



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



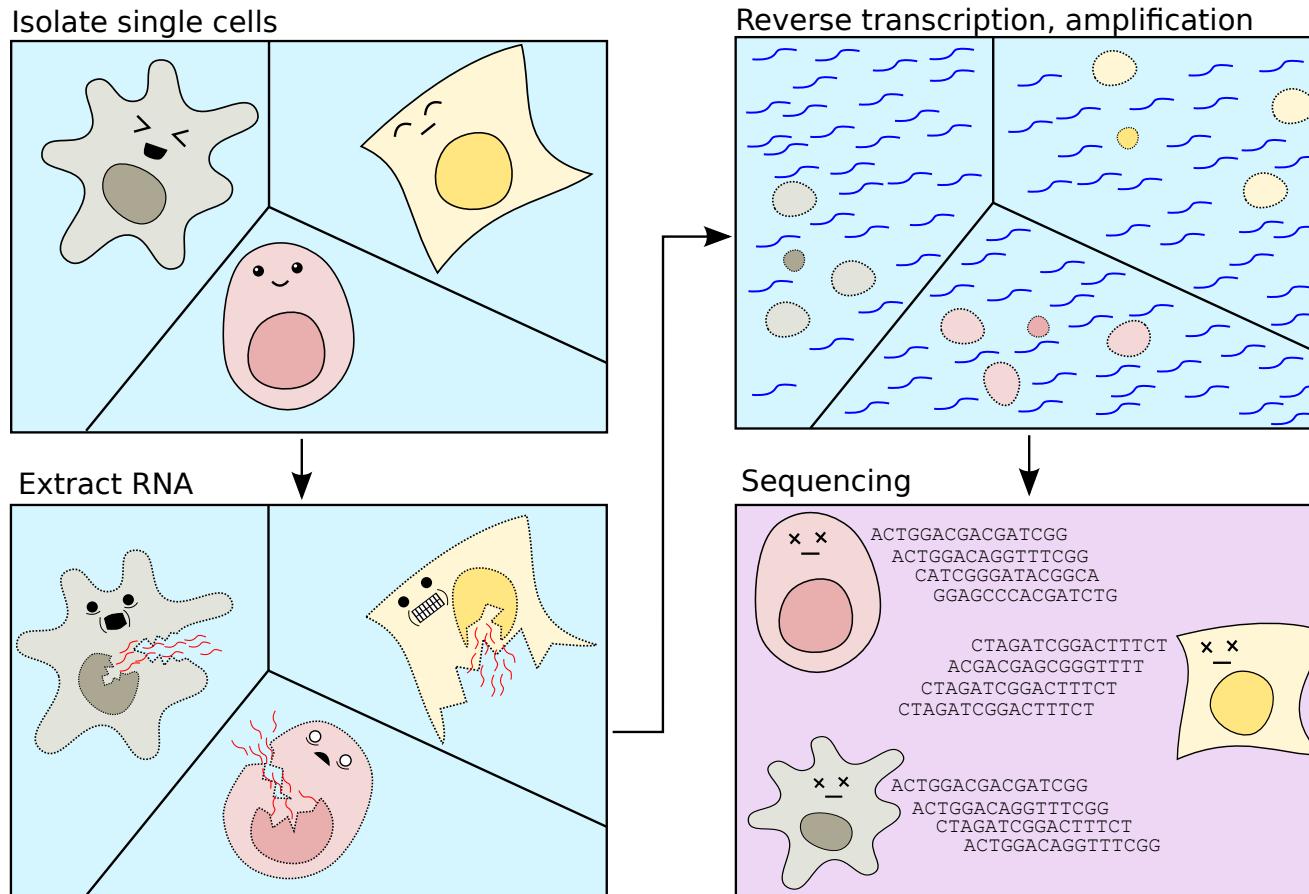
University of
Zurich UZH

10
01
101

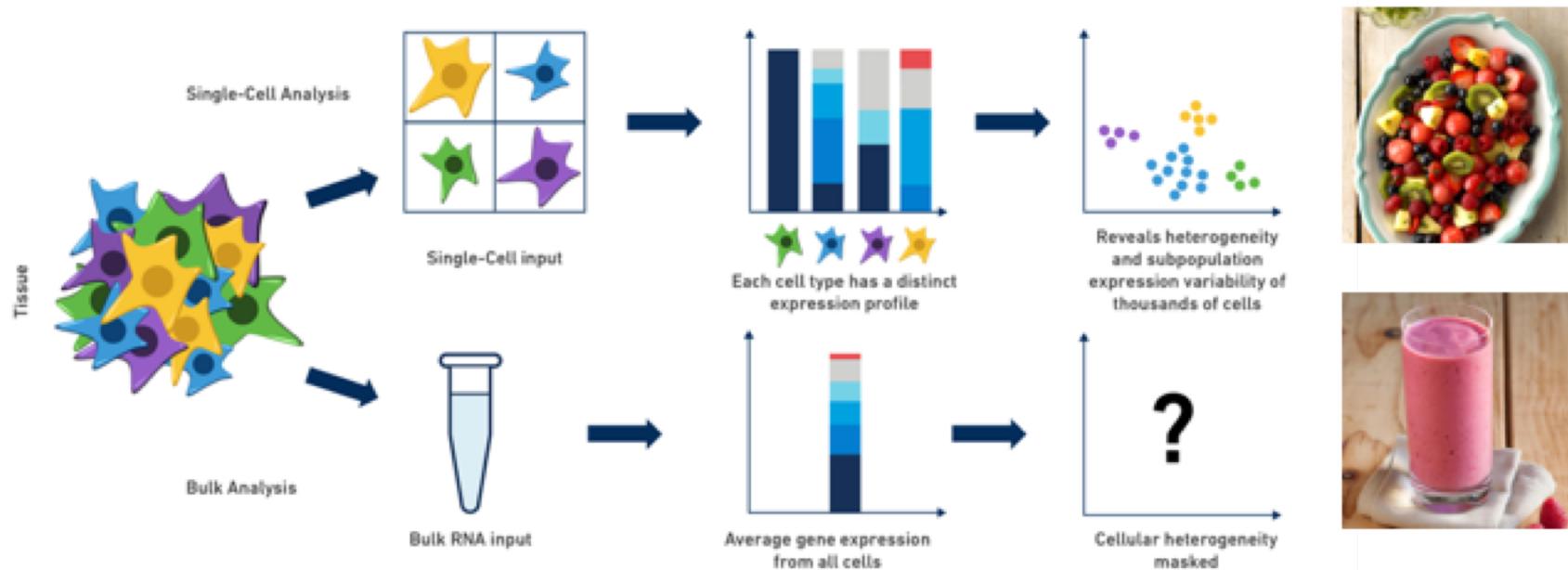
functional genomics center zurich

010 01
101 10
010 01
01 1

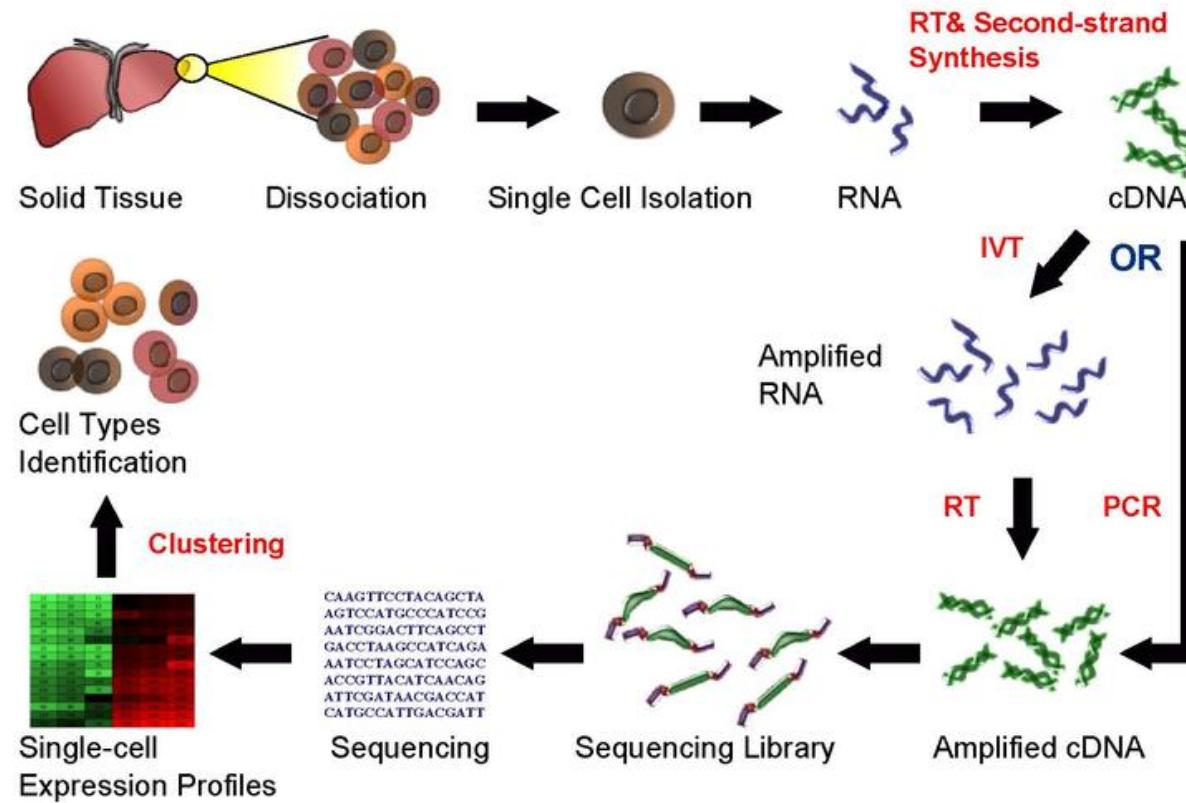
f g c z
10 0
01 1



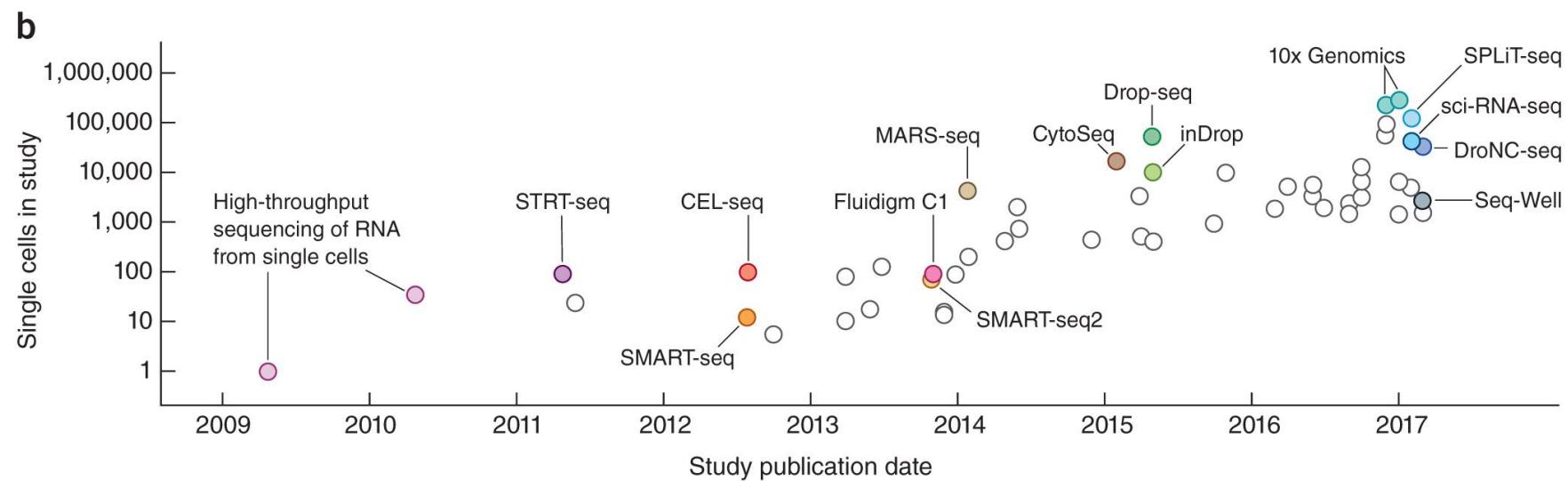
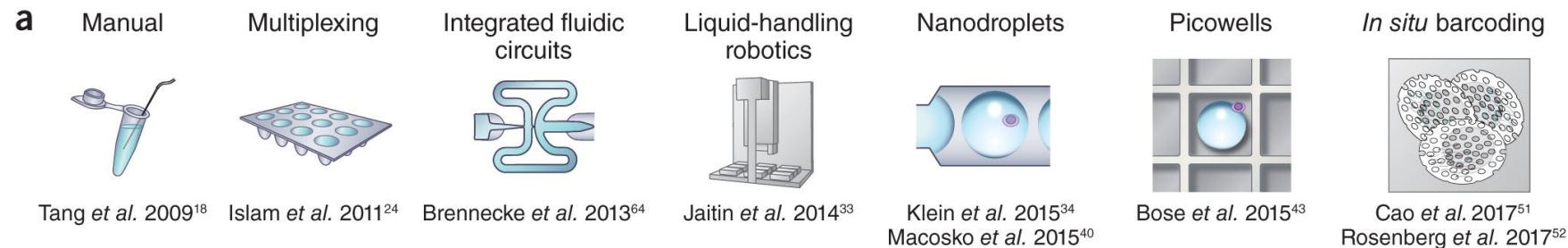
Adapted from Aaron Lun



Single Cell RNA Sequencing Workflow



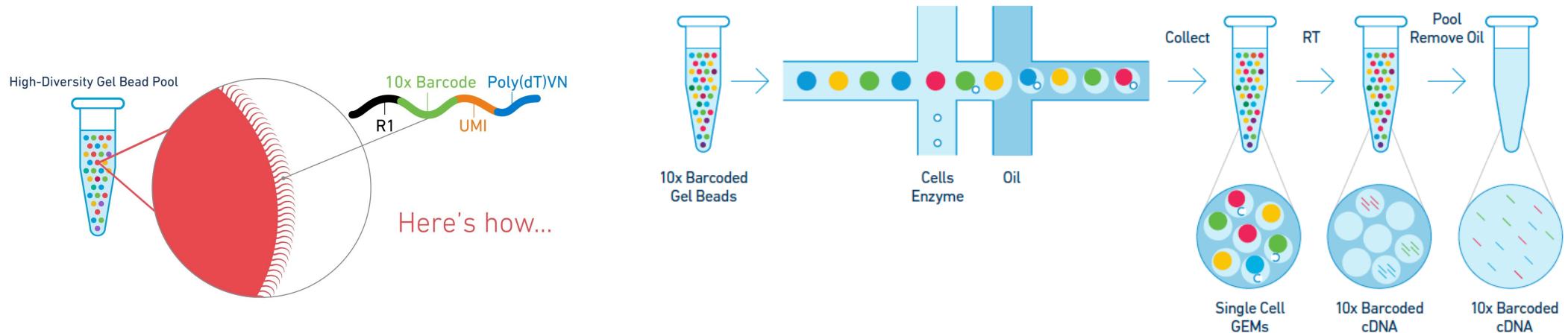
Technologies



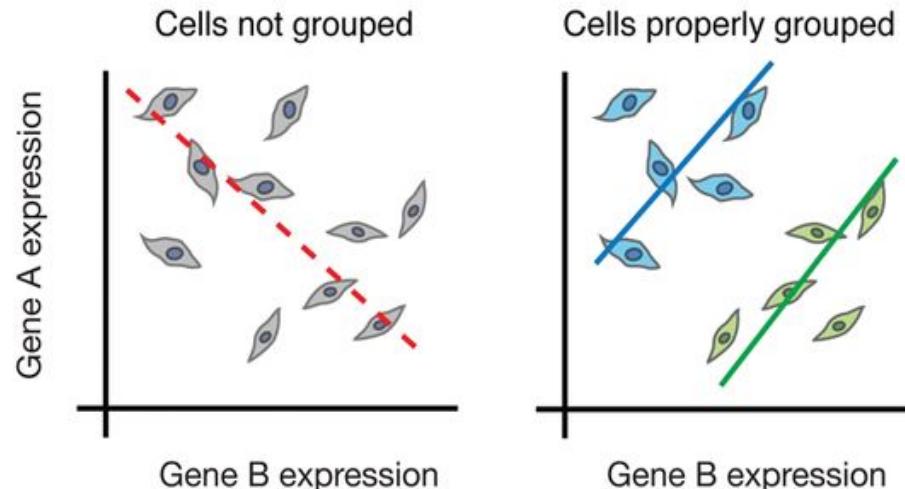


Nanodroplets systems – e.g. 10X Genomics

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



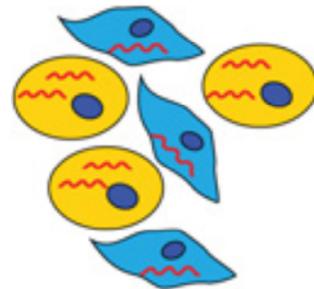
Single cell measurements preserve crucial information that is lost by bulk measurements



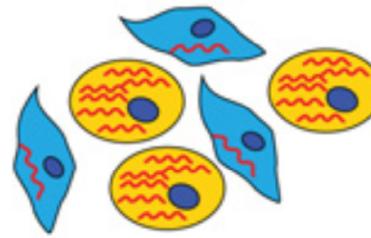
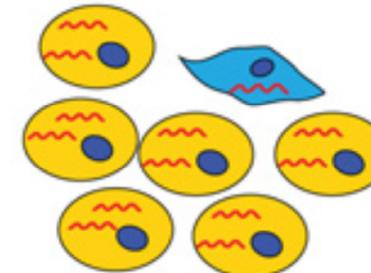
Simpson's Paradox describes the misleading effects that arise when averaging signals from multiple individuals

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

Control cells

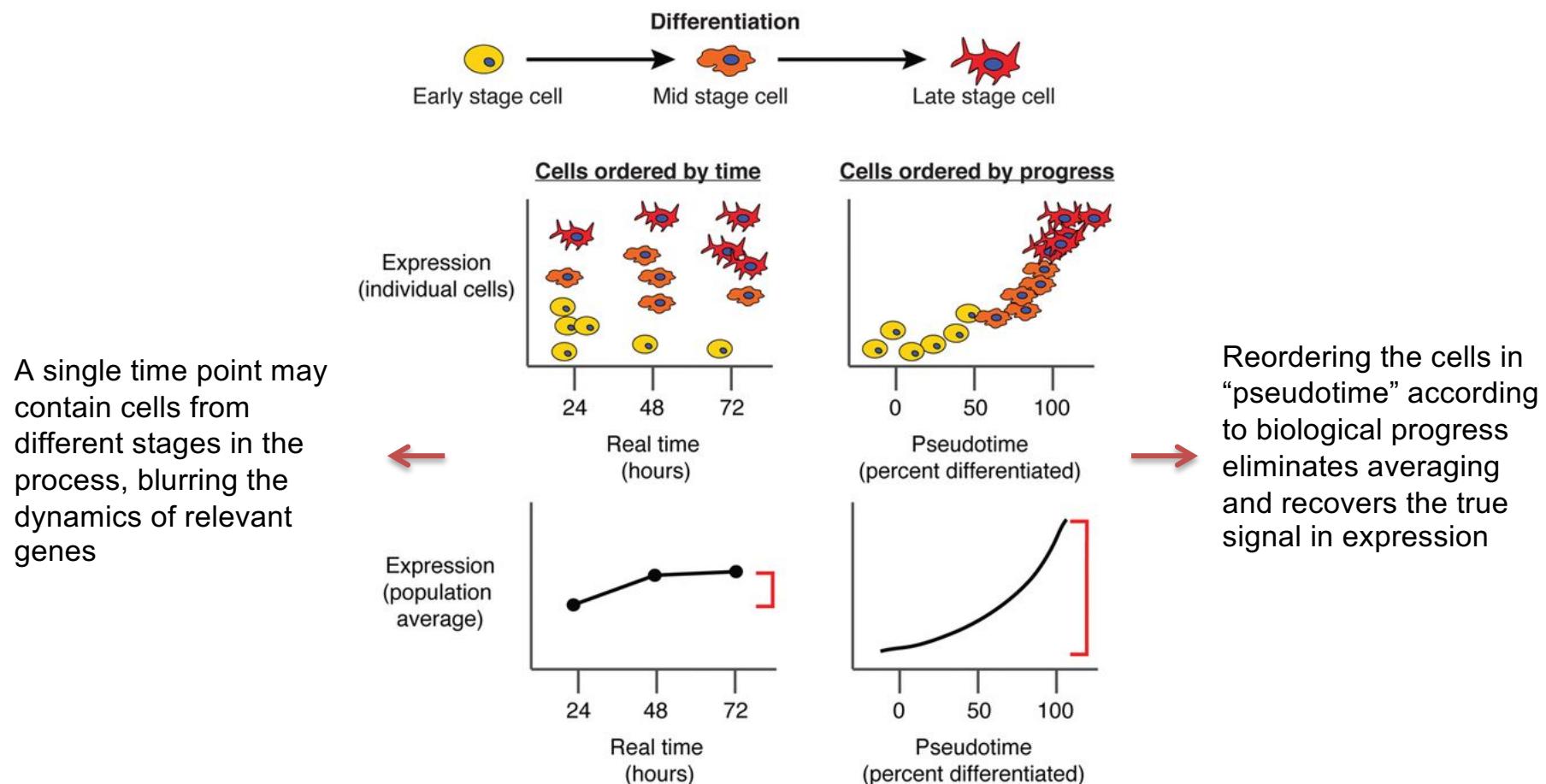


Perturb cells

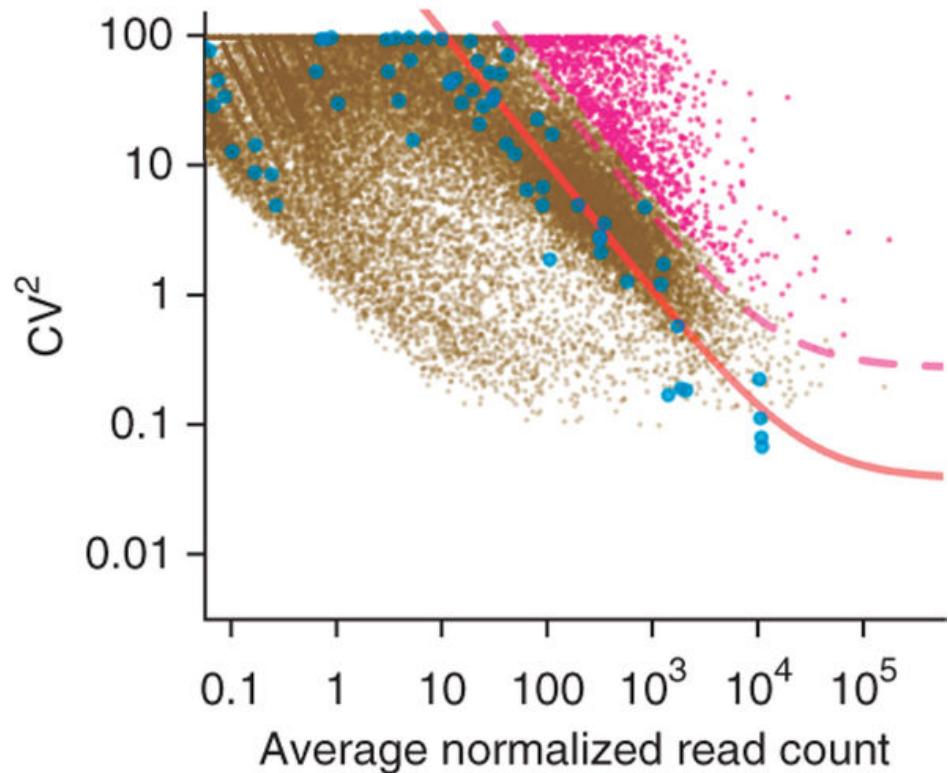
Change in
regulationChange in
composition



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



Variability of scRNA-seq counts

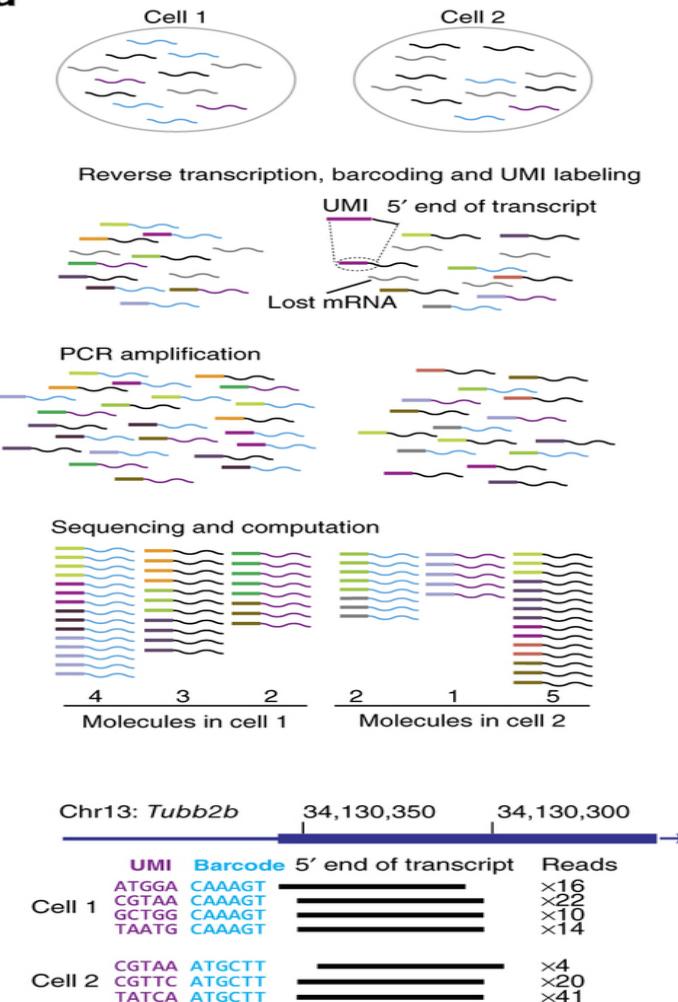


- Blue: spike-ins
- Pink: genes with potentially differential expression between cells

Brennecke et al. *Nature Methods* 2013

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

Unique Molecular Identifiers (UMIs)

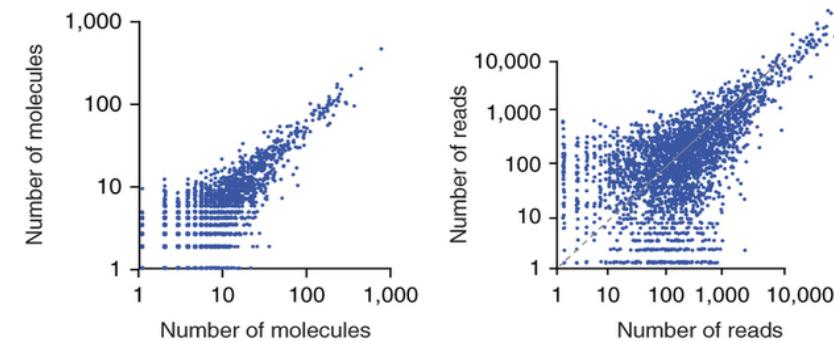
a


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.

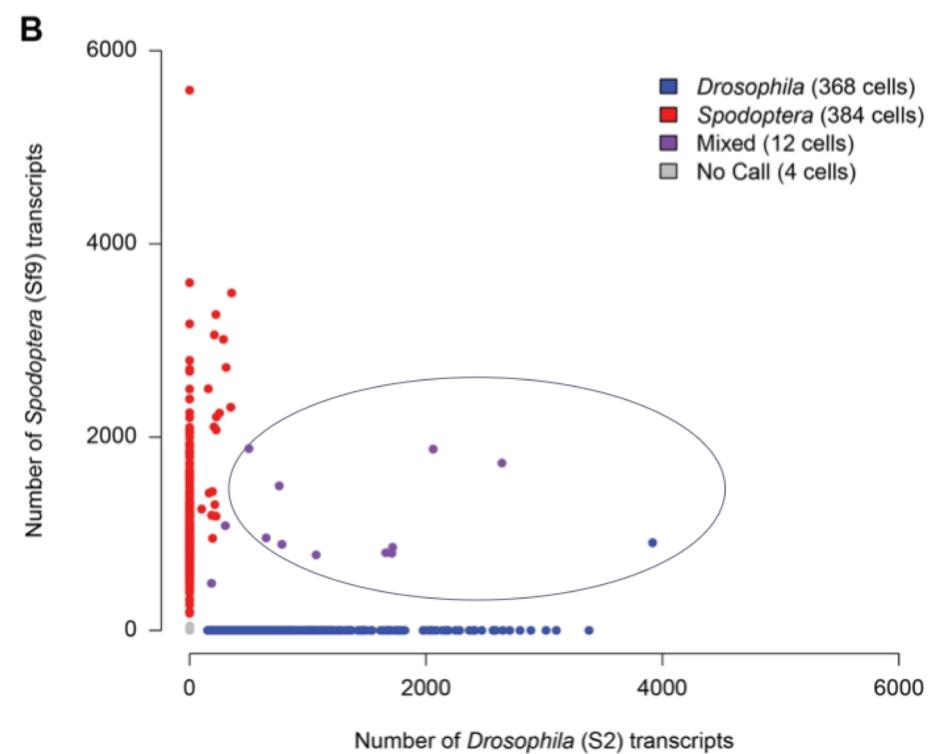
UMIs are only available for 3'-tagging protocols, e.g. the 10X.

Islam et al. Nature Methods 2014



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



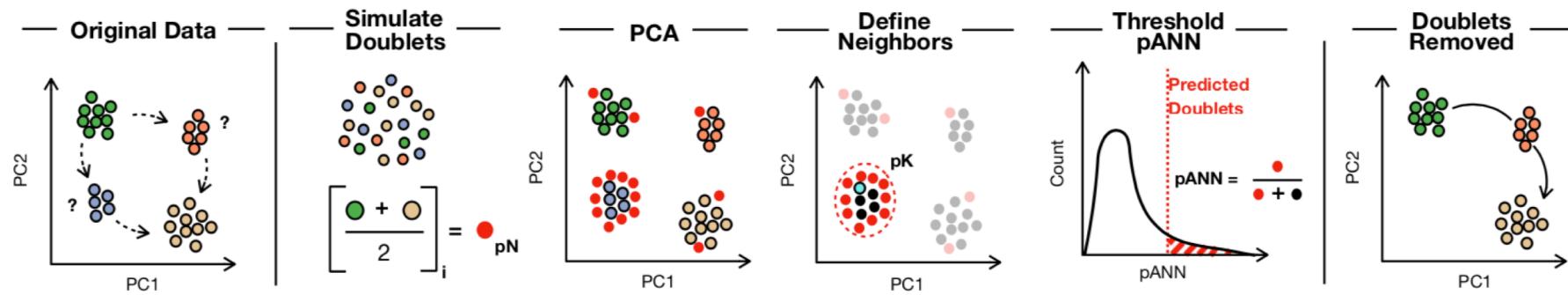


10

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

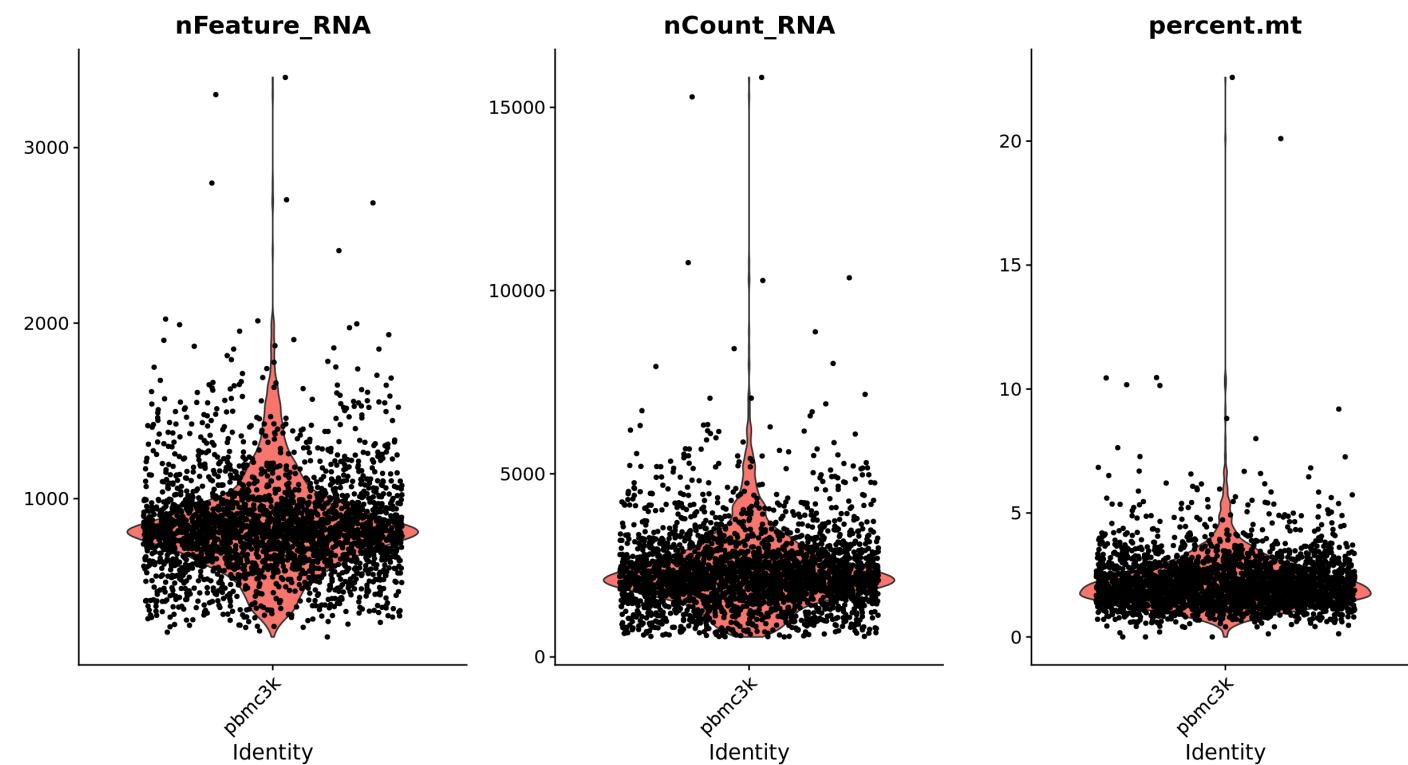
10
01
101

functional genomics center zurich

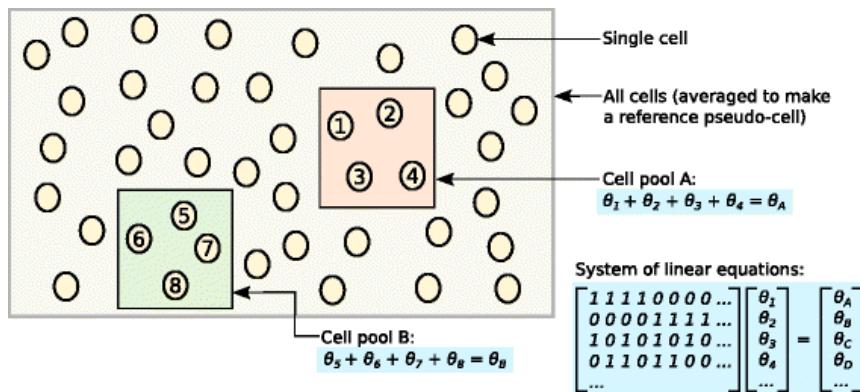
010 01
101 10
010 01
010 01

f g c z
10 10
01 01

QC Filtering Plots



scran: Normalization using pools of cells



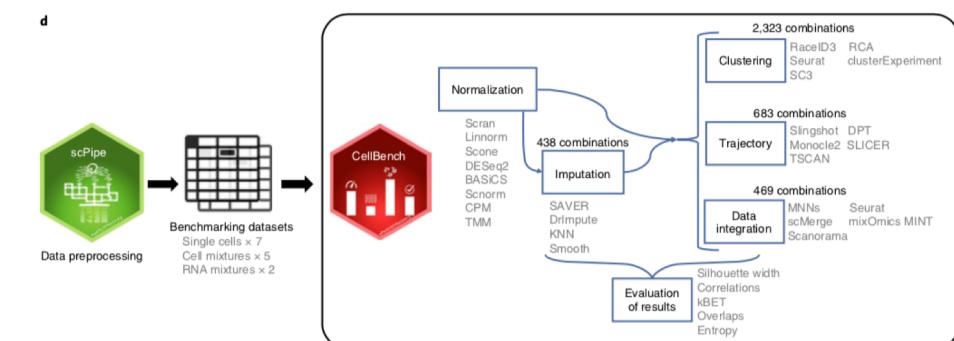
- Define a pool of cells
- Sum expression values across all cells in the pool
- Normalize the cell pool against an average reference, using the summed expression values
- Repeat this for many different pools of cells to construct a linear system
- Deconvolute the pool-based size factors to their cell-based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

Normalization: Performance comparison papers



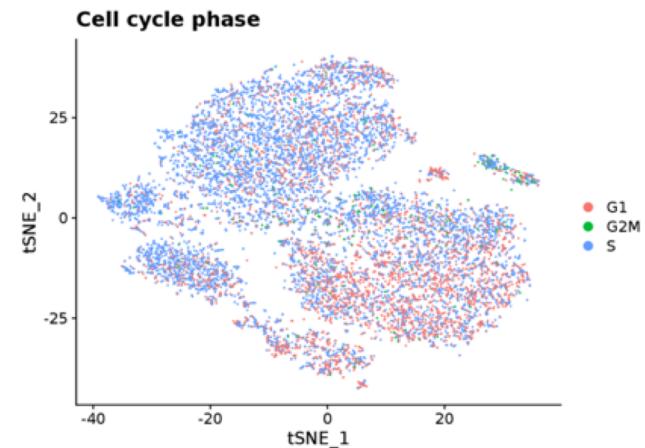
Benchmarking single cell RNA-sequencing analysis pipelines using mixture control experiments

Luyi Tian^{ID}^{1,2*}, Xueyi Dong^{1,3}, Saskia Freytag^{1,4}, Kim-Anh Lê Cao^{ID}⁵, Shian Su¹, Abolfazl JalalAbadi¹⁵, Daniela Amann-Zalcenstein^{1,2}, Tom S. Weber^{ID}^{1,2}, Azadeh Seidi⁶, Jafar S. Jabbari⁶, Shalin H. Naik^{ID}^{1,2} and Matthew E. Ritchie^{ID}^{1,2*}



Normalization: Cell Cycle causes unwanted variation

- Requires normalization considering cell cycle as latent variable



Dimensionality Reduction

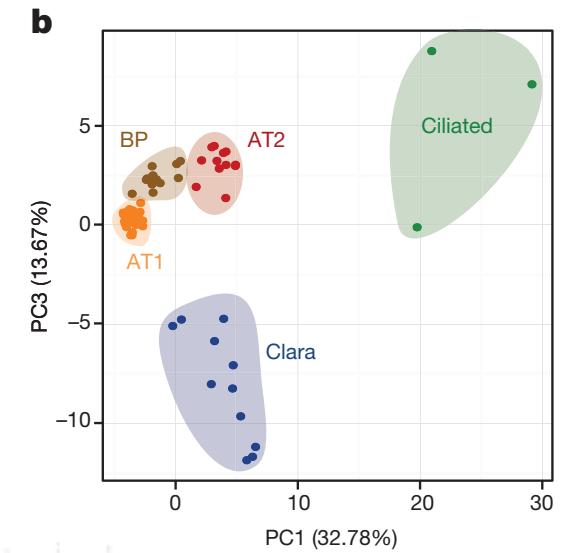
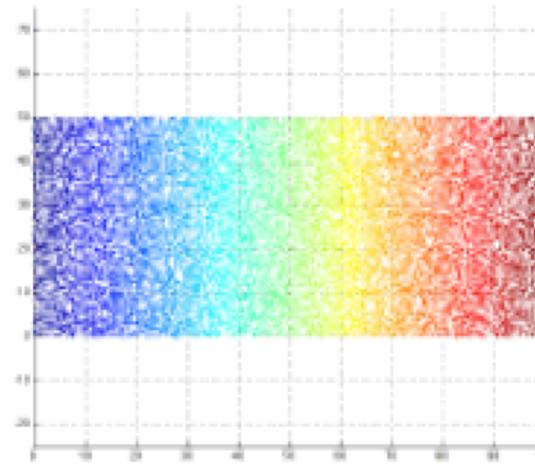
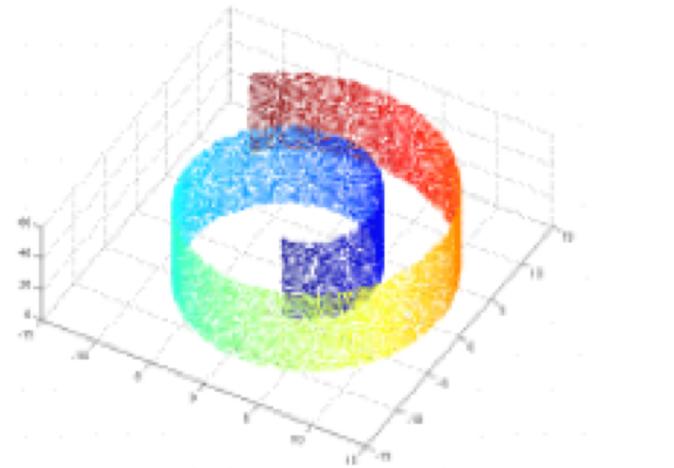
- Removes redundancy in data
 - Makes subsequent analyses more efficient and robust
 - clustering
 - classification
 - cell characterization
 - Many features and high dimensionality make classification more erroneous

Dimensionality Reduction Methods

- Matrix Factorization
 - PCA
 - MDS – multi-dimensional scaling
 -
 - Graph-based
 - t-SNE
 - UMAP
 - Autoencoder (based on matrix factorization)

PCA

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

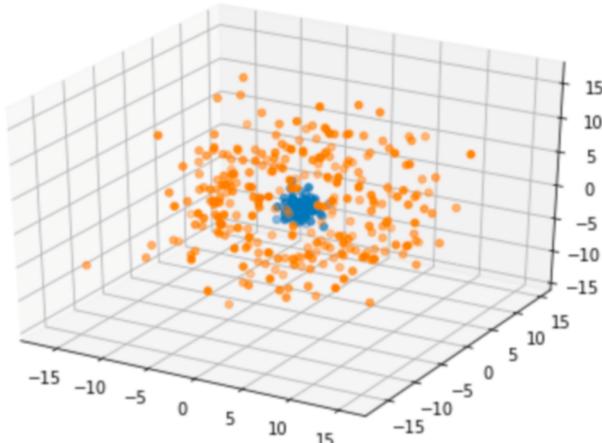




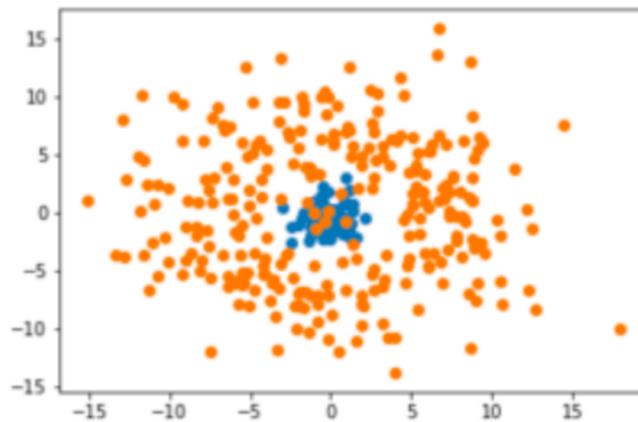
PCA vs t-SNE

original space:

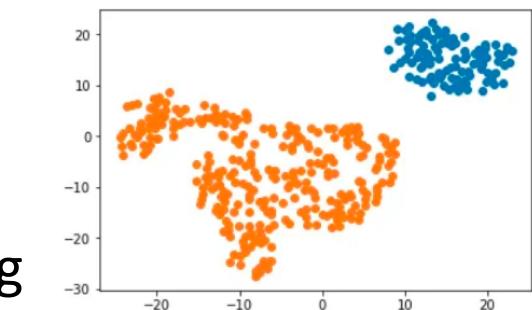
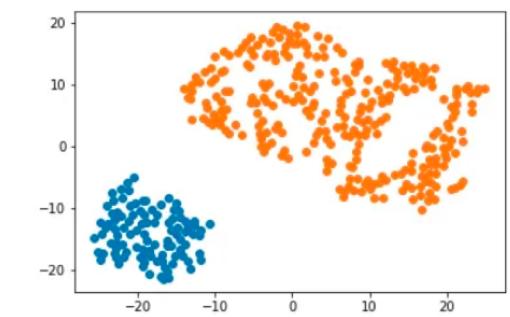
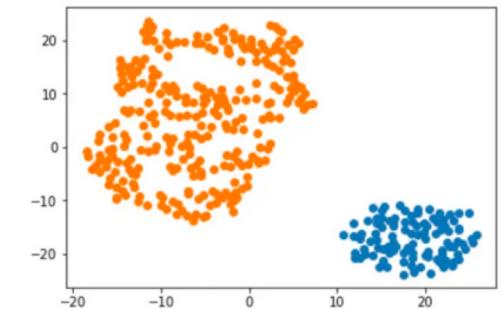
- blue center
- orange hull



PCA



t-SNE



- 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.



10
01
101



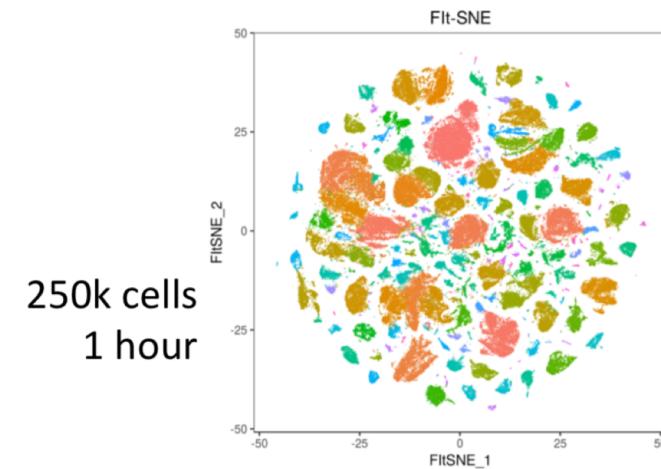
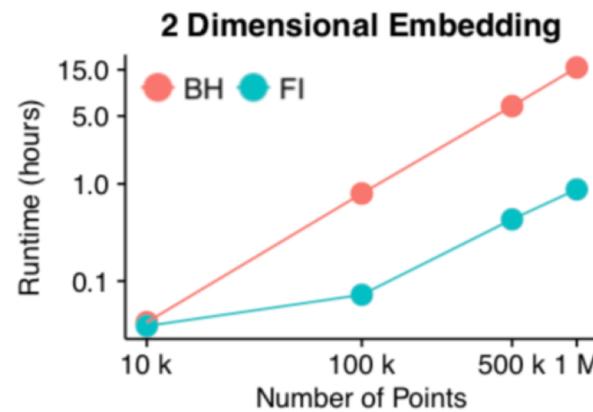
t-SNE

- 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*





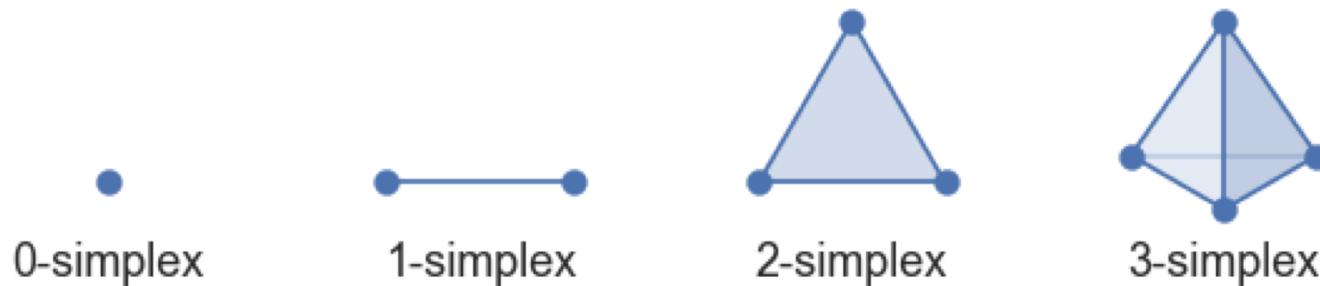
University of
Zurich UZH

10
01
101

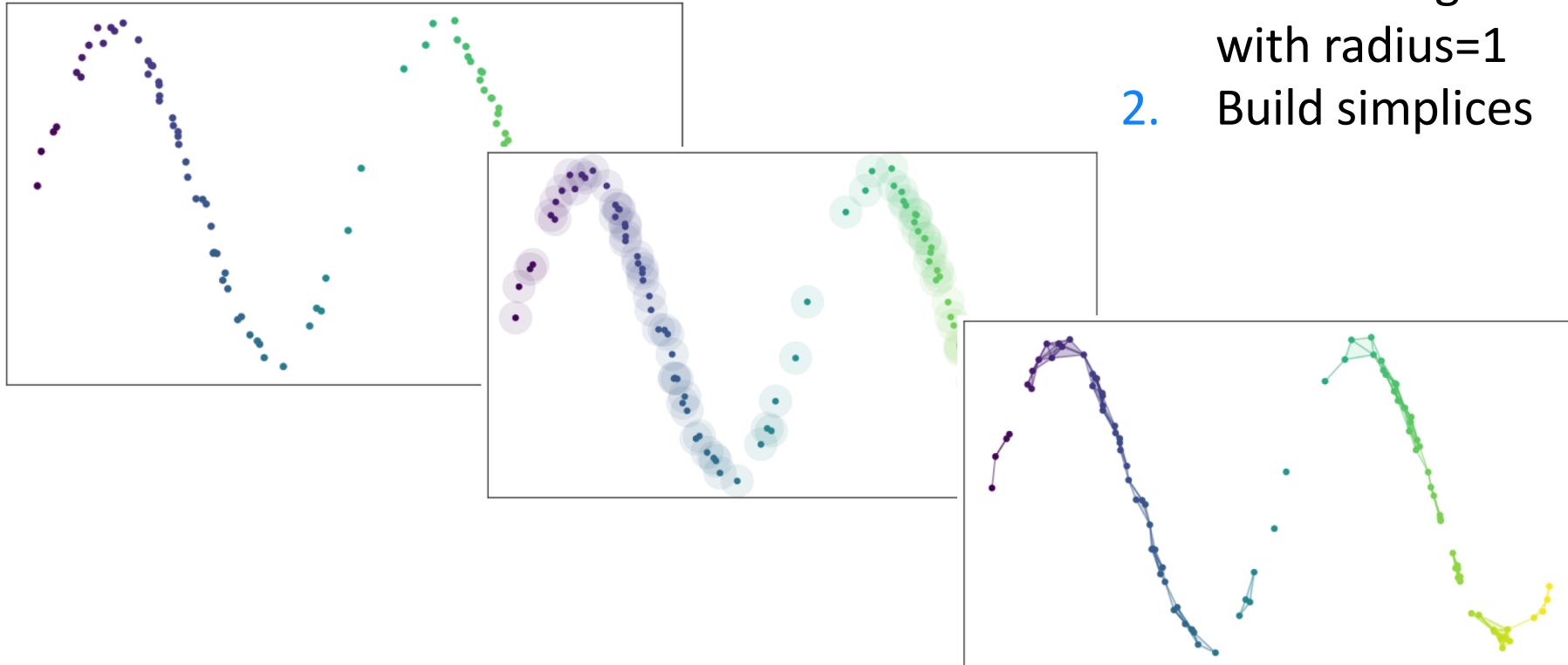


UMAP

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



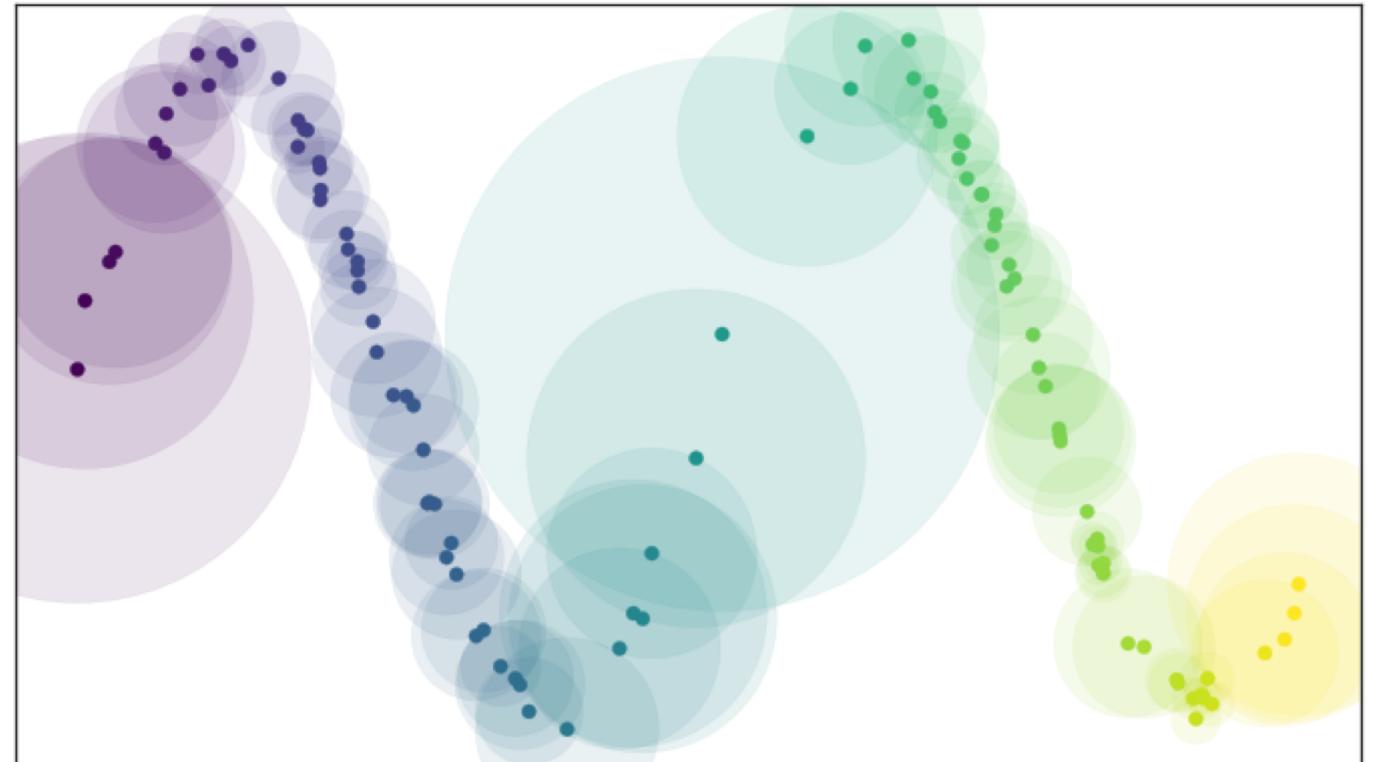
UMAP steps



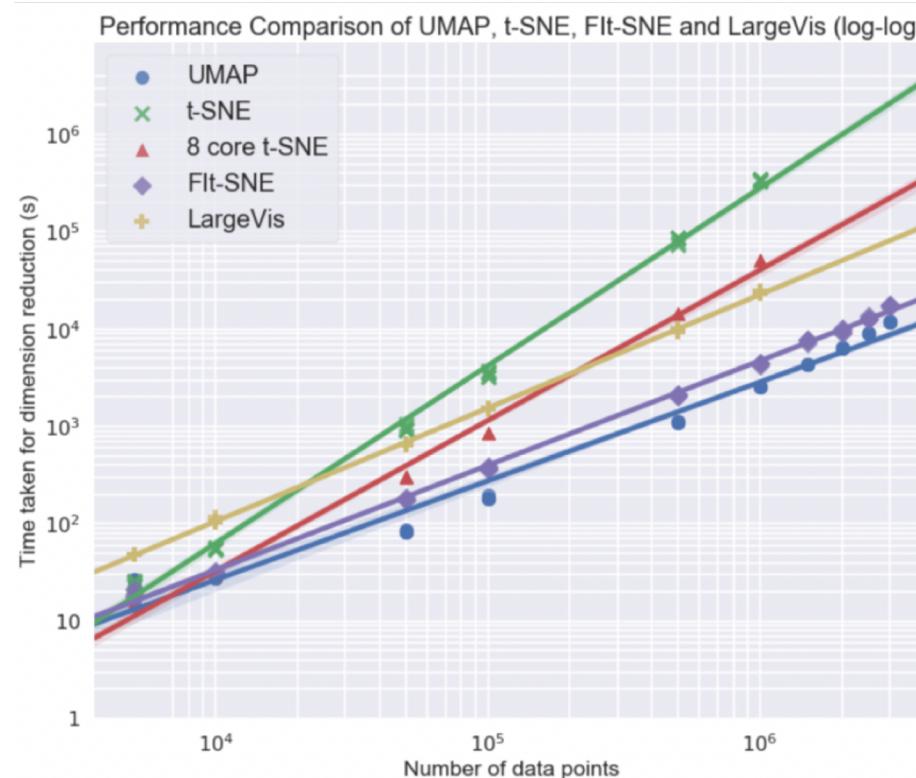


UMAP

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



UMAP computing time



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

Clustering

A systematic performance evaluation of clustering methods for single-cell RNA-seq data

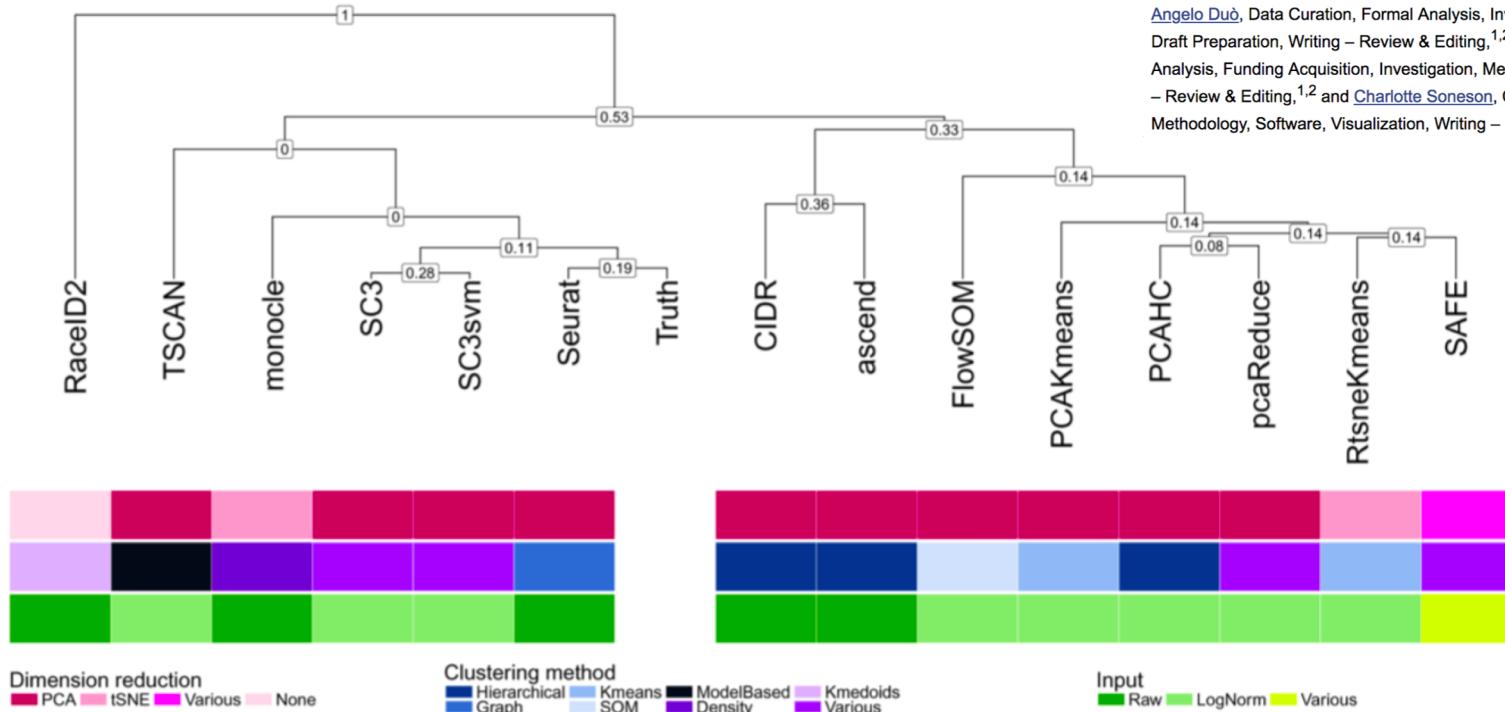


Figure 4. Clustering of the methods based on the average similarity of their partitions across data sets, for the true number of clusters.

Numbers on internal nodes indicate the fraction of dendograms from individual data sets where a particular subcluster was found.