



Single Cell RNA-seq

Preprocessing and QC

Hubert Rehrauer

(with slides from Ge Tan)



University of
Zurich^{UZH}

ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich



Outline

- Biological need, technical solution, noise characteristics
- QC/Issues: mitochondrial genes, ribosomal genes, few genes, emptyDrops, doublets, dropouts
- normalization
- dimensionality reduction
 - Matrix factorization
 - graph-based (t-SNE, UMAP)
 - Autoencoder
- batch correction, batch integration
- clustering



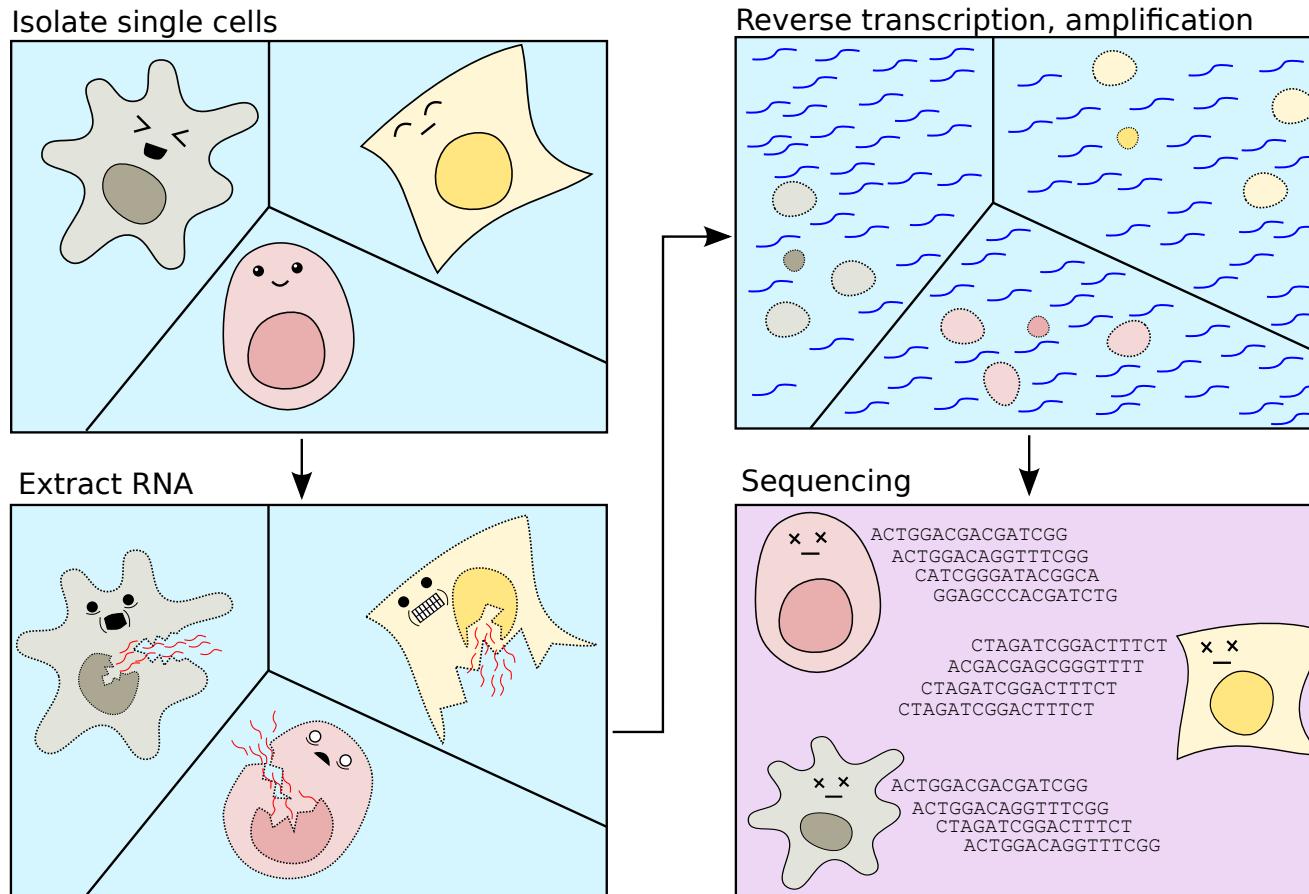
University of
Zurich UZH

10
01
101

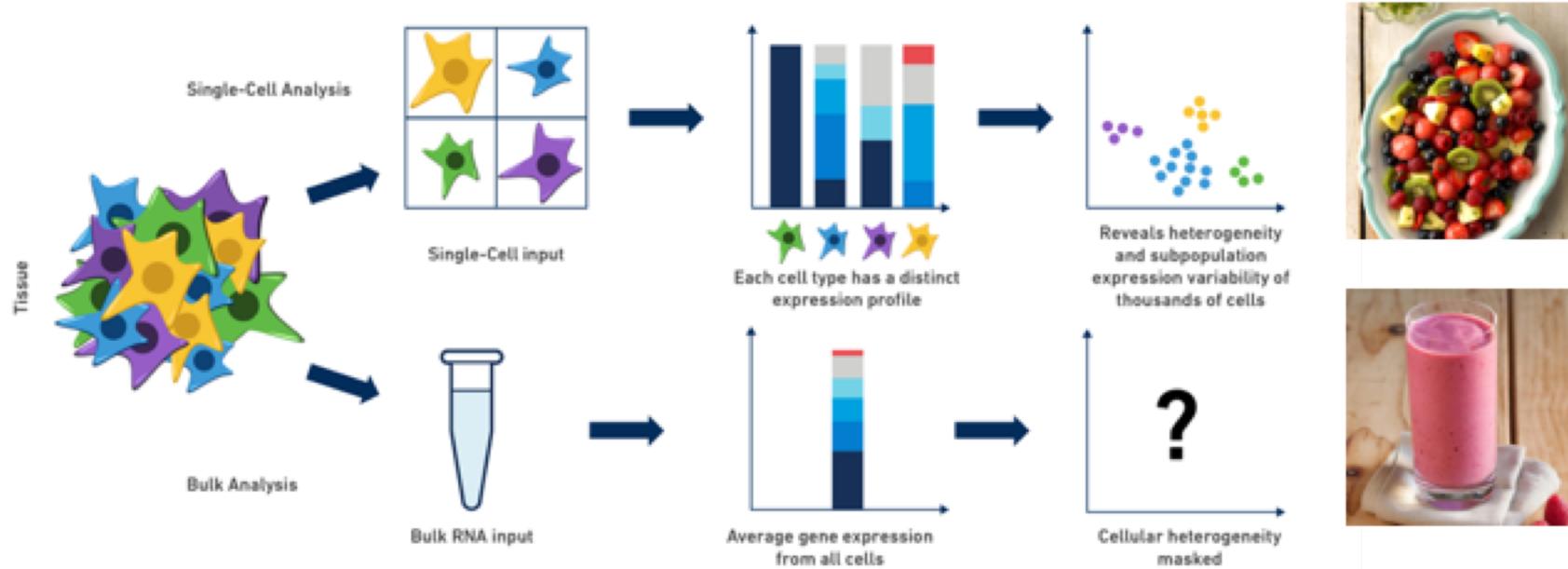
functional genomics center zurich

010 01
101 10
010 01
01 1

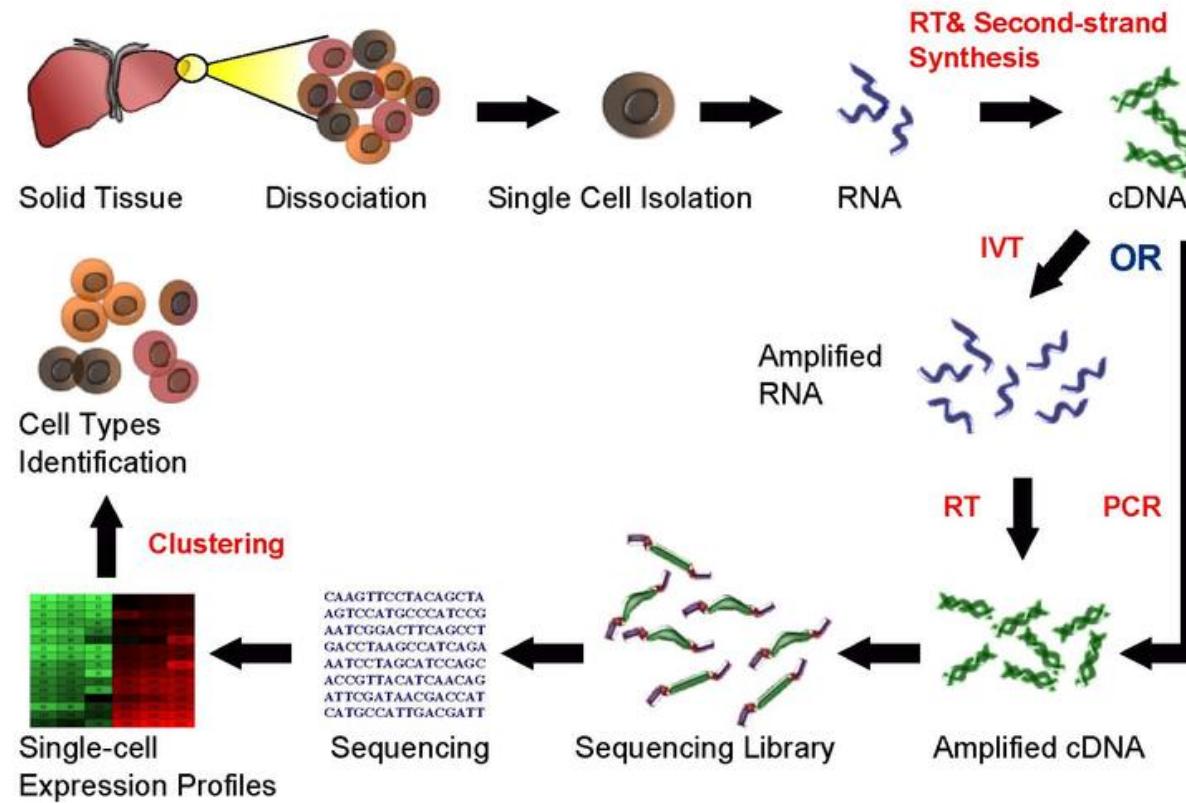
f g c z
10 0
01 1



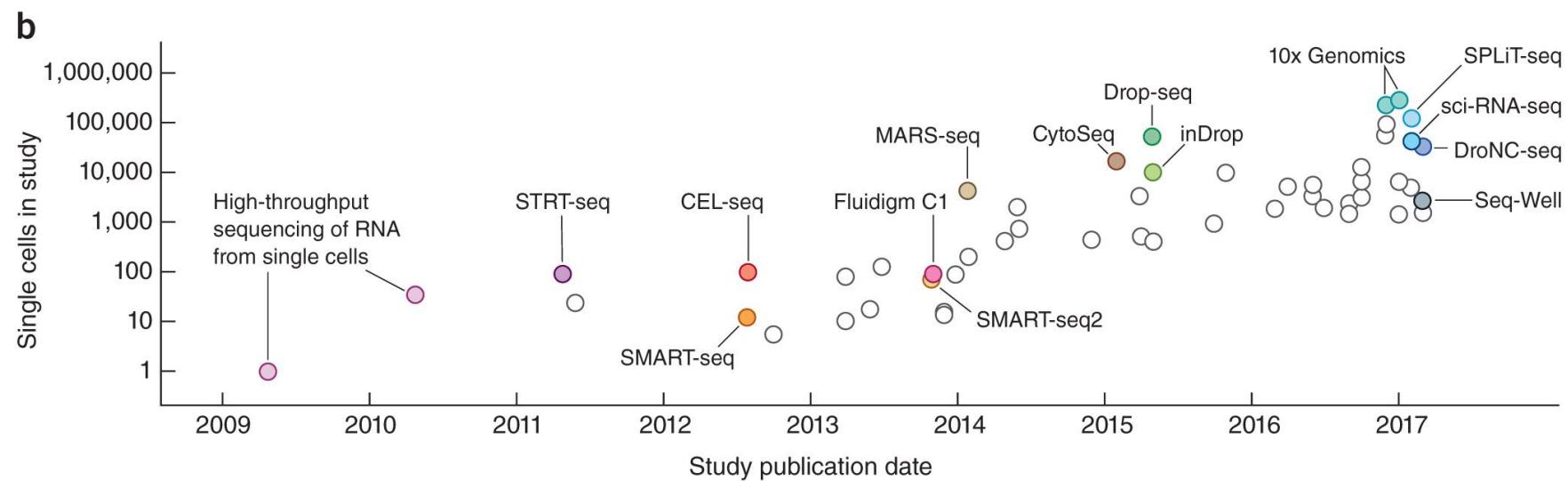
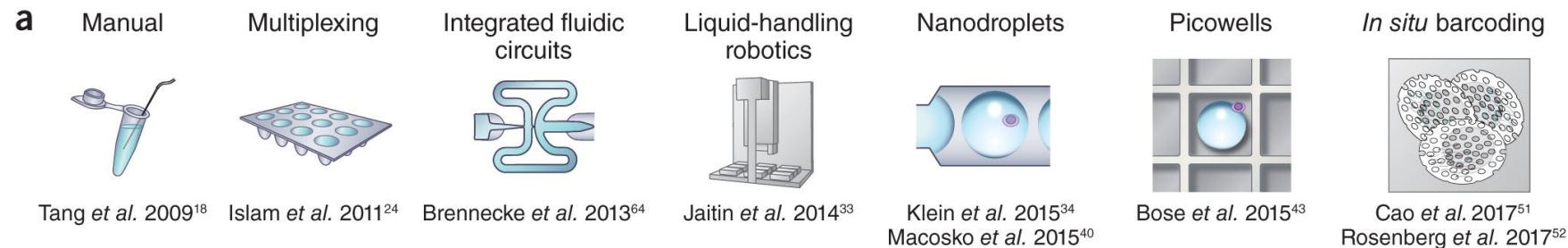
Adapted from Aaron Lun



Single Cell RNA Sequencing Workflow



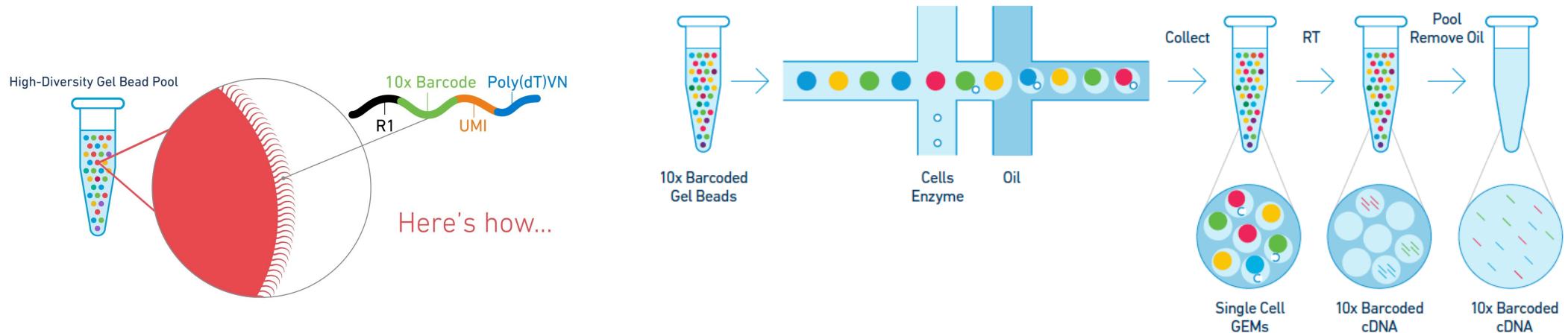
Technologies

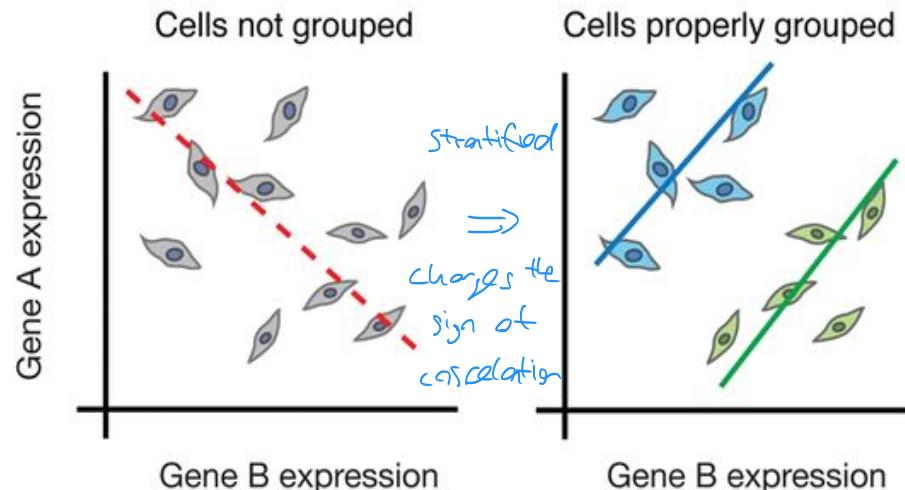




Nanodroplets systems – e.g. 10X Genomics

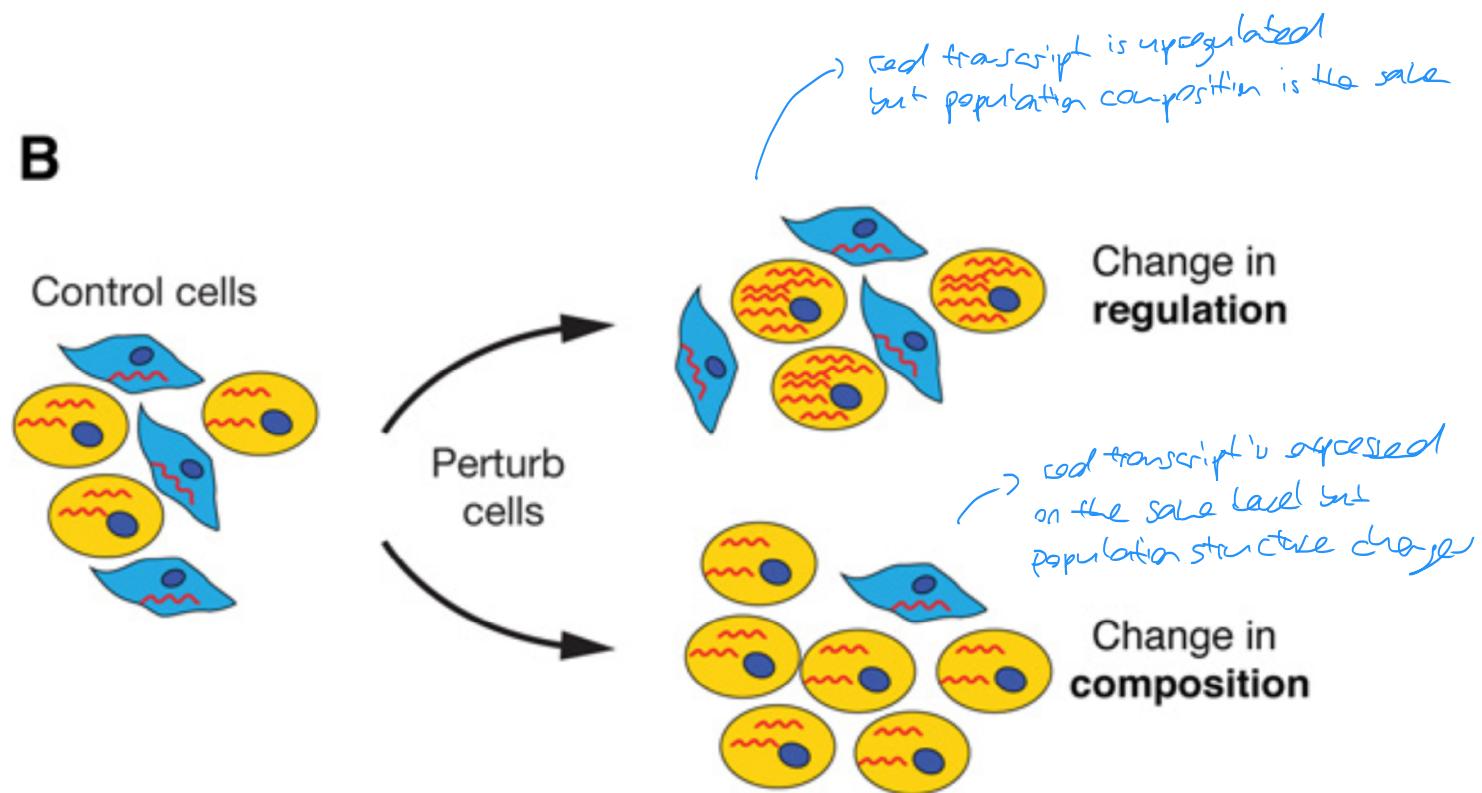
- High throughput (up to 80,000 cells per sample).
- Current standard in the field.
- Restricted to certain cell-types / sizes.





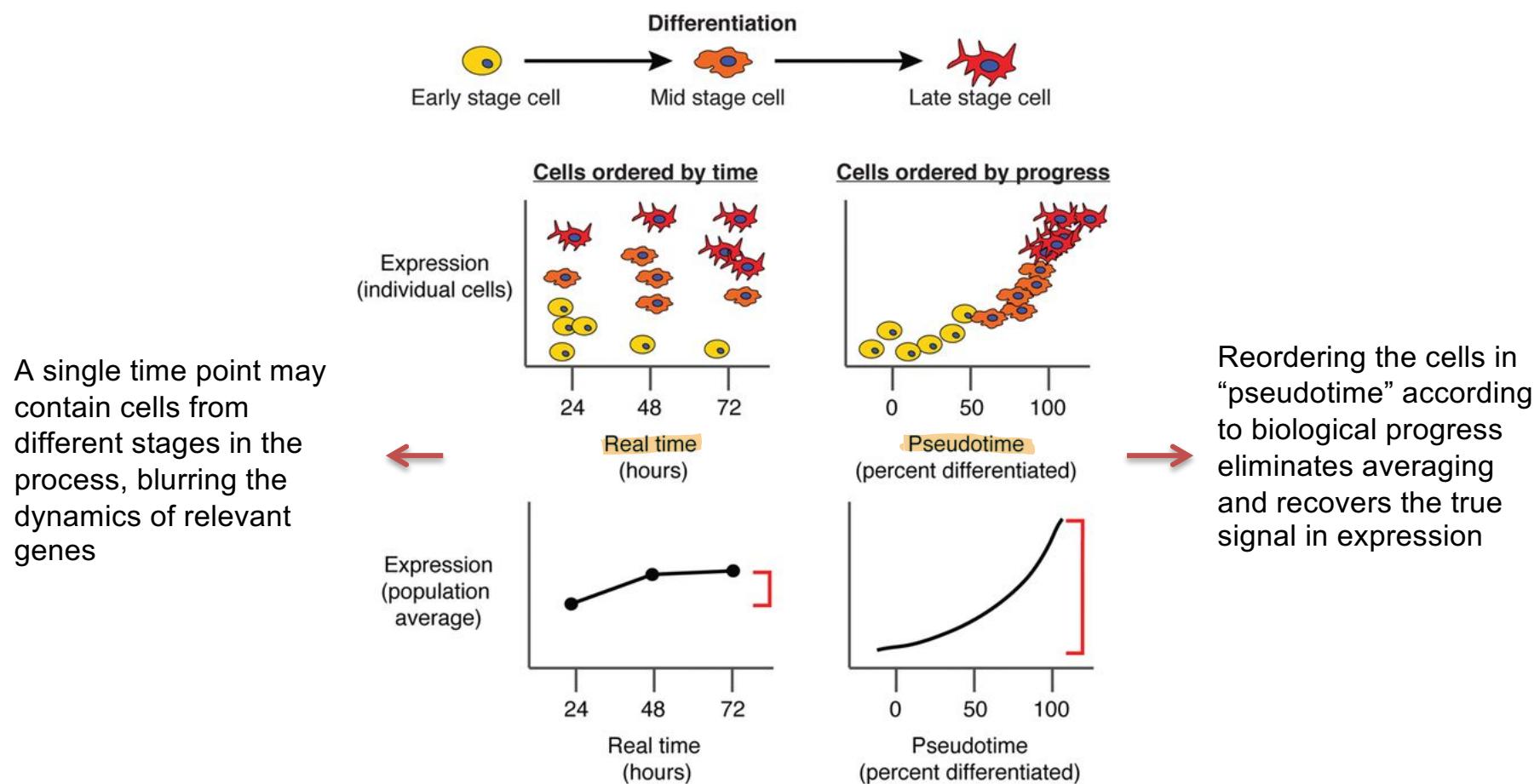
Simpson's Paradox describes the misleading effects that arise when averaging signals from multiple individuals

Trapnell Genome Res. 2015

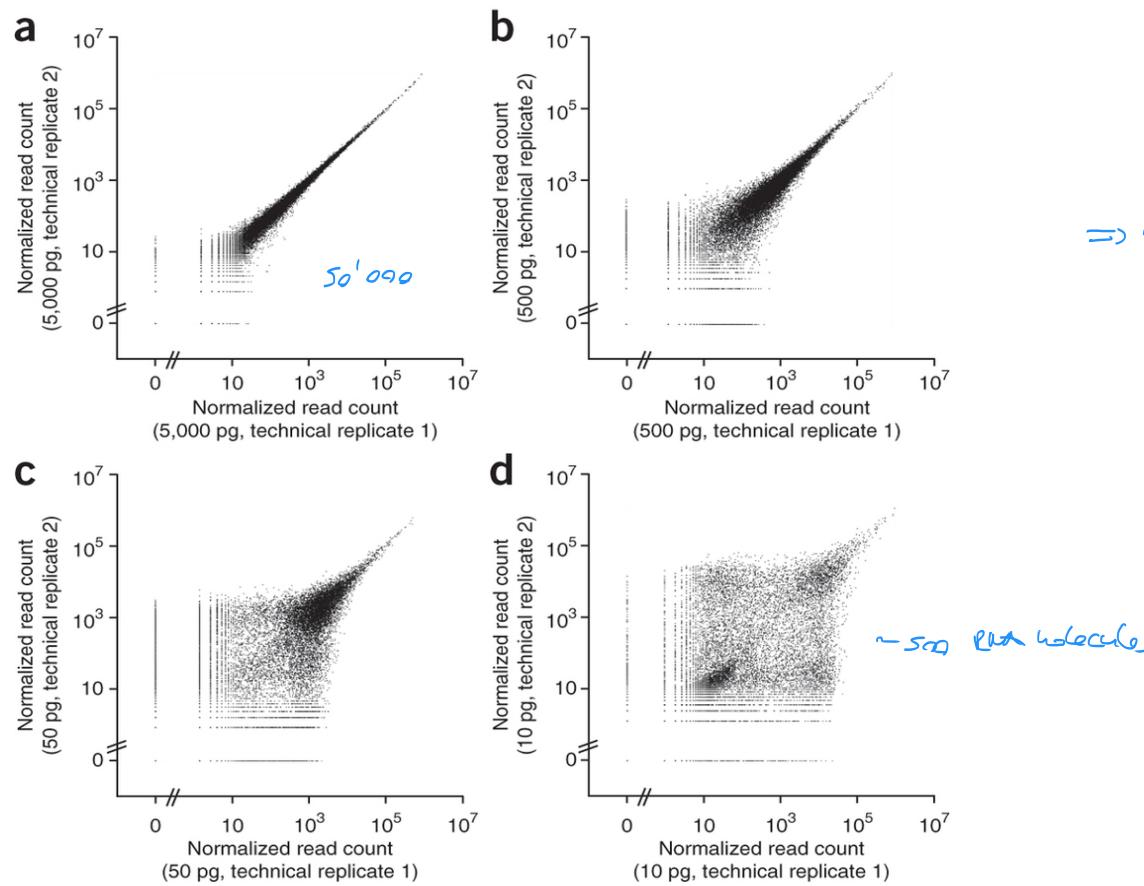
10
01
101010 01
101 10
010 01
010 010
1
0
1
0
0
1
1**B**



Time series experiments are affected by averaging when cells proceed through a biological process in an unsynchronized manner



Challenge in scRNAseq: high technical noise

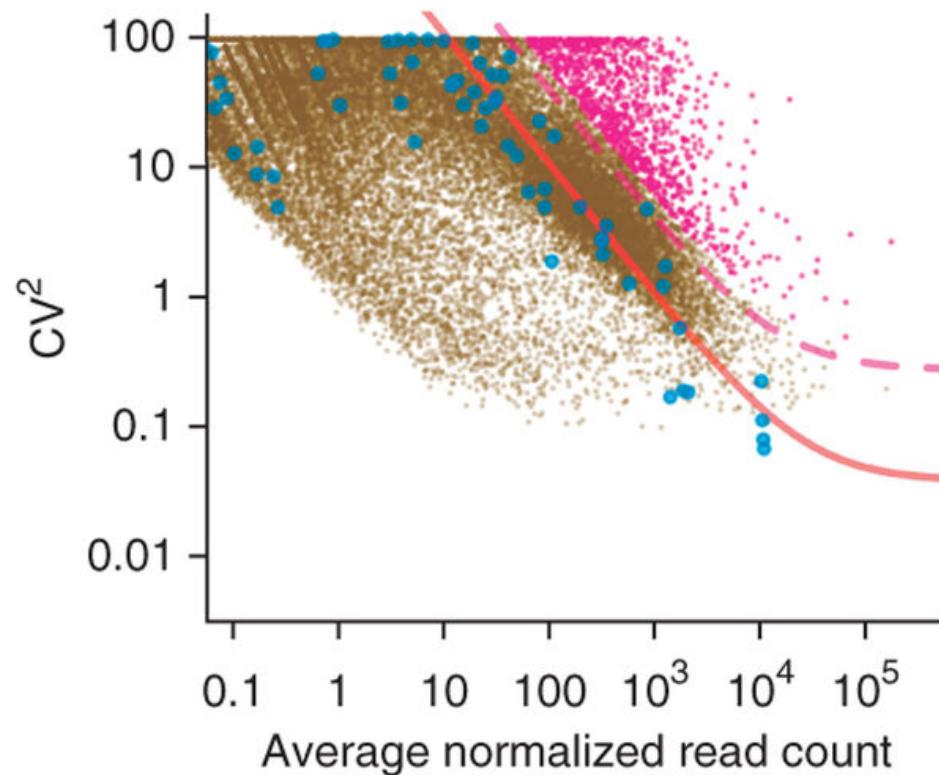


- low input means high amplification and high technical noise

\Rightarrow more noise in low abundance species

Brennecke et al. Nature Methods 2013

Variability of scRNA-seq counts

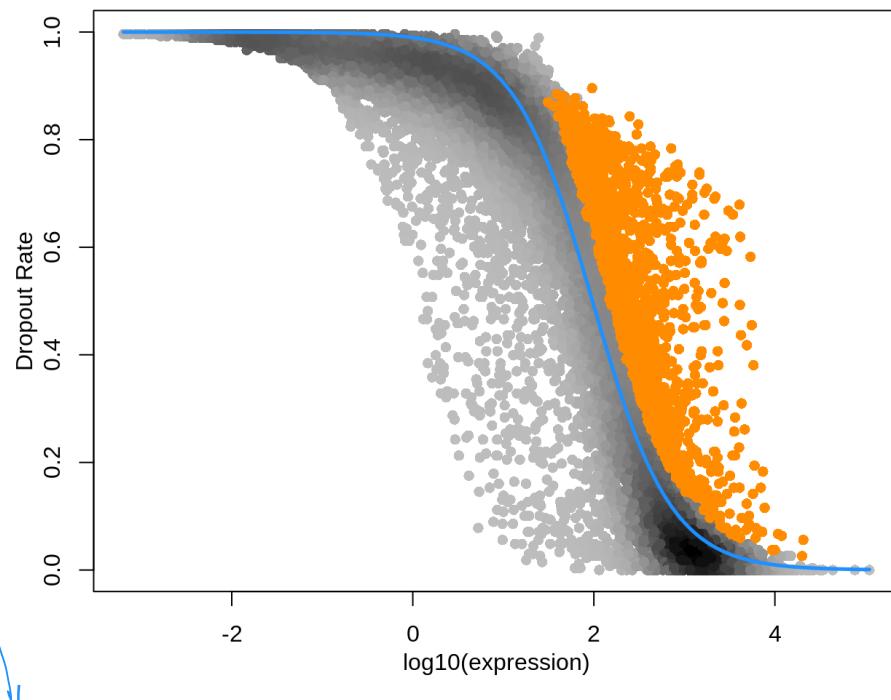


- Blue: spike-ins
 - Pink: genes with potentially differential expression between cells

Brennecke et al. *Nature Methods* 2013

$$CV = \text{st.dev}/\text{mean}$$

Dropouts: Genes not expressed

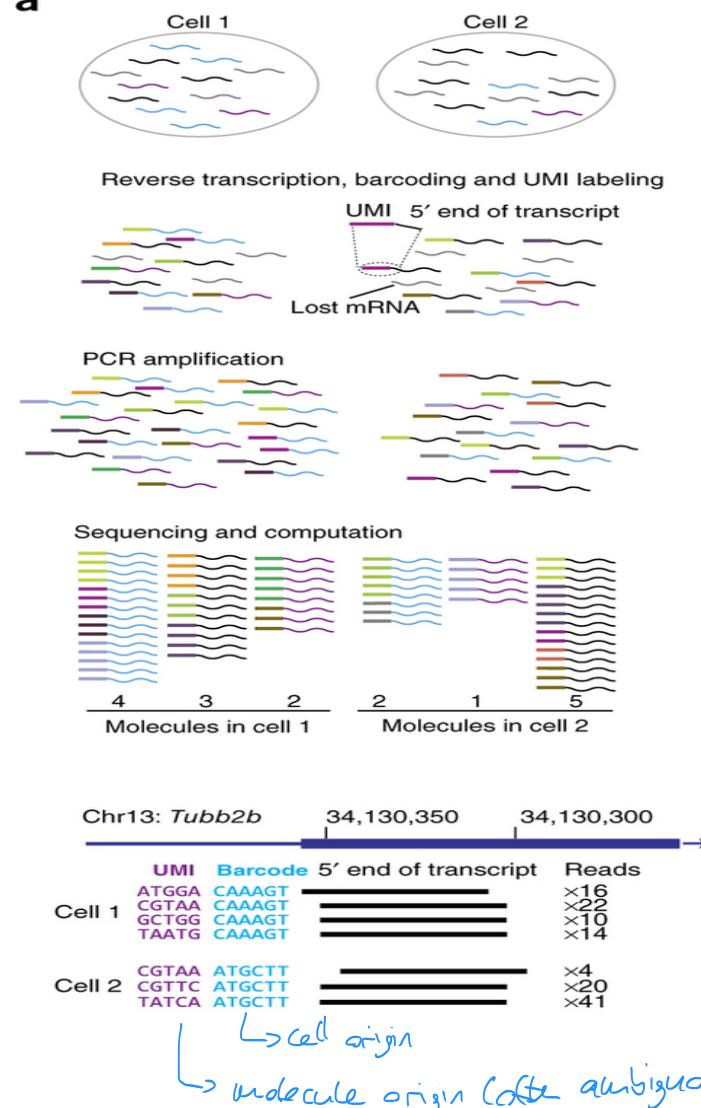


fraction of the cell in
which the individual gene is not
detected.

- Dropout Rate: Fraction of cells where gene is not detected
- genes with low average expression have high dropout rate
- 1000 – 10'000 genes detected per cell (depends on single cell protocol and sequencing depth)
- > 15'000 genes detect in bulk tissue
- Reasons for dropout
 - stochastic expression bursts
 - not enough reads in sequencing – lost by the sampling process
 - library prep failed to capture the transcript

Unique Molecular Identifiers (UMIs)

a

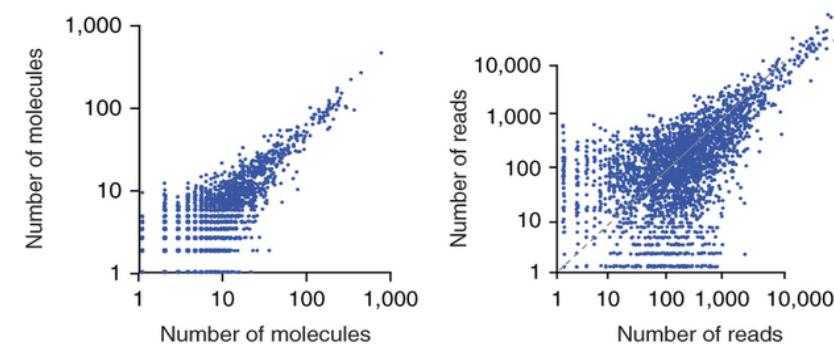


PCR duplicates can be identified

Two reads are considered duplicates if they have the same UMI, the same cell barcode and map to the same gene.

UMIs are only available for 3'-tagging protocols, e.g. the 10X.

Islam et al. Nature Methods 2014



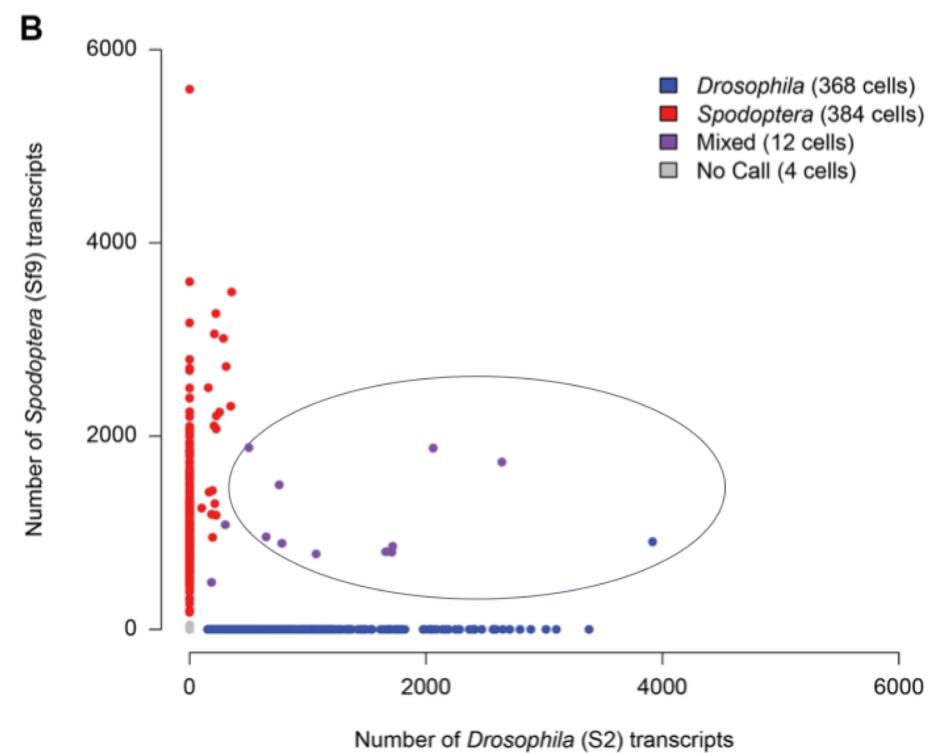
3'-AAAAA ~~~~~ ~246 s' ← not targeted
 UUUU - TTTTT ~~~~~
~~~~~  
 20nt  $\Rightarrow$  only 20nt are sequenced

$$UHI - HIT = \underbrace{\phantom{0000}}_{2001t}$$

⇒ only 200 nt are sequenced / read  
thus the UMI protocol can only be applied  
to 3'-end sequencing of RNA transcripts

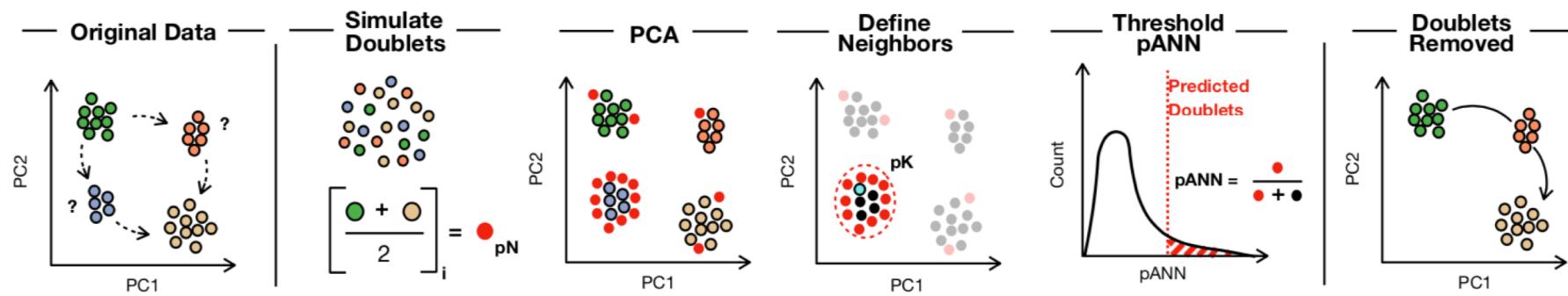
## Quality Control: Doublets

- Barcode collisions: not enough barcodes
- Technical doublets: two cells in the same droplet (10X specs: +1% per 1000 cells)
- Biological doublets: two cells sticking tightly together and form a unit; need to do nuclear single-cell RNA-seq
- Test datasets for doublets consist of cells from two species



10  
01  
101

# Doublet Detection



- see also:
  - <https://github.com/plger/scDblFinder>



## Other Quality Control Metrics

- high content of mitochondrial RNA
  - apoptotic cell
- high content of ribosomal RNA
  - failed library prep
- few reads sequenced or few genes detected
  - **empty drop**
  - failed library prep
  - alternative explanation: cell with low transcriptional activity



University of  
Zurich UZH

ETH zürich

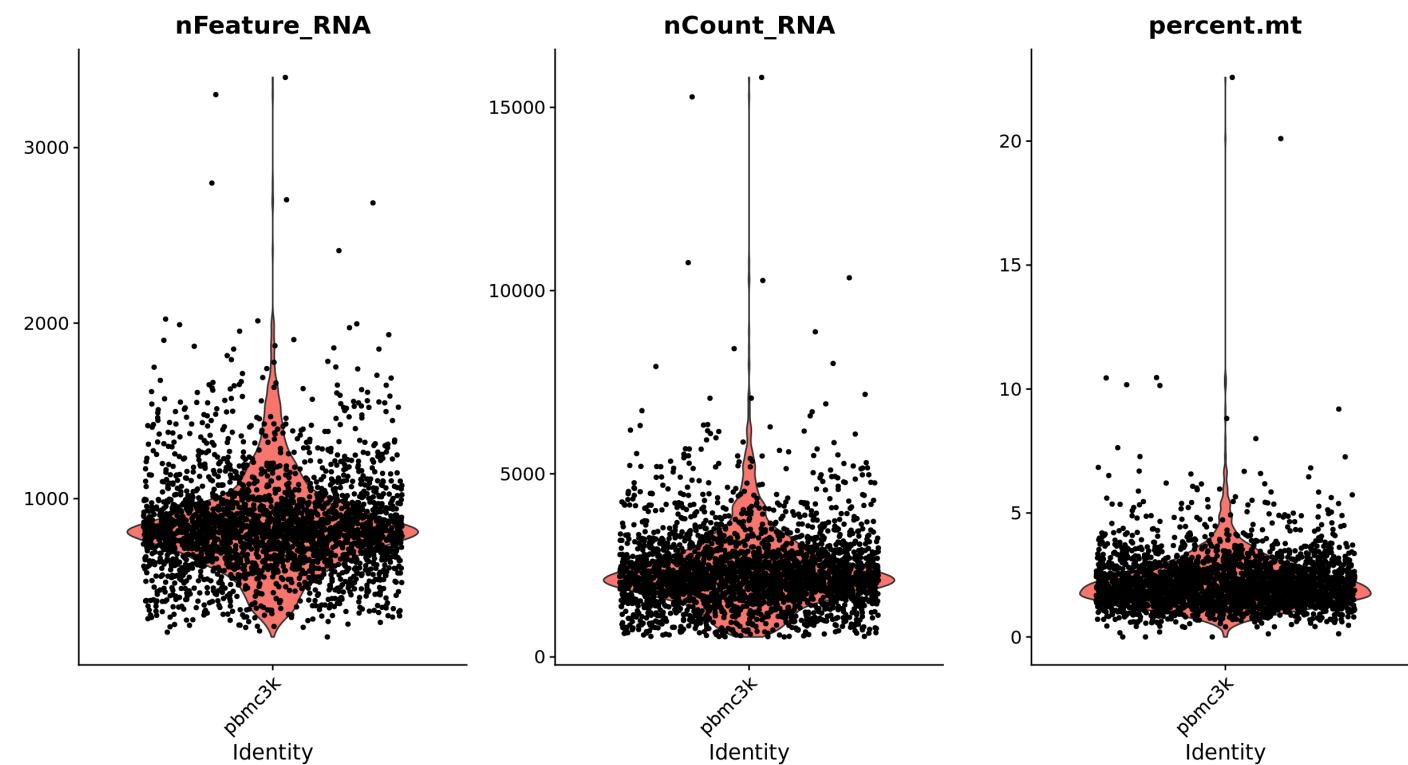
10  
01  
101

functional genomics center zurich

010 01  
101 10  
010 01  
010 01

f g c z  
10 10  
01 01

## QC Filtering Plots

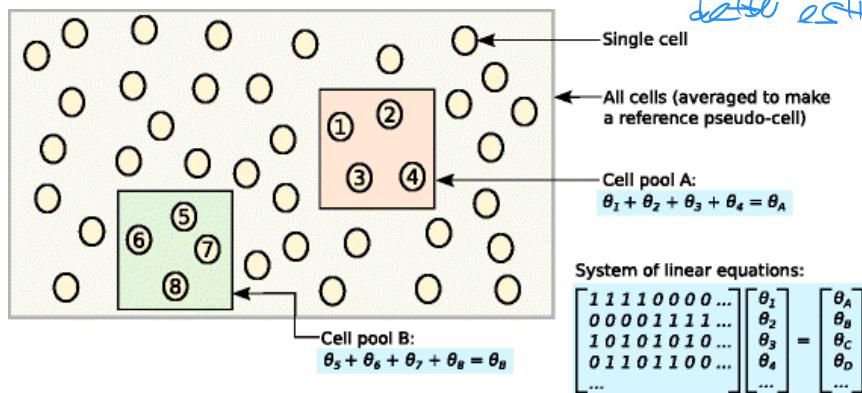


## scran: Normalization using pools of cells

⇒ In bulk RNAseq we normalize by column (i.e condition) (we see expression for next genes)

and since we have bulk most of the values are non-zero

⇒ In single cell RNAseq we have a large fraction of genes with no detected expression, thus it is difficult to estimate a column-wise value for normalization. ⇒ pool cells to get better estimate.



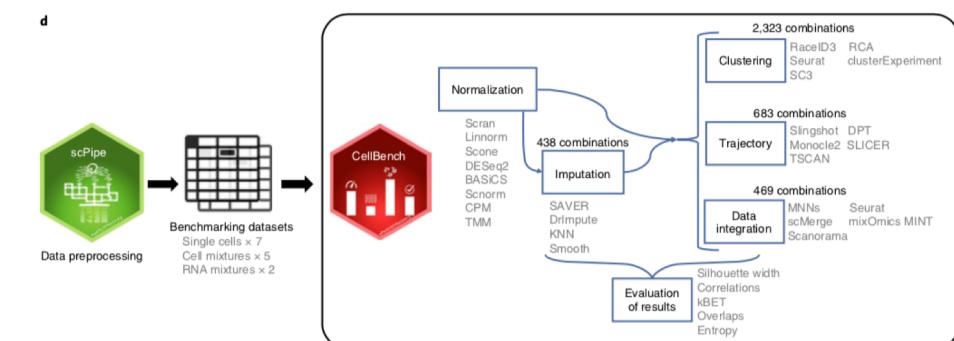
- Define a pool of cells
- Sum expression values across all cells in the pool
- Normalize the cell pool against an average reference, using the summed expression values
- Repeat this for many different pools of cells to construct a linear system
- Deconvolute the pool-based size factors to their cell-based counterparts

## Normalization: Performance comparison papers



# Benchmarking single cell RNA-sequencing analysis pipelines using mixture control experiments

Luyi Tian<sup>ID</sup><sup>1,2\*</sup>, Xueyi Dong<sup>1,3</sup>, Saskia Freytag<sup>1,4</sup>, Kim-Anh Lê Cao<sup>ID</sup><sup>5</sup>, Shian Su<sup>1</sup>, Abolfazl JalalAbadi<sup>15</sup>, Daniela Amann-Zalcenstein<sup>1,2</sup>, Tom S. Weber<sup>ID</sup><sup>1,2</sup>, Azadeh Seidi<sup>6</sup>, Jafar S. Jabbari<sup>6</sup>, Shalin H. Naik<sup>ID</sup><sup>1,2</sup> and Matthew E. Ritchie<sup>ID</sup><sup>1,2\*</sup>

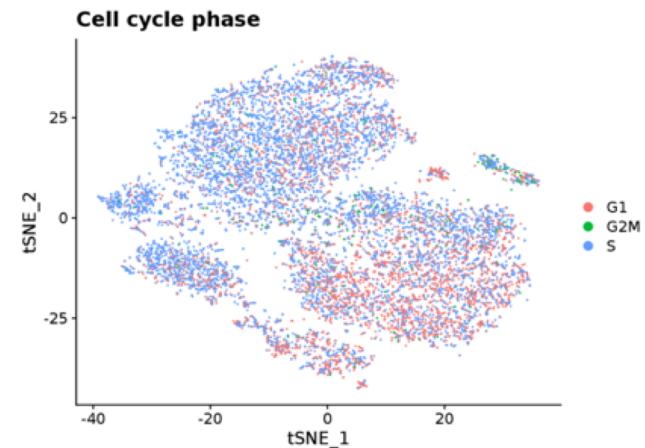




## Normalization: Cell Cycle causes unwanted variation

- Requires normalization considering cell cycle as latent variable

→ cell cycle influences data





10

01

101

010

01

101

10

010

01

10

0

01

1

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

01

0

01

1

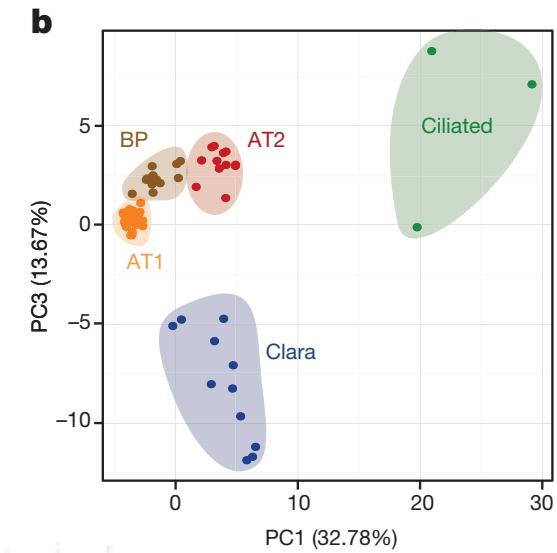
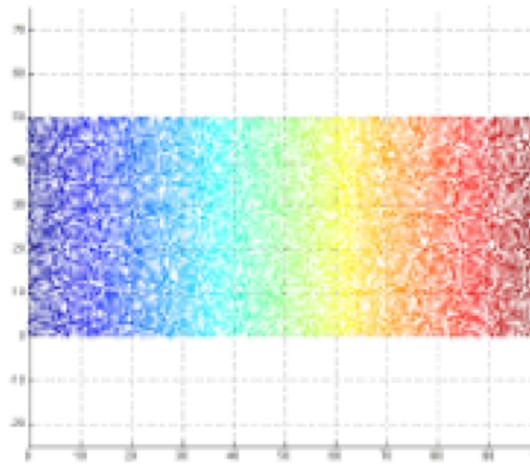
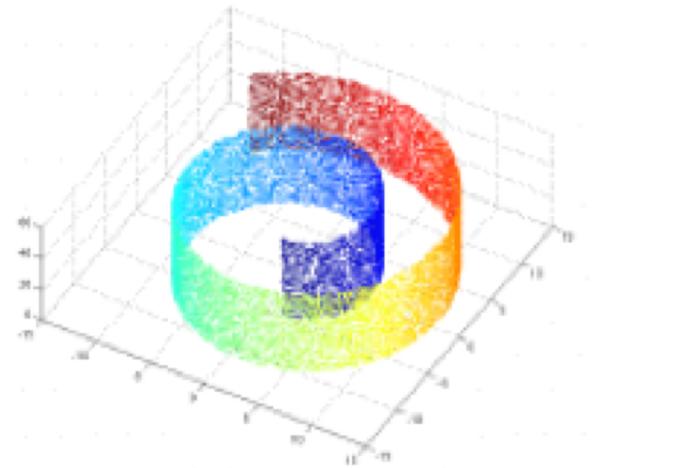
01

0&lt;/



# PCA

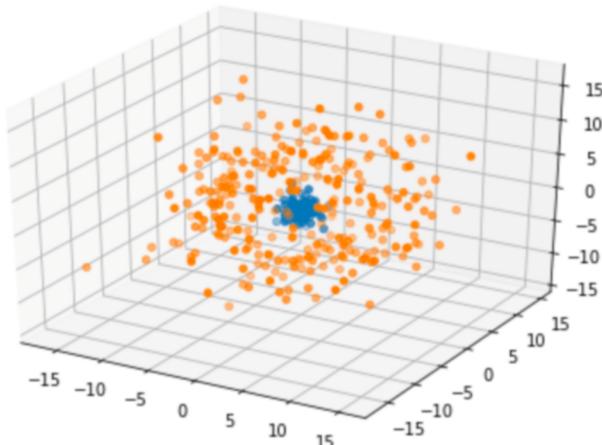
- PCA is simple and efficient
- But can not cope with complex structures because it is a linear projection.



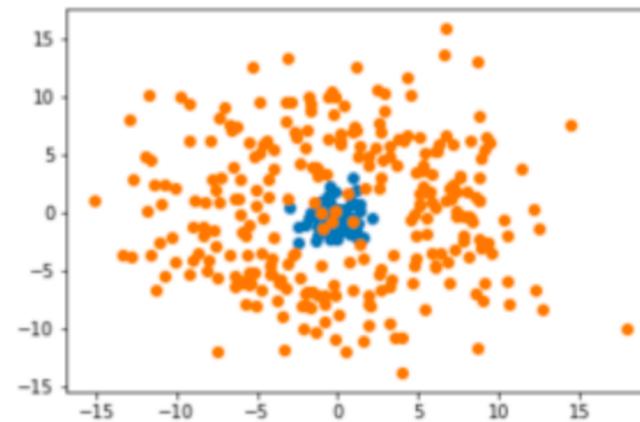
# PCA vs t-SNE

original space:

- blue center
  - orange hull

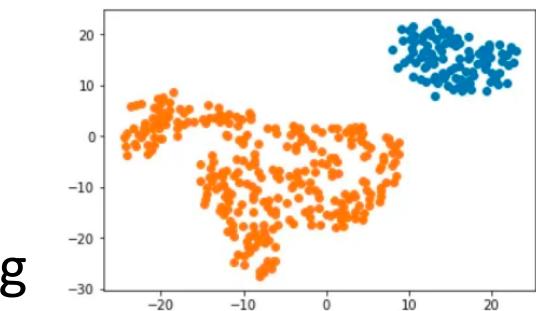
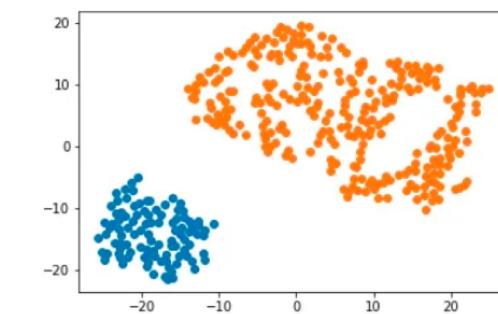
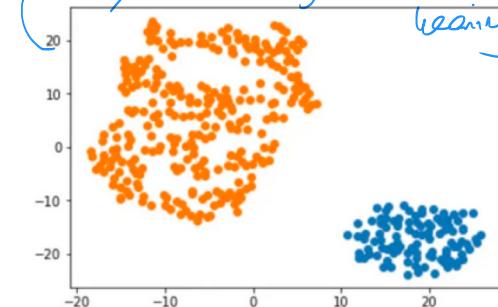


PCA



## t-SNE

→ only local neighborhood has learning



- t-SNE is not deterministic
  - only local neighborhood has a meaning



## t-SNE

- Step 1: In the high-dimensional space, create a probability distribution that dictates the relationships between various neighboring points
- Step 2: Recreate a low dimensional space that follows that probability distribution as best as possible.
- the “t” in t-SNE comes from the t-distribution, which is the distribution used in Step 2. The “S” and “N” (“stochastic” and “neighbor”) come from the fact that it uses a probability distribution across neighboring points.

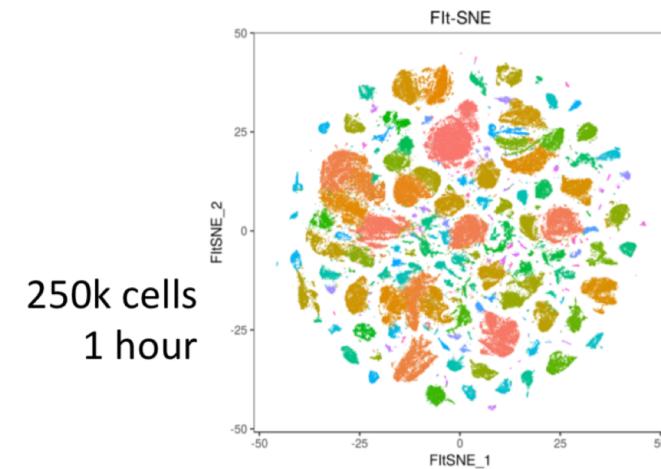
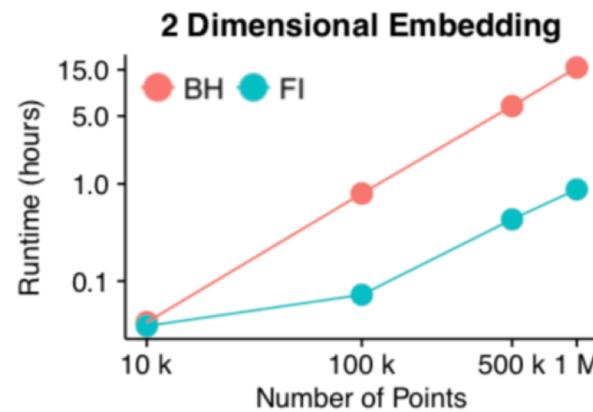


## t-SNE

- Finds a low-dimensional representation of high-dimensional data
  - preserve distances to neighbouring cells
  - non-linear, different transformations on different regions
- Powerful, but need to fiddle with random seed and perplexity
  - 5-50 usually; default value for 10X: 30
  - Often on PC space, but not mandatory
- Nice blog about t-SNE: <https://distill.pub/2016/misread-tsne/>
- Implementation: Rtsne or much faster Flt-SNE

- Fast Fourier Transform-accelerated Interpolation-based t-SNE -  $O(n)$

Linderman et al (2017) *BioRxiv*





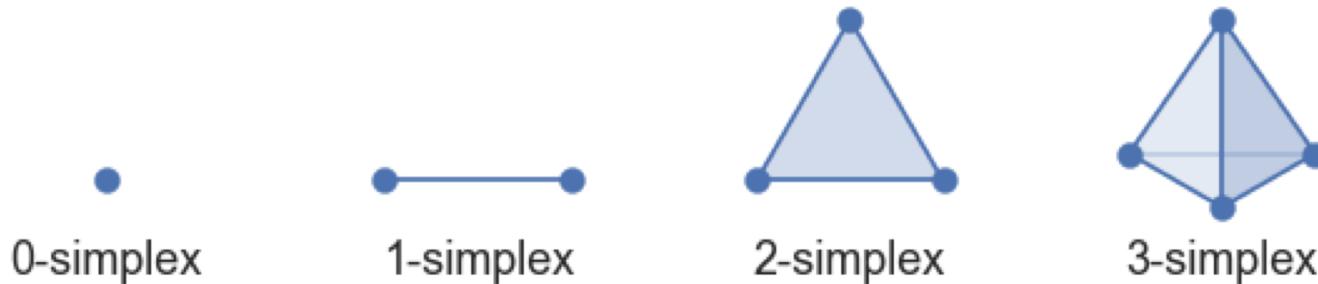
University of  
Zurich UZH

10  
01  
101

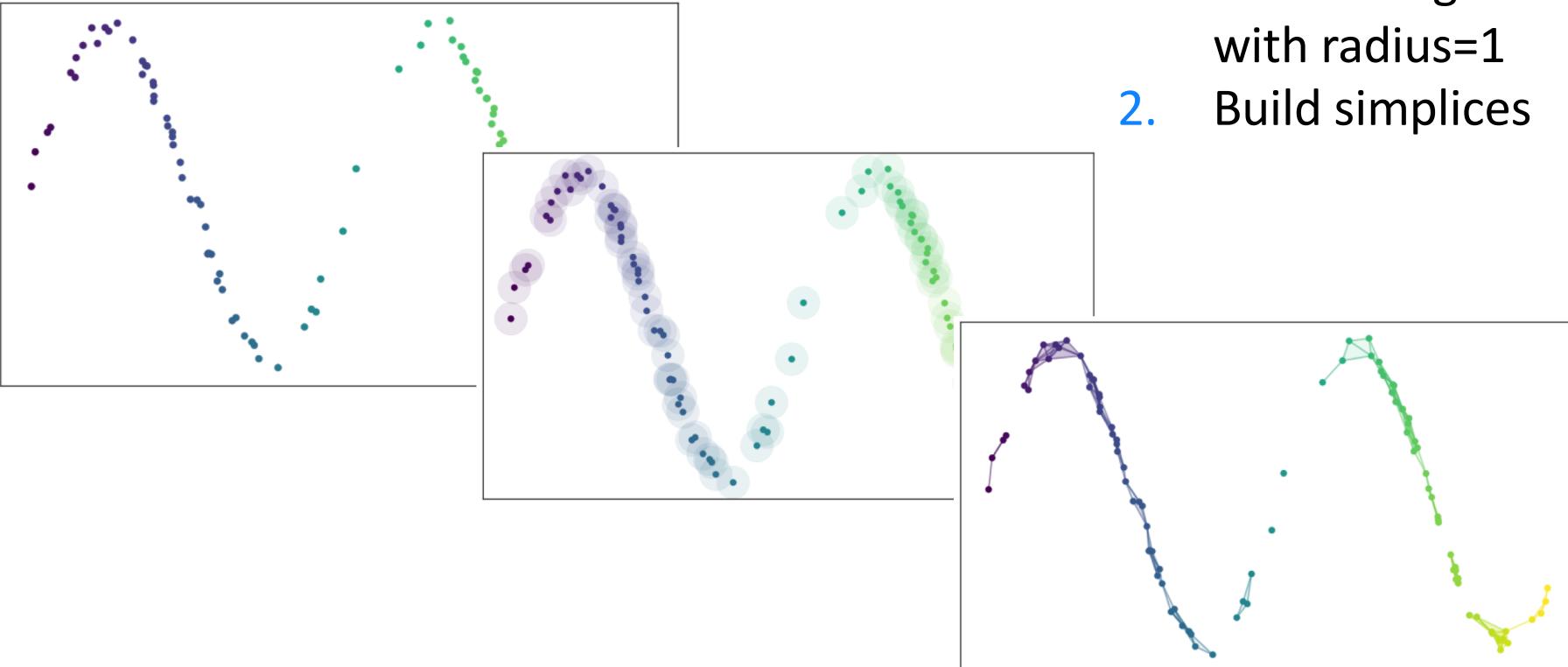


# UMAP

- Uniform Manifold Approximation and Projection
- Approach: Find for each point the neighbors and build simplices



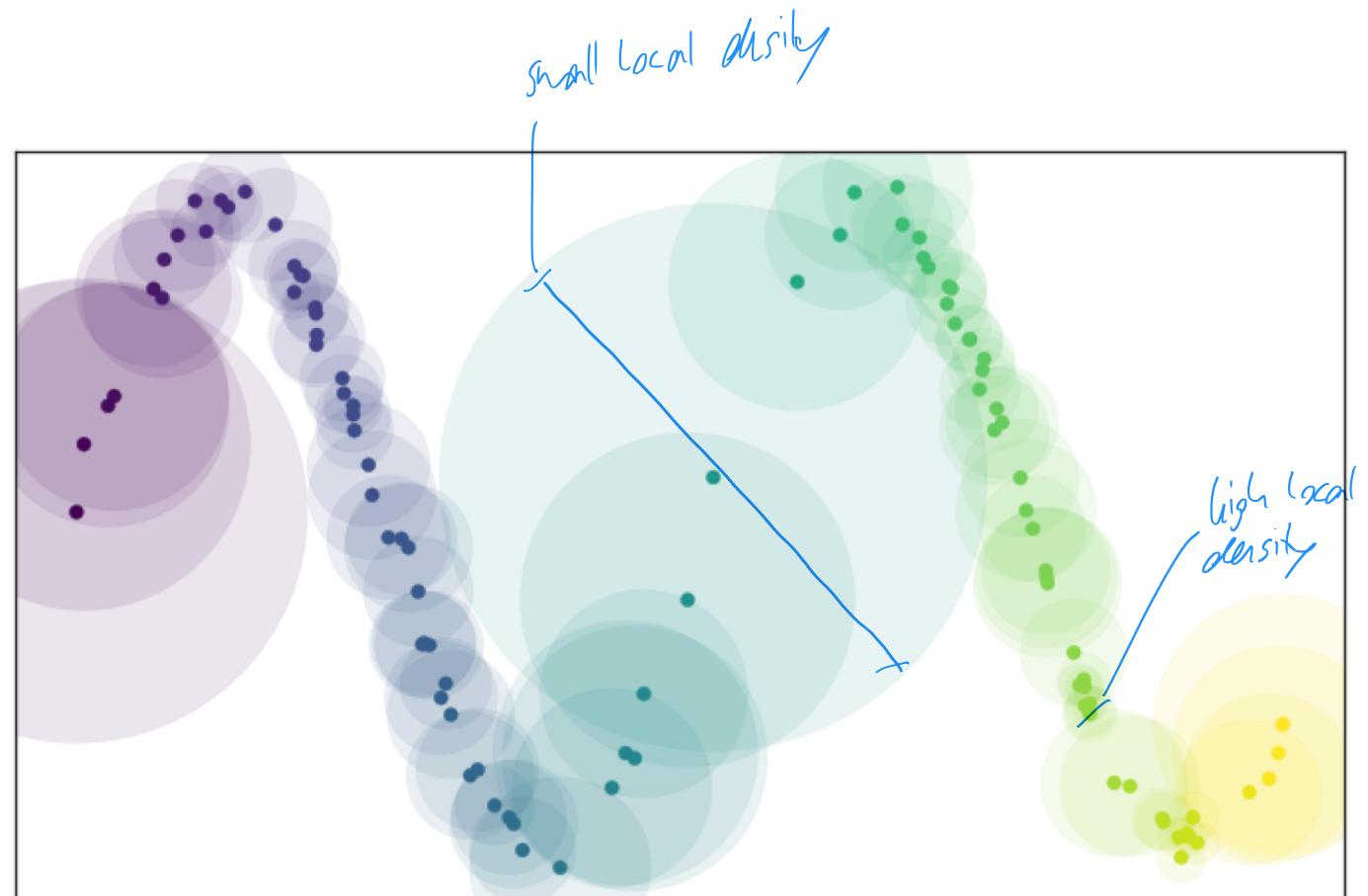
## UMAP steps





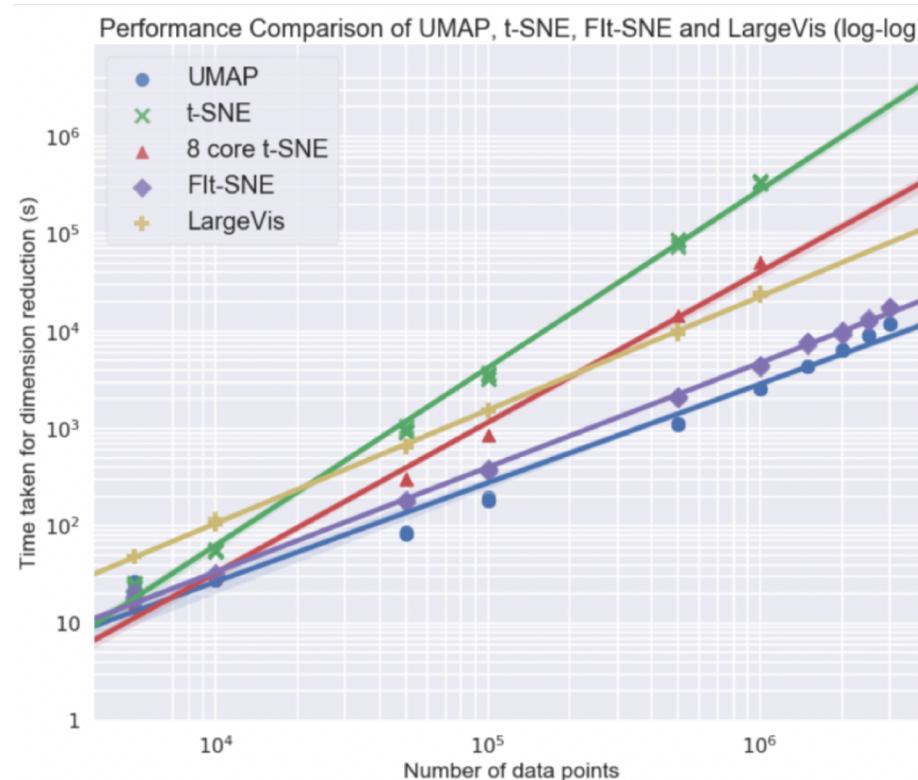
## UMAP

- All balls have radius=1 using a locally varying metric
- Avoids Clumping



10  
01  
101

## UMAP computing time



- UMAP is faster and scales better than classical t-SNE
- comparable to Flt-SNE

↳ Fairly

