# A beginner's guide to scRNAseq analysis

Florian Heyl

In cooperation with

nfdi

# Questions

**(1) Laboratory:**
- What is the difference between bulk and single cell sequencing?
- How do you perform a single cell experiment?
- What confounders do I have?

**(2) Bioinformatics:**
- What data do I get?
- What should I check for my single cell data?
- What can I do with my data?

**(3) Products:**
- Why should I automatize my data analysis?
- What is important for a workflow?
- What open challenges do we still have?

GHGA THE GERMAN HUMAN GENOME- PHENOME ARCHIVE
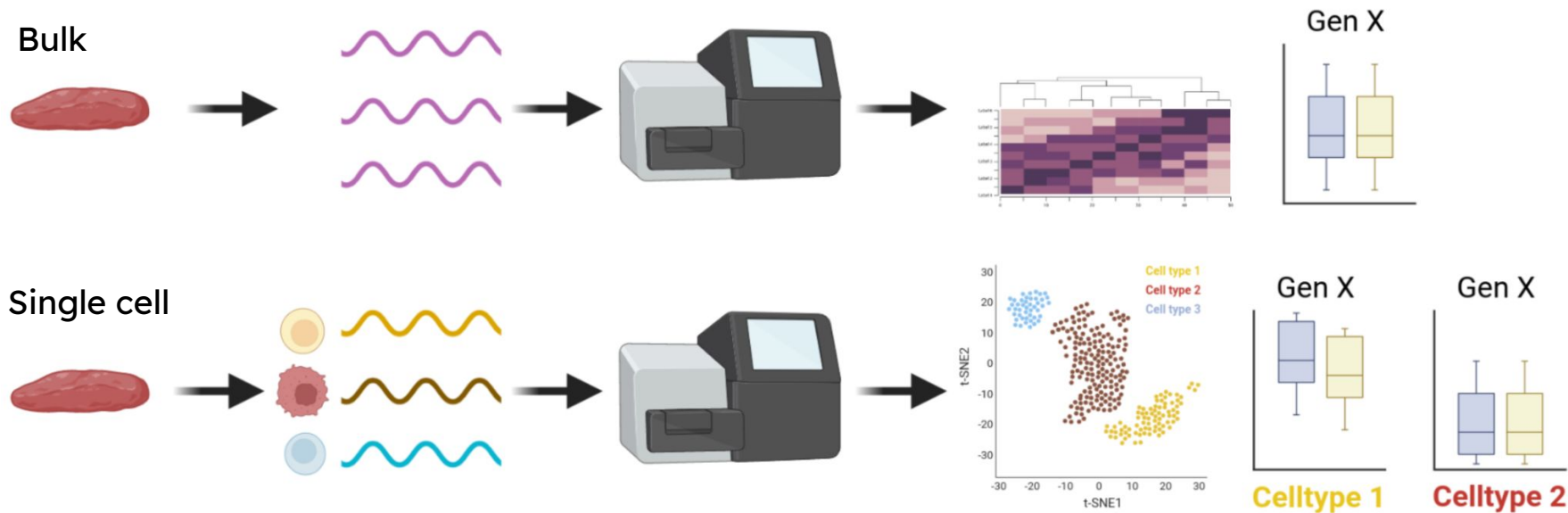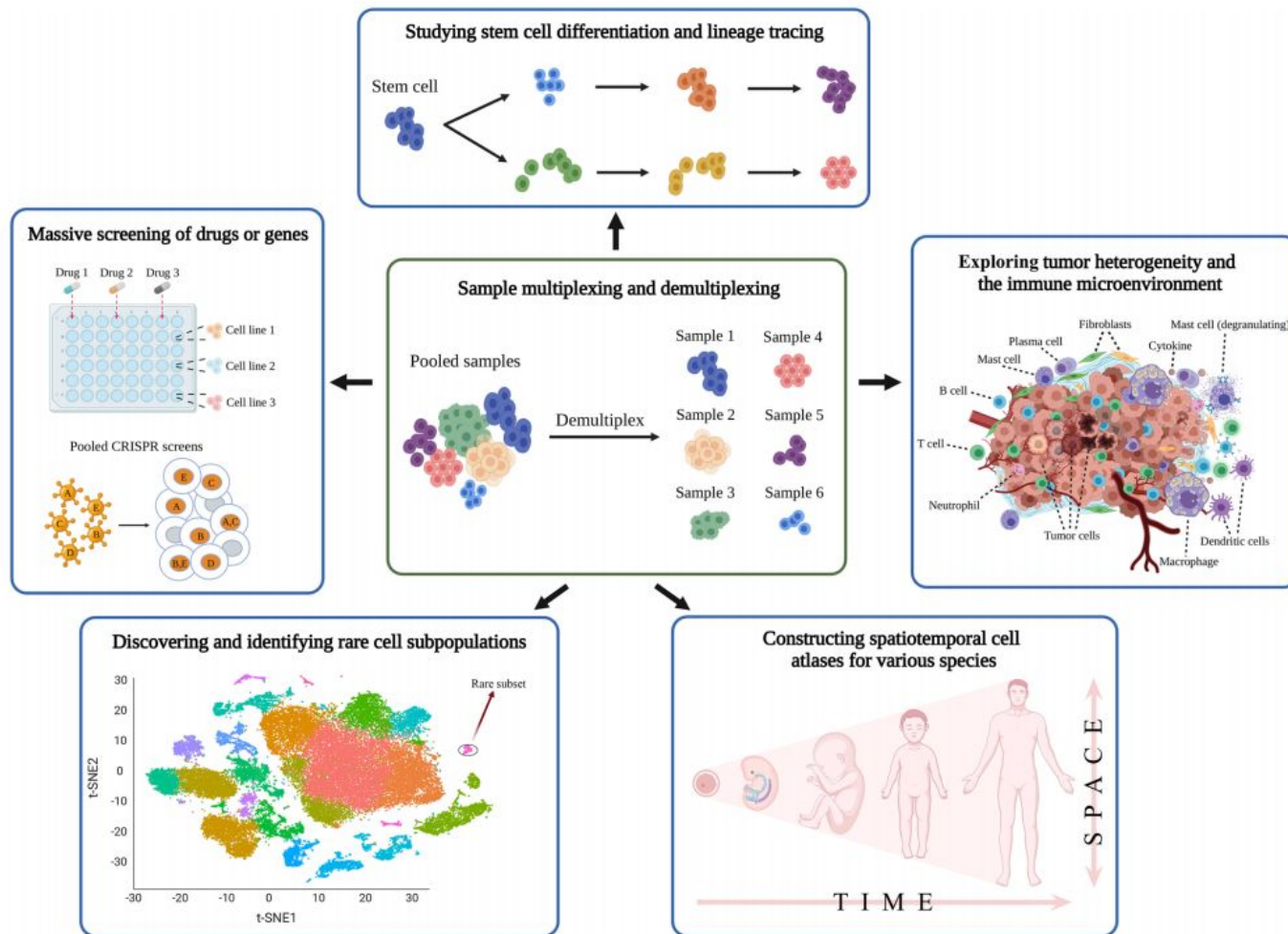
# Slides will be available

Take home message

Literature

(1) What is single cell sequencing?

**Bulk seq.:** Average-based expression profile
**Single cell seq.:** Cell level expression profile
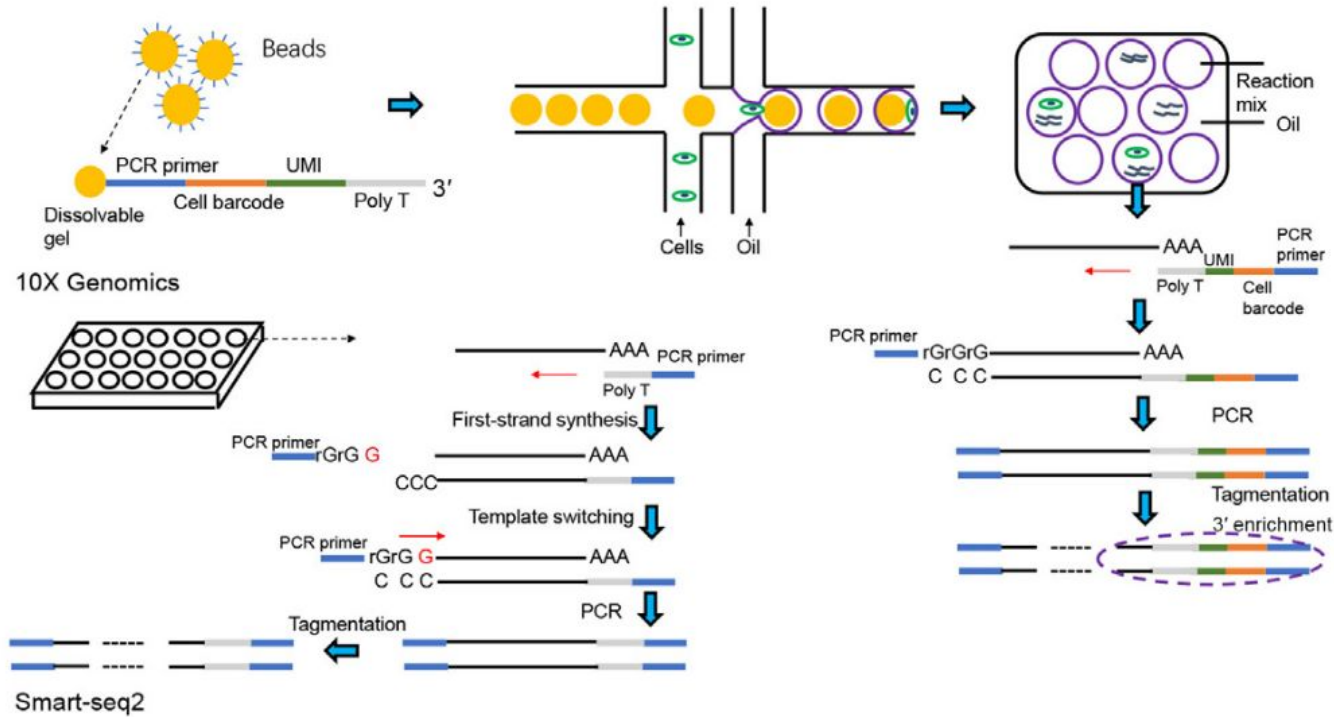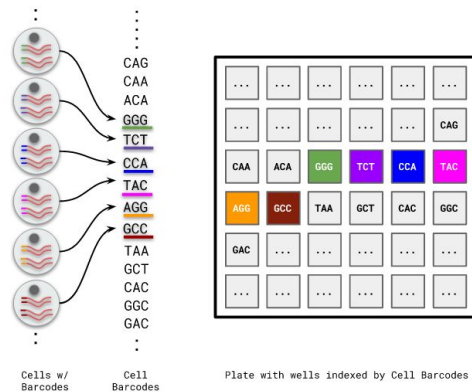
Yulong Zhang et al. (2022)

**Plate-based:** cells into wells on a plate
**Droplet-based:** each cell in its own microfluidic droplet
Each cell is a sample which cannot be replicated.

Jeanette Baran-Gale et al. (2018)
Malte D. Luecken & Fabian J. Theis (2019)
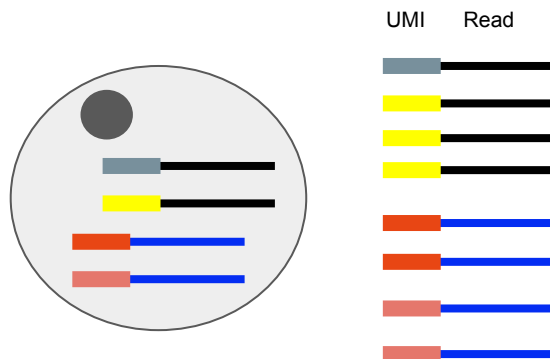Xiliang Wang et al. (2021)

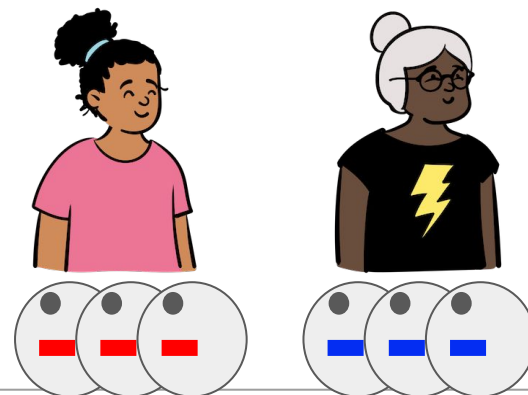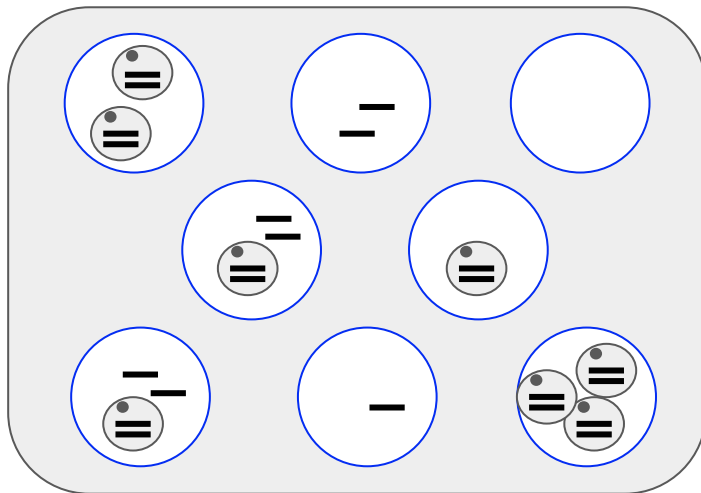| Cell Barcodes | Unique Molecular Identifier (UMI) | "Donor" (Multiplex) Barcode |
|---|---|---|
| To connect read (e.g., RNA) to a cell. | To reduce amplification bias (keep unique reads). In example: each gene (black & blue) has just two reads. | To connect read (e.g., RNA) to a donor (e.g., patient). |

Mehmet Tekman / Galaxy Training material (An introduction to scRNA-seq data analysis)
Eric Vallabh Minikel (2012) How PCR duplicates arise
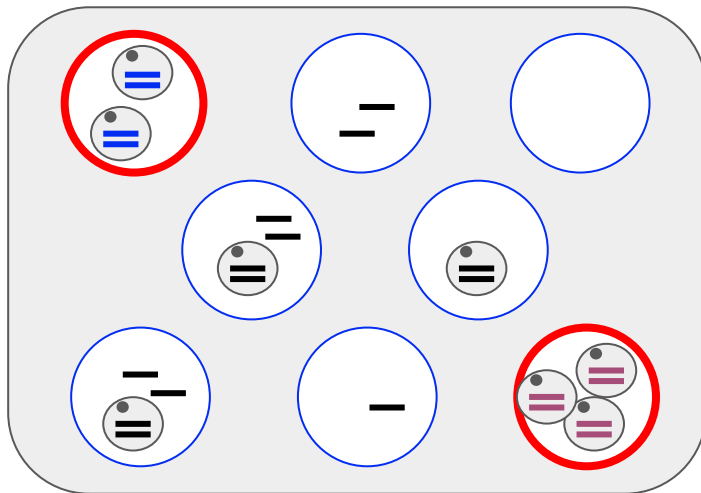Yulong Zhang et al. (2022)

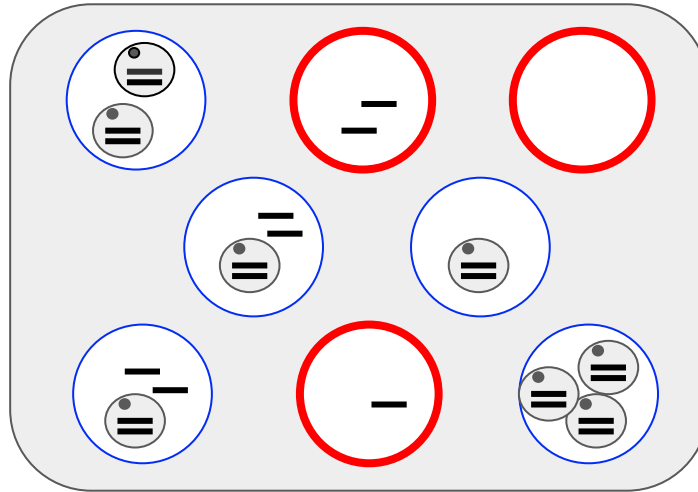# Confounders

# Problem 1: Doublets/Multiplets



**Two or more cells are sequenced together**
A high (e.g., RNA) count or number of detected regions is the result.

Malte D. Luecken & Fabian J. Theis (2019)
Samuel L. Wolock et al. (2019)
Tallulah S. Andrews et al. (2021)
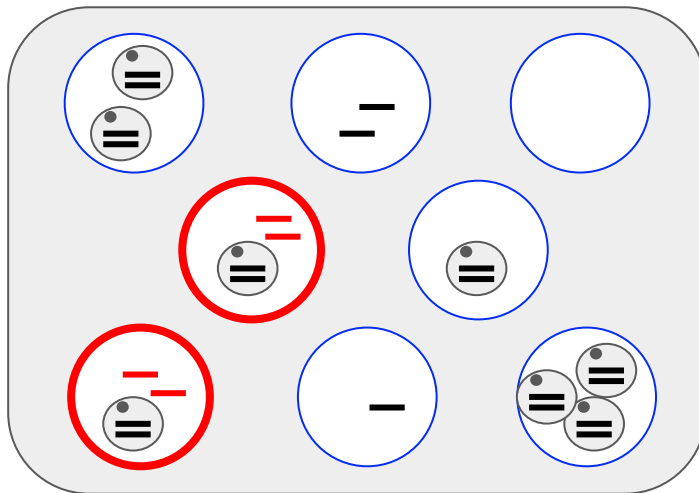
# Problem 2: Empty droplets/wells



**"Broken" cells or no cell can be collected**
A low (RNA) count, few detected genes, and a high fraction of mitochondrial counts can be the result.

Malte D. Luecken & Fabian J. Theis (2019)
Aaron T. L. Lun et al. (2019)
Tallulah S. Andrews et al. (2021)
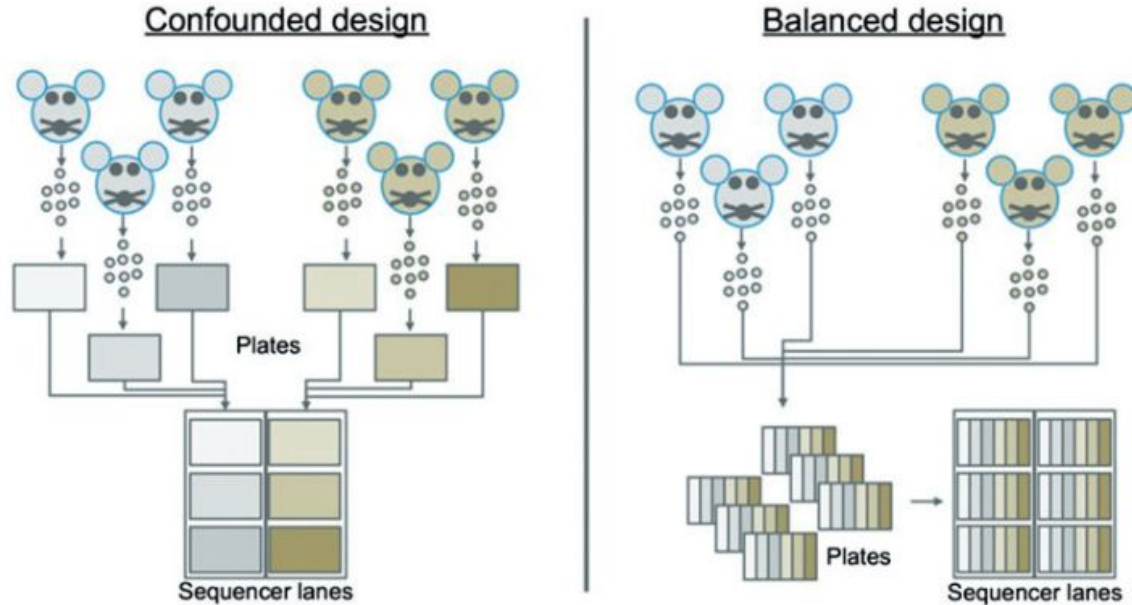
# Problem 3: Ambient RNA



**Counts that do not originate from the true barcoded cell, but from other lysed cells.**
Can lead to overrepresentation of some cell clusters (spurious clusters), higher cluster overlaps, higher gene coverage.

Malte D. Luecken & Fabian J. Theis (2019)
Shiyi Yang et al. (2020)
Tallulah S. Andrews et al. (2021)
Emre Caglayan et al. (2022)
Stephen J. Fleming et al. (2022)
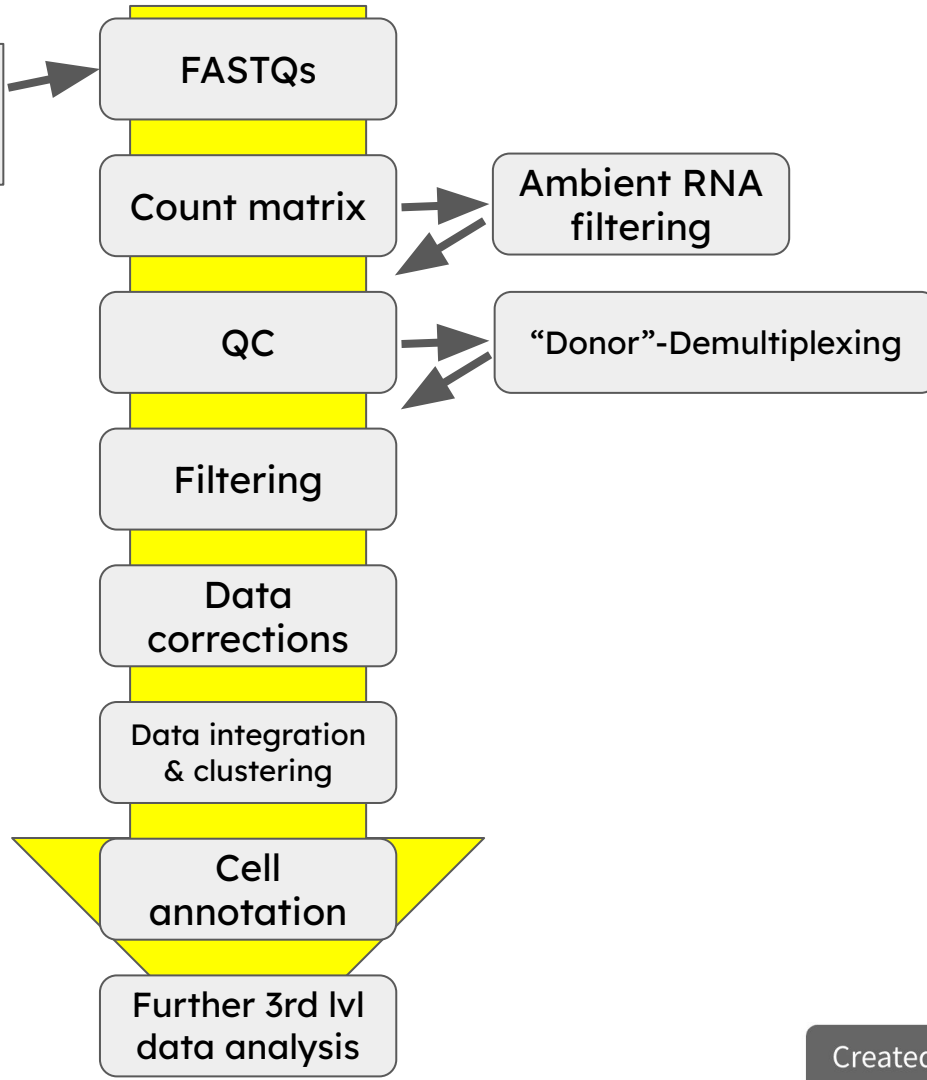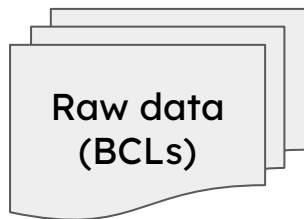
# Problem 4: Batch effects



**Sequencing your experiments in batches**
Might lead to spurious results (e.g., clusters or correlations)

Jeanette Baran-Gale et al. (2018)
Malte D. Luecken & Fabian J. Theis (2019)
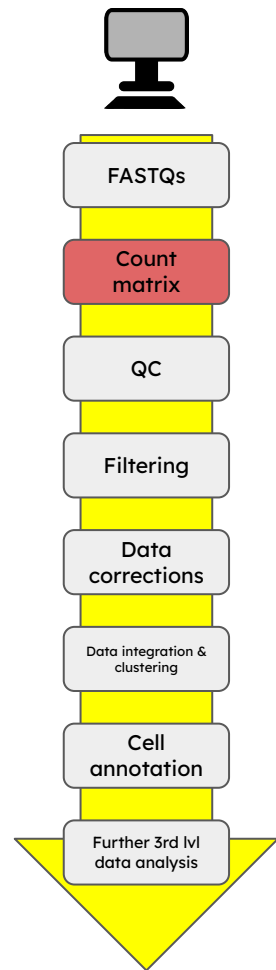Tallulah S. Andrews et al. (2021)

# (2) How can I use single cell sequencing data?

Raw data (BCLs)

FASTQs

Count matrix

Ambient RNA filtering

QC

"Donor"-Demultiplexing

Filtering

Data corrections

Data integration & clustering
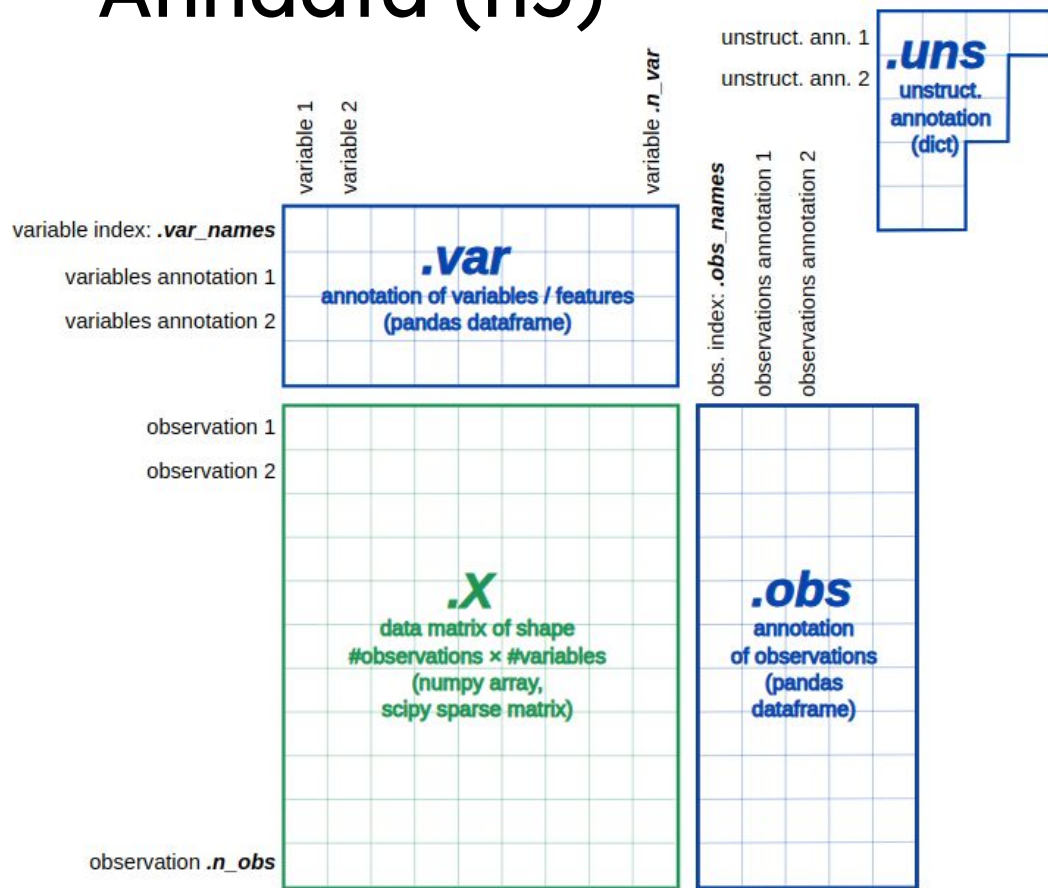
Cell annotation

Further 3rd lvl data analysis

# Cellranger

1. Read trimming
2. Genome/Transcriptome alignment
3. MAPQ adjustment
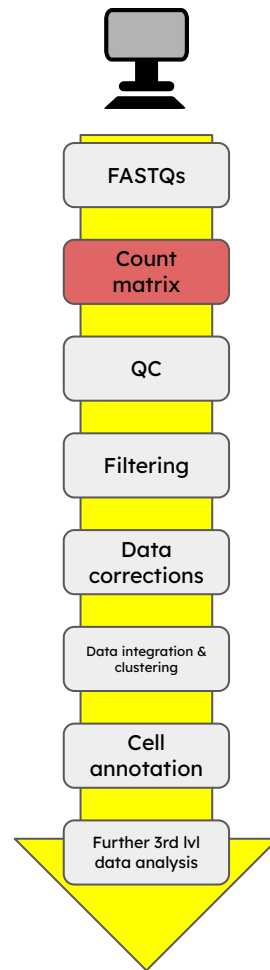4. 10x barcode correction
5. UMI counting
6. Calling cell barcodes

**10x Genomics** (Gene Expression Algorithms Overview)
Ralf Schulze Brüning et al. (2022)

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

# Anndata (h5)



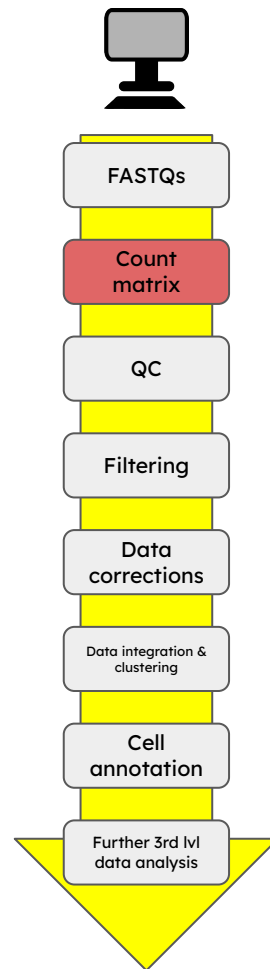Bérénice Batut et al. / Galaxy Training material (Clustering 3K PBMCs with Scanpy)

# Anndata (h5)

- Cell labels
- QC
- Sample IDs

Gene

| | A | B | C | D | E |
|---|---|---|---|---|---|
| | 1 | 0 | 0 | 5 | 0 |
| | 0 | 2 | 0 | 4 | 4 |
| | 1 | 1 | 0 | 0 | 0 |
| | 0 | 0 | 3 | 0 | 3 |
| | 0 | 0 | 0 | 2 | 0 |

Cell

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

Bérénice Batut et al. / Galaxy Training material (Clustering 3K PBMCs with Scanpy)

# Quality control of cells / barcodes



Malte D. Luecken & Fabian J. Theis (2019)
Tallulah S. Andrews et al. (2021)
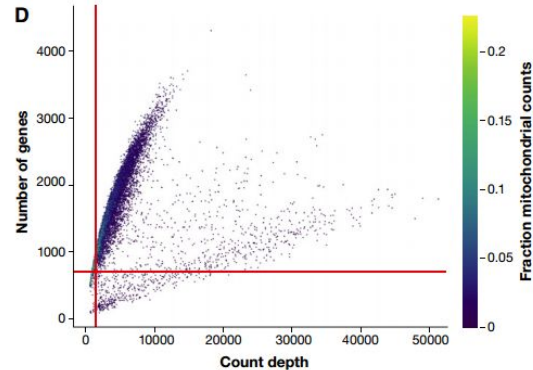Rui Hong et al. (2022)

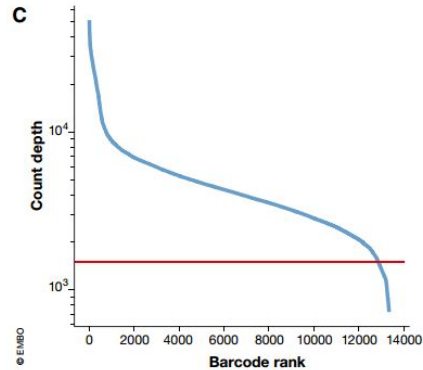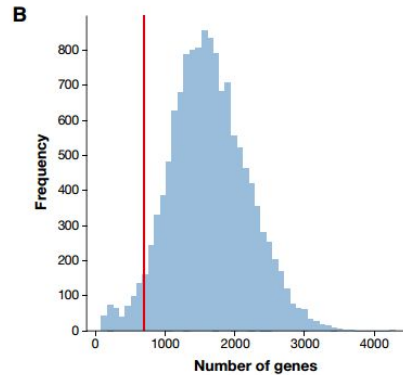FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

# Quality control of experiment



10x Genomics (Web summary)

# Quality control of experiment



galaxy_dataset_3c1f0a62-119f-46ed-b18b-d65a8a57ac98.dat

Lucille Delisle et al. / Galaxy (ATAC-Seq data analysis)

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

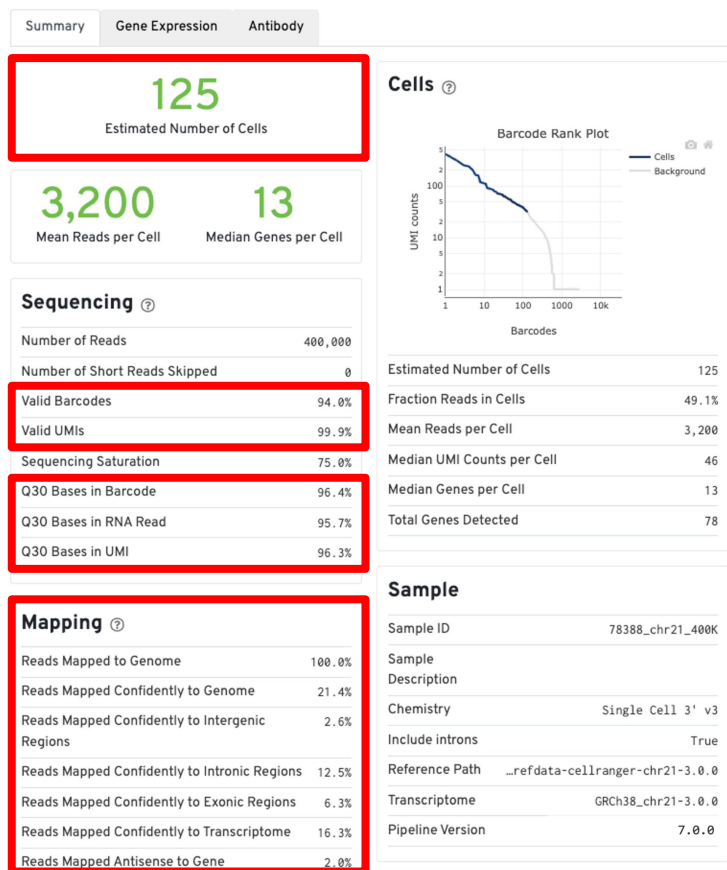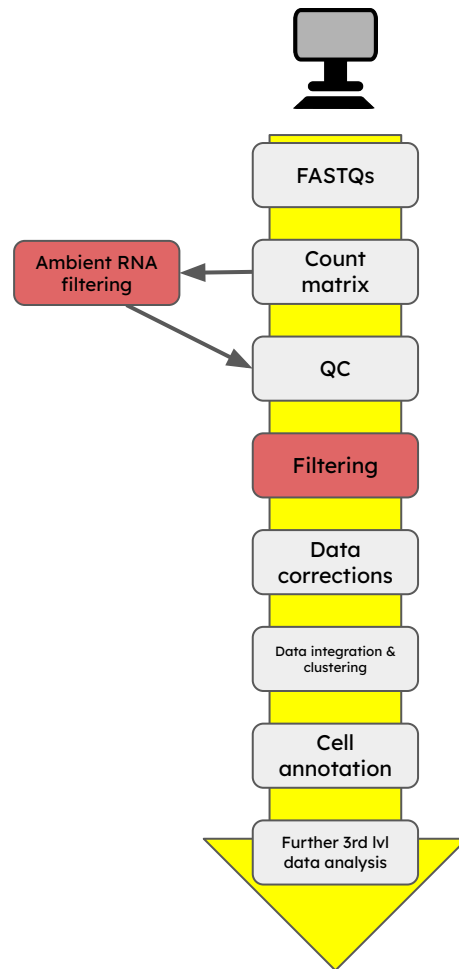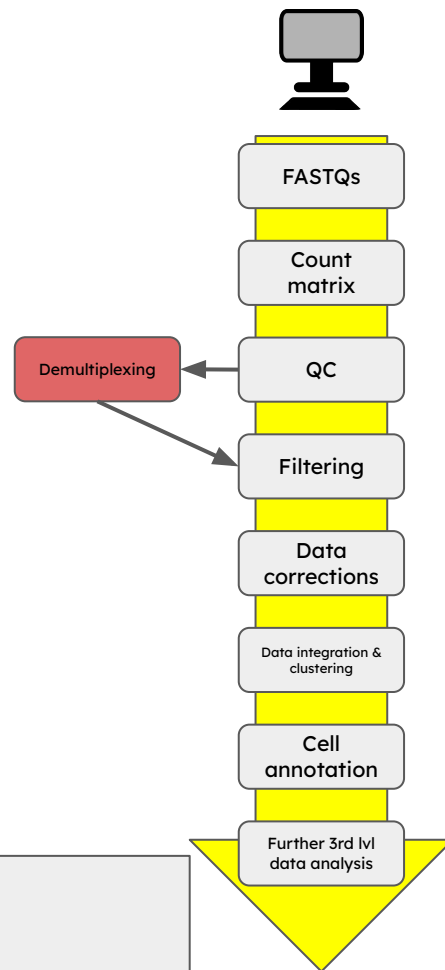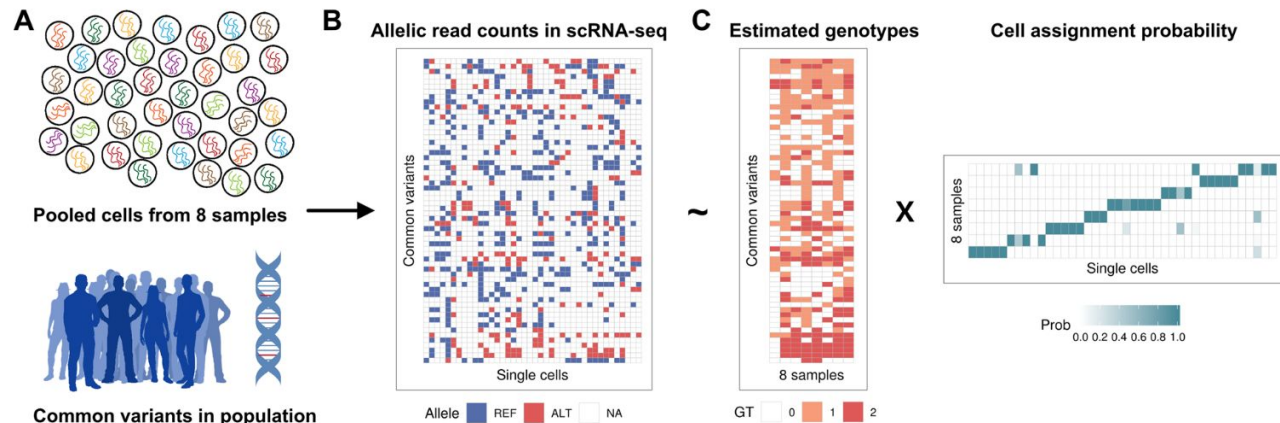Cell annotation

Further 3rd lvl data analysis

1. Emptyplets/Douplets/Multiplets
2. Ambient RNA
3. Experimental important QC measures (e.g., fragment sizes scATAC)

# "Donor"-Demultiplexing



A. Pooled cells from 8 samples / Common variants in population
B. Allelic read counts in scRNA-seq — Common variants / Single cells — Allele: REF, ALT, NA
C. Estimated genotypes — Common variants / 8 samples — GT: 0, 1, 2 × Cell assignment probability — 8 samples / Single cells — Prob 0.0 0.2 0.4 0.6 0.8 1.0

Pipeline: FASTQs → Count matrix → QC → Demultiplexing → Filtering → Data corrections → Data integration & clustering → Cell annotation → Further 3rd lvl data analysis

**Demultiplexing approach depends on experimental design.**

Yuanhua Huang et al. (2019)
Xianjie Huang and Yuanhua Huang (2021)
Yulong Zhang et al. (2022)
Joseph F. Cardiello et al. (2022)

# Normalization



Cell

| | | | | |
|---|---|---|---|---|
| 1 | 0 | 0 | 5 | 0 |
| 0 | 2 | 0 | 4 | 4 |
| 1 | 1 | 0 | 0 | 0 |
| 0 | 0 | 3 | 0 | 3 |
| 0 | 0 | 0 | 2 | 0 |

Size factors

| |
|---|
| 1.5 |
| 2.0 |
| 1.3 |
| 2.3 |
| 4.1 |

Malte D. Luecken & Fabian J. Theis (2019)
Nicholas Lytal et al. (2020)
Tallulah S. Andrews et al. (2021)

Brenda Marquina-Sanchez et al. (2020)
Xin Wang et al. (2021)

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

# Spike-In Normalization

**Reference Cell**



**Reference RNA**



FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

Malte D. Luecken & Fabian J. Theis (2019)
Nicholas Lytal et al. (2020)
Tallulah S. Andrews et al. (2021)

Brenda Marquina-Sanchez et al. (2020)
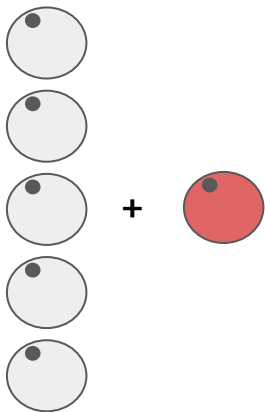Xin Wang et al. (2021)

# Batch correction



No batch correction — Batch correction

Malte D. Luecken & Fabian J. Theis (2019)
Tallulah S. Andrews et al. (2021)

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis
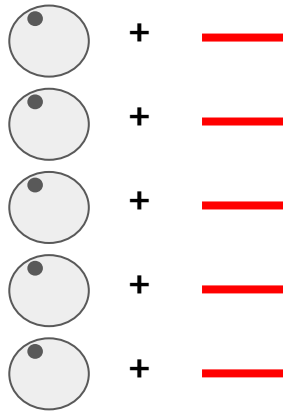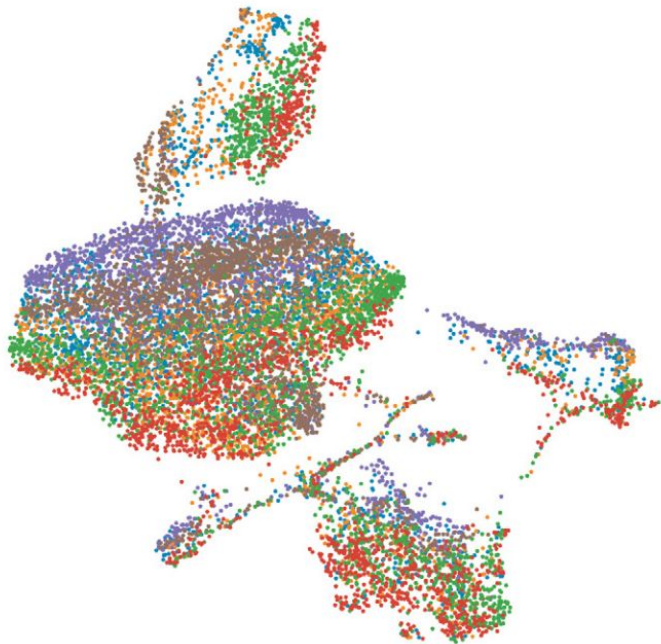
# Imputation and smoothing   ! Caution !

Gen

|   | A | B | C | D | E |
|---|---|---|---|---|---|
|   | 1 | 0 | 0 | 5 | 0 |
| Cell | 0 | 2 | 0 | 4 | 4 |
|   | 1 | 1 | 0 | 0 | 0 |
|   | 0 | 0 | 3 | 0 | 3 |
|   | 0 | 0 | 0 | 2 | 0 |

|   |   |   |   |   |
|---|---|---|---|---|
| 1 | 1 | 2 | 5 | 2 |
| 0 | 2 | 1 | 4 | 4 |
| 1 | 1 | 2 | 3 | 2 |
| 0 | 0 | 3 | 1 | 3 |
| 0 | 1 | 1 | 2 | 2 |

**Zeros result from either (a) true expression or (b) technical variance.**

Malte D. Luecken & Fabian J. Theis (2019)
Wenpin Hou et al. (2020)
Tallulah S. Andrews et al. (2021)

FASTQs

Count matrix

QC

Filtering

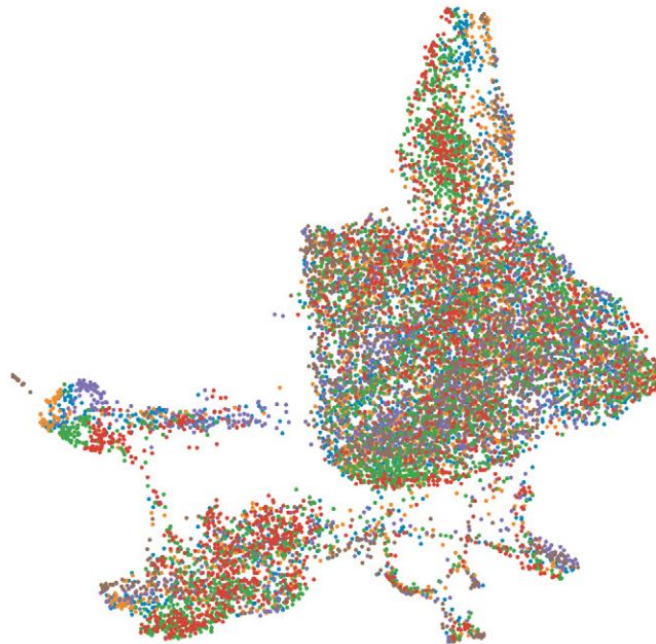Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

# Cell cycle removal   ! Caution !
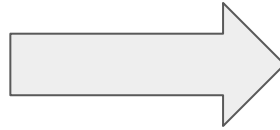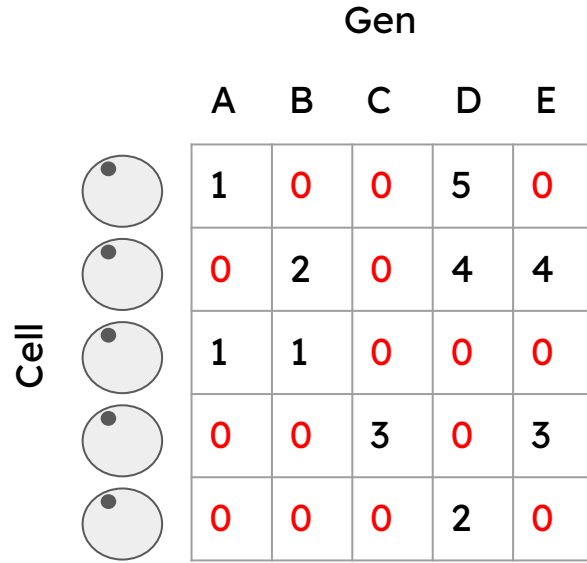


FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

Malte D. Luecken & Fabian J. Theis (2019)
Daniel Schwabe et al. (2020)
Tallulah S. Andrews et al. (2021)
Jiajia Liu et al. (2021)

# Clustering



Malte D. Luecken & Fabian J. Theis (2019)
Ren Qi et al. (2020)
Tallulah S. Andrews et al. (2021)

FASTQs

Count matrix

QC

Filtering

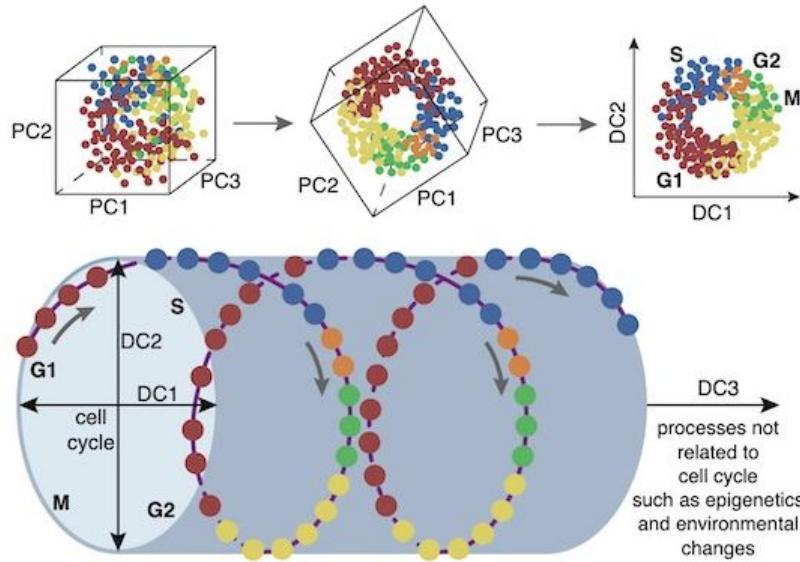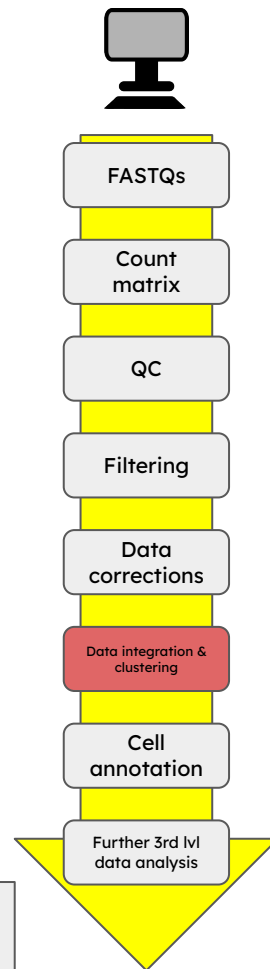Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

# Cell annotation



Tuft cells

Goblet cells

EEC

Paneth cells

TA

EP (early)

Stem cells

EP (late)

Enterocytes

Malte D. Luecken & Fabian J. Theis (2019)

FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering

Cell annotation

Further 3rd lvl data analysis

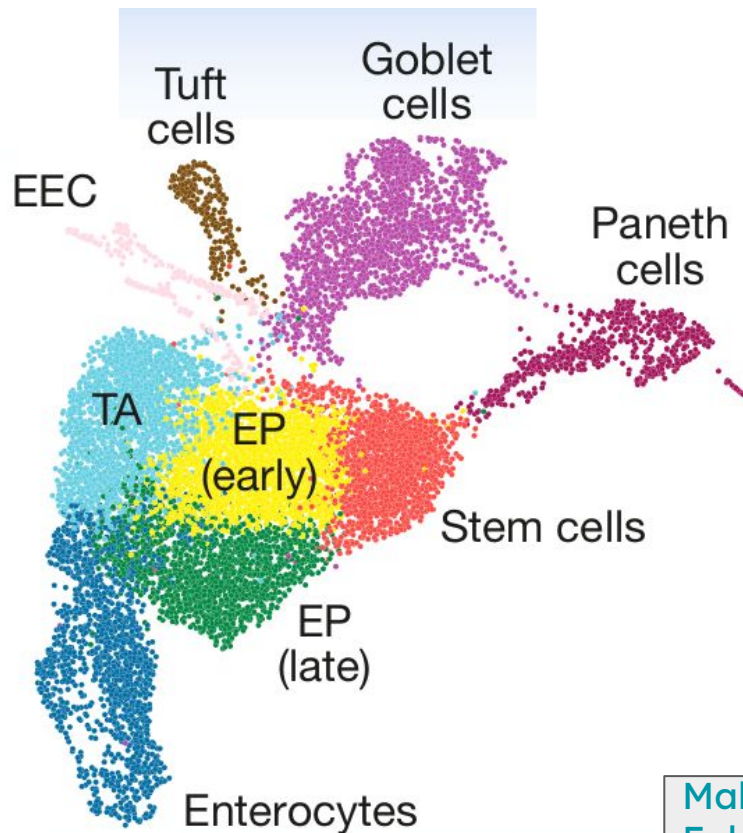# Compositional Analysis



Sean Simmons (2022)

# Trajectory Analysis



Shijie C. Zheng et al. (2022)

# Differential Expression Analysis



Samarendra Das et al. (2022)

# Gene Set Analysis



Farhad Maleki et al. (2020)



FASTQs

Count matrix

QC

Filtering

Data corrections

Data integration & clustering
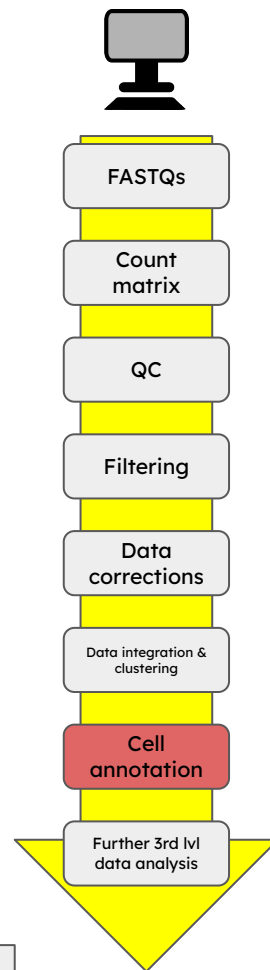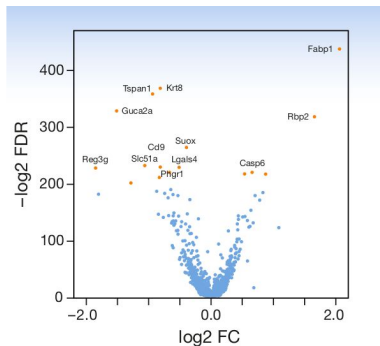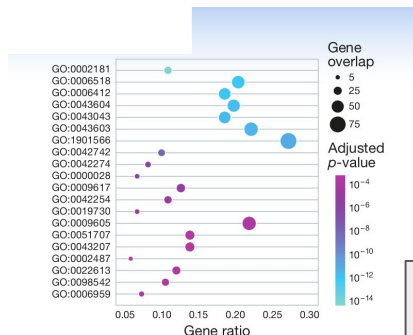
Cell annotation

Further 3rd lvl data analysis

# (3) What tools can I use?

# We need a **standard** for single cell data to make it FAIR

GHGA
THE GERMAN HUMAN GENOME-PHENOME ARCHIVE

**WHY?**

# nf-core/scrnaseq

- nf-core based
- scrnaseq version 2.0 released in June (DSL2)

- Protocols: SmartSeq2, 10xChromium, Drop-Seq

- 4 tools generating count matrix

# What does nf-core 🍎 provide?

- Documentation

- CI Testing

- Stable Releases

- Packaged software

- Portable and reproducible

- Cloud-ready

# Future direction

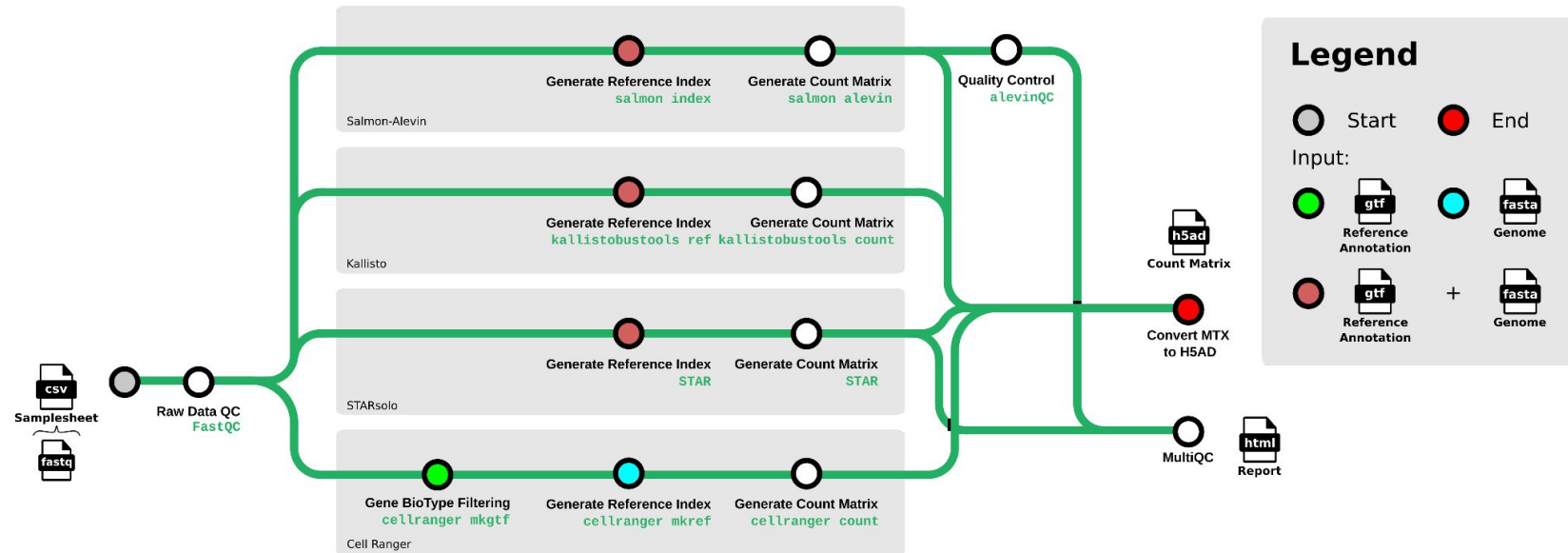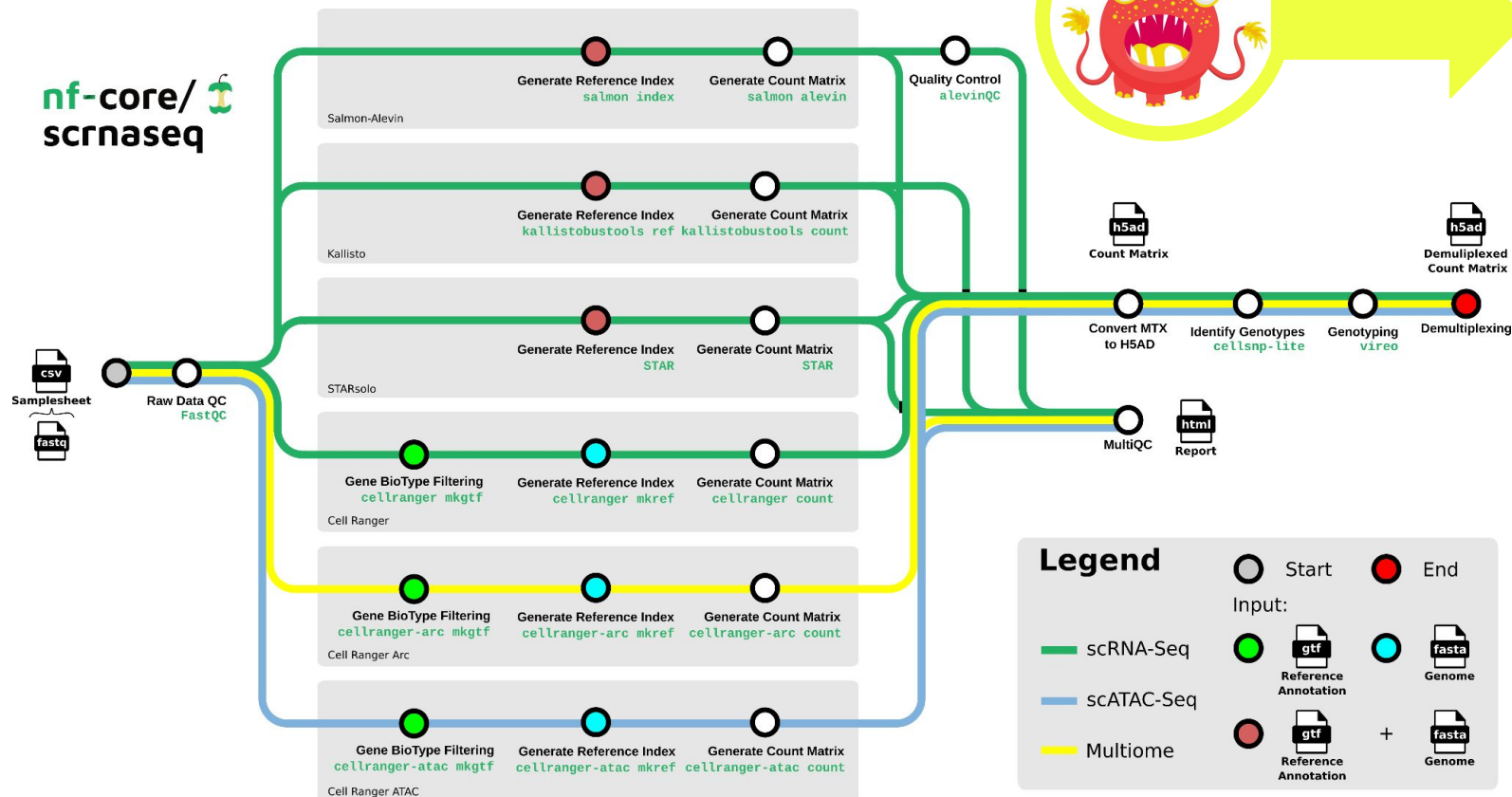# scflow for further analysis

# A benchmark "set" helps to cement a standard for GHGA

**WHY?**

GHGA
THE
GERMAN
HUMAN
GENOME-
PHENOME
ARCHIVE

- Clear Tasks
- Easily accessible
- Quantitative metrics
- Runtime
- CI/CD + Continuous ranking (CR) = CI/CD/CR



Open Problems in Single-Cell Analysis

https://openproblems.bio/#about



Batch 1
Batch 2

**Batch integration graph**

Removing batch effects while preserving biological variation (graph output)

Christopher Lance et al. (2021)

# Tools – Literature

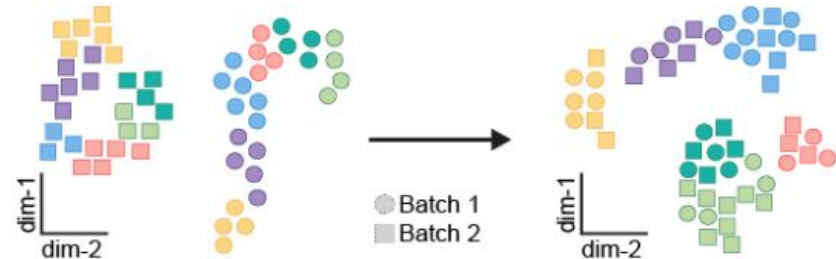| | | |
|---|---|---|
| Guidelines | <ul><li>Jiajia Liu et al. (2021)</li><li>Malte D. Luecken & Fabian J. Theis (2019)</li><li>Galaxy Training Material (Single cell)</li><li>Tallulah S. Andrews et al. (2021)</li><li>Christopher Lance et al. (2021)</li><li>Rui Hong et al. (2022)</li><li>Sean Davis (https://github.com/seandavi/awesome-single-cell)</li><li>Ren Qi et al. (2020)</li><li>Yulong Zhang et al. (2022)</li></ul> | <ul><li>Samarendra Das et al. (2022)</li><li>Farhad Maleki et al. (2020)</li></ul> |
| Software | <ul><li>Cellranger: QC, count matrix, etc.</li><li>scverse (scanpy, muon, scvi-tools, …): QC, normalization, data integration, clustering, cell annotation, etc.</li><li>Scrublet: doublet/multiplet filtering</li><li>Cellsnp-lite: read pileup & genotyping</li><li>Vireo: genotyping & demultiplexing</li><li>ArchR: scATAC QC, etc.</li></ul> | <ul><li>CellBender: emptyplet filtering</li><li>MultiQC: QC</li><li>Seurat: QC, filtering, etc.</li><li>nf-core/scrnaseq: QC, count matrix, etc.</li><li>nf-core/scflow: QC, clustering, etc.</li></ul> |
| Benchmarks | <ul><li>Luyi Tian et al. (2019)</li><li>Malte D. Luecken et al. (2022)</li><li>Cody N. Heiser et al. (2021)</li><li>Ralf Schulze Brüning et al. (2022)</li><li>Huidong Chen et al. (2019)</li><li>Wenpin Hou et al. (2020)</li></ul> | <ul><li>Sean Simmons (2022)</li></ul> |