

scRNA-seq : visualization

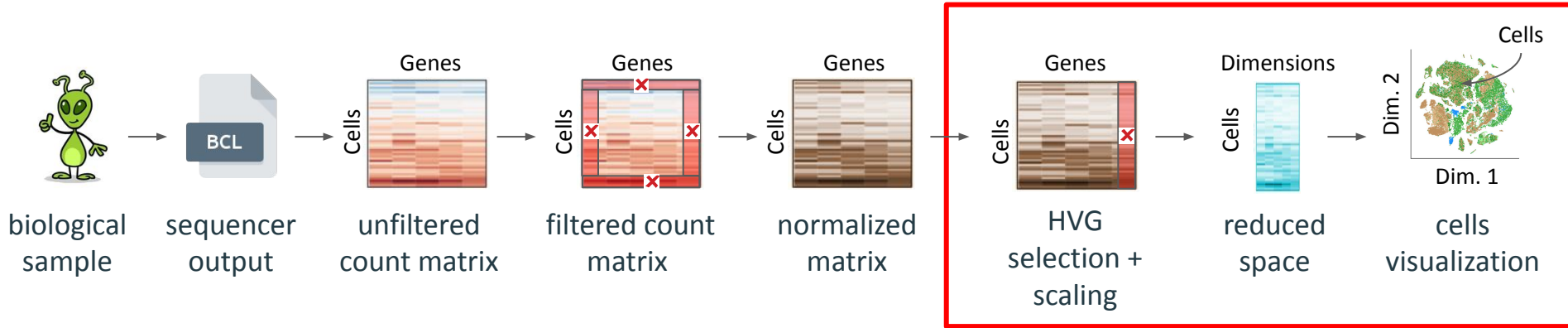
Bastien Job, Gustave Roussy, Villejuif

Lilia Younsi, Institut Cochin

Nathalie Lehmann, Institut Pasteur, Paris

Audrey Onfroy, Institut Mondor, Créteil

scRNA-Seq pipeline overview

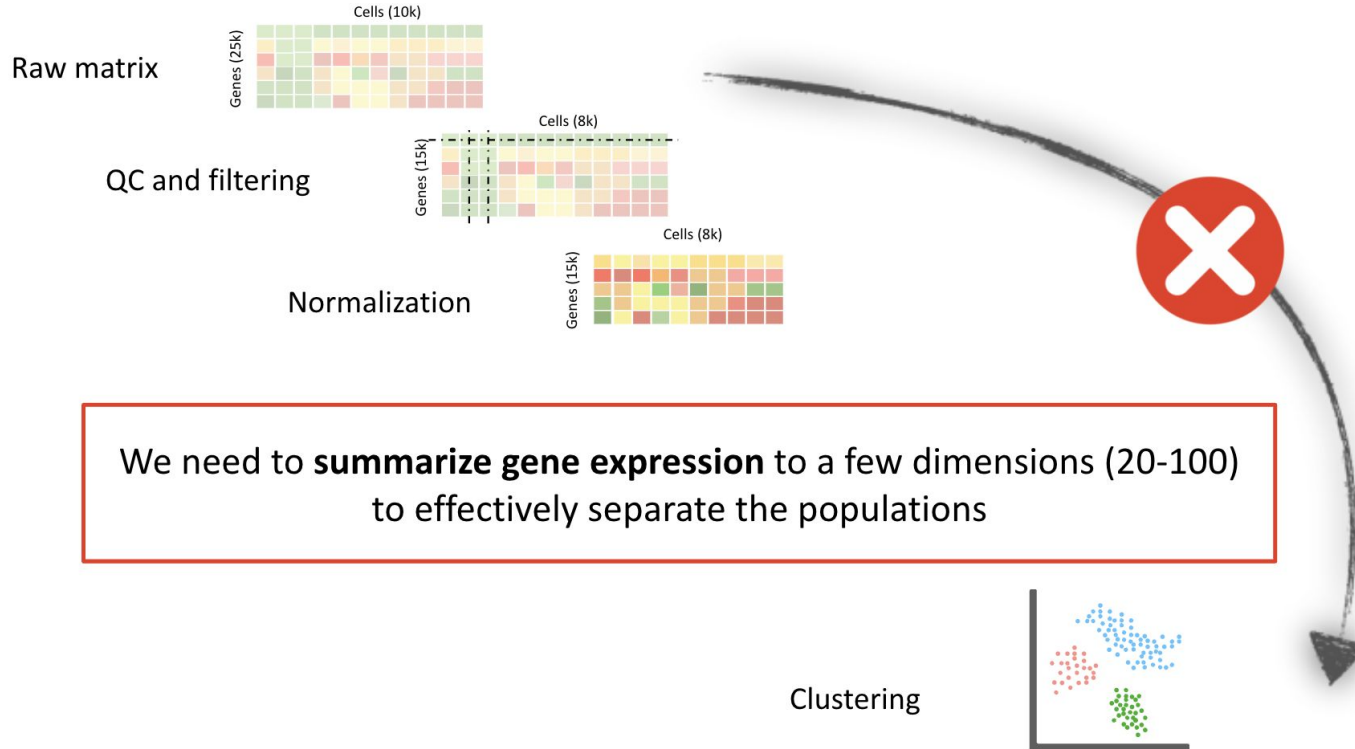


We want a visual summary of thousands cells' gene expression.

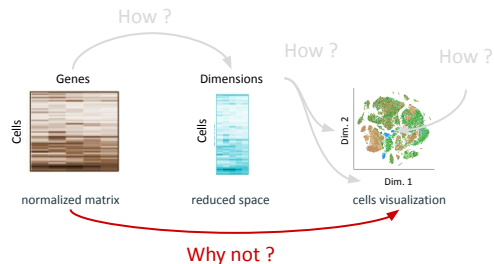
How do we get to data visualization and clustering ?



How do we get to data visualization and clustering ?



Why an intermediary step is necessary ?



scRNA-Seq data are sparse

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

Dense Matrix

1	2	31	2	9	7	34	22	11	5
11	92	4	3	2	2	3	3	2	1
3	9	13	8	21	17	4	2	1	4
8	32	1	2	34	18	7	78	10	7
9	22	3	9	8	71	12	22	17	3
13	21	21	9	2	47	1	81	21	9
21	12	53	12	91	24	81	8	91	2
61	8	33	82	19	87	16	3	1	55
54	4	78	24	18	11	4	2	99	5
13	22	32	42	9	15	9	22	1	21

Sparse Matrix

1	.	3	.	9	.	3	.	.	.
11	.	4	2	1
.	.	1	.	.	.	4	.	1	.
8	.	.	.	3	1
.	.	.	9	.	.	1	.	17	.
13	21	.	9	2	47	1	81	21	9
.
.	.	.	.	19	8	16	.	.	55
54	4	.	.	.	11
.	.	2	22	.	21

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

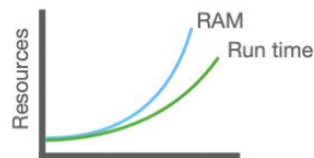


```
prop(expr_mat == 0)
```

Data are noisy

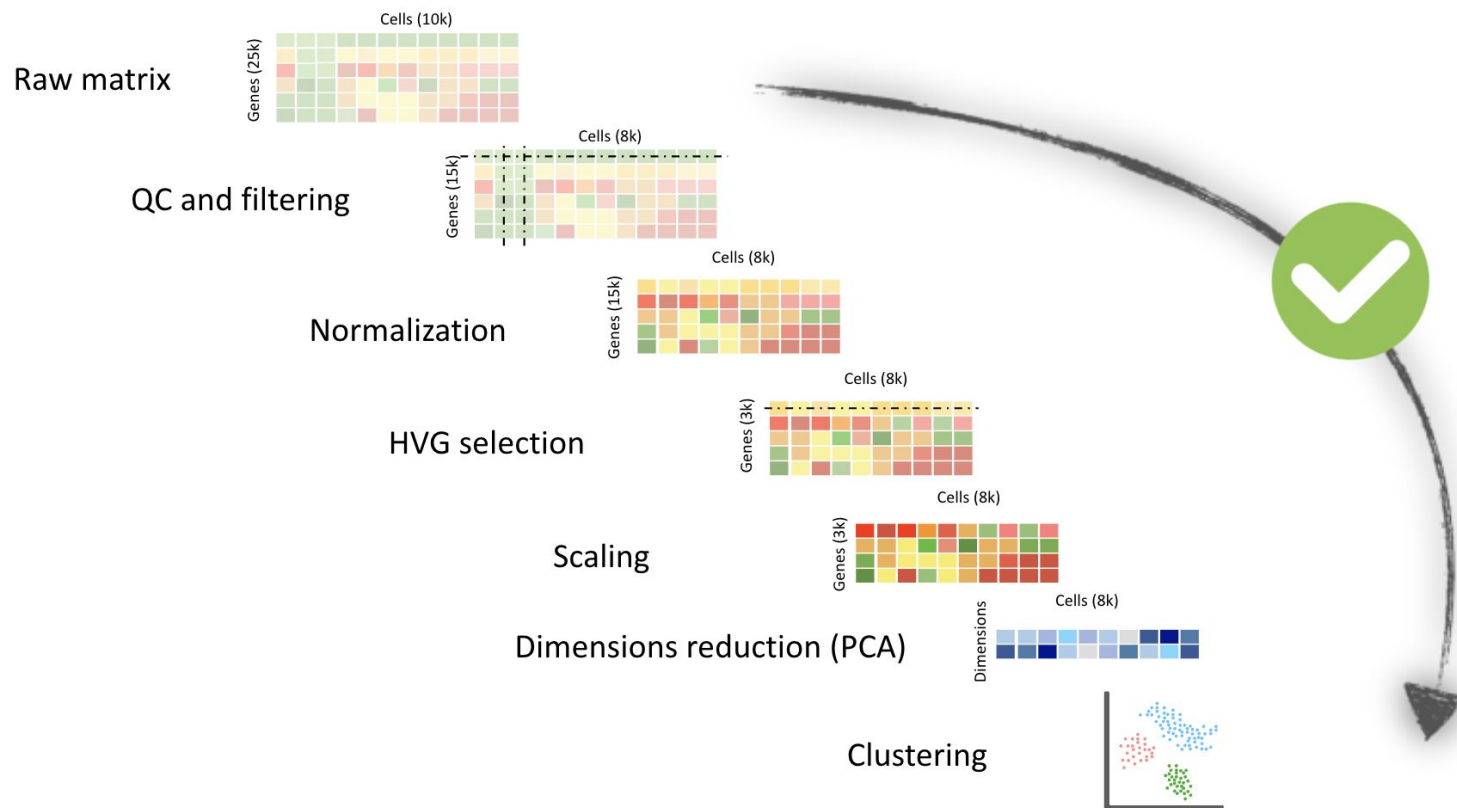
Some genes are more informative than some other.
There is **biological / technical noise** in gene expression.

Computational time and resources

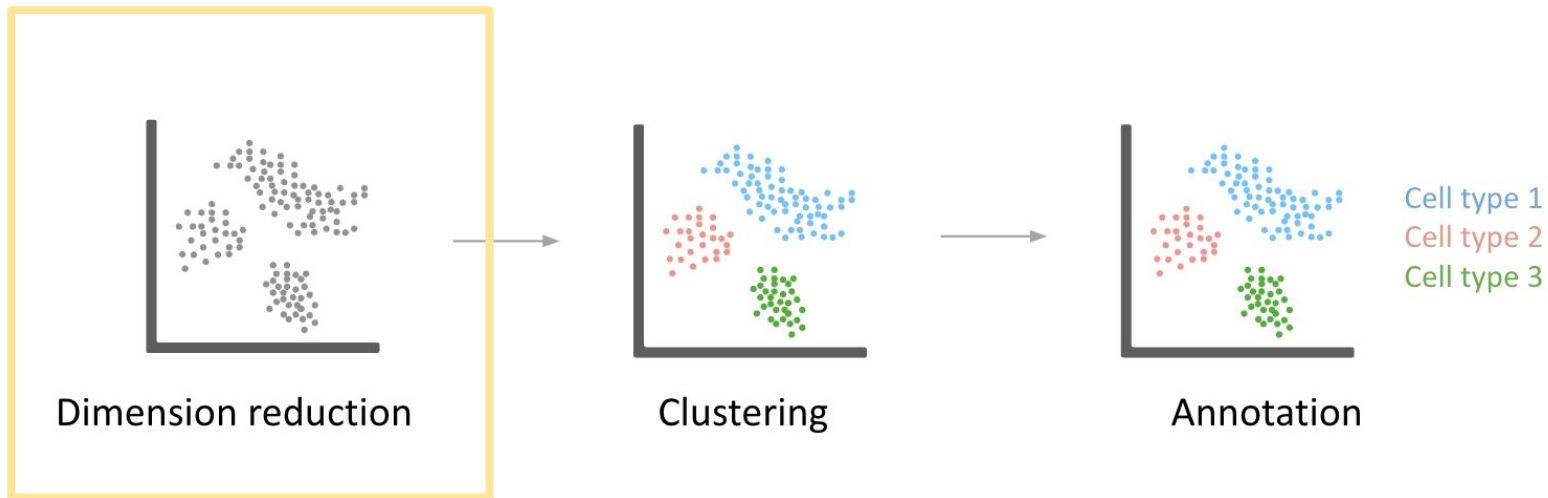


We will summarize genes expression in few dimensions, before building the 2D projection. 5

The right way to get to data visualization and clustering

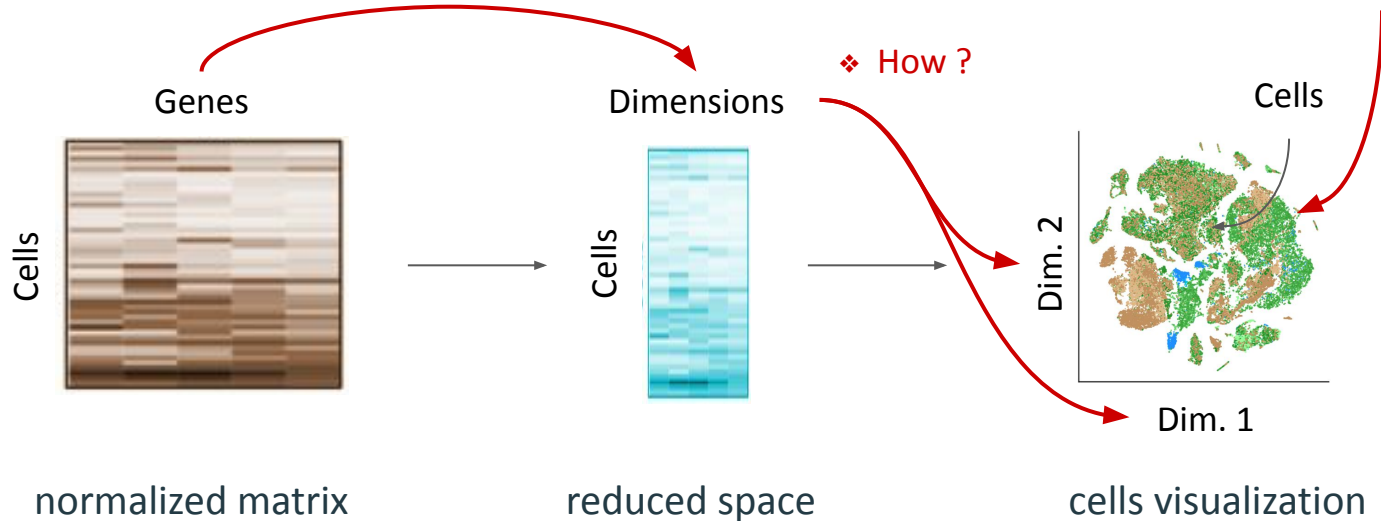


Our analyses goals



Challenges

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

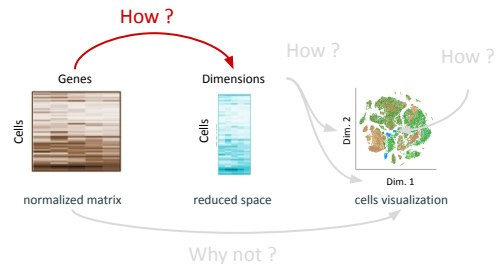


- ❖ How to identify cell populations ?

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

Dimensionality reduction

Overview

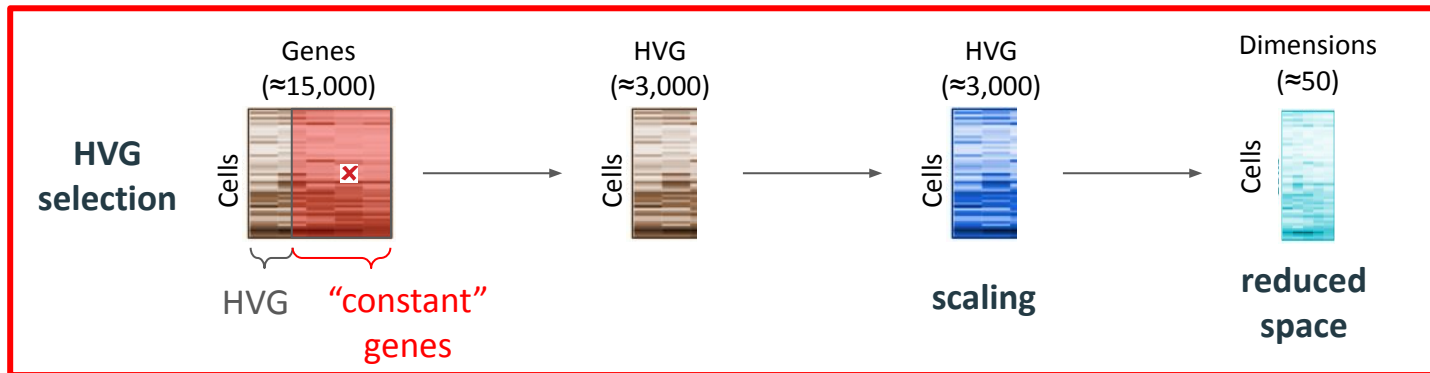


Commonly used **dimensionality reduction methods**

- **PCA** Principal **C**omponent **A**nalysis
- **BFA** Binary **F**actor **A**nalysis
- **ICA** Independent **C**omponent **A**nalysis
- **LSI** Latent **S**emantic **I**ndexing
- **LDA** Linear **D**iscriminant **A**nalysis
- ...

Important parameters

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



Dimensionality reduction

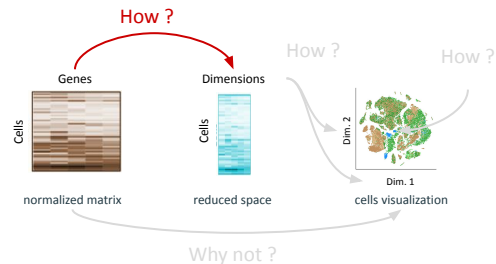
PCA introduction

1000 wine bottles



Features that will vary from one bottle to another :

- Acidity
- Tannins
- Alcohol level
- Aroma
- Color
- Clarity
- Color intensity
- Freshness (acidity driven)
- ...



Which features explain the big differences between my bottles ?

How can I sum up this data ?

Dimensionality reduction

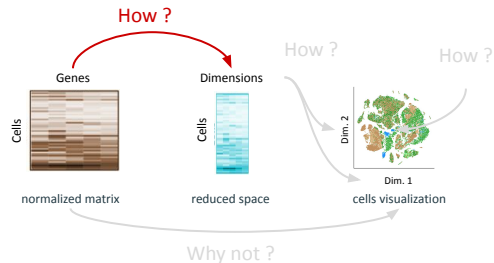
PCA introduction

1000 wine bottles



Features that will vary from one bottle to another :

- Acidity
- Tannins
- Alcohol level
- Aroma
- Color
- Clarity
- Color intensity
- Freshness (acidity driven)
- ...



TOO MUCH INFO !

Impossible to compare all these variables 1 to 1 for all bottles without getting lost.

Dimensionality reduction

PCA introduction

1000 wine bottles



Visual Inspection

This chart will help you build your mental repertoire for identifying wines by hue and intensity. You will find a useful key for blind tasting and assessing quality.

Hue in Red Wine

Algebraic or not, wine is affected by several factors including pH level. Wines with red and blue hues tend to have a lower pH than wines with blue hues.

Advice for Viewing Wine

Look at wine under bright, efficient, natural lighting over a white background for best results. This is best observed at the point where the wine meets the glass.

WINE FOLLY

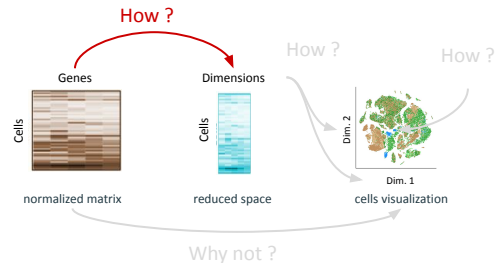
© 2019 Wine Folly Ltd. Made in Seattle, WA, USA. Learn more at www.winefolly.com

Features that will vary from one bottle to another :

- Acidity
- Tannins
- Alcohol level
- Aroma
- Color
- Clarity
- Color intensity
- Freshness (acidity driven)
- ...

REDUNDANT INFORMATION

Acidity ⇔ Freshness



Dimensionality reduction

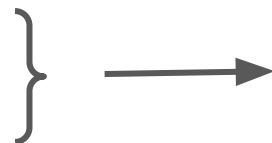
PCA introduction

1000 wine bottles

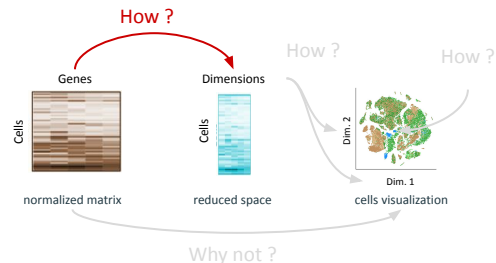


Features that will vary from one bottle to another :

- Acidity
- Tannins
- Alcohol level
- Aroma
- Color
- Clarity
- Color intensity
- Freshness (acidity driven)
- ...



Features can be combined in one dimension : “**robe du vin**” (or *wine apprerance*).



Dimensionality reduction

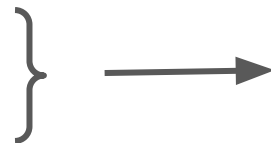
PCA introduction

1000 wine bottles



Features

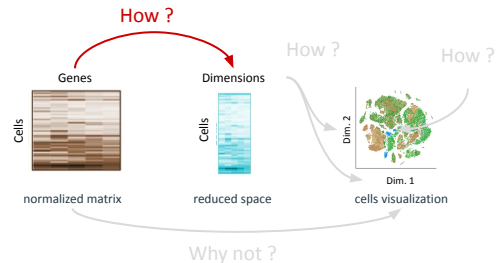
- Acidity
- Tannins
- Alcohol level
- Aroma
- **Color**
- **Clarity**
- **Color intensity**
- Freshness (acidity driven)
- ...



Features can be combined in one dimension : “**robe du vin**” (or *wine apperance*).



Principal Component (PC)



Visual Inspection

This chart will help you build your mental repertoire for identifying wines by hue and intensity. You will find a useful key for blind tasting and assessing quality.

Hue in Red Wine

Aliphatic or malic acid is critical for several factors including pH levels. Wines with red-based hues tend to have a lower pH than wines with blue-based hues.

Advice for Viewing Wine

Look at wine under bright, indirect, natural lighting over a white background for best results. This is best observed at the point where the wine meets the glass.

WINE FOLLY

© 2019 Wine Folly LLC. Made in Seattle, WA, USA. Learn more at www.winefolly.com

Dimensionality reduction

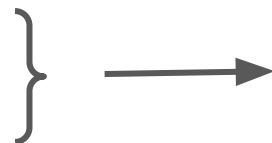
PCA introduction

1000 wine bottles



Features

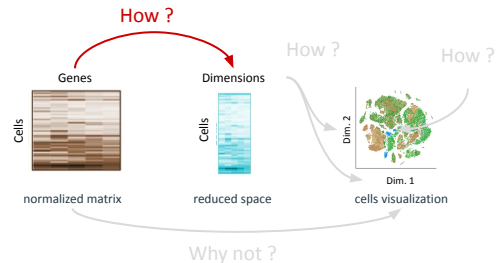
- Acidity
- Tannins
- Alcohol level
- Aroma
- Color
- Clarity
- Color intensity
- Freshness (acidity driven)
- ...



Features can be combined in one dimension : “**robe du vin**” (or *wine apperance*).



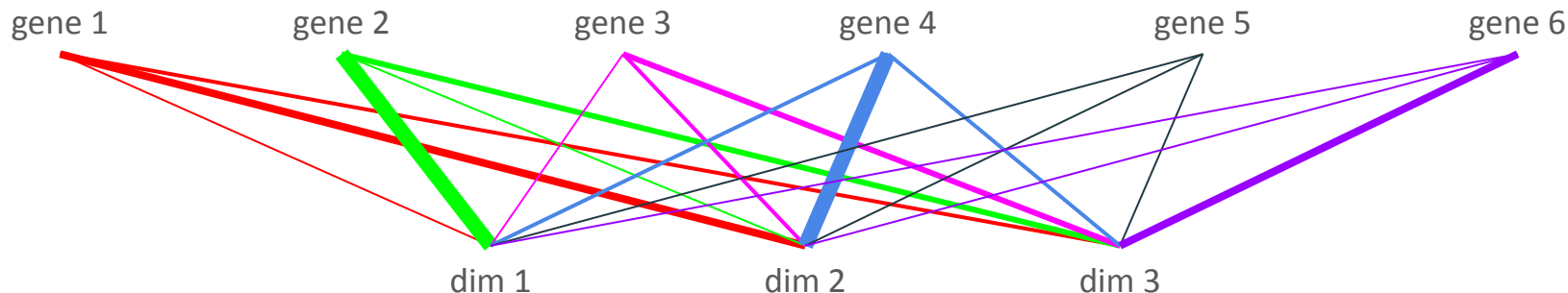
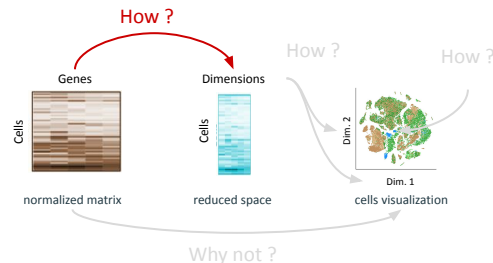
Principal Component (PC)



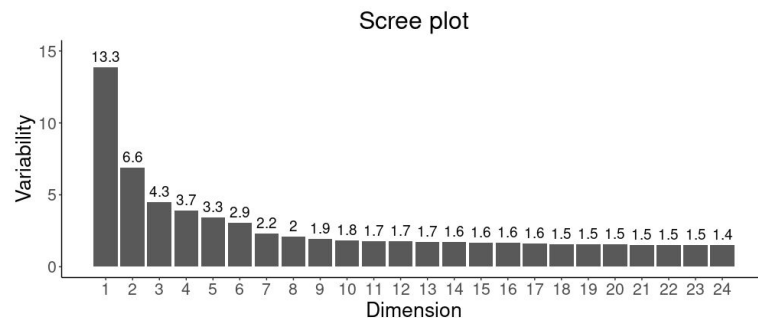
Dimensionality reduction

Principal Component Analysis - principle

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



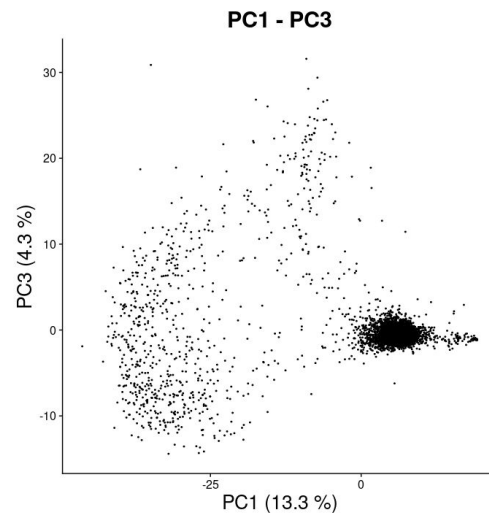
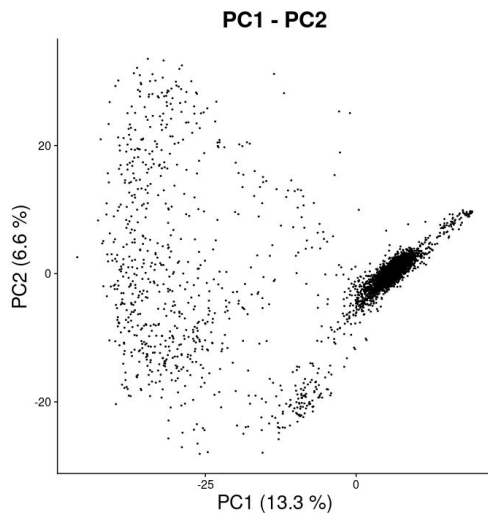
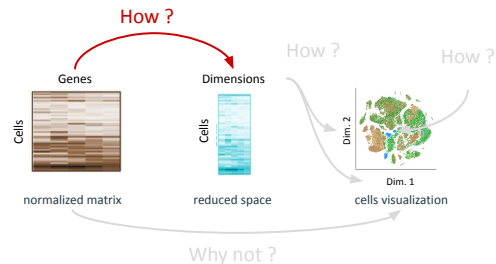
- Output : **Z** ($\approx 50 - 100$) dimensions “Principal Component”
- Each PC summarizes a certain amount of the input data variability
 - First PC recapitulates the most part of information
 - Last PC can be considered as noise



Dimensionality reduction

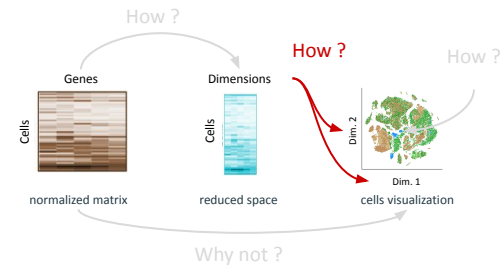
Principal Component Analysis - visualization

- 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 this reduced space to build a 2D graphical representation.

2D space for cells visualization

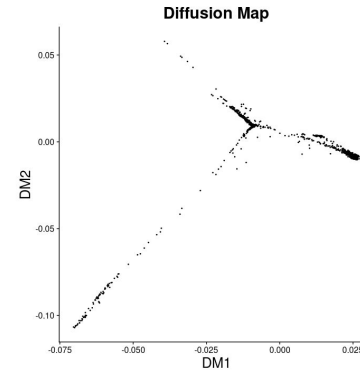
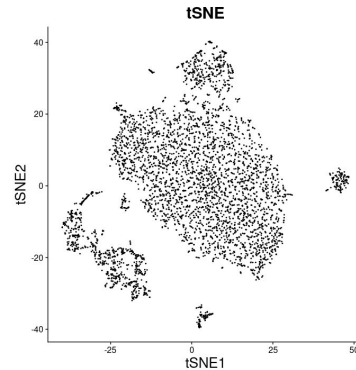
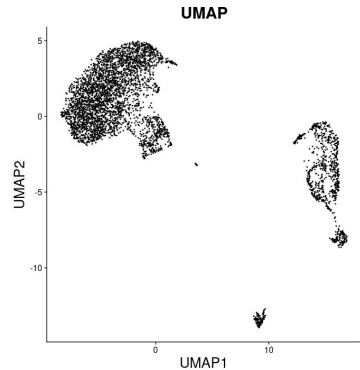


Commonly used 2D spaces

- UMAP
- tSNE
- Diffusion Map
- ...

Important parameters

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



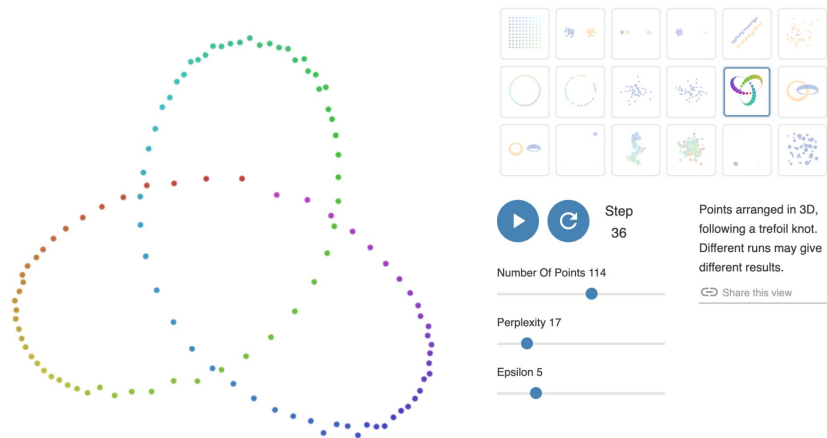
The same cells can be represented using **different 2D spaces**.

Do not make too 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

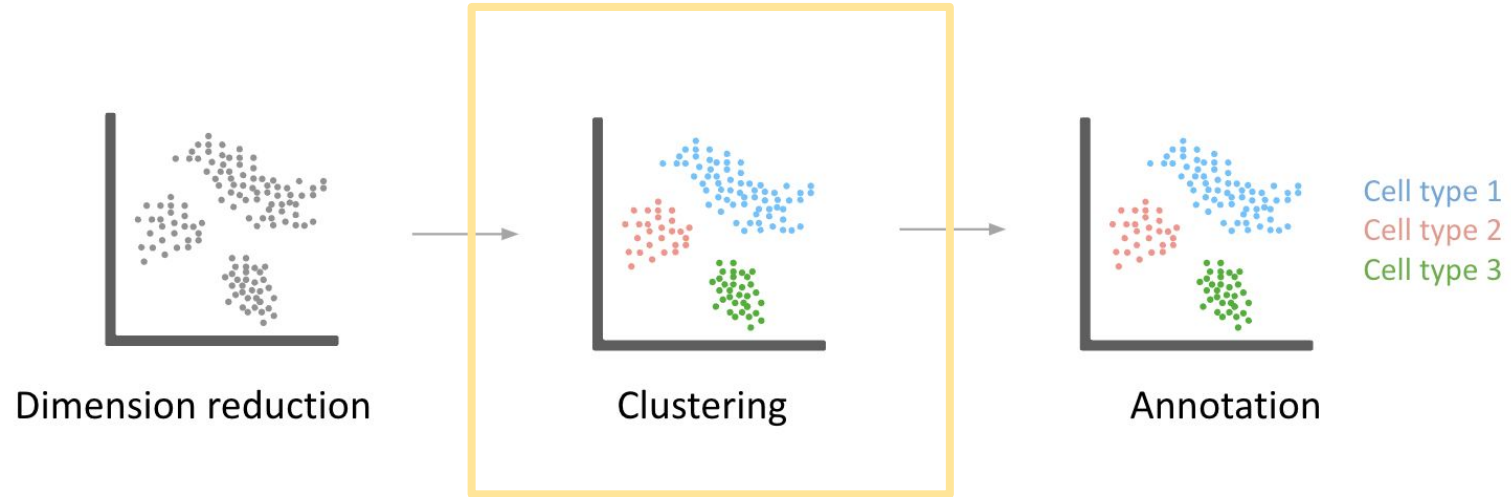
How to Use t-SNE Effectively

Although extremely useful for visualizing high-dimensional data, t-SNE plots can sometimes be mysterious or misleading. By exploring how it behaves in simple cases, we can learn to use it more effectively.



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

Our analyses goals



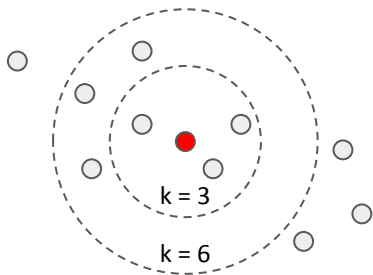
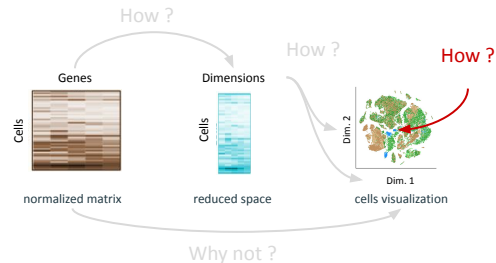
Clustering

Commonly used methods

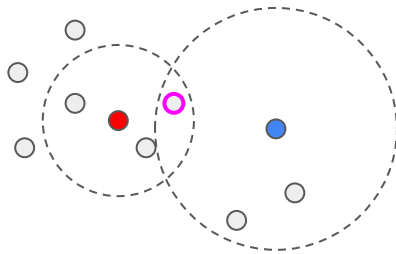
- Louvain clustering
- Leiden clustering
- k-means
- ...

Important parameters

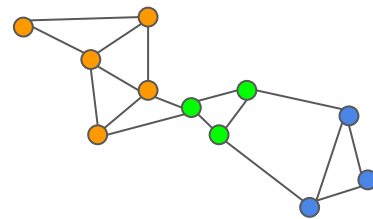
- **input information** : number of dimensions
- cells **neighborhood** parameters : number of neighbors, distance measurement method, **resolution**...



k-nearest neighbors
(kNN)



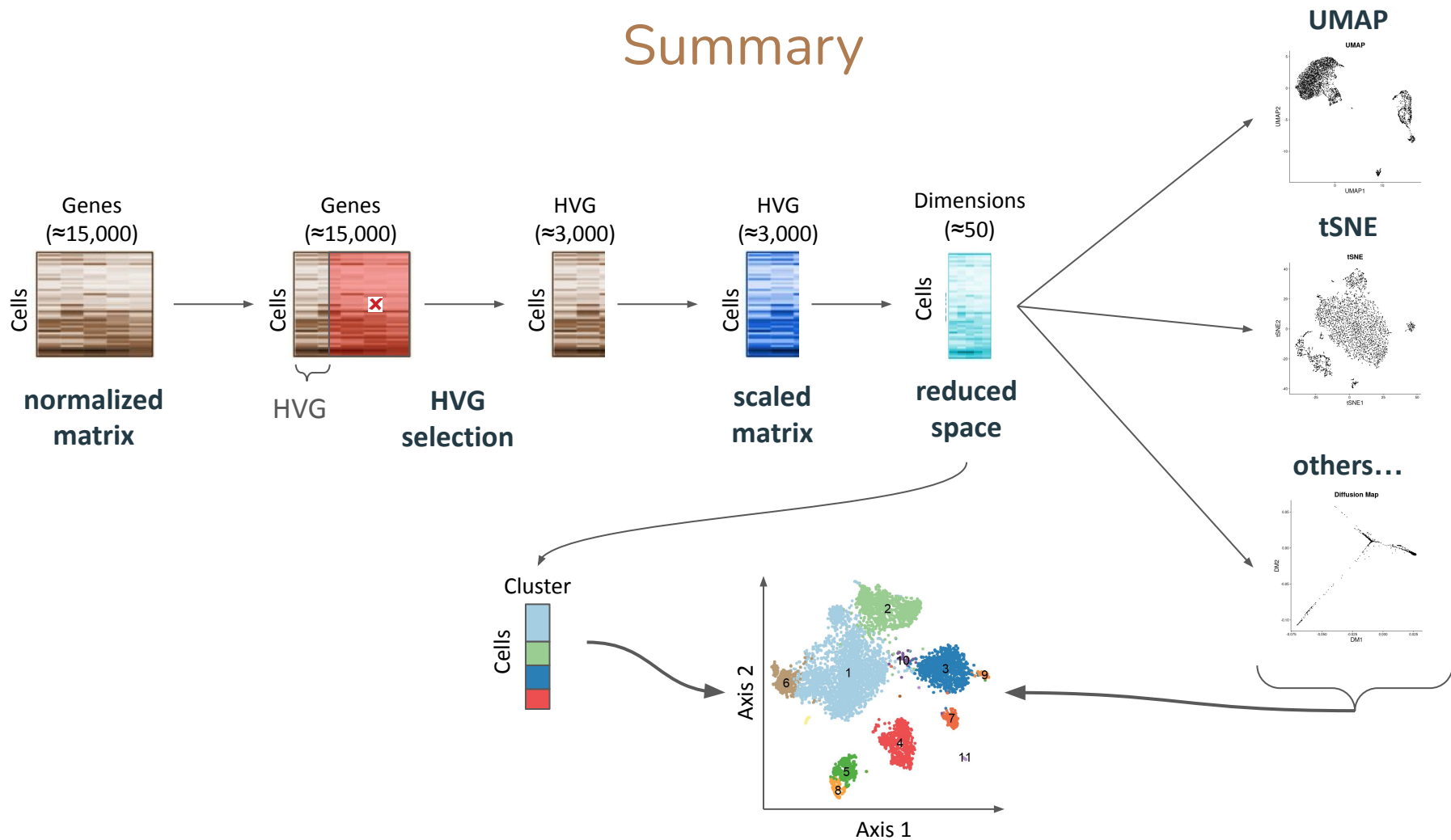
shared nearest neighbors
(SNN)



clustering
(from SNN graph)

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

Summary



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 global transcriptomic profile
- **tSNE** is suited to visualize sub cell types and their local transcriptomic particularity
- **Diffusion Map** is suited to visualize cell differentiation data
- The **resolution** impacts the number of clusters : not enough clusters / not biologically interpretable clusters

Advice :

1. Make the analysis with all default settings :

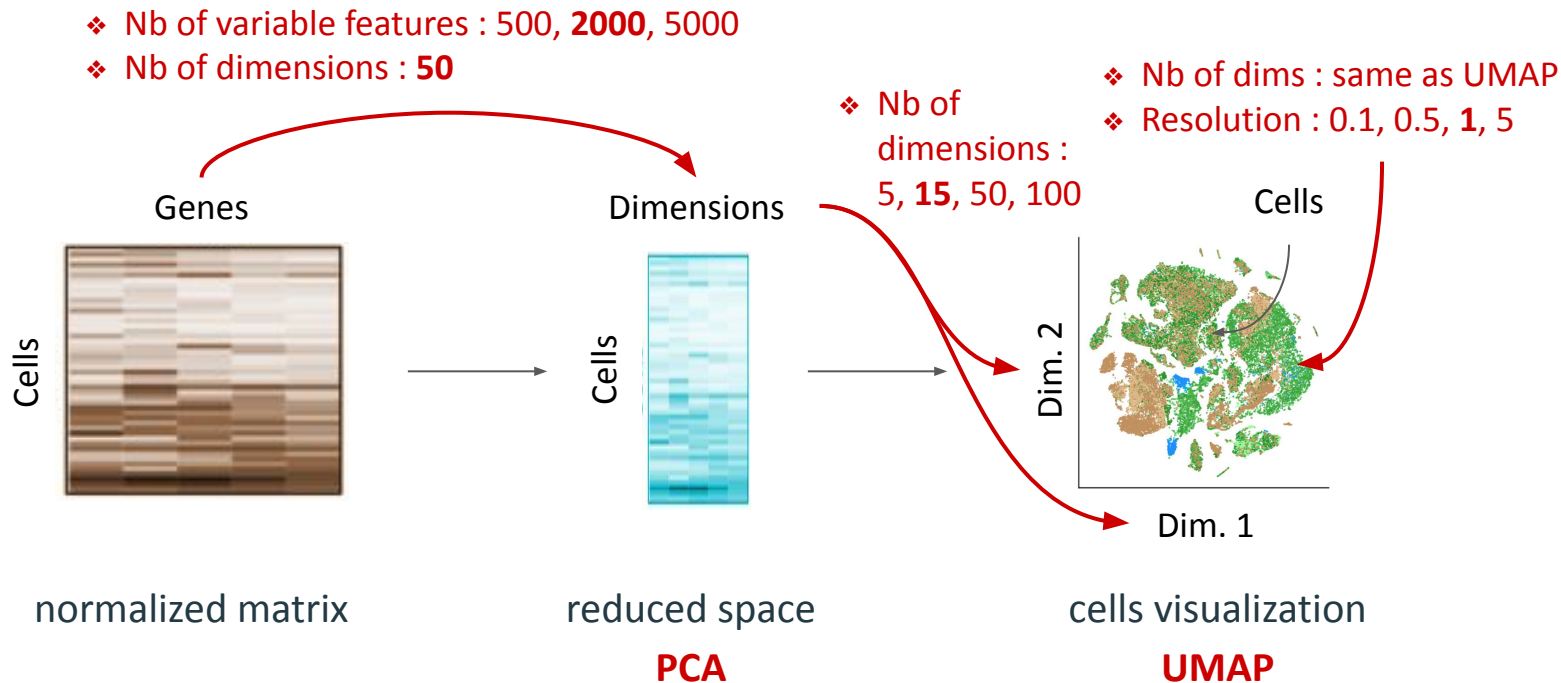
- **2000** HVG
- **15** PC to generate a UMAP (or tSNE)
- Resolution **1** for the clustering

2. Identify your cell populations

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

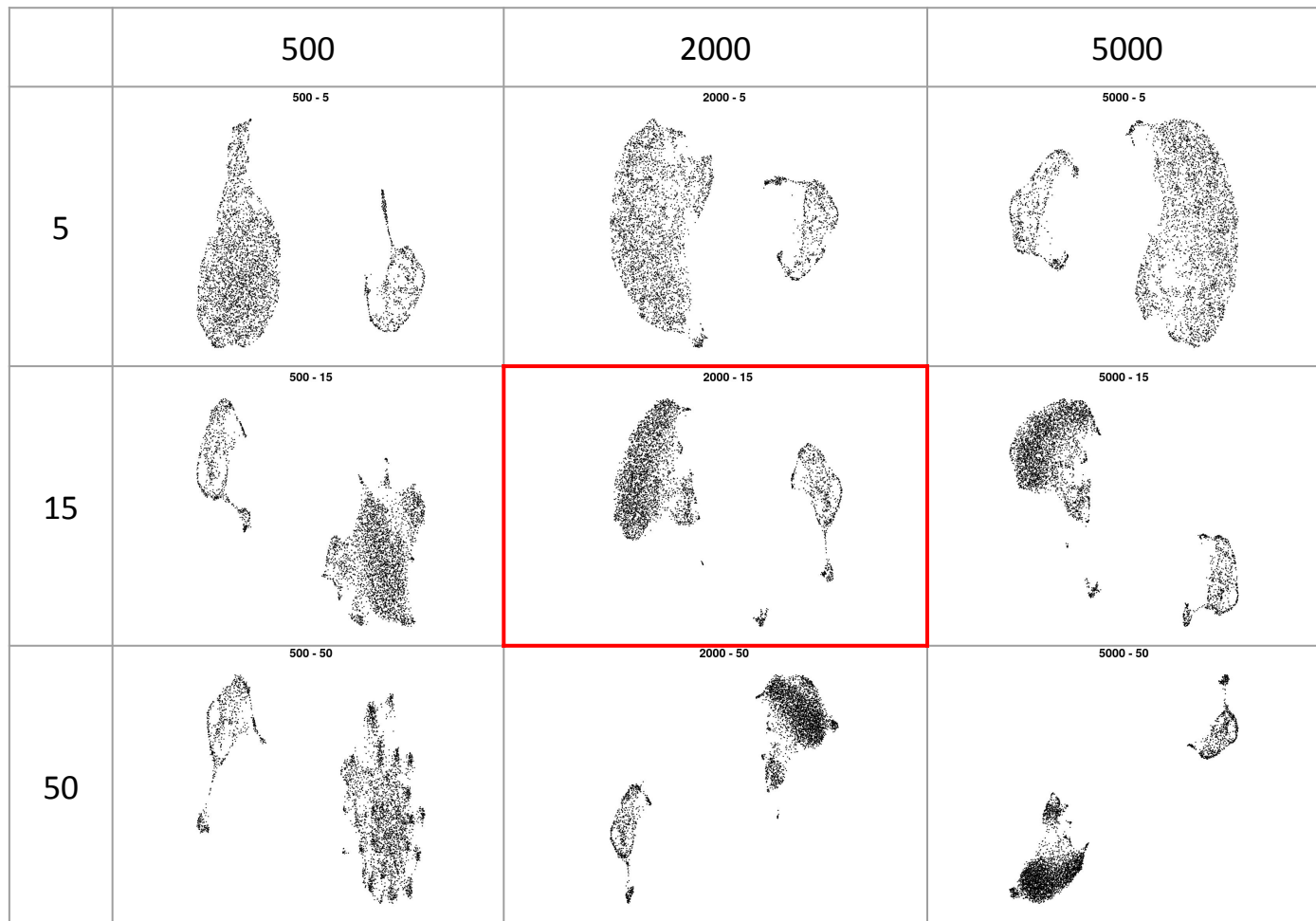
The goal is to generate a quick representation for your cells. Run your favorite analyses and represent results on the representation. Do not make too many interpretations from the 2D representation itself.

Let's go to practice



Number of PC (/50) to make the UMAP

Number of variable features



Resolution

0.1

0.5

1

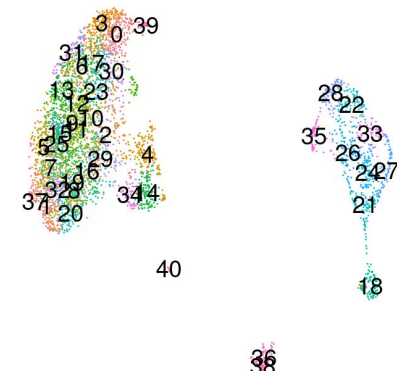
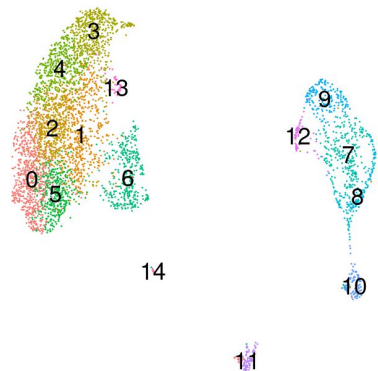
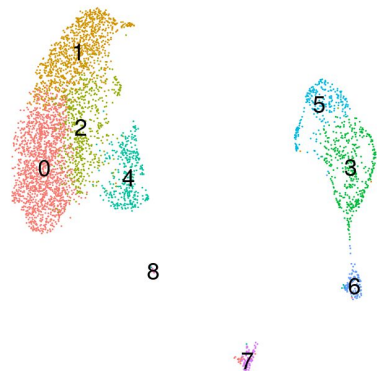
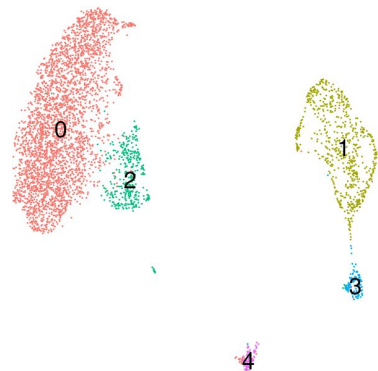
5

Resolution : 0.1

Resolution : 0.5

Resolution : 1

Resolution : 5



Choosing resolution wisely

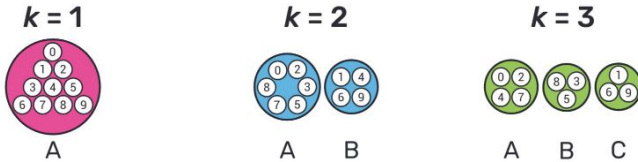


- Too low resolution → losing information of populations
- Too high resolution → **overclustering**

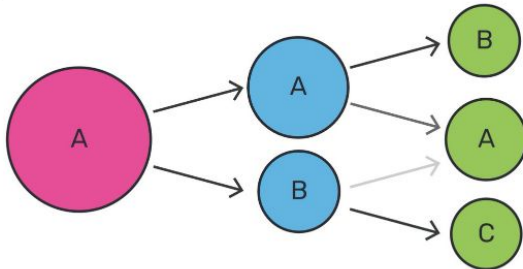
→ **Clustering trees** can be used to help choose the optimal resolution !

Step 1 - Clustering at multiple resolutions

```
res <- c(0.1, 1.2, 0.1)
```



Step 2 - Building and visualizing the clustering tree



Example :

