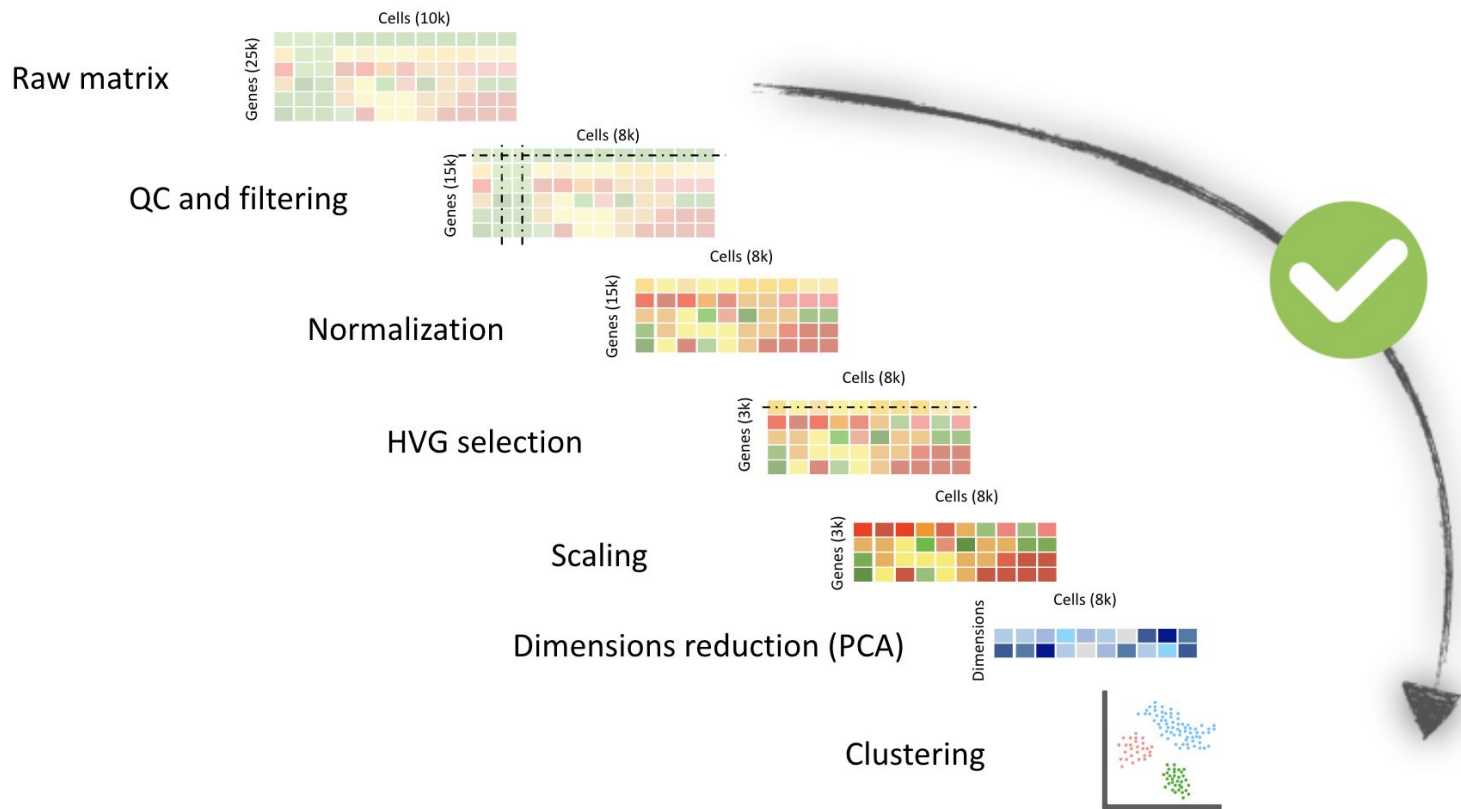We want a visual summary of thousands cells' gene expression.

# Dimensionality reduction

## Overview



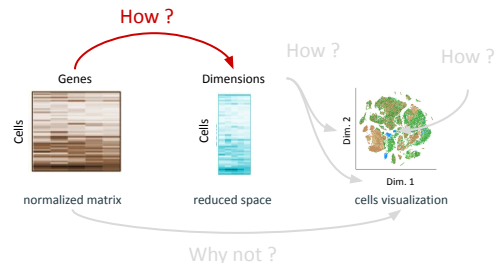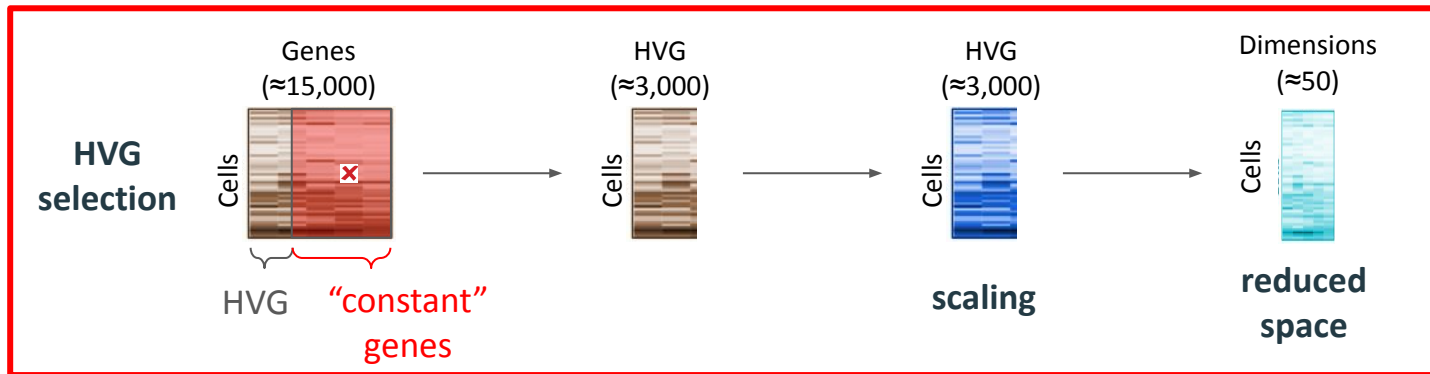## Commonly used **dimensionality reduction methods**

- **PCA**     **P**rincipal **C**omponent **A**nalysis
- **BFA**     **B**inary **F**actor **A**nalysis
- **ICA**     **I**ndependent **C**omponent **A**nalysis
- **LSI**     **L**atent **S**emantic **I**ndexing
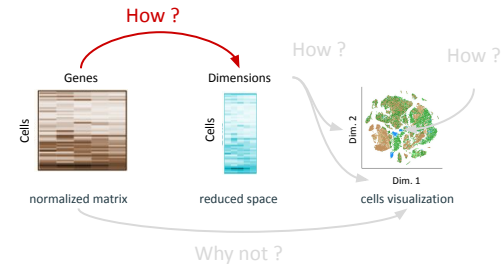- **LDA**     **L**inear **D**iscriminant **A**nalysis
- …

## Important parameters

- **information** : number of <u>variable</u> genes (HVG)
- number of **dimensions** to generate (signal / noise)
- *randomness : random seed*
- *convergence criteria*
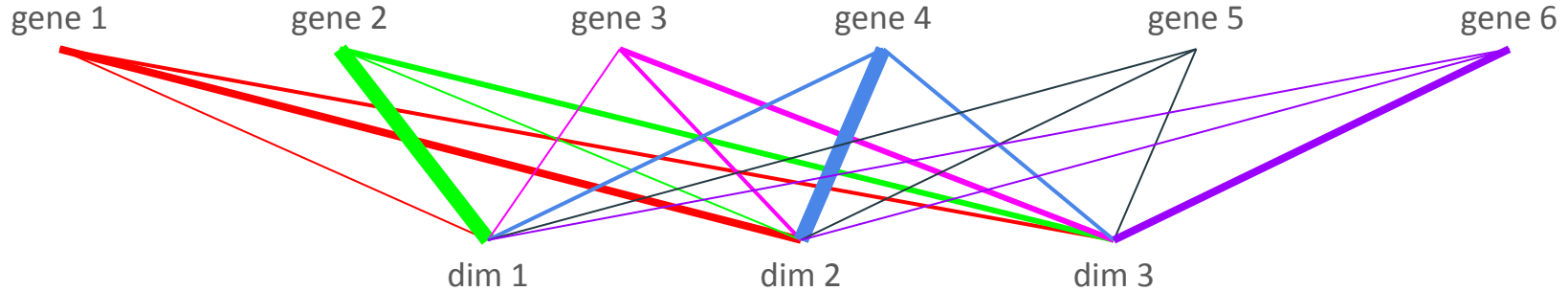


9

# Dimensionality reduction

## Principal Component Analysis - principle



- Input : **X** (≈ 2 000 - 5 000) HVG with **scaled** expression levels
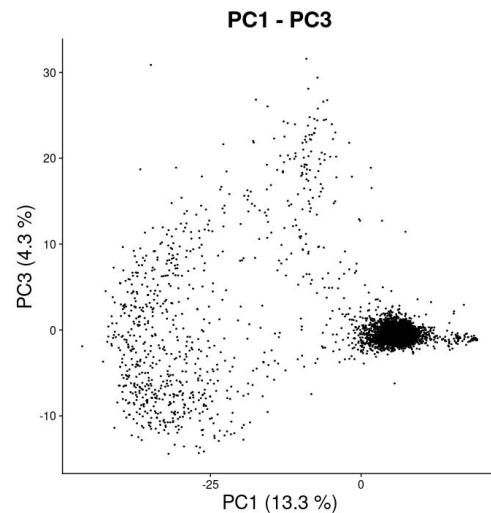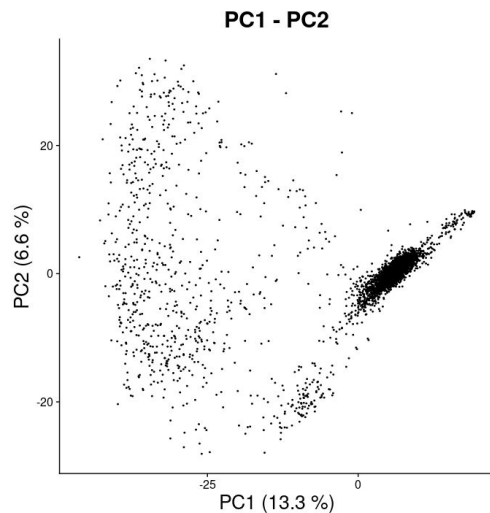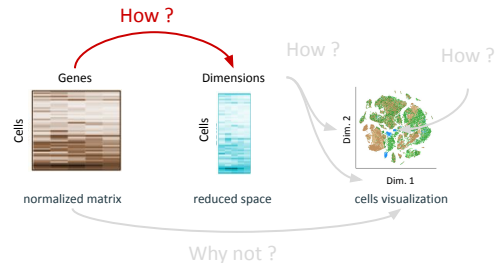- Goal : Group genes by dimensions when they have similar expression across cells



- Output : **Z** (≈ 50 - 100) dimensions "Principal Component"
- Each PC summarizes a certain amount of the input data variability
  - First PC recapitulates the most part of information
  - Last PC can be considered as noise
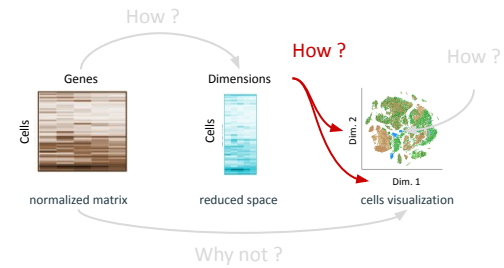


Scree plot

# Dimensionality reduction
## Principal Component Analysis - visualization



- Input : **X** most variable genes
- Goal : Group genes by dimensions when they have similar expression across cells
- Output : **Z** dimensions "Principal Component"
- Each PC summarizes a certain amount of the input data variability



Now, we will use the reduced space to make a 2D representation.

# 2D space for cells visualization



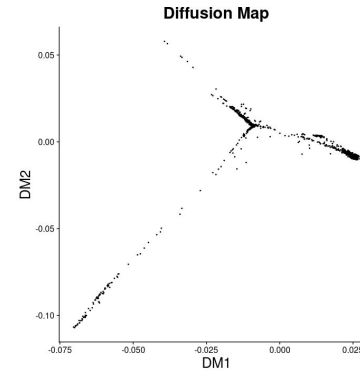## Commonly used 2D space

- **UMAP**
- **tSNE**
- **Diffusion Map**
- …

## Important parameters

- **input information** : number of dimensions
- cells **neighborhood** : number of neighbors, perplexity, distance method, …



The same cells can be represented using **different 2D spaces**.
Do not make to many interpretations from the 2D space, it is an **over-simplified representation** of cells.

# There are an infinite way to represent our data into 2D



https://distill.pub/2016/misread-tsne/

# Our analyses goals



Dimension reduction — Clustering — Annotation

Cell type 1
Cell type 2
Cell type 3

# Clustering

How ?
Genes
Dimensions
How ?
How ?
Cells
Cells
Dim. 2
Dim. 1
normalized matrix
reduced space
cells visualization
Why not ?

## Commonly used methods

- **Louvain clustering**
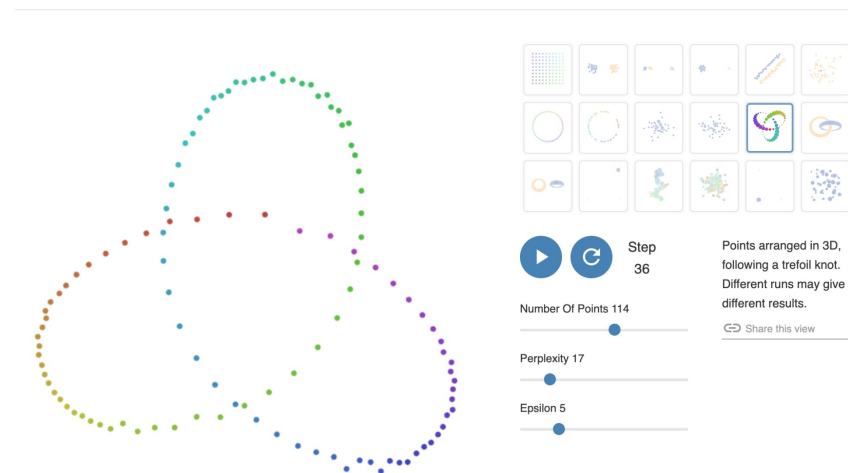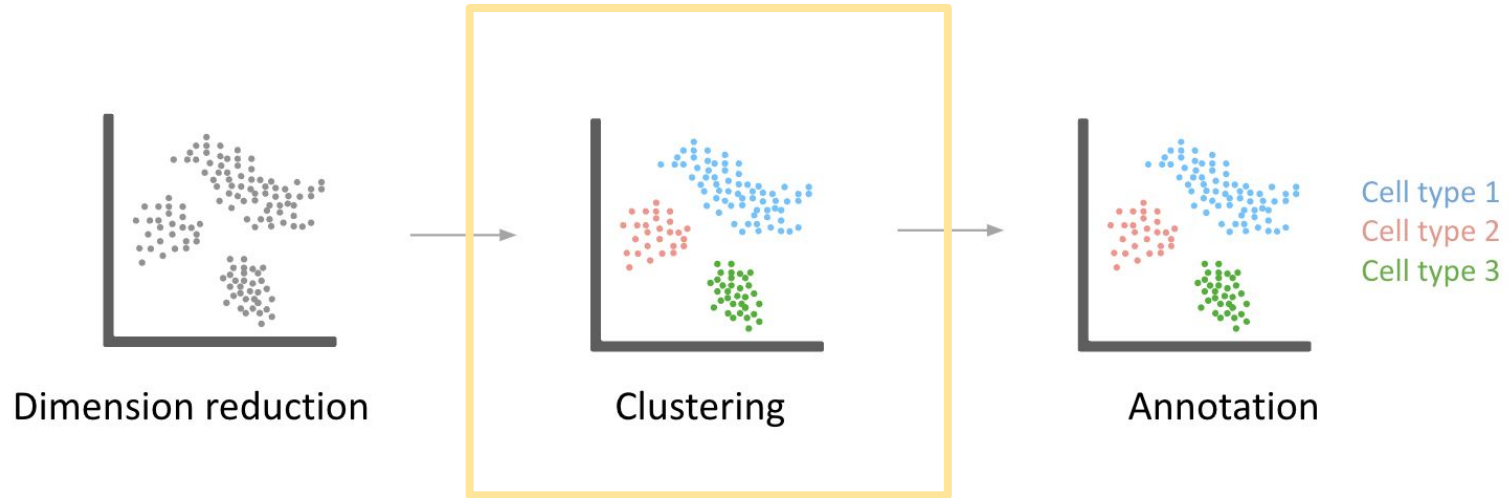- Leiden clustering
- k-means
- …

## Important parameters

- **input information** : number of dimensions
- cells **neighborhood** parameters : number of neighbors, distance measurement method, **resolution**…



**k-nearest neighbors**
(kNN)

**shared nearest neighbors**
(SNN)

**clustering**
(from SNN graph)

Clustering is made on expression matrix or reduced space, <u>not</u> on the 2D projection.
The 2D projection is not a clustering. A clustering is an **annotation**.

# Summary



Genes (≈15,000) — **normalized matrix** → Genes (≈15,000) — HVG / **HVG selection** → HVG (≈3,000) → HVG (≈3,000) — **scaled matrix** → Dimensions (≈50) — **reduced space**

Cells

**UMAP**

**tSNE**

**others…**

Cluster — Cells

Axis 2 / Axis 1

# Take Home Messages

- The **number of variable genes** impact the PCA, thus the 2D space. It depends on the expected number of cell populations in the dataset.

- Number of **dimensions** = amount of information (not enough < - - > noisy data)

- **UMAP** is suited to visualize several cell types and their <u>global</u> transcriptomic profile
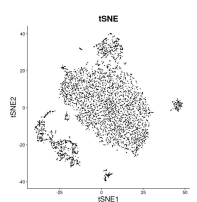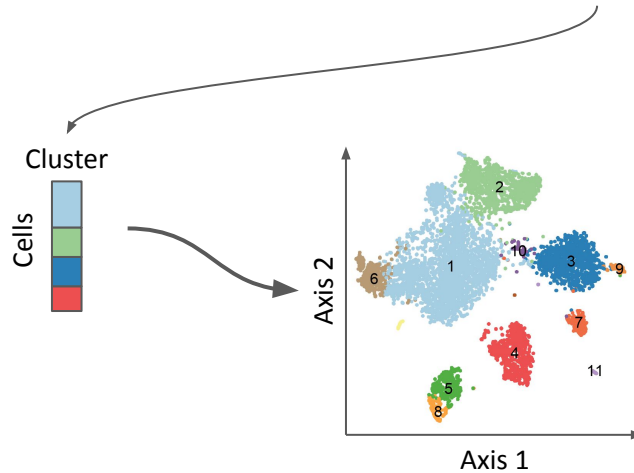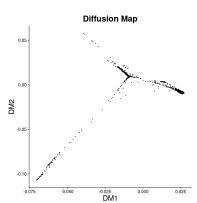
- **tSNE** is suited to visualize sub cell types and their <u>local</u> transcriptomic particularity

- **Diffusion Map** is suited to visualize cell <u>differentiation</u> data

- The **resolution** impacts the number of clusters : not enough clusters / not biologically interpretable clusters

Advice :

1. Make the analysis with all default settings :

   - **2000** HVG
   - **15** PC to generate a UMAP (or tSNE)
   - Resolution **1** for the clustering

   The goal is to generate a quick representation for your cells. Run your favorite analyses and represent results on the representation. Do not make to many interpretations from the 2D representation itself.
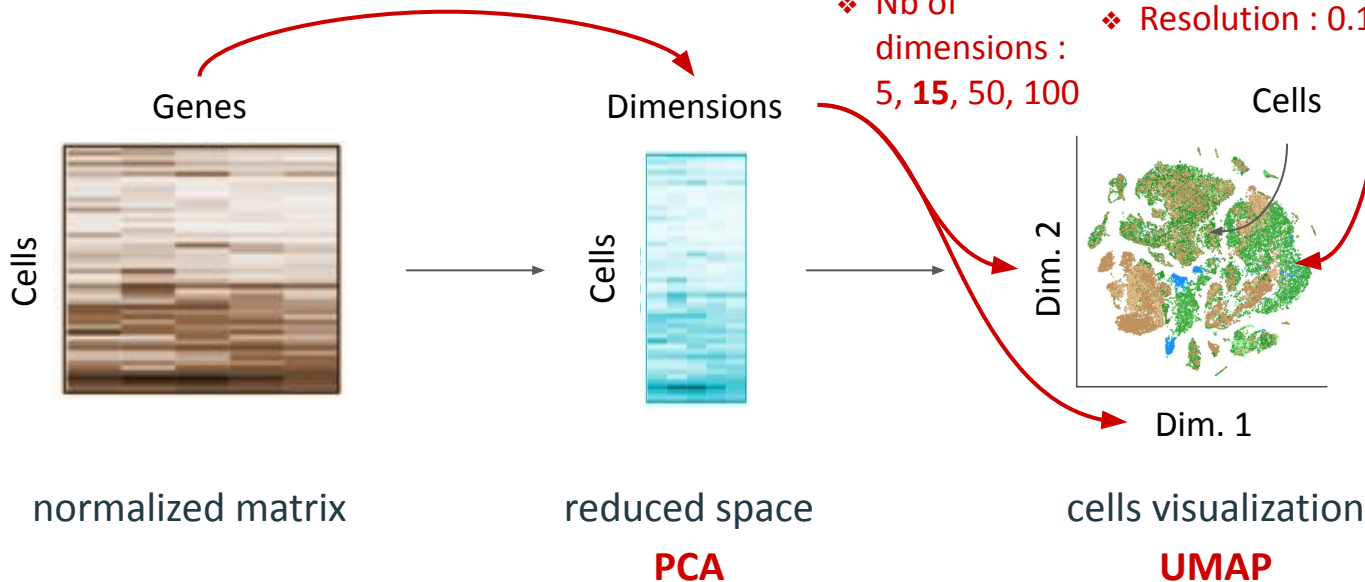
2. Identify your cell populations

3. Change the settings to make the representation showing what you identified

# Let's go to practice



❖ Nb of variable features : 500, **2000**, 5000
❖ Nb of dimensions : **50**

❖ Nb of dimensions : 5, **15**, 50, 100

❖ Nb of dims : same as UMAP
❖ Resolution : 0.1, 0.5, **1**, 5

Genes

Dimensions

Cells

Cells

Cells

Dim. 2

Dim. 1

normalized matrix

reduced space

**PCA**

cells visualization

**UMAP**

# Number of variable features

|   | 500 | 2000 | 5000 |
|---|-----|------|------|
| **5** | 500 - 5 | 2000 - 5 | 5000 - 5 |
| **15** | 500 - 15 | 2000 - 15 | 5000 - 15 |
| **50** | 500 - 50 | 2000 - 50 | 5000 - 50 |

Number of PC (/50) to make the UMAP

# Resolution

| 0.1 | 0.5 | 1 | 5 |
|---|---|---|---|