



University of  
Zurich<sup>UZH</sup>

10  
01  
101

101 1  
010 0  
0101 10

functional genomics center zurich

010 01  
101 10  
010 01

01 1  
0  
10 0  
01 1

# Single Cell RNA-seq: Characteristics, Preprocessing and QC

Hubert Rehrauer



University of  
Zurich<sup>UZH</sup>

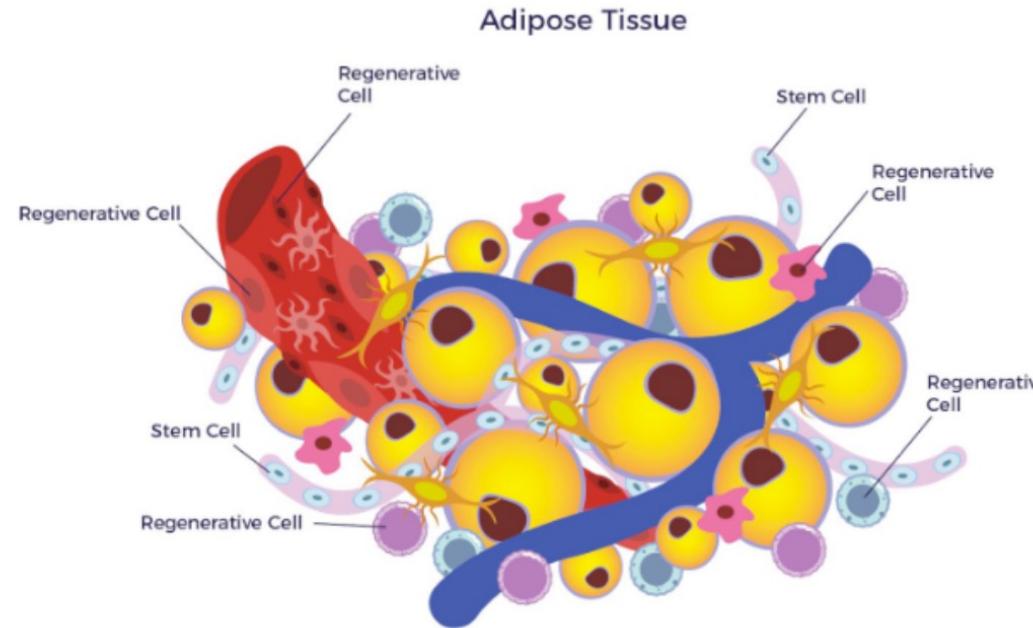
**ETH**

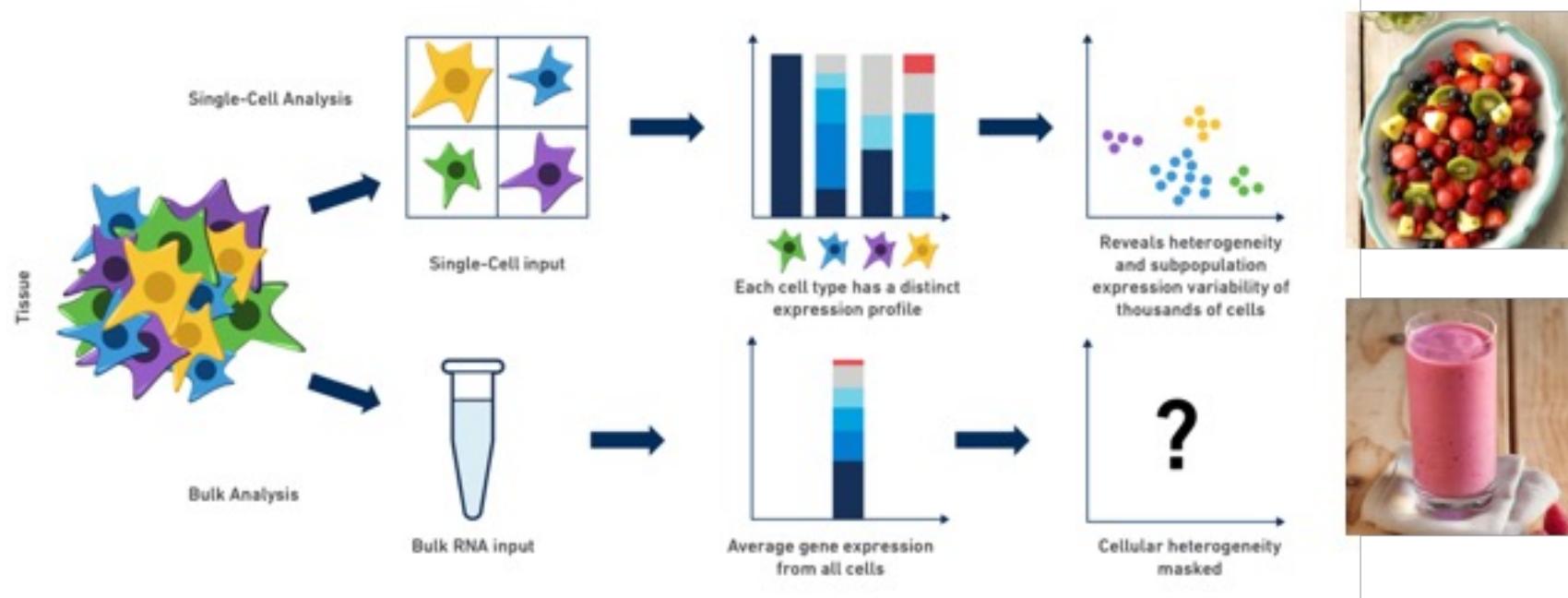
Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich

10  
01  
101101 1  
010 0  
0101 1001 1  
10 0  
101 10  
010 01  
010 01

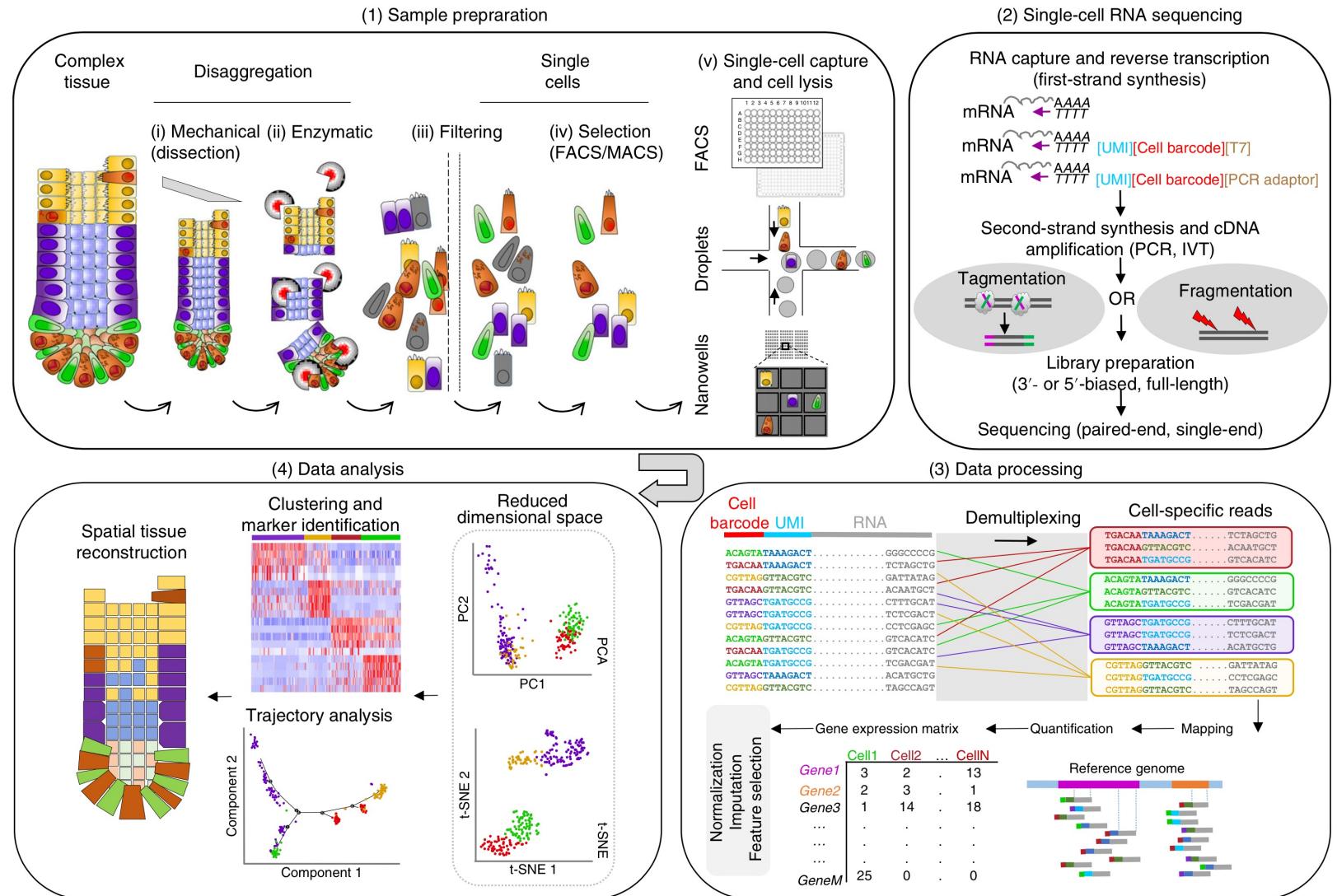
f g c z + 01 1

# Tissues are heterogeneous





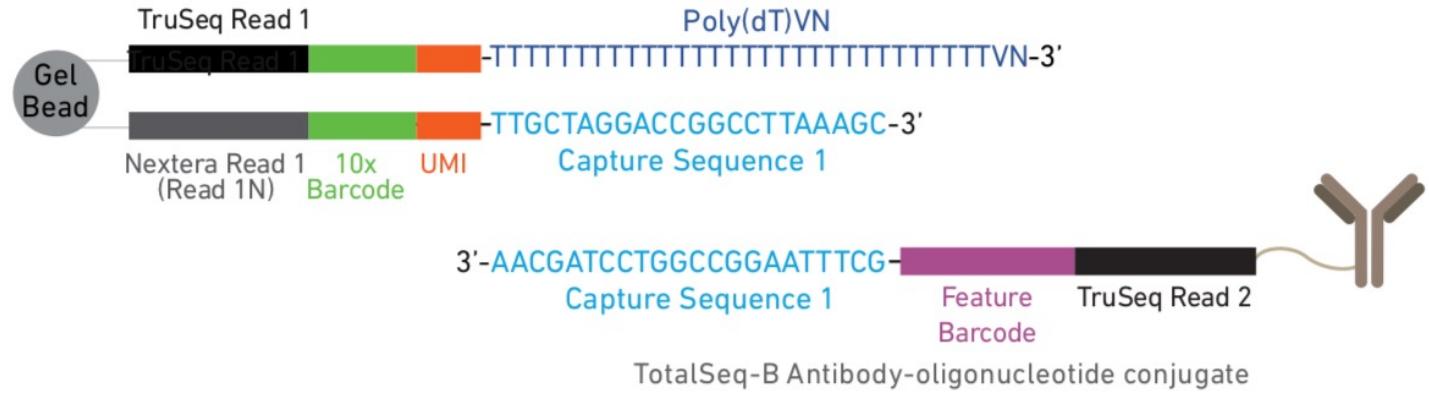
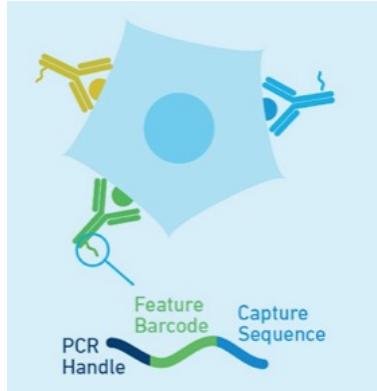
# scRNA-seq





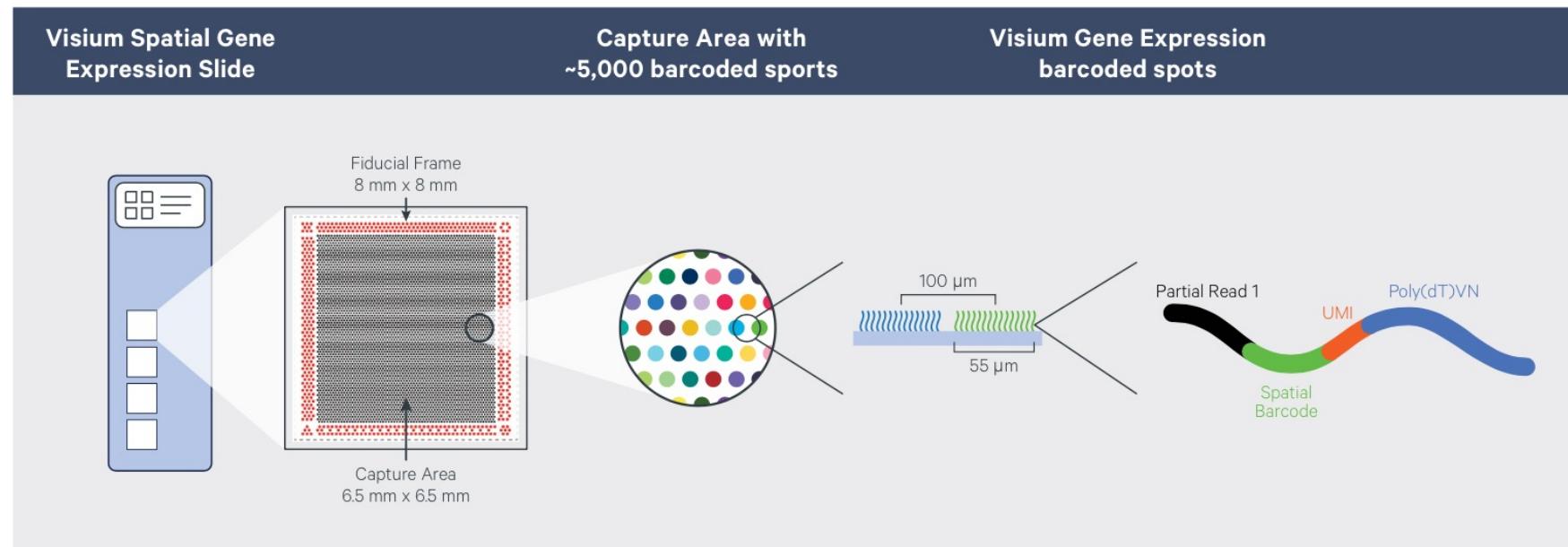
# Cell hashing (feature barcoding)

CITE-seq

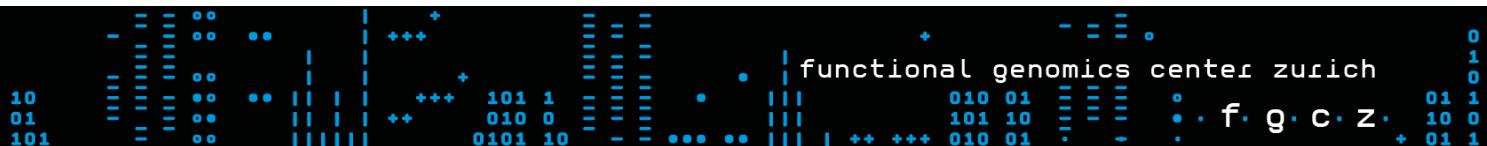


- Can be used
  - to combine cells from different pools in one experiment
  - detect presence of surface markers of cells
  - detect doublets

# scRNA-seq with location information (spatial transcriptomics)

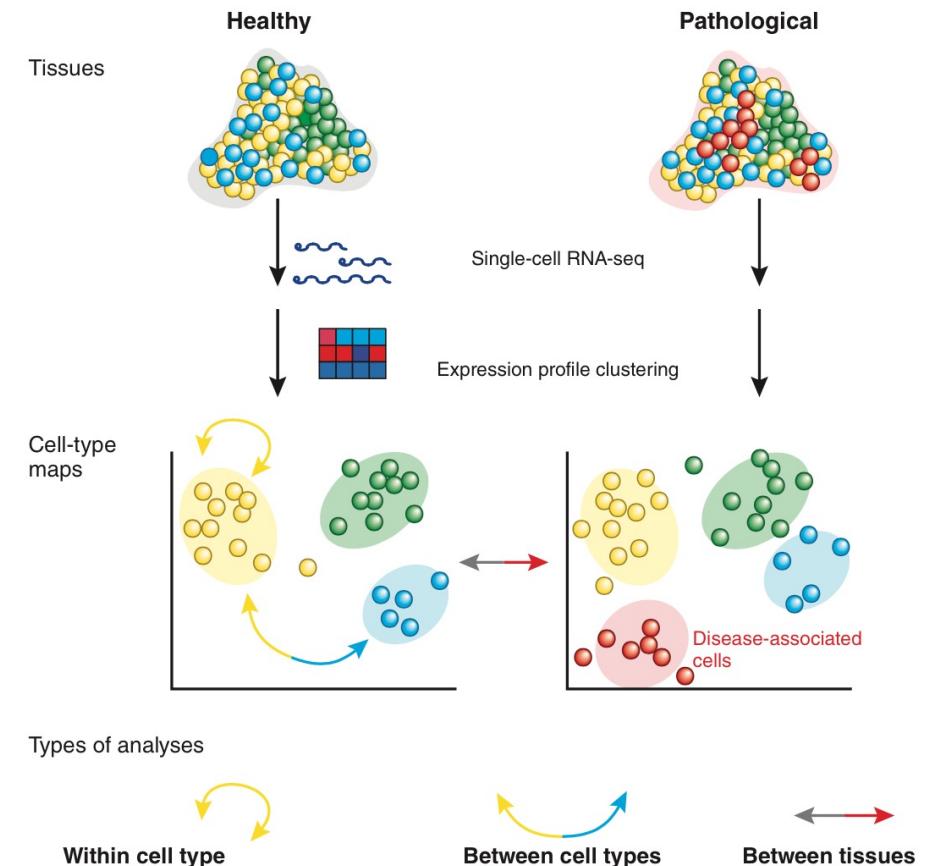


- can be combined with imaging the tissue slice
  - each spot is characterized by expression profile + image



## scRNA-seq vs bulk RNA-seq

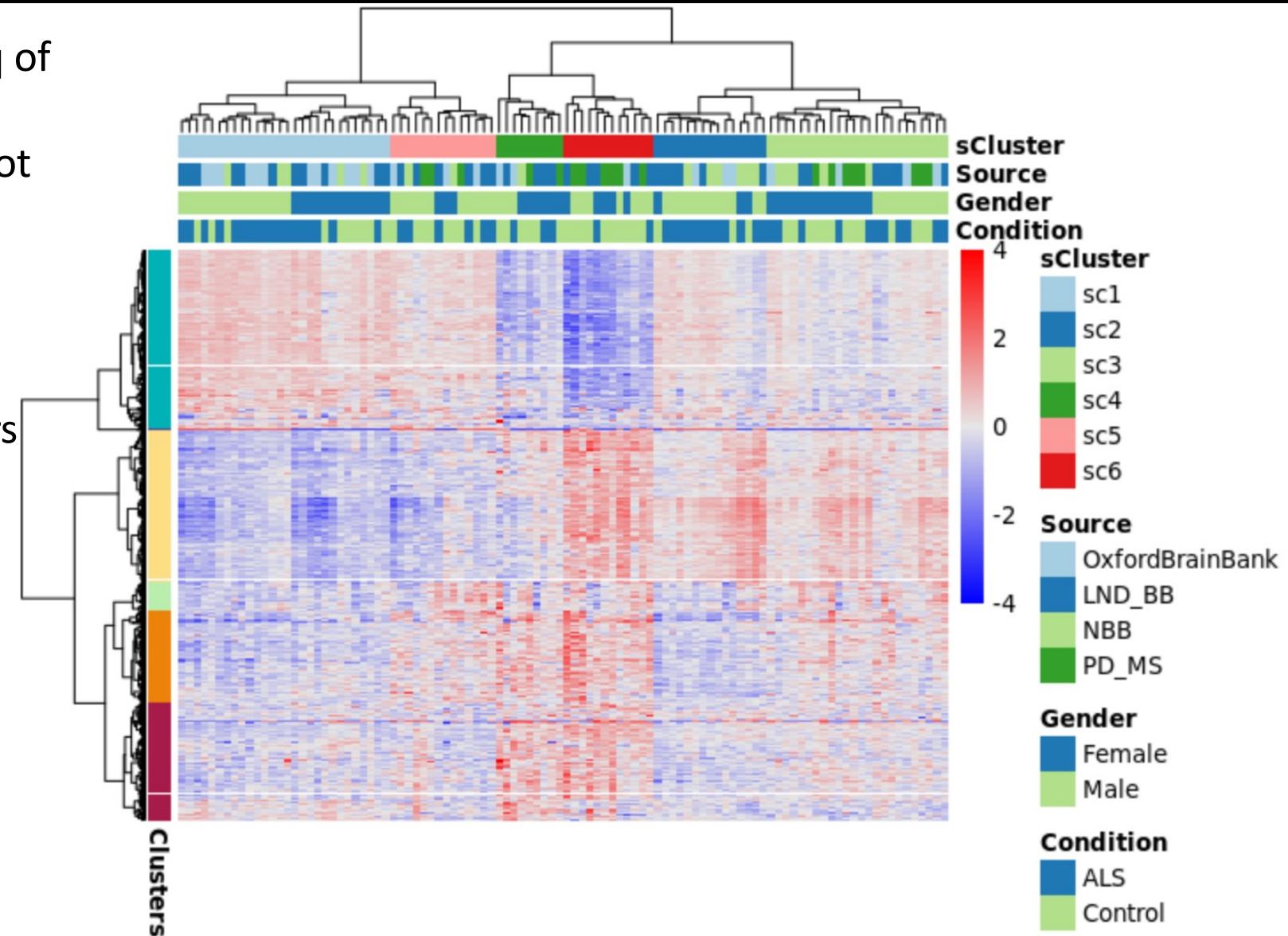
- bulk RNA-seq
  - lyse entire tissue → expression vector is average of all cells
  - sort cells by surface marker → expression profile is average of all cells exhibiting that surface marker
- cell sorting can affect cell viability and can thus affect mRNA measurements



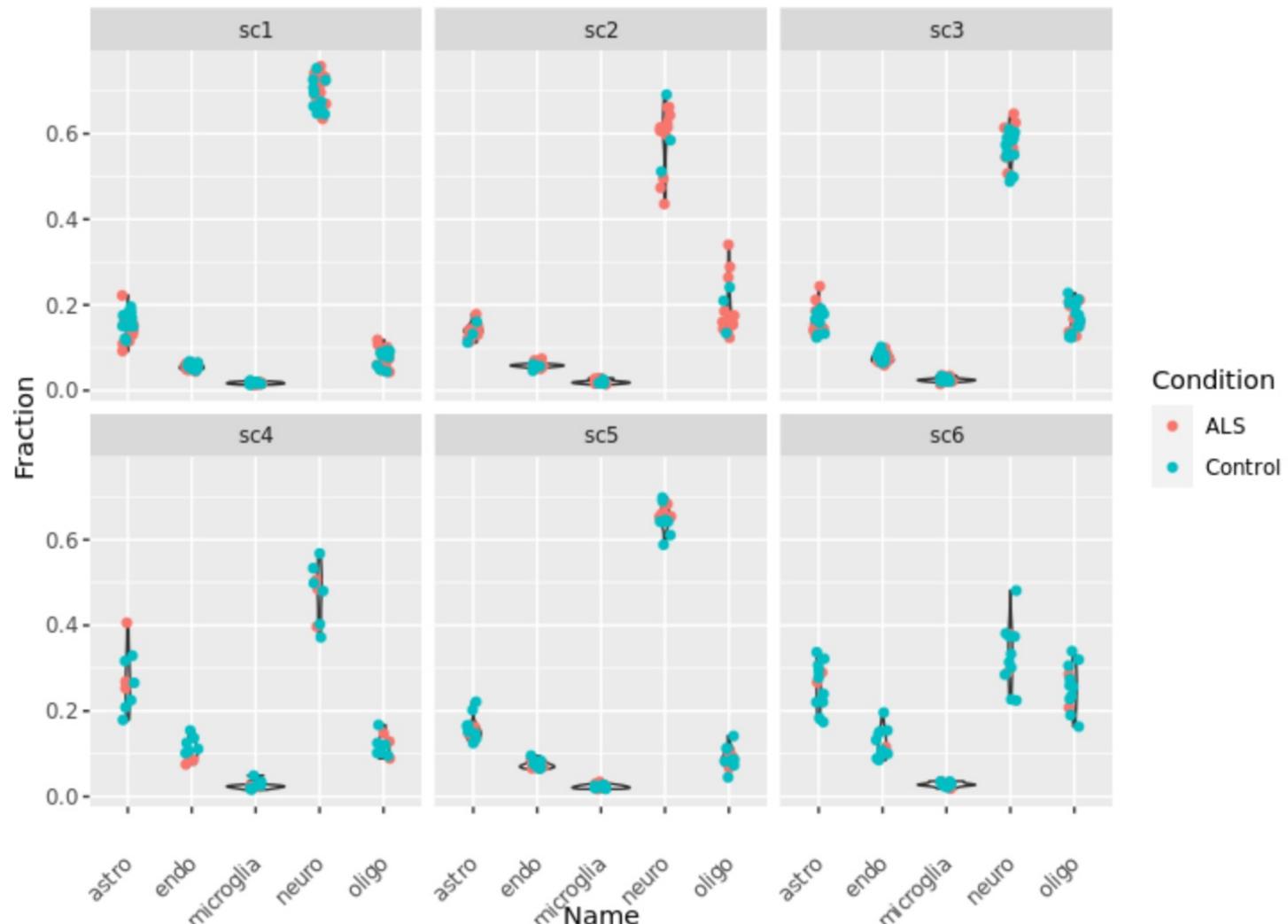


Example: bulk RNA-seq of brain tissue

- Sample clusters do not correlate with experimental factors
- Gene clusters can be associated to
  - cell type markers
  - gender

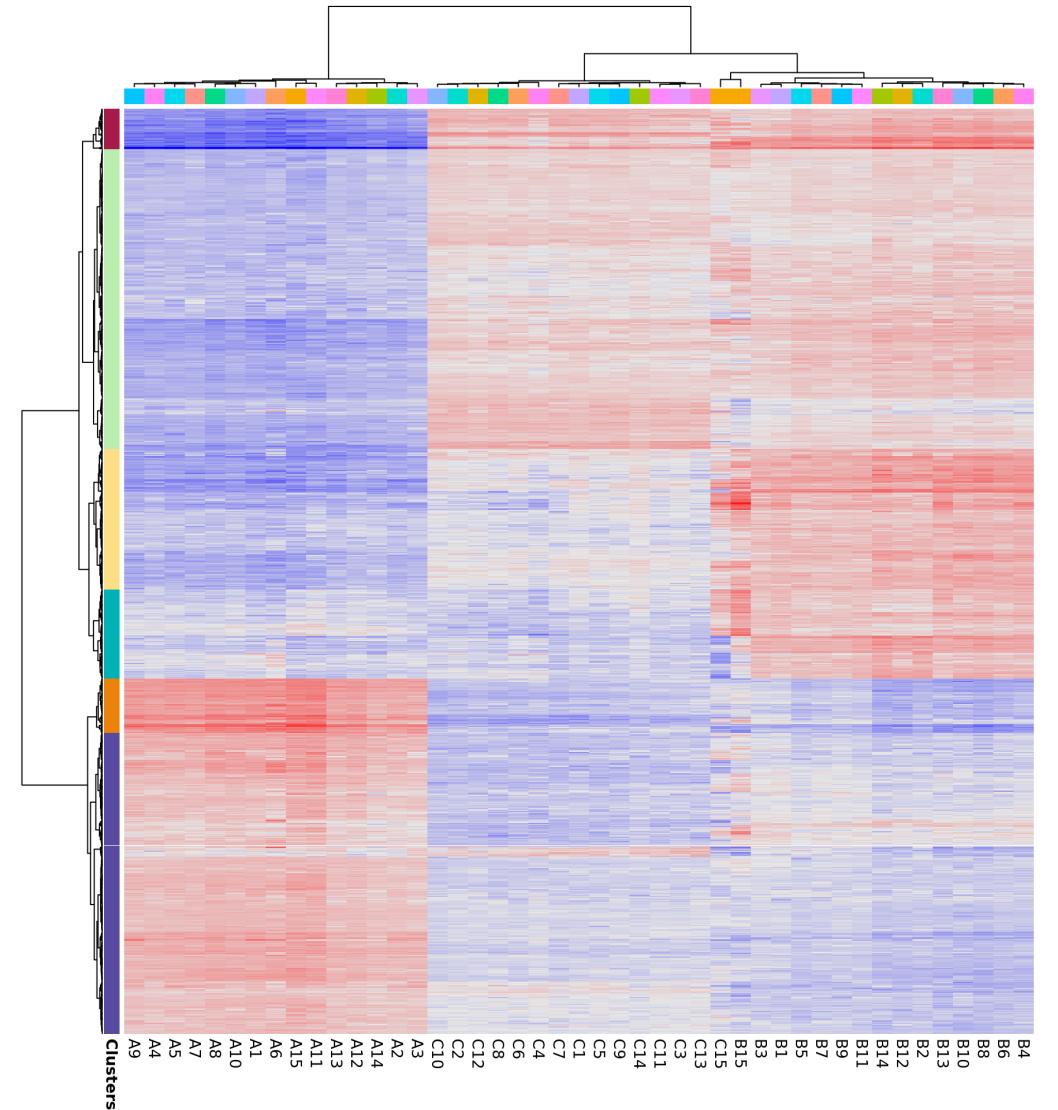


# Cell-type deconvolution



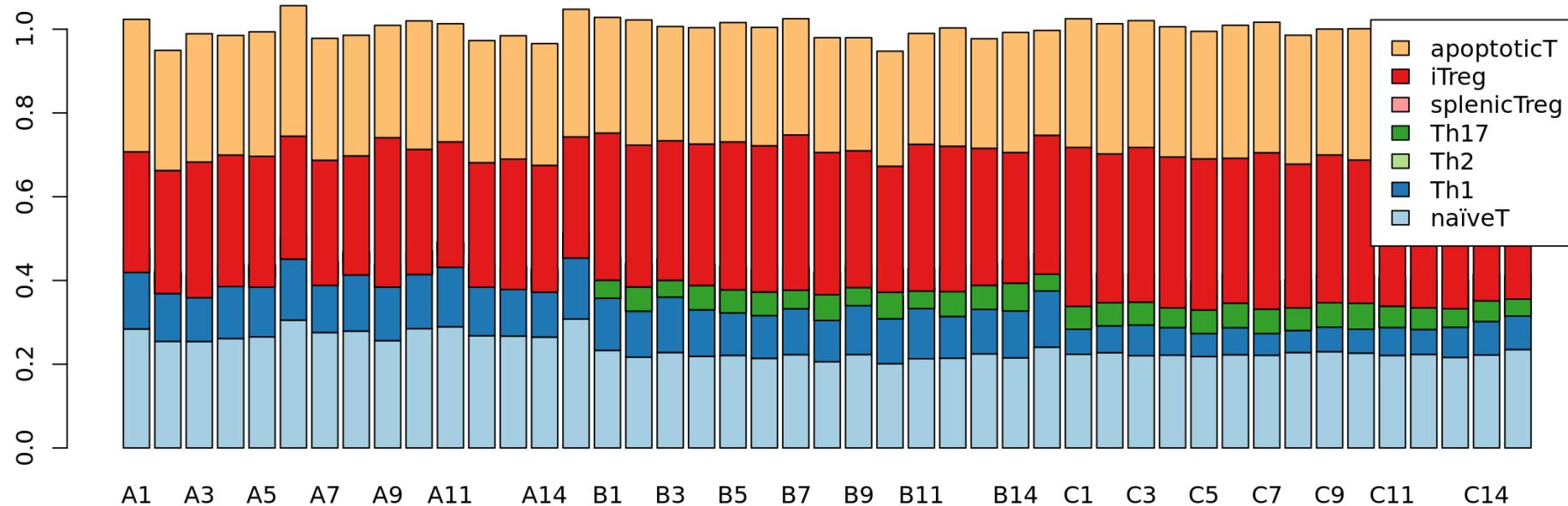
## Batch effect: Treated T-cells

- bulk RNA-seq of T-cells treated with different drugs
- Samples cluster according to batch



## Batch effect: Treated T-cells

- Batch A samples do not have the helper T cells Th17



# Data characteristics

- bulk RNA-seq:
  - 10k to 10Mio cells → billions of mRNA molecules → 100s of billions of fragments after fragmentation and amplification → sequencing samples ~20 Mio reads out of a much larger pool of fragments
  - 10-20k different genes detected in a sample
- single cell RNA-seq
  - 1 cell → 200k mRNA molecules
  - 70-90% of mRNA molecules are typically lost during library prep; precise numbers are unknown
  - library prep **overamplifies deliberately**, i.e. sequencing is expected to capture multiple reads that originate from the same mRNA fragments; without strong amplification the detection sensitivity of current technologies is too low
  - protocols include Unique Molecular Identifiers (UMIs) that allow to identify and deduplicate reads
  - 200 – 8000 different genes detected in a cell
  - transcription happens in **transcriptional bursts** → considerable cell-to-cell variability

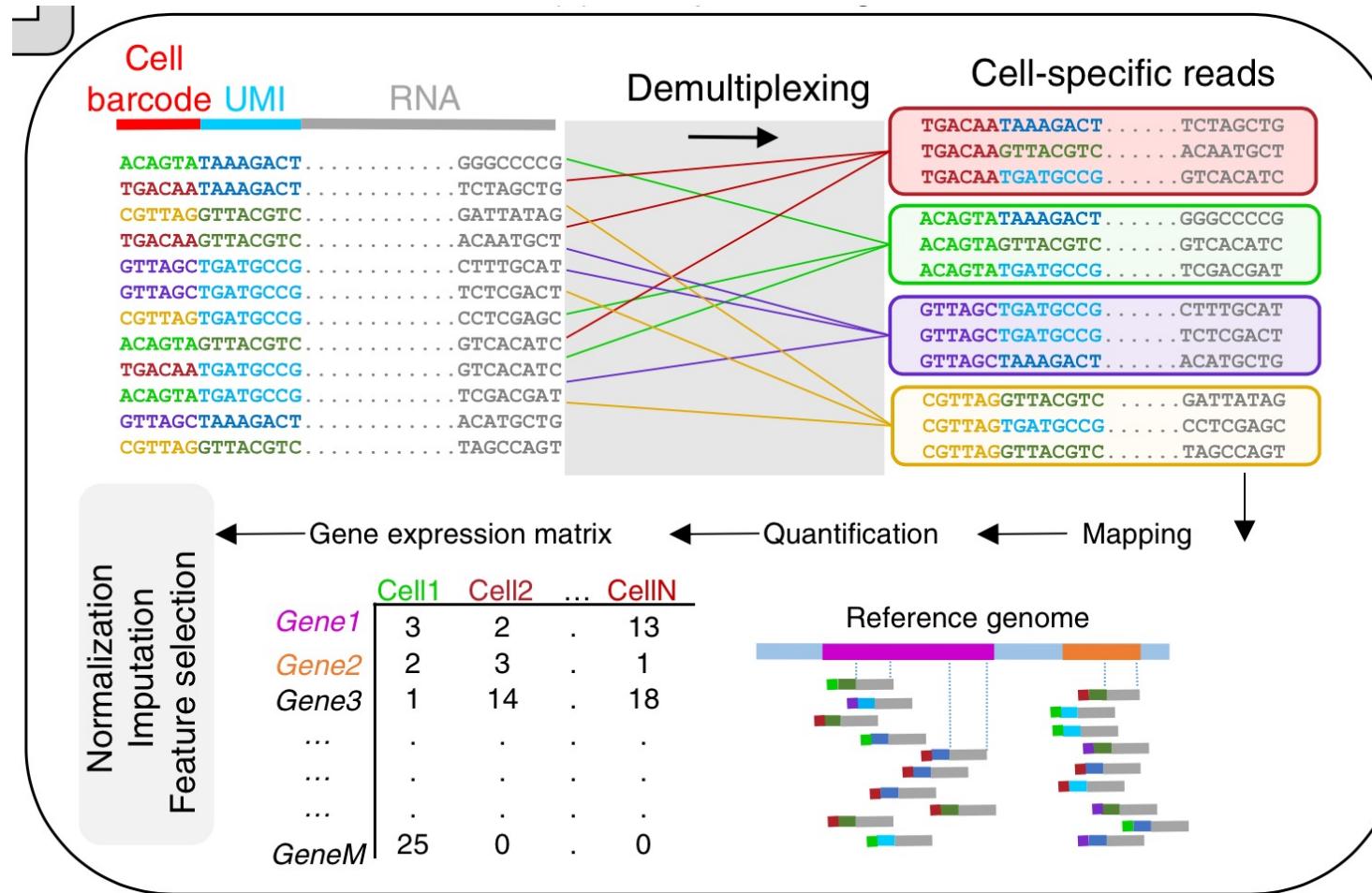


## UMIs and 3'-tagging

- The use of UMIs implies that only one end of the transcript can be captured! Expression quantification is only at the **gene level** not at the isoform level
- bulk RNA-seq:
  - typically coverage of entire transcript body
- single cell RNA-seq
  - transcript end tagging and UMI deduplication
  - or:
    - whole transcript coverage but increased variability caused by clonal reads



# Generating the single cell count matrix



- Expression matrix cells vs genes
- Measurement does not capture the identity of the cells
- Identity of cells needs to be derived from gene expression
- Up to 95% of the expression matrix may be zero

# Filtering of cells/barcodes

## Reasons for low quality cells

- Empty drops (ambient RNA)
  - Sequencing depth
  - Low quality bead
  - Non-viable or apoptotic cell
  - Nucleus only
  - Doublet
  - ...



## Empty drops

- The expression matrix contains also columns/barcodes that correspond to empty drops
  - Reads in empty drops correspond to free-floating RNA in the suspension of single cells
  - Empty drops are expected to have few reads and the reads in empty drops all come from the **same distribution of ambient RNA**
- Implemented in Bioconductor package *DropletUtils* and 10X Genomics CellRanger software

Method | [Open Access](#) | Published: 22 March 2019

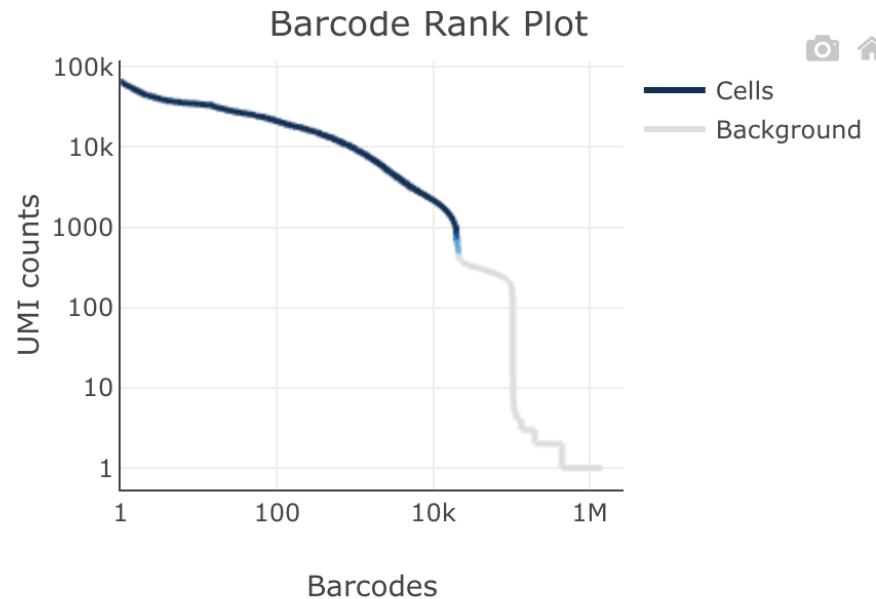
### EmptyDrops: distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data

Aaron T. L. Lun [✉](#), Samantha Riesenfeld, Tallulah Andrews, The Phuong Dao, Tomas Gomes, participants in the 1st Human Cell Atlas Jamboree & John C. Marioni [✉](#)

[Genome Biology](#) **20**, Article number: 63 (2019) | [Cite this article](#)



## EmptyDrops algorithm



- In a typical single-cell 10X experiment >1Mio barcodes are detected
- These are error-corrected barcodes; error-correction is done by matching barcode reads against *known-good* barcodes from 10X
- Vast majority do not represent cells but are empty drops capturing few ambient mRNA molecules

# EmptyDrops algorithm

The algorithm has two key steps:

1. Select x% of the barcodes with the highest total UMI counts → primary mode of cells with high RNA content cells.
  2. Select y% of the barcodes with lowest total UMI counts  
→ assume they represent empty cells and build ambient RNA model  
→ call barcodes that disagree with ambient model as valid cells

Barcodes are called cells if they

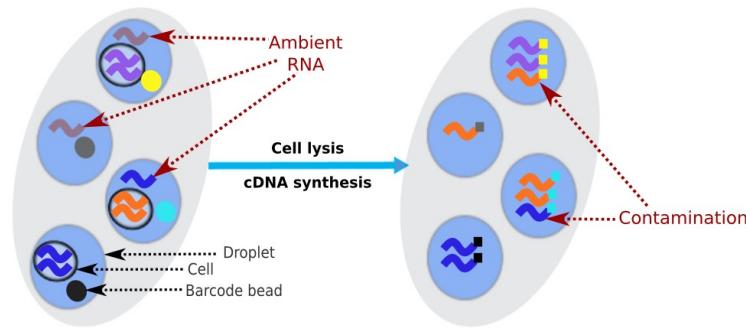
- have either sufficient reads (above knee-point)
  - if they are significantly different from background distribution



## Ambient RNA

- EmptyDrops removes barcodes that do contain only ambient RNA
- Ambient RNA still may affect expression of true cells by an additive contribution
- This may be relevant for cells that have low transcriptional activity
- Ambient RNA contributions may be removed using Bayesian approach.

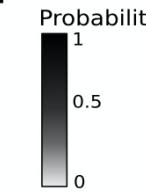
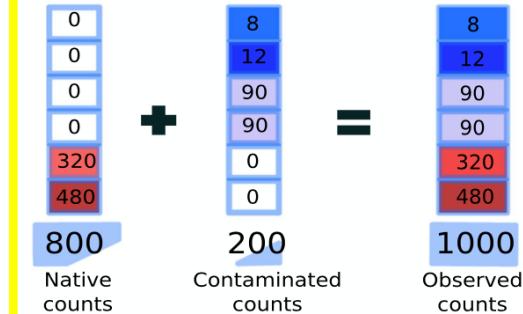
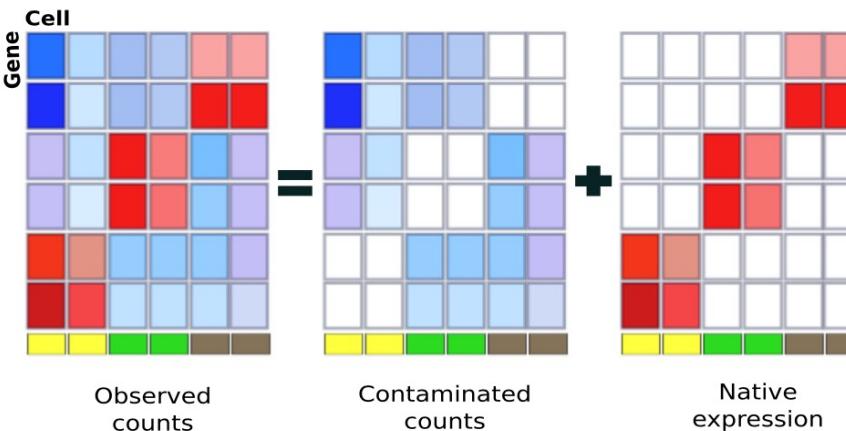
# Ambient RNA

**A**

**B**
**Expression distribution**

|     |     |     |
|-----|-----|-----|
|     |     | 0.4 |
|     |     | 0.6 |
|     | 0.5 |     |
|     | 0.5 |     |
| 0.4 |     |     |
| 0.6 |     |     |

**Contamination distribution**

|      |     |      |
|------|-----|------|
| 0.04 | 0.2 |      |
| 0.06 | 0.3 |      |
| 0.45 |     | 0.25 |
| 0.45 |     | 0.25 |
|      | 0.2 | 0.2  |
|      | 0.3 | 0.3  |

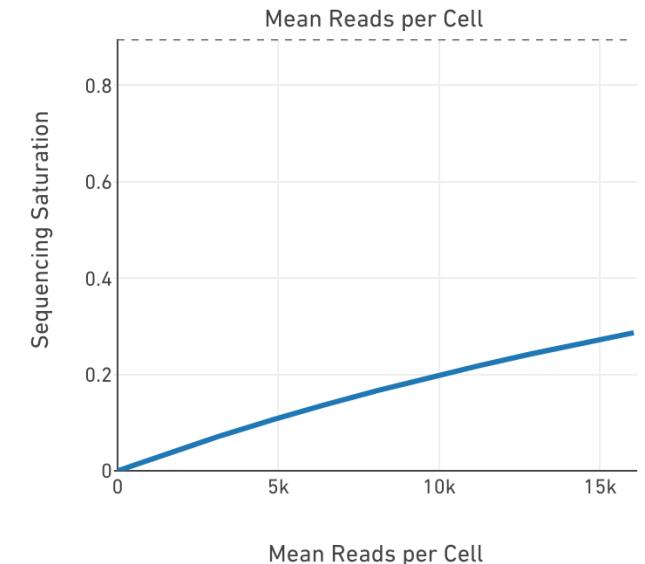
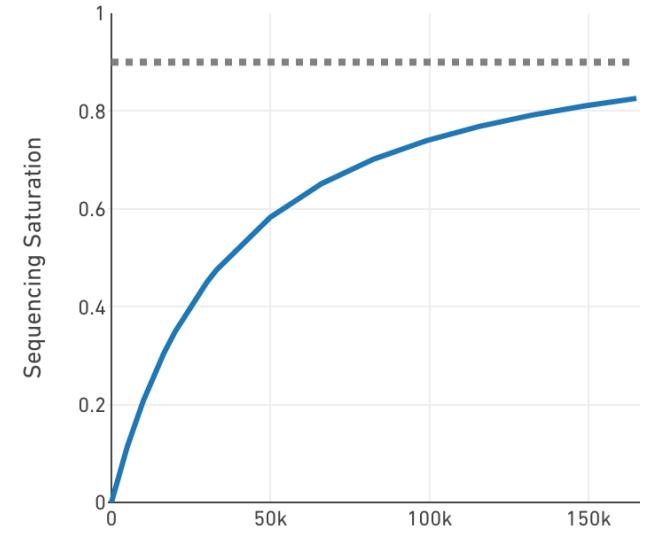

**Example cell from cluster 1**

**C**


- Ambient RNA originates from damaged or lysed cells
- Protocols include washing steps to remove ambient RNA but some ambient RNA may remain



# Sequencing Depth

- Based on resampling one can estimate how many genes would be additionally detected if more reads were sequenced
- Typically, 50'000 reads per cell are targeted, since the yield in terms of detected cells may vary the actual number of reads per cell can be very different





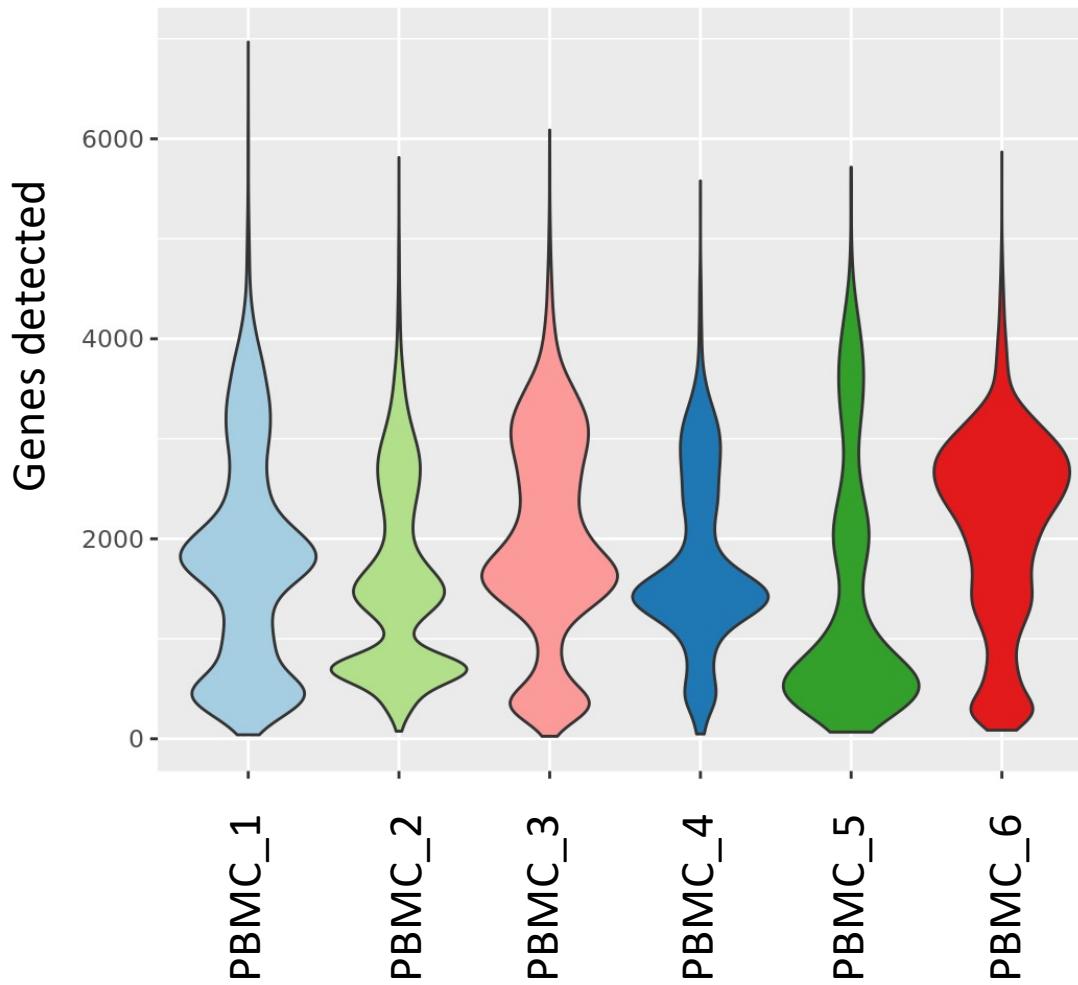
10

01

101

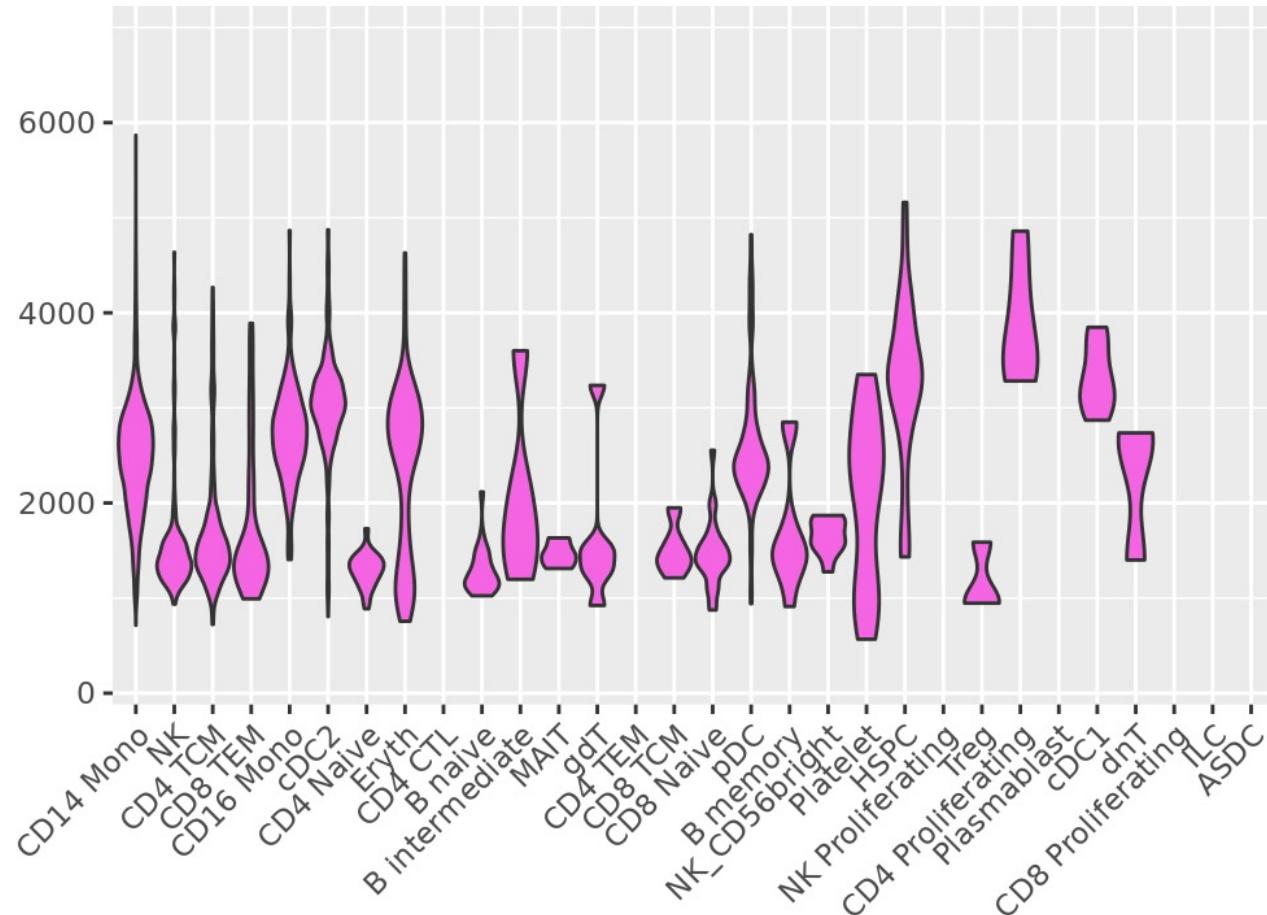


## Filtering example: Peripheral Blood Mononuclear Cells (PBMCs)



- distribution
  - is multi-modal
  - differs between samples

## Filtering example: PBMCs



- PBMCs consist of different cell types with varying overall transcriptional activity
- Risk: cell types with few genes detected may be filtered ab



## Quality control: Low quality bead

- Remove barcodes with few reads sequenced
- Remove barcodes with few genes detected
- Risk: might remove cells with low transcriptional activity

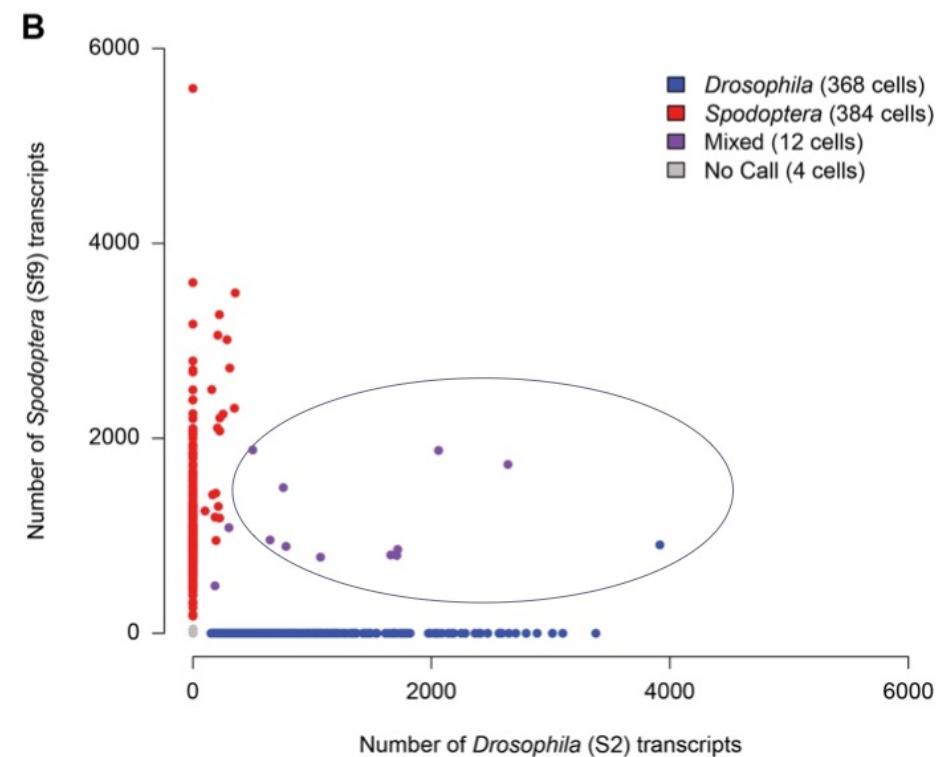
10  
01  
101101 1  
010 0  
0101 10functional genomics center zurich  
010 01  
101 10  
010 010  
1  
0  
0  
1  
1  
0  
0  
1  
1

## Quality control: Low quality cell

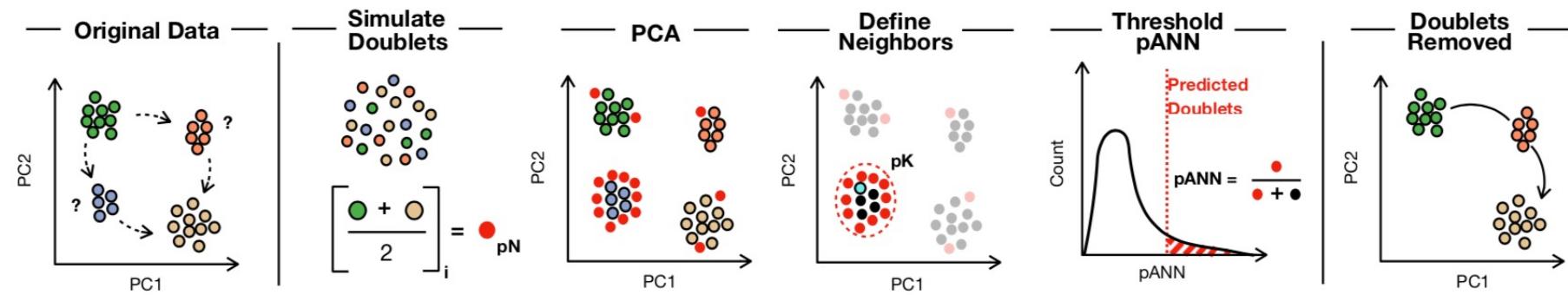
- Remove barcodes with
  - high fraction of mitochondrial genes (in apoptotic or partially lysed cells)
  - zero mitochondrial genes (cells that have lost the cytoplasm?)
- Potential confounder: number of mitochondria differs between cells
  - red blood cells: no mitochondria
  - liver cells: up to 2000 mitochondria

## 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 single-nucleus RNA-seq
- Test datasets for doublets consist of cells from two species



# Doublet Detection



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

<https://github.com/chris-mcginnis-ucsf/DoubletFinder>

## Quality control thresholds

- Quality metrics, like the number of reads or genes detected per cell, do depend on factors like
  - cell type
  - sequencing depth
  - overall quality
  - ...
- As a consequence thresholds for these metrics need to be adapted to individual samples and cell types



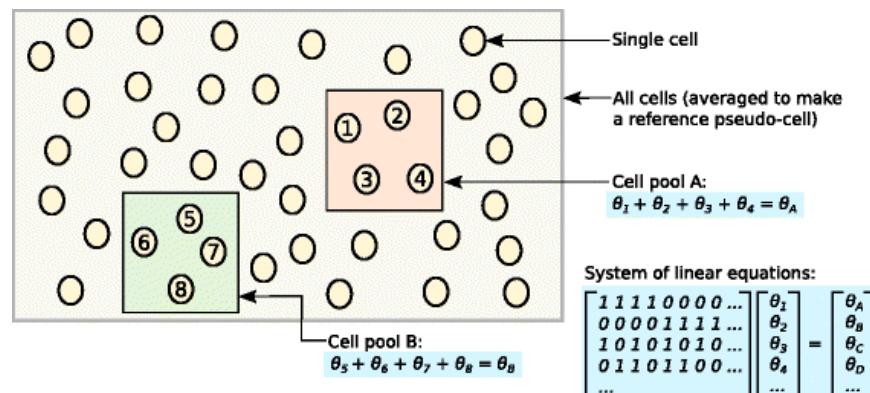
# Normalization

- Simple straightforward normalization:
  - scale by total number of counts per cell
  - transform the result with  $\log(x+1)$
- Expression is at Log-scales
- Avoids that zero counts become negative infinite at log scale



functional genomics center zurich  
01 1  
0 0  
010 01  
101 10  
010 01  
01 1  
10 0  
01 1

## scran: Normalization using pools of cells

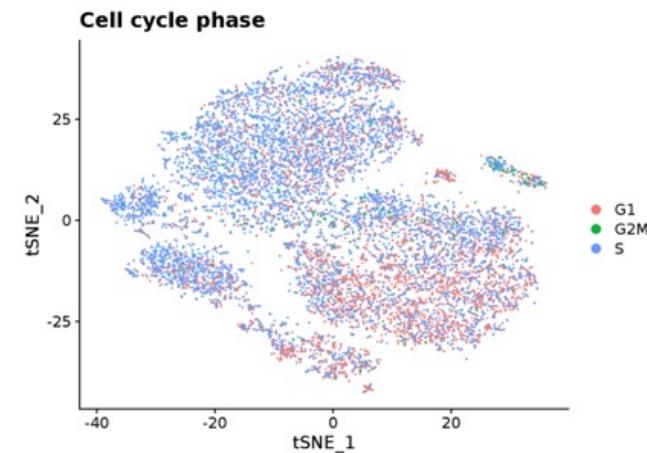


- Define a pool of cells
- Sum expression values across all cells in the pool
- Normalize the cell pool against an average reference, using the summed expression values
- Repeat this for many different pools of cells to construct a linear system
- Deconvolute the pool-based size factors to their cell-based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



# Normalization: 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 (high dimensionality) make classification more erroneous

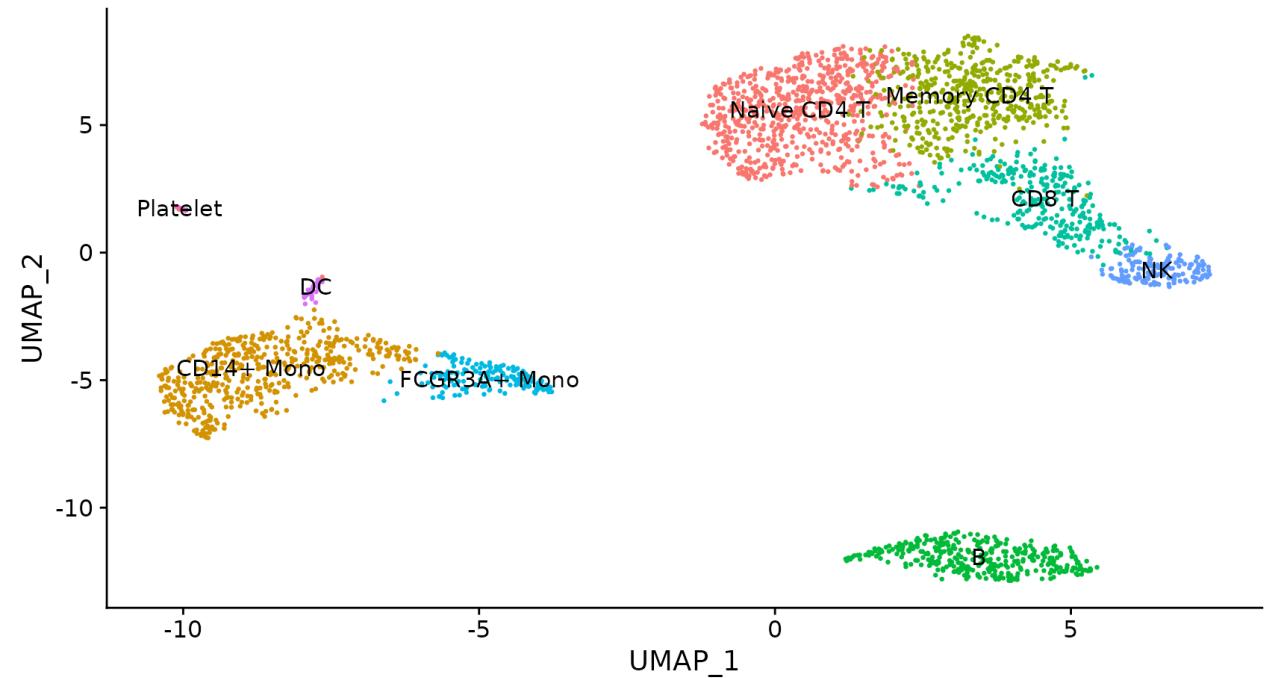


# Dimensionality Reduction Methods

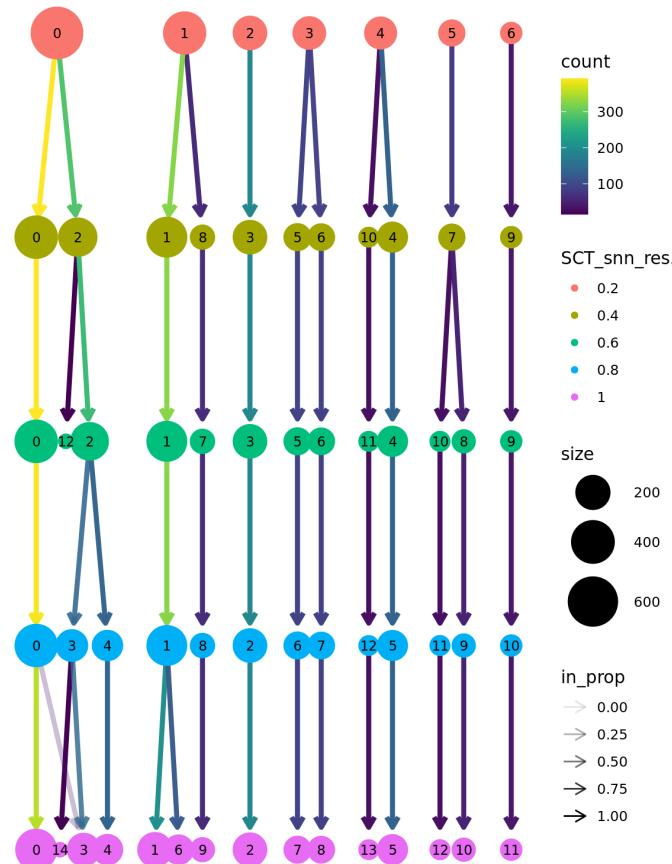
- Matrix Factorization
  - PCA
  - MDS – multi-dimensional scaling
  - ....
- Graph-based
  - t-SNE
  - UMAP
- Autoencoder (based on matrix factorization)

## Two-dimensional representation

- cells cluster according to **cell types**
- **cell type** is the major cause for variation in sc data



# Cluster resolution should match cell types of interest



- resolution of clustering controls cell type granularity
- Often there is no perfect clustering
- If in doubt, overclustering is preferred
  - it is easy to handle a cell type split in two clusters
  - it is cumbersome to have a cluster that contains two cell types



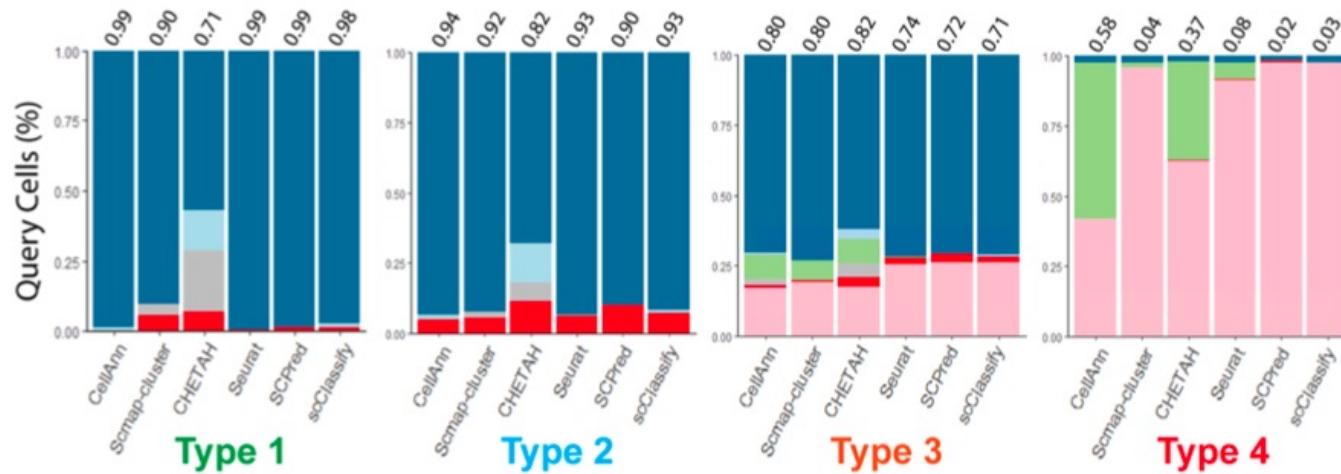
# Cell type annotation tools

|             | marker genes  | reference expression profiles   |
|-------------|---|---|
| per cell    | <ul style="list-style-type: none"><li>• <b>AUCell</b></li><li>• scROSHI</li></ul> | <ul style="list-style-type: none"><li>• <b>SingleR</b></li><li>• scArches</li><li>• <b>Azimuth</b></li><li>• scmap-cell</li></ul> |
| per cluster | <ul style="list-style-type: none"><li>• <b>enrichR</b></li></ul>                  | <ul style="list-style-type: none"><li>• <b>SingleR</b></li><li>• FR-Match</li><li>• CellAnn</li><li>• scmap-cluster</li></ul>     |

## Per-cell vs per-cluster labeling

- generally they should agree for the major cell-types
- discrepancies occur at sub-types
- often only a small percentage of cells (<10%) has single-cell labels that diverge from the cluster label
- major discrepancies might be explained by
  - cluster resolution inappropriate
  - reference data / marker gene set inappropriate
  - subtype of cells
- discrepancies suggest manual resolution

# Reference-based Annotation



Type 1:  
Query Ref

Type 2:  
Ref  
Query

Type 3:  
Query  
Ref

Type 4:  
Query Ref

- generally good performance
- situation where query and reference do not overlap are of a concern



## Algorithms:

- **CellAnn:** compute **correlations** of cluster averages, analyze correlations to define thresholds, compute rank-test to define sub-types
- **Seurat:** integrate data in UMAP space and assign nearby cell type
- **scmap-cell:** approximate **nearest neighbors**
- **SingleR:** identify variable genes in the reference set, **correlate** individual cells with the reference using the variable genes

# Single Cell Reference Data

- <https://atlas.brain-map.org/>
  - <https://www.humancellatlas.org/>
  - <https://tabula-sapiens-portal.ds.czbiohub.org/>
  - <https://www.flycellatlas.org/>
  - <https://azimuth.hubmapconsortium.org/references/>
  - Single Cell Expression Atlas: <https://www.ebi.ac.uk/gxa/sc/home>
    - individual studies
  - Chan Zuckerberg: <https://cellxgene.cziscience.com/collections>
  - BROAD:
    - [https://singlecell.broadinstitute.org/single\\_cell](https://singlecell.broadinstitute.org/single_cell)
  - ...

*Generally: The more similar to your system, the better*

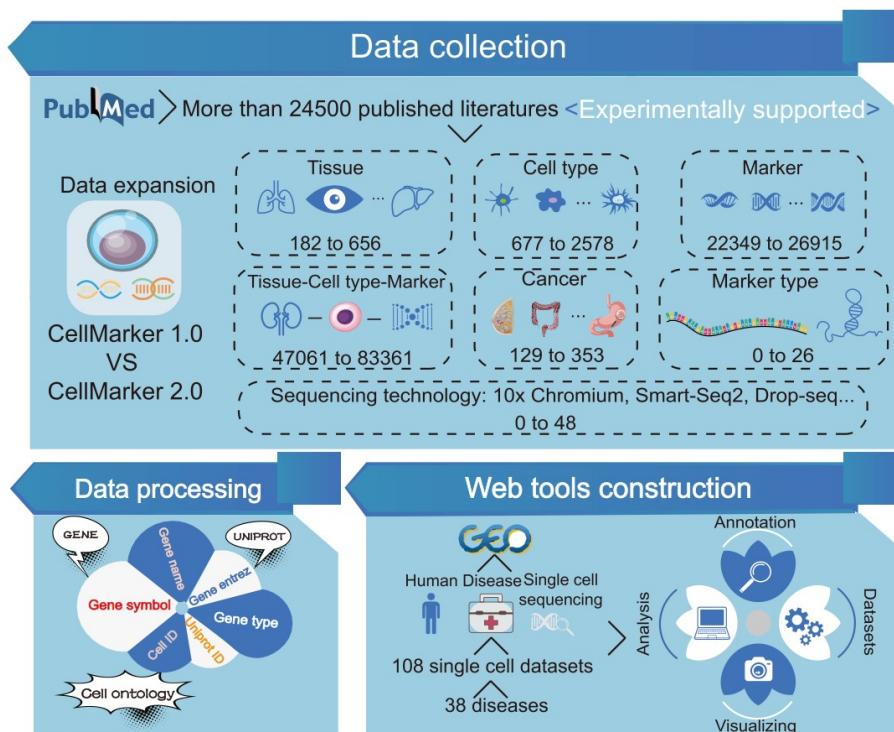


# Celltype Marker-Based Annotation

## Tools:

- **AUCell**: Area und the curve to estimate marker enrichment
- **enrichR**: compute cluster-markers, subsequently compute fuzzy enrichment among cell-type markers
- **scROSHI**: rank test to identify if cell-type markers are higher expressed than others

# CellMarker 2.0



# Welcome to CellMarker 2.0



# EnrichR

- database of databases
- diverse source of gene lists
- hosts also gene sets specific to cell types
- Gene sets are defined for human/mouse
- EnrichR databases for model organisms exist but do not include cell type specific gene sets

Description No description available (512 genes)

**Descartes Cell Types and Tissue 2021**

Lymphoid cells in Adrenal  
Lymphoid cells in Intestine  
Lymphoid cells in Placenta  
Lymphoid cells in Lung  
Lymphoid cells in Kidney

**CCLE Proteomics 2020**

CMK HAEMATOPOIETIC AND LYMPHOID TISSUE  
SUDHL4 HAEMATOPOIETIC AND LYMPHOID TISSUE  
KYSE70 OESOPHAGUS Tissue Px18  
U937 HAEMATOPOIETIC AND LYMPHOID TISSUE  
KARPAS299 HAEMATOPOIETIC AND LYMPHOID TISSUE

**Tabula Sapiens**

Kidney-nk Cell  
Lung-nk Cell  
Bone Marrow-nk Cell  
Spleen-nk Cell  
Thymus-nk Cell

**HuBMAP ASCTplusB augmented 2022**

Effector Memory CD8-positive, Alpha-Beta T  
Mature Natural Killer cell - Blood  
CD8 T Cell memory - Blood  
Mature CD8 T cell - Blood  
CD16-negative CD56-bright Natural Killer ce

**CellMarker Augmented 2021**

Natural Killer T (NKT) cell:Fetal Kidney  
CD4+ Cytotoxic T cell:Liver  
Effector CD8+ Memory T (Tem) cell:Peripheral  
Natural Killer cell:Splenic Red Pulp  
Natural Killer cell:Lung

**PanglaoDB Augmented 2021**

NK Cells  
T Cells  
Natural Killer T Cells  
T Memory Cells  
Decidual Cells

**Tabula Muris**

Natural Killer Cell Fat CL:0000623  
Natural Killer Cell Lung CL:0000623  
Unknown Brain Non-Microglia  
Natural Killer Cell Marrow CL:0000623  
Natural Killer Cell Liver CL:0000623

**Human Gene Atlas**

CD56+ NKCells  
CD8+ Tcells  
CD4+ Tcells  
PrefrontalCortex  
WholeBlood

**Azimuth Cell Types 2021**

Natural Killer CL0000623  
CD56-dim Natural Killer 2 CL0000939  
CD56-dim Natural Killer 3 CL0000939  
CD56-dim Natural Killer CL0000939  
Natural Killer Cell CL0000623

**ProteomicsDB 2020**

Hematopoietic PBMC BTO:0001008 234 repl  
Hematopoietic PBMC BTO:0001008 234 repl

**Mouse Gene Atlas**

NK cells  
follicular B-cells  
spleen  
mast cells  
lymph nodes

**ARCS4 Tissues**

NATURAL KILLER CELLS  
LYMPHOCYTE  
CD4+ T CELL  
BLOOD PBMC  
PERIPHERAL BLOOD



University of  
Zurich UZH

10  
01  
101

functional genomics center zurich

01 1  
10  
01 1  
101 10  
010 01  
010 01  
10 0  
01 1

## Azimuth

- Next to reference scRNA-seq data, Azimuth hosts also corresponding cell type markers

▼ celltype.l2 (default)

| Label          | Expanded Label      | OBO Ontology ID | Markers   |
|----------------|---------------------|-----------------|---|
| B intermediate | Intermediate B cell | mature B cell   | MS4A1, TNFRSF13B, IGHM, IGHD, AIM2, CD79A, LINC01857, RALGPS2, BANK1, CD79B |
| B memory       | Memory B cell       | memory B cell   | MS4A1, COCH, AIM2, BANK1, SSPN, CD79A, TEX9, RALGPS2, TNFRSF13C, LINC01781  |
| B naive        | Naive B cell        | naive B cell    | IGHM, IGHD, CD79A, IL4R, MS4A1, CXCR4, BTG1, TCL1A, CD79B, YBX3             |



## Summary

- Quality control is an essential step in single cell processing
- Cell-type annotation typically requires manual curation
- Major cell-type are well automatically detected, subtypes however are often misclustered and misassigned in automatic workflows