

Genome Informatics 2022

Lesson 5 - Single cell RNA & Spatial transcriptomics analysis

Why Single Cell study?

- Hidden variation in gene expression
- Regulatory process of biotechnological or medical relevance
- Relationship between cellular processes and external stimuli

Why Single Cell study?

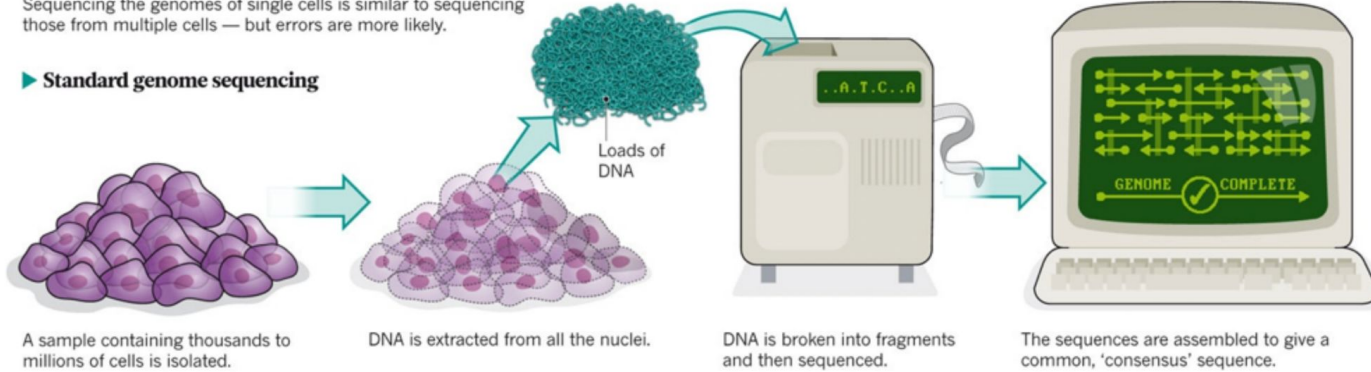
- Developmental biology
 - Discover more complicated mechanisms in cellular development
 - Confirm the distinct gene expression signatures across different cell types
 - Identify functional differences among the same cell cell type
- Cancer biology
 - Find evidence for models of cancer
 - Infer timing of mutations and the drivers
 - Evaluate effectiveness of targeted therapy
- Microbiology
 - Discover low-abundance species that are are difficult to culture in vitro
 - Monitor transcriptional gene activation mechanisms for functional annotation

Bulk RNA sequencing vs Single cell

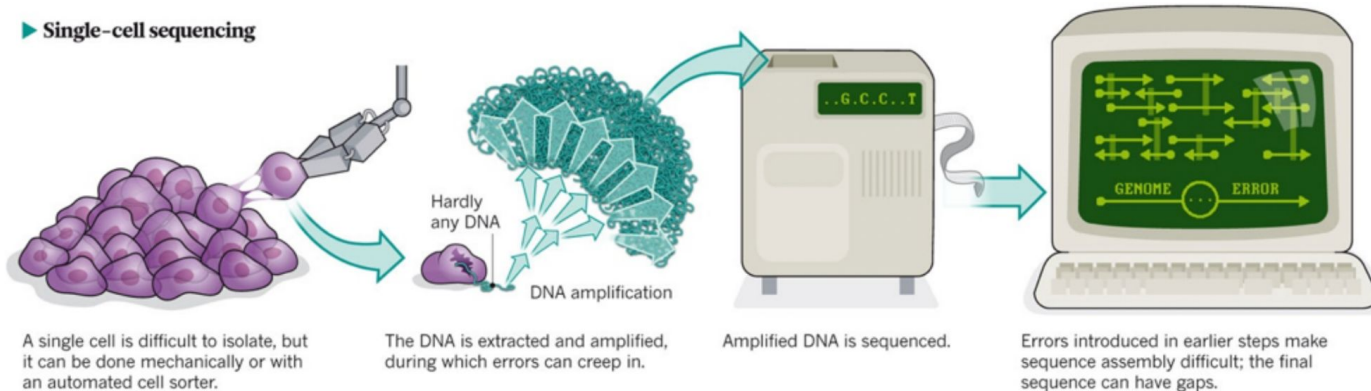
ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

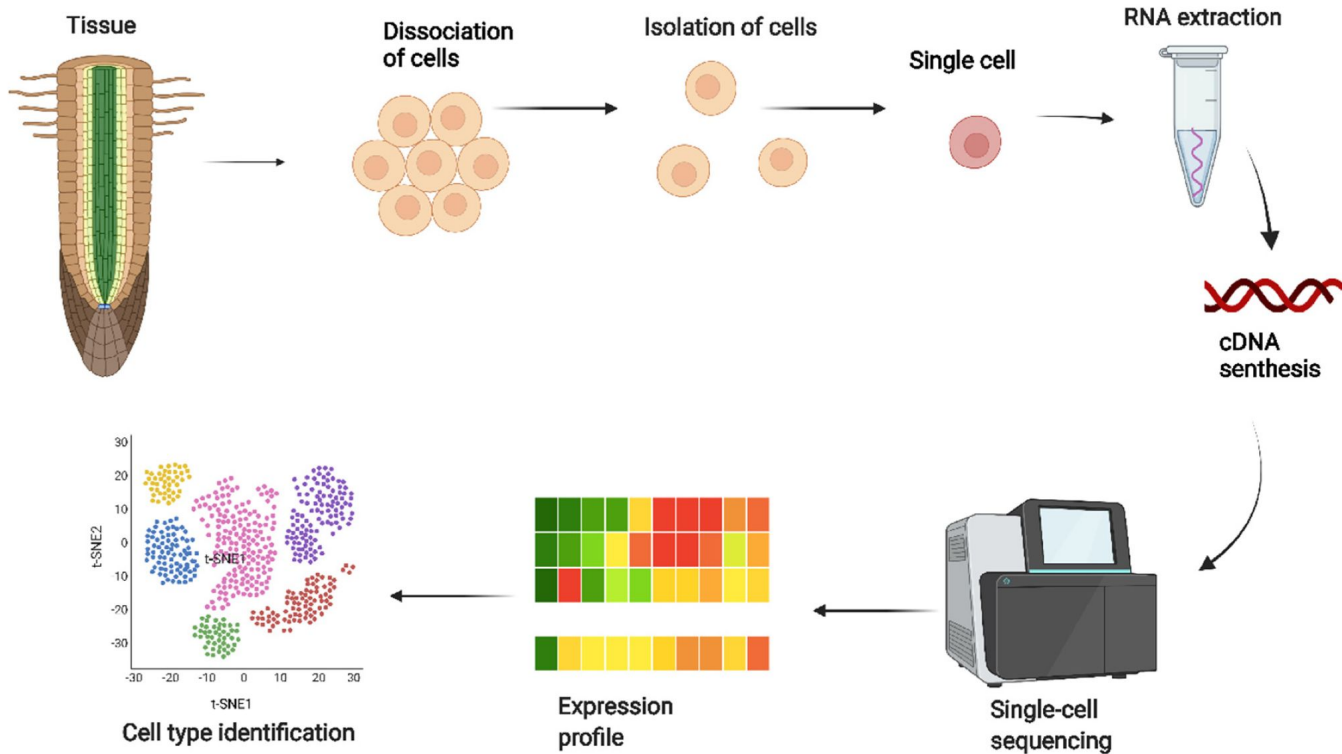
► Standard genome sequencing



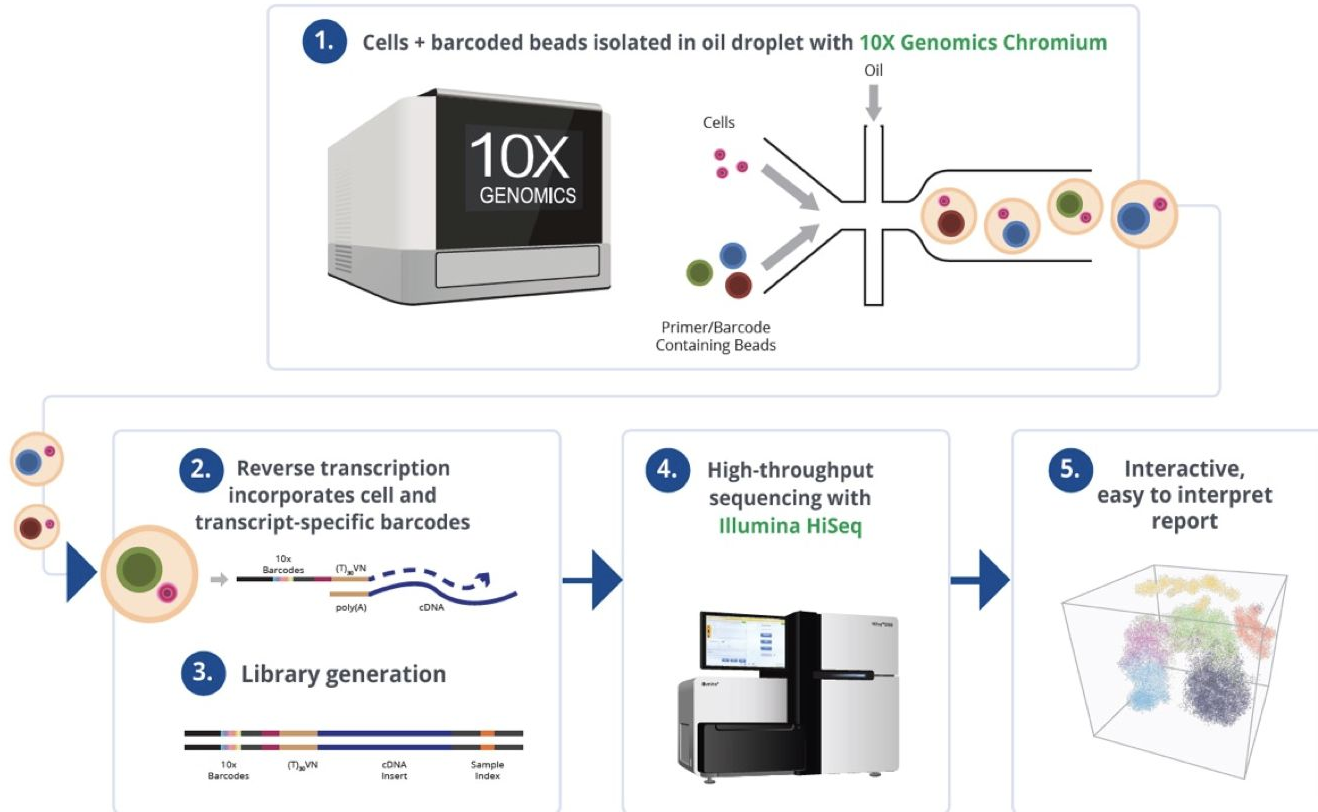
► Single-cell sequencing



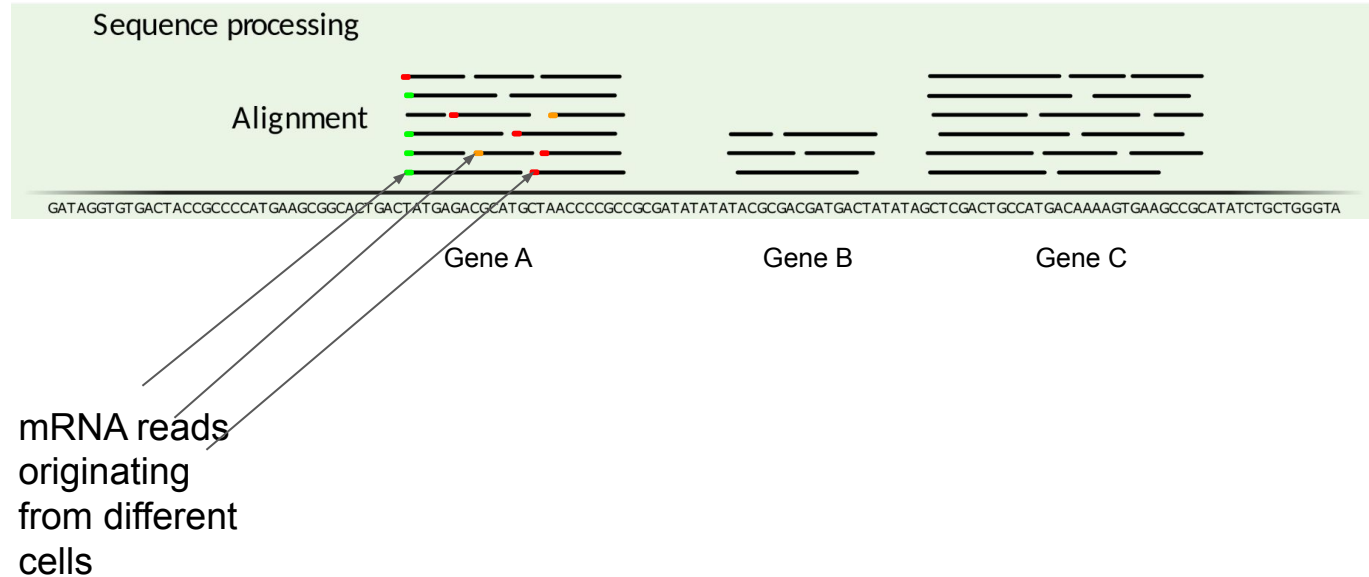
Single cell sequencing



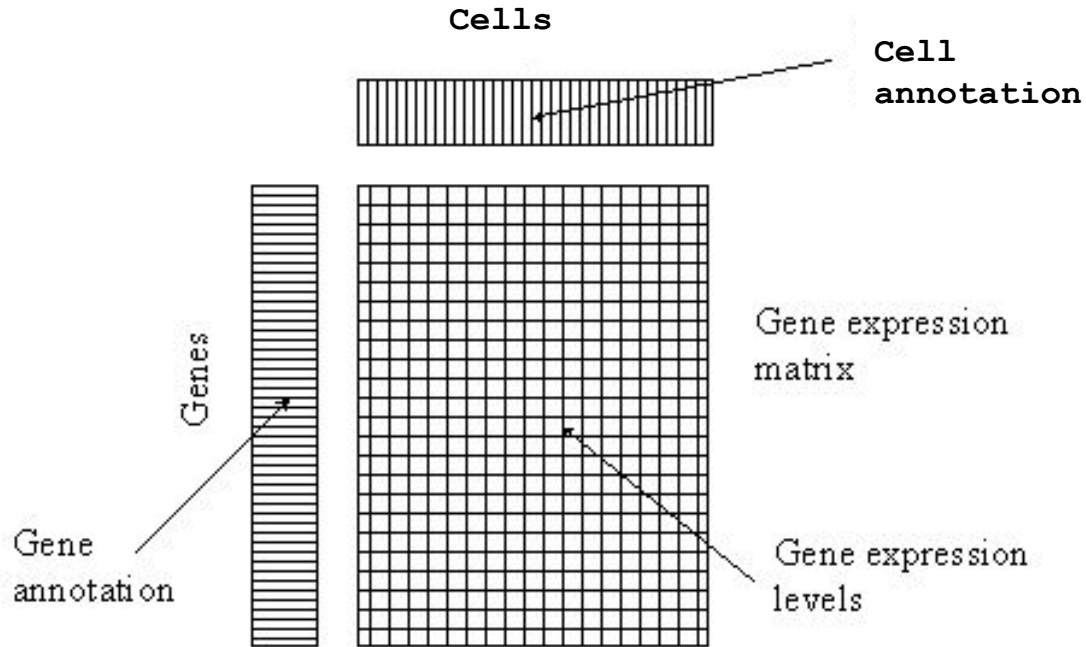
Single Cell RNA-seq: Easy as 1,2, 3, ... 5



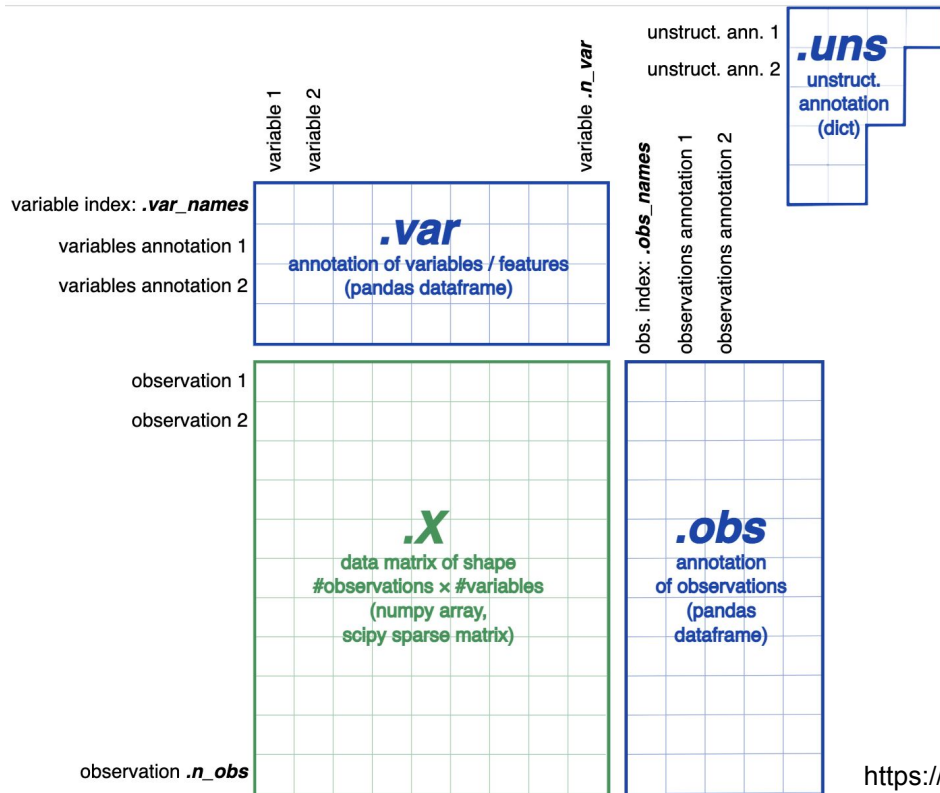
Single cell sequencing alignment and gene count



Cell-gene matrix

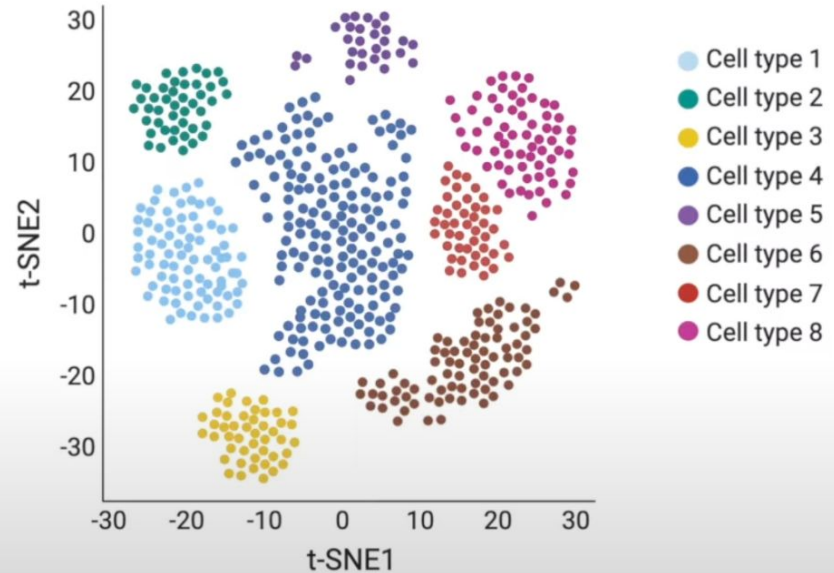


Annotated data object



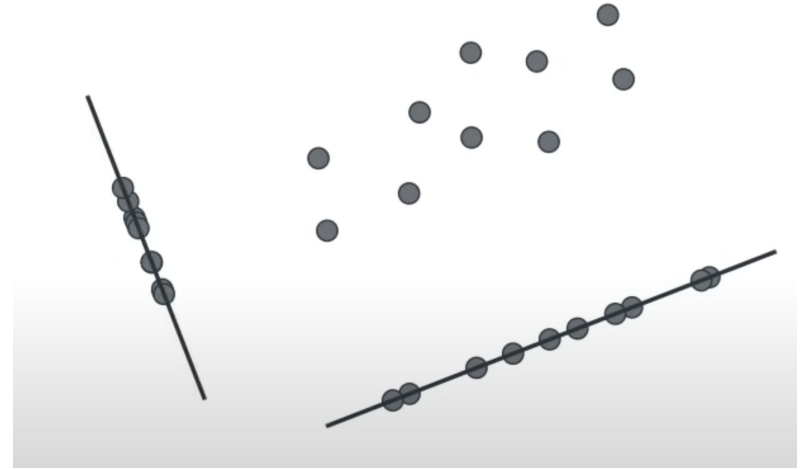
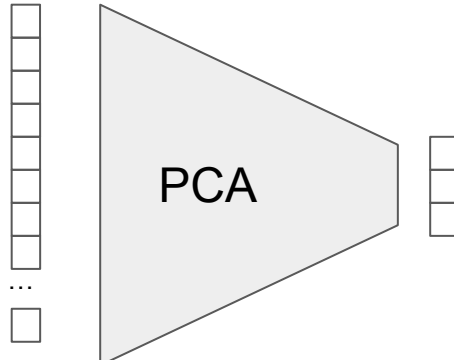
Cell-gene matrix

	Cell 1	cell 2	cell 3	...	cell X
Gene a	3	5	6	...	3
gene b	3	5	3	...	2
gene c	5	6	5	...	4
gene d	5	6	7	...	8
...	
gene z	7	8	4	...	3



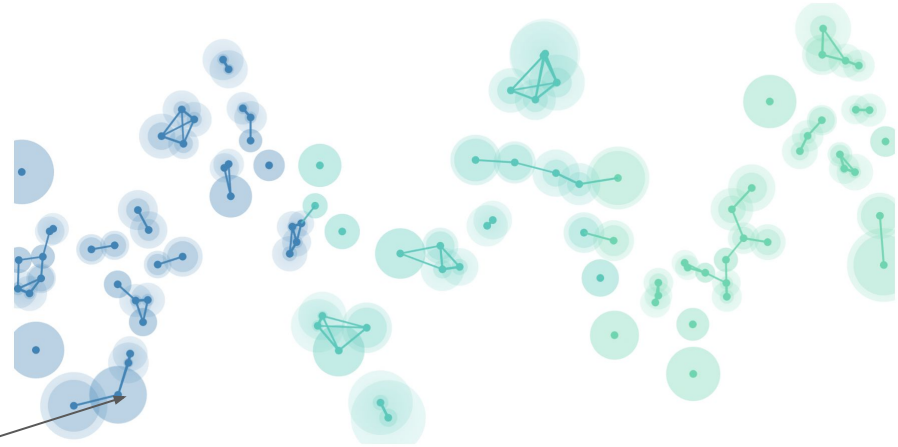
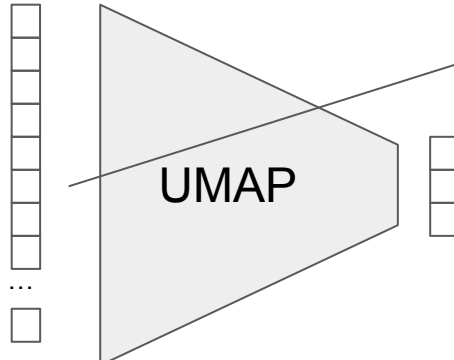
Latent (low-dimensional) representation of data

- Principal Component analysis - linearly transforming the data into a new coordinate system where (most of) the variation in the data can be described with fewer dimensions than the initial data

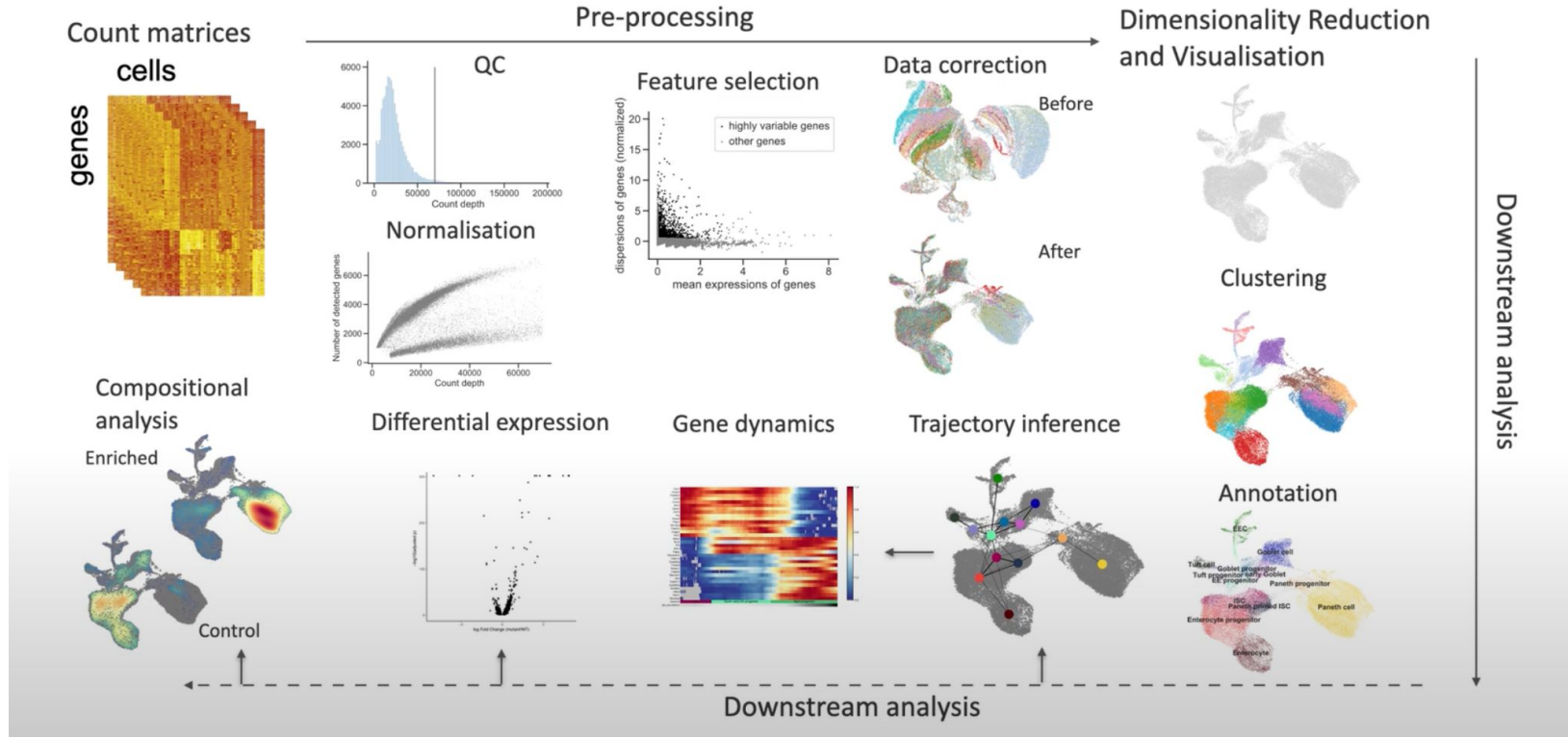


Latent (low-dimensional) representation of data

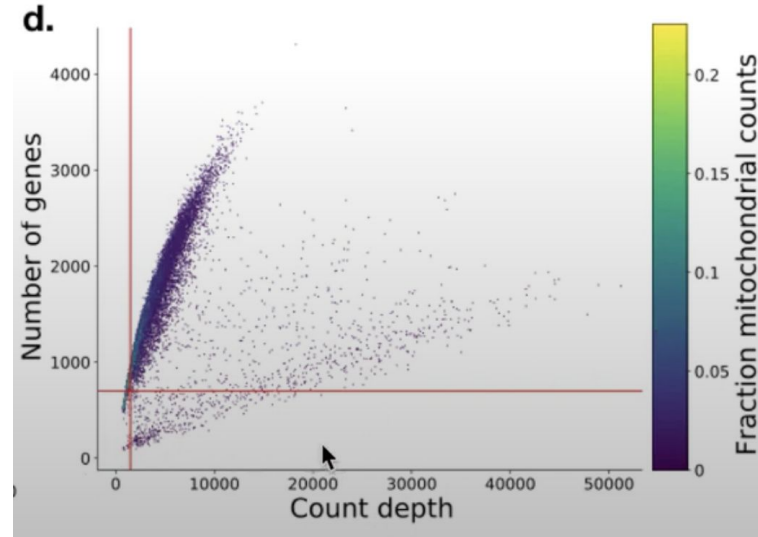
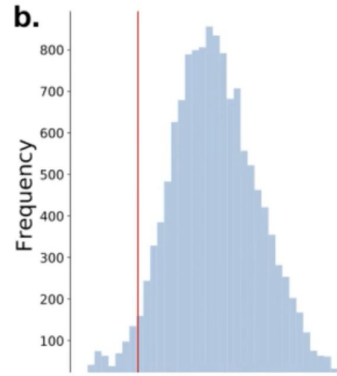
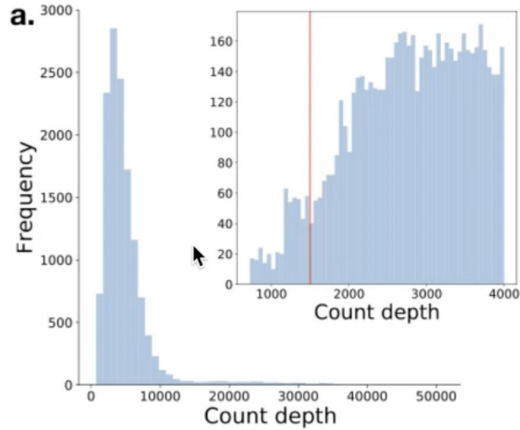
- Uniform Manifold Approximation and Projection (UMAP) tends to better preserve the global structure of the data when projecting from high to low dimensions



Single-cell RNA downstream analysis workflow

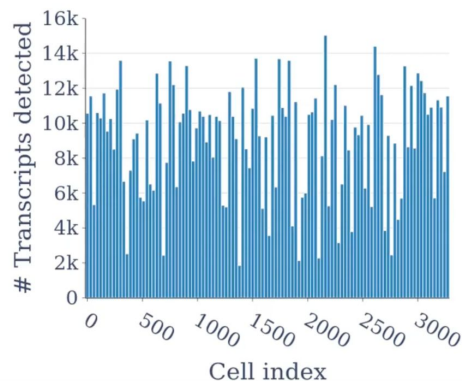


Preprocessing

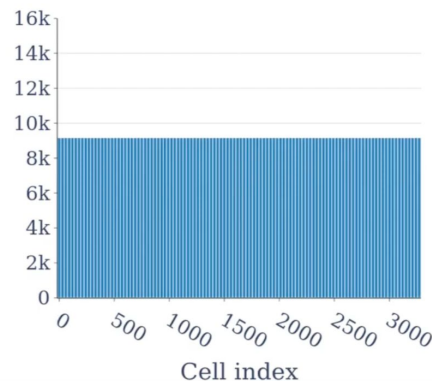


Normalisation

- Gene length might affect the number of captured reads
- Scaling
- Transformation
 - Log
 - Square root
 - Pearson residual (scTeansform)



Scaling



Raw data

	Cell Type A	Cell Type B	Δ
Gene 1	1	2	1
Gene 2	100	200	100

Log₂ transform

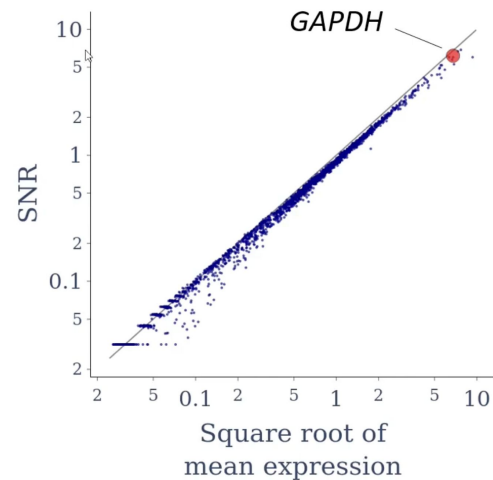
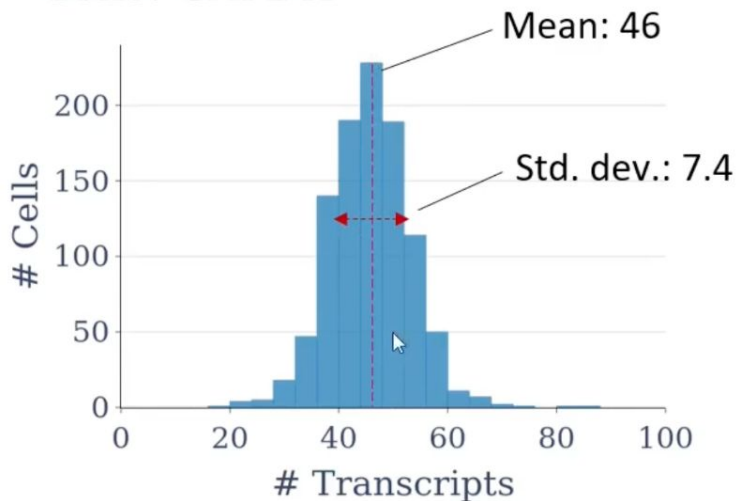
Cell Type A	Cell Type B	Δ
0	1	1
6.64	7.64	1

Square root transform

Cell Type A	Cell Type B	Δ
1	1.41	0.41
10	14.1	4.1

Normalisation

Gene: *GAPDH*



⇒ Let's quantify the measurement accuracy for *GAPDH* using the **signal-to-noise ratio (SNR)**:

$$SNR = \frac{\overset{\text{mean}}{\mu}}{\underset{\text{standard deviation}}{\sigma}} = \frac{46}{7.4} = 6.2$$

Normalisation - Pearson residual

- Simple transformations $\longrightarrow y_{ij} = f(x_{ij})$
 - Log transform
 - Square root transform

- Pearson residuals $\longrightarrow y_{ij} = w_j * x_{ij}$ $w_i = \frac{1}{\sqrt{\mu_i}}$

\Rightarrow Instead of transforming each measurement individually, Pearson residuals apply a weight to **all** measurements of a gene.

\Rightarrow This makes it so that each gene contributes to the analysis according to **how much evidence** there is that it is non-uniformly expressed.

\Rightarrow This favors genes that are expressed in **only a small fraction of cells**.

Pearson residuals

	Cell Type A (50%)	Cell Type B (50%)	Δ
Gene 1	0.816	1.63	0.814
Gene 2	8.16	16.3	8.14

Raw data

	Cell Type A (50%)	Cell Type B (50%)	
		Subtype 1 (48%)	Subtype 2 (2%)
Gene 1	0	8	8
Gene 2	0	0	4.5

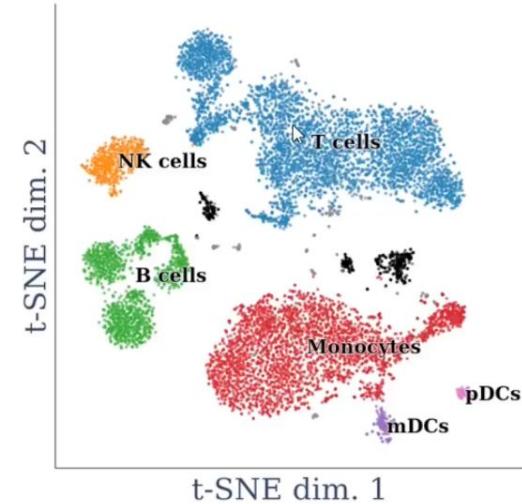
Pearson residuals

	Cell Type A (50%)	Cell Type B (50%)	
		Subtype 1 (48%)	Subtype 2 (2%)
	0	4	4
	0	0	15

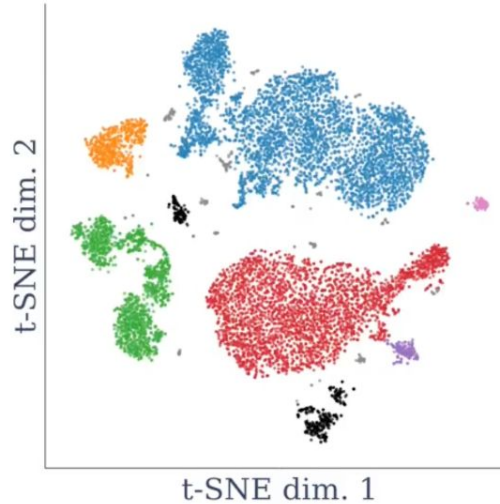
Normalisation - Log, square root and Pearson residual

A real-world comparison

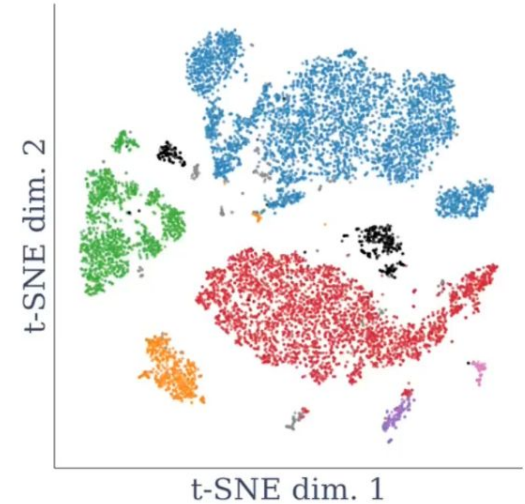
$y = \ln(x + 1)$ (Log)



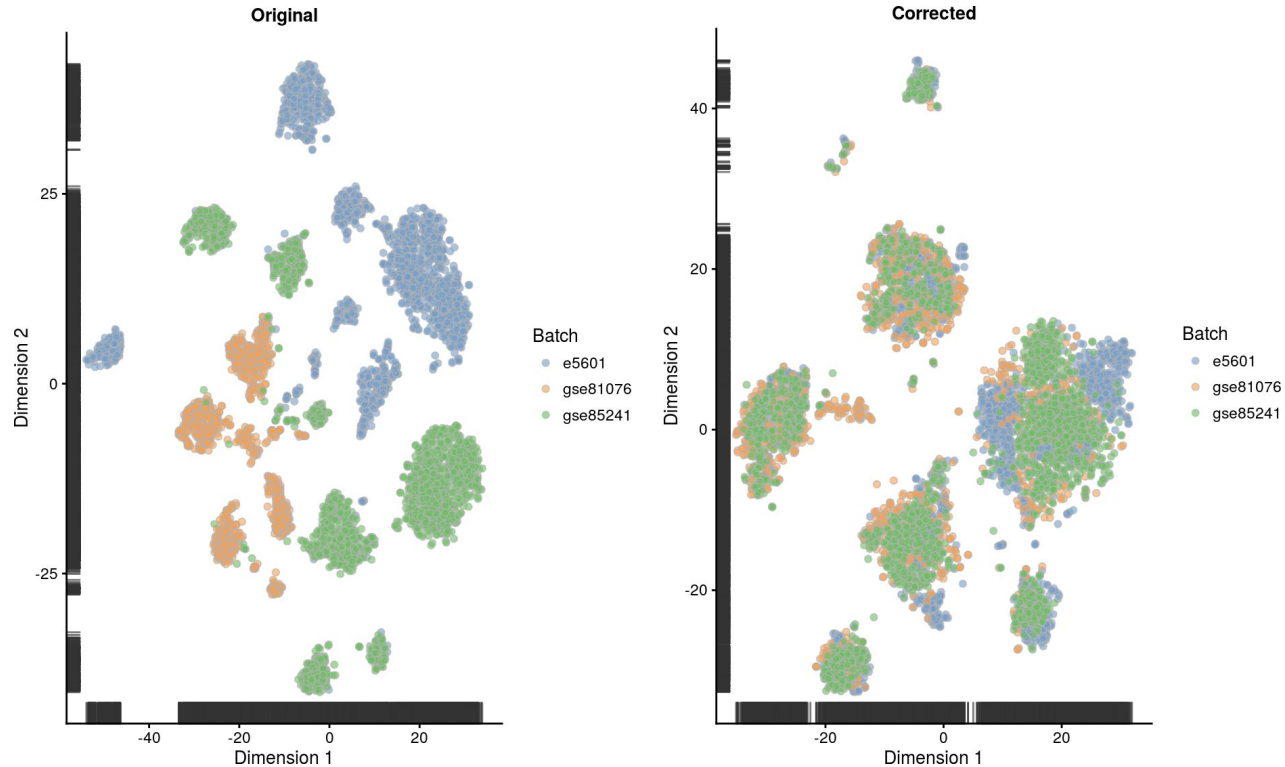
$y = \sqrt{x} + \sqrt{x + 1}$ (Freeman–Tukey)



Pearson residuals

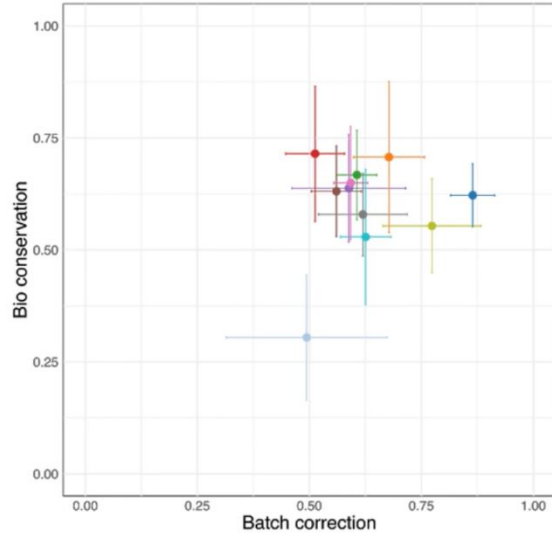


Batch effects correction











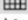




t-SNE plots of the pancreas datasets, before and after MNN correction. Each point represents a cell and is coloured by the batch of origin.


Batch effects correction




- BBKNN (graph, HVG, unscaled)
- Scanorama (embed, HVG, scaled)
- scVI (embed, HVG, unscaled)
- Conos (graph, HVG, unscaled)
- Scanorama (gene, HVG, scaled)
- ComBat (gene, HVG, unscaled)
- MNN (gene, HVG, scaled)
- Harmony (embed, HVG, unscaled)
- Seurat v3 (gene, HVG, unscaled)
- trVAE (embed, HVG, unscaled)
- LIGER (embed, HVG, unscaled)

Method					RNA					Simulations		Usability	Scalability	
Rank	Name	Output	Features	Scaling	Pancreas	Lung	Immune (human)	Immune (human/mouse)	Mouse brain	Sim 1	Sim 2	Usability	Time	Memory
1	BBKNN	 HVG	-	1	1		3	1				3	2	2
2	Scanorama	 HVG	+	3	2	1	1			1				
3	scVI	 HVG	-				2	3						1
4	Conos	 HVG	-		3	2					3			
5	Scanorama	 HVG	+							3				
6	ComBat	 HVG	-					2					1	
7	MNN	 HVG	+											3
8	Harmony	 HVG	-			3					2	1	3	
9	Seurat v3	 HVG	-	2						2	1	2		
10	trVAE	 HVG	-											
11	Unintegrated	 FULL	-											
12	LIGER	 HVG	-											

 gene

 embed


 graph


+


 scaled


-

 unscaled









↑ 1

12 ↓

Output

- gene
- embed
- graph

Scaling

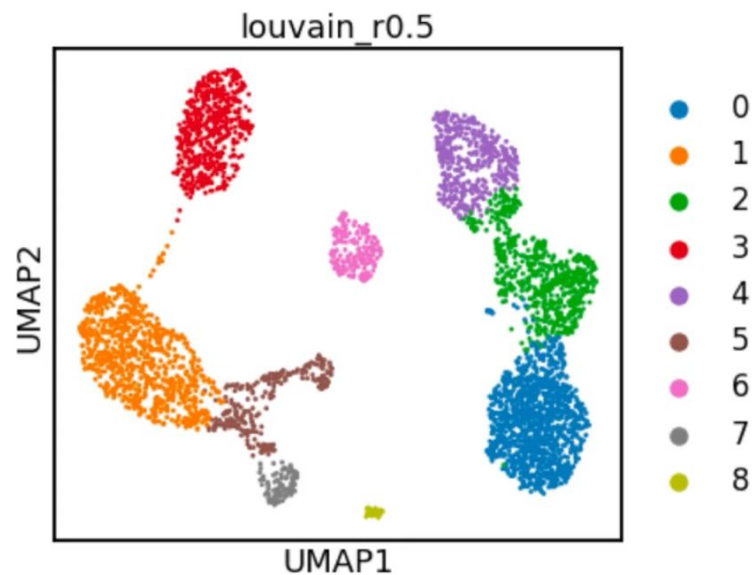
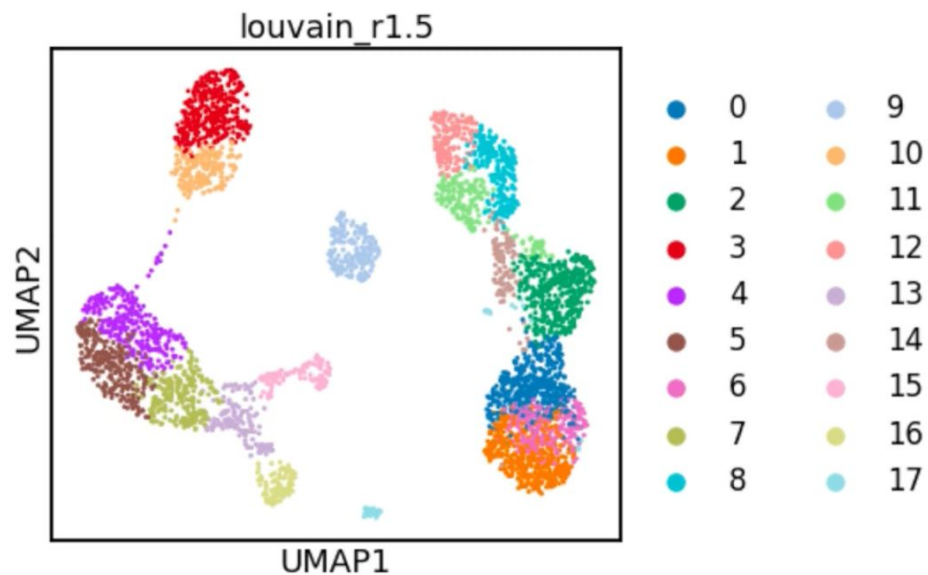
- +
-

Ranking



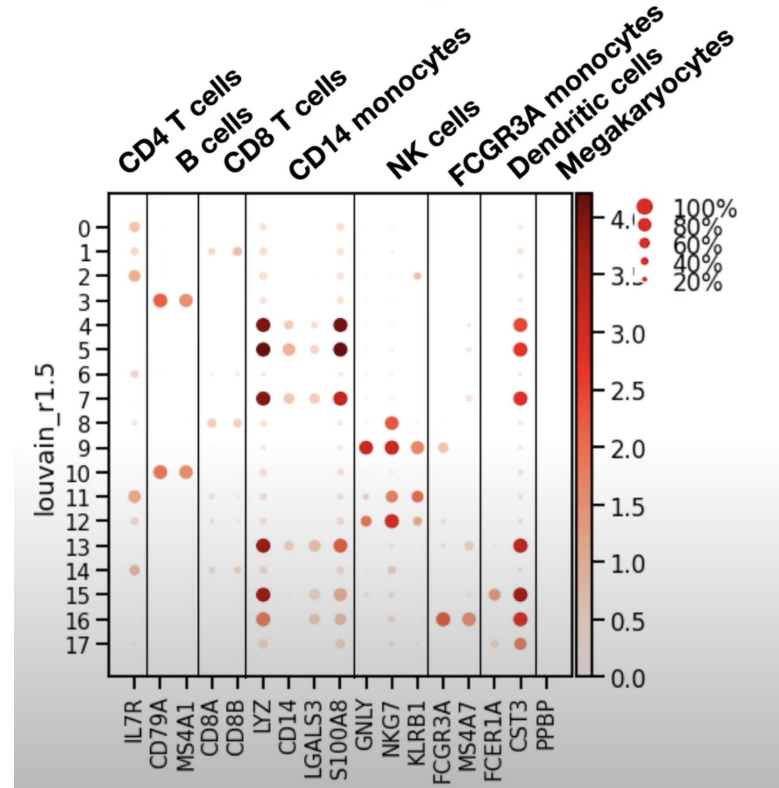
Clustering

Louvain clustering



Clustering and cell type annotation

- Marker genes - genes with statistically significant variation in the specific cluster comparing to the rest of the cells

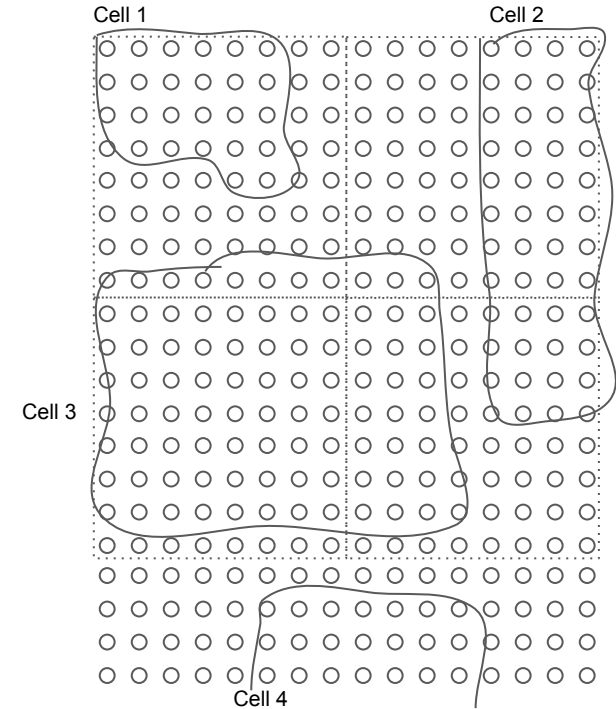


Scanpy library

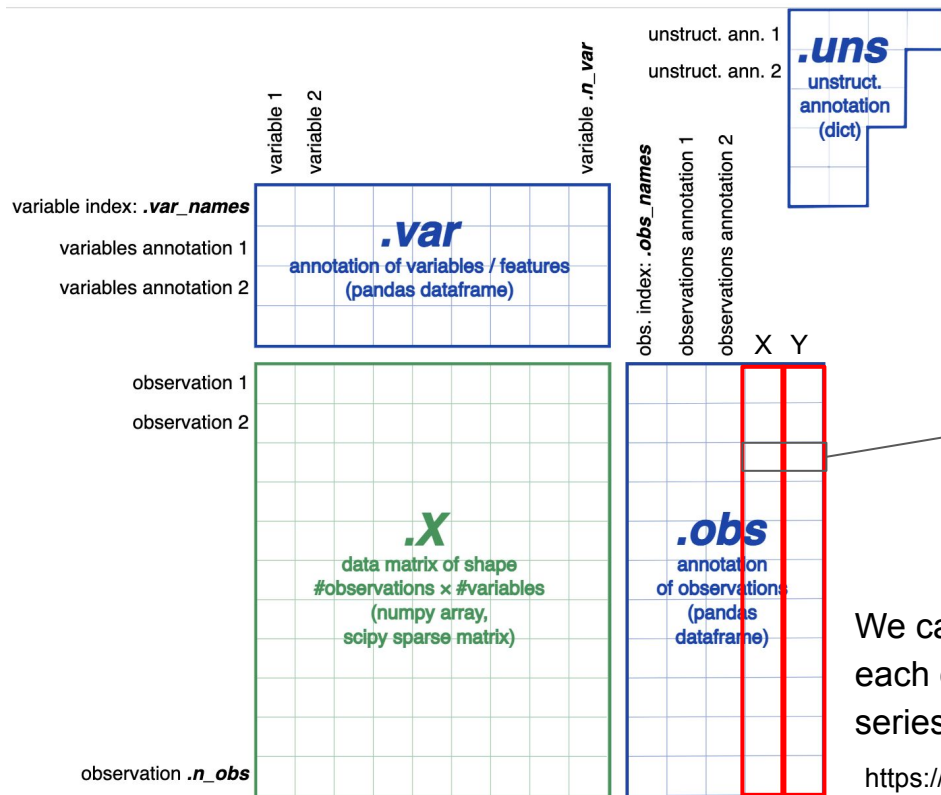
- Scalable toolkit for analyzing single-cell gene expression data
- Annotated data object
- Tutorial: [Preprocessing and clustering 3k PBMCs](#)

Spatial transcriptomics

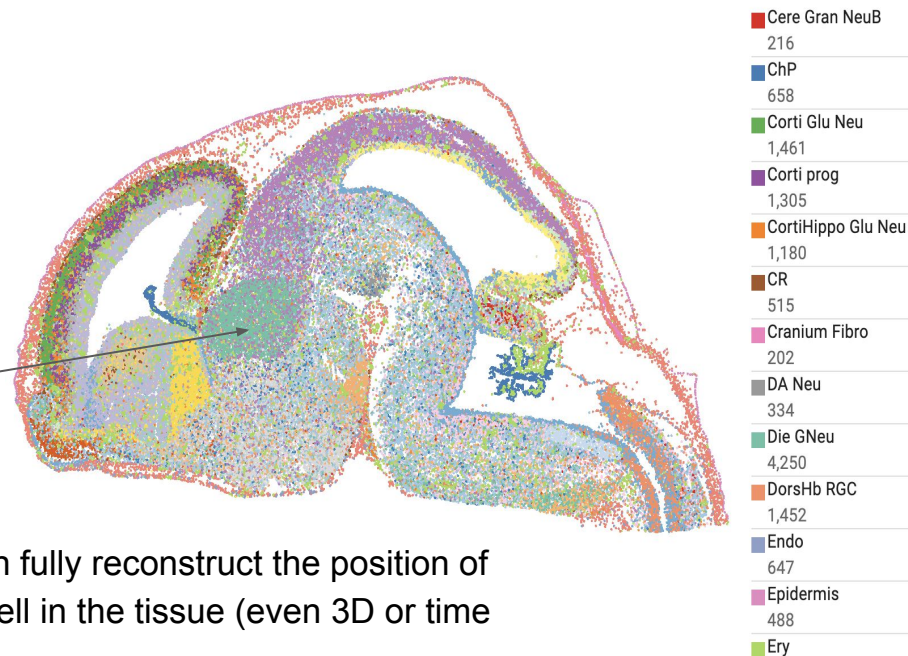
- What if we obtain for every cell beside expression its spatial coordinates
- Each spot is 250 nm in diameter and the center-to-center distance between neighboring spots is 500 nm



Annotated data object



Mouse embryo brain, day 16

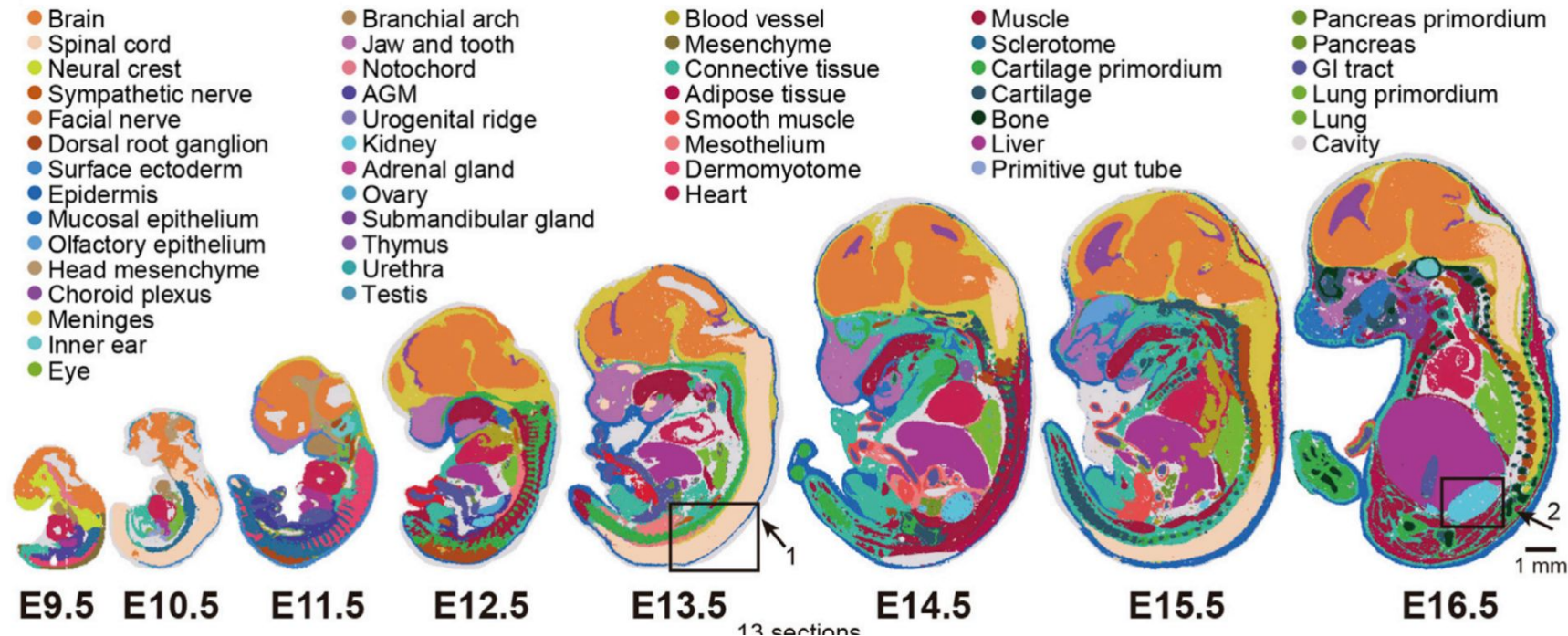


We can fully reconstruct the position of each cell in the tissue (even 3D or time series)

<https://anndata-tutorials.readthedocs.io/en/latest/getting-started.html>

Stereo-Seq: Spatiotemporal transcriptomic atlas of mouse organogenesis

A



Spatial transcriptomics data analysis

- [Stereopy](#) library
- [Quick start tutorial](#)