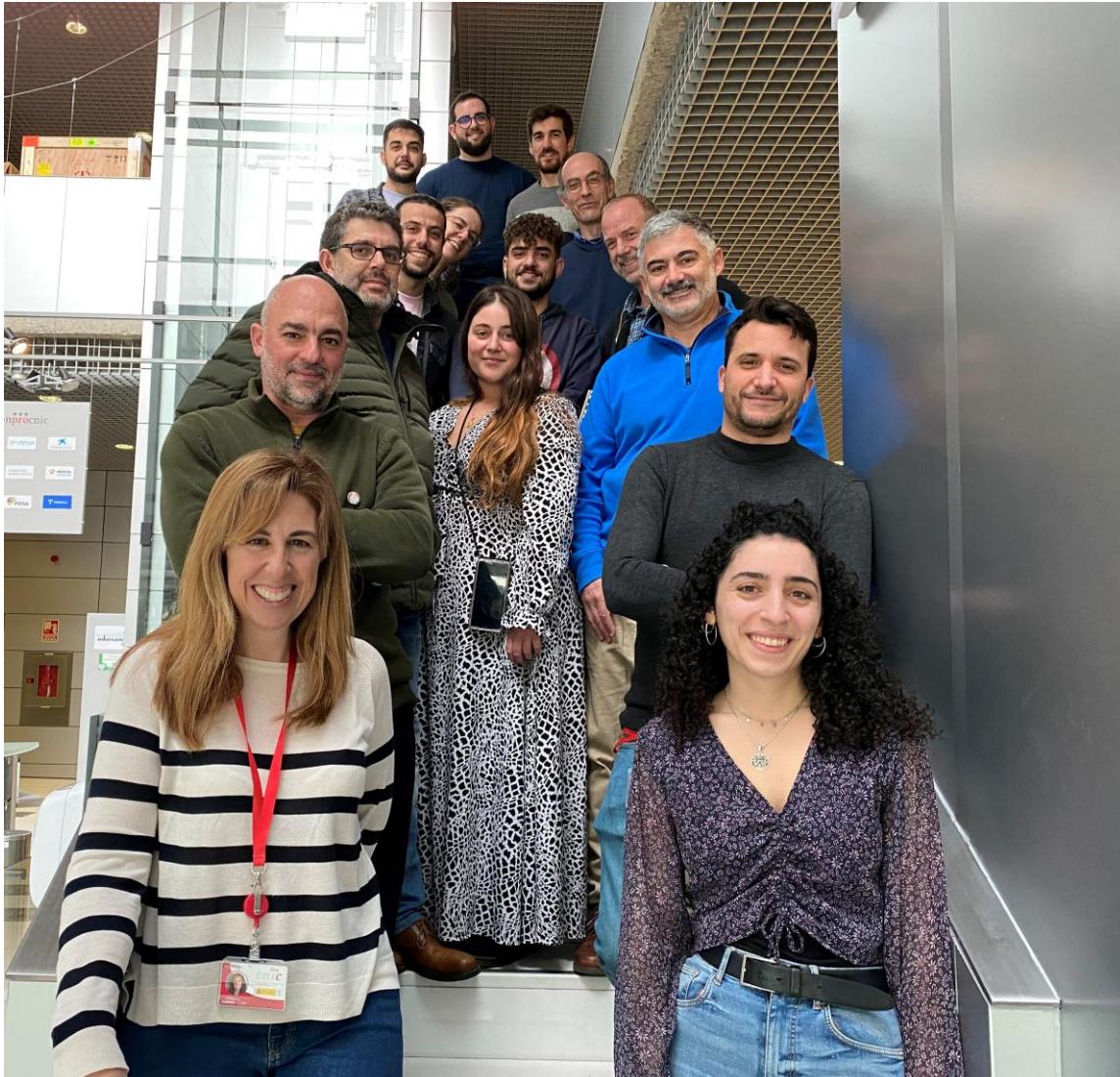




# Single Cell RNA-Seq

Carlos Torroja  
Bioinformatics Unit

- Experimental Design Considerations:
  - Replicates
  - Conditions
- Platforms: MARS-Seq, 10X, Rhapsody
- Libraries:
  - 3' UMI
  - 5' UMI
  - VDJ
  - ATAC-Seq
  - CITE-Seq
- Spatial Single Cell: From Imaging to Transcripts
- How the data looks like
- Data Processing: QC -> Filtering -> Normalization -> Clustering -> DR
- Data Analysis: Markers, Differential Expression -> Functional Analysis
- Cell Type Identification
- Cell-Cell Communication
- Trajectory/Pseudotime Analysis, Interactome, Regulome



## CNIC Bioinfo Unit

Dr. Fátima Sánchez Cabo  
Dr. Fernando Martínez  
Dr. Carlos Torroja  
Dr. Manuel J. Gómez  
Jorge de la Barrera  
Dr Daniel Jiménez  
Lucía Sánchez García  
Juan Carlos Silla

## Embedded Bioinformatician

Jose Luis Cabrera (JLP Lab)  
Marina Rosa Moreno

## PhD Students

Beatriz de las Mercedes (JJF Lab)  
Víctor Jiménez (MAP Lab)  
Jon Sicilia (AH Lab)  
Alvaro Serrano (AR Lab)  
Inés Rivero García (MT Lab)  
Diego Mañanes (DS Lab)

## Master Students

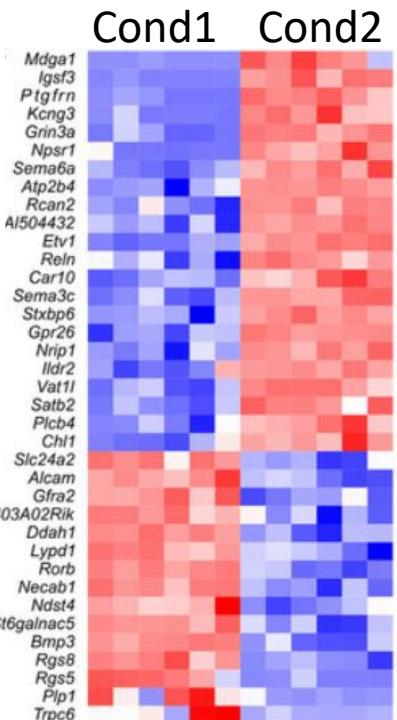
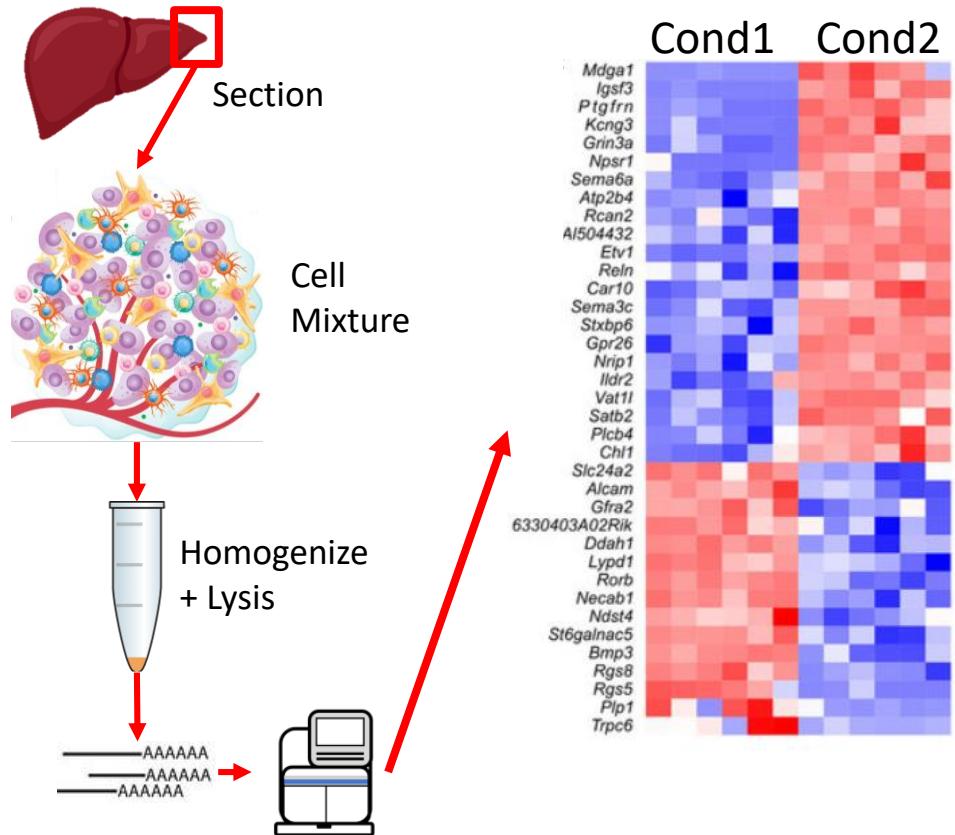
Pablo Beltrán López (UAM)  
Juan Ignacio Álarez Arenas (UAM)

<https://www.cnic.es/nextcloud/s/beoMoMQbzR59FXb>

[https://drive.google.com/drive/folders/1XT6mZ2H0rqzsPUGMpZzhZ\\_9Ls\\_gDcw6?usp=sharing](https://drive.google.com/drive/folders/1XT6mZ2H0rqzsPUGMpZzhZ_9Ls_gDcw6?usp=sharing)

---

## Typical Bulk RNA-Seq experiment



**What is the origin of these differences?**

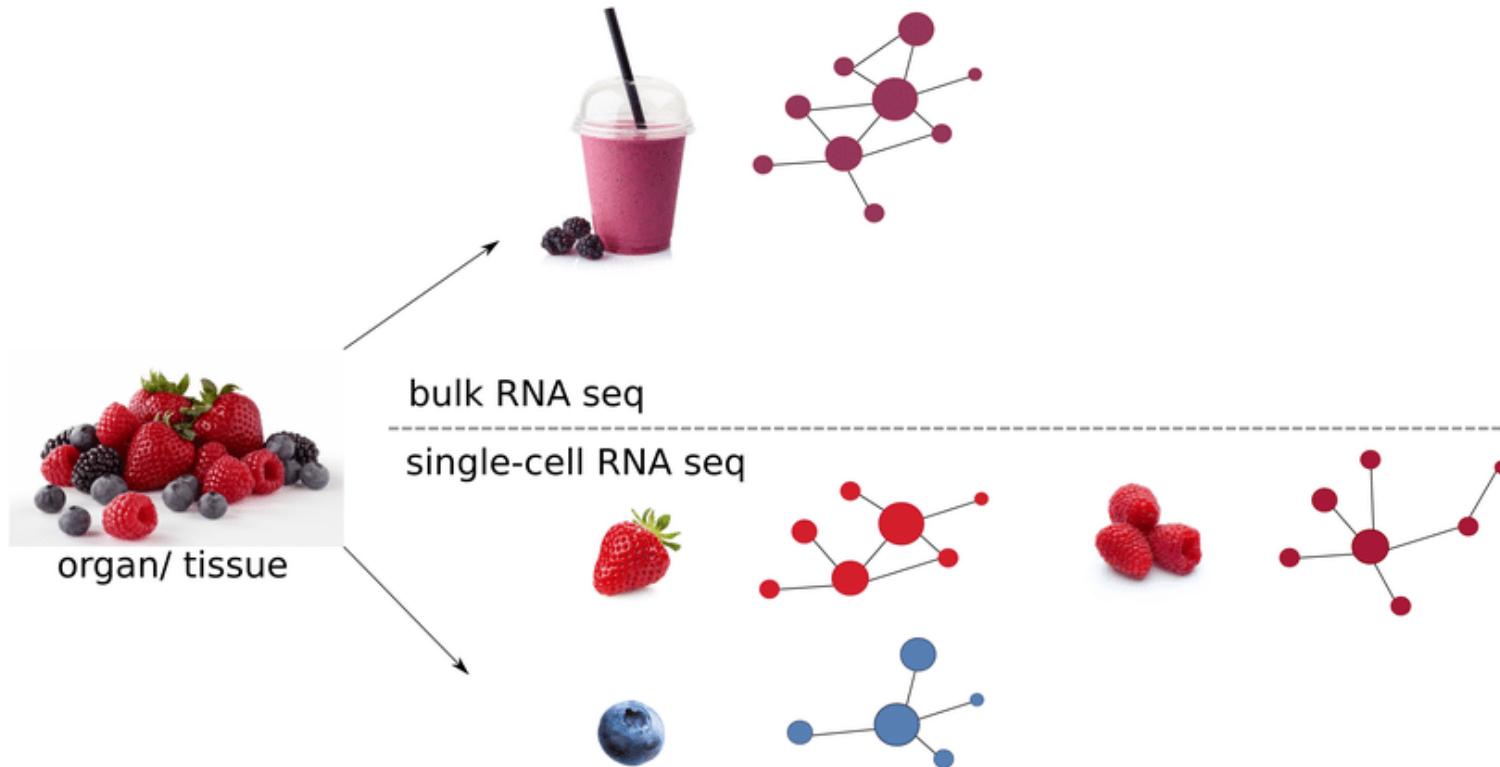
Have all cells changed their physiological status in the same way?

Are there more cells of one type?

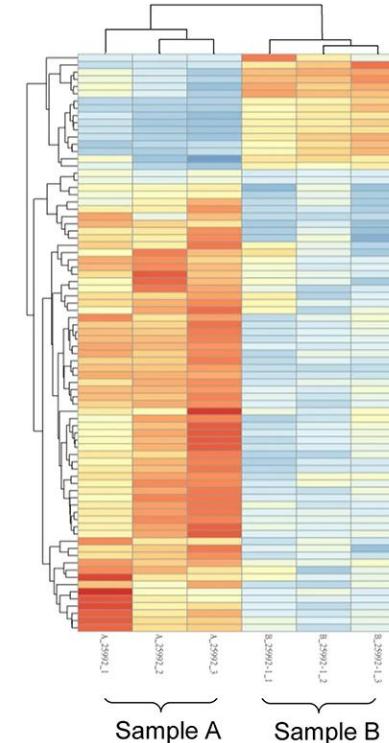
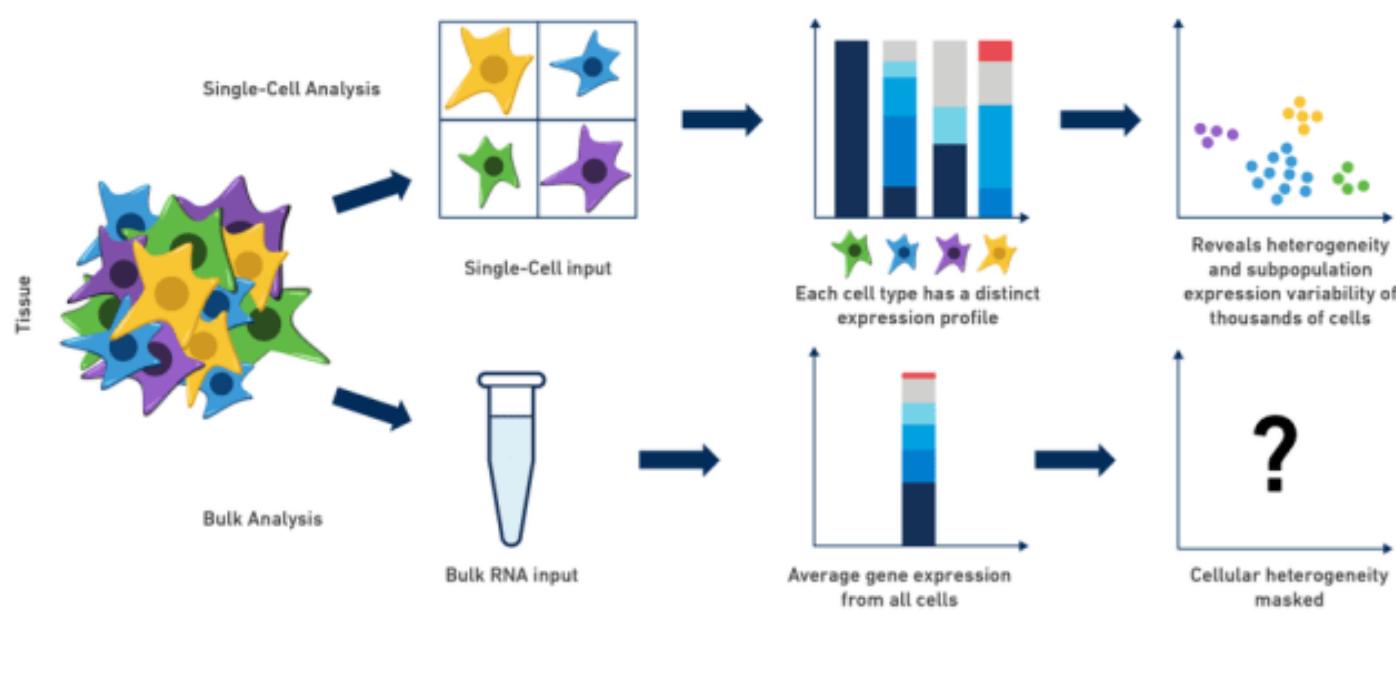
Is there a particular abundant cell that has changed?

Is there a new cell type in the tissue?

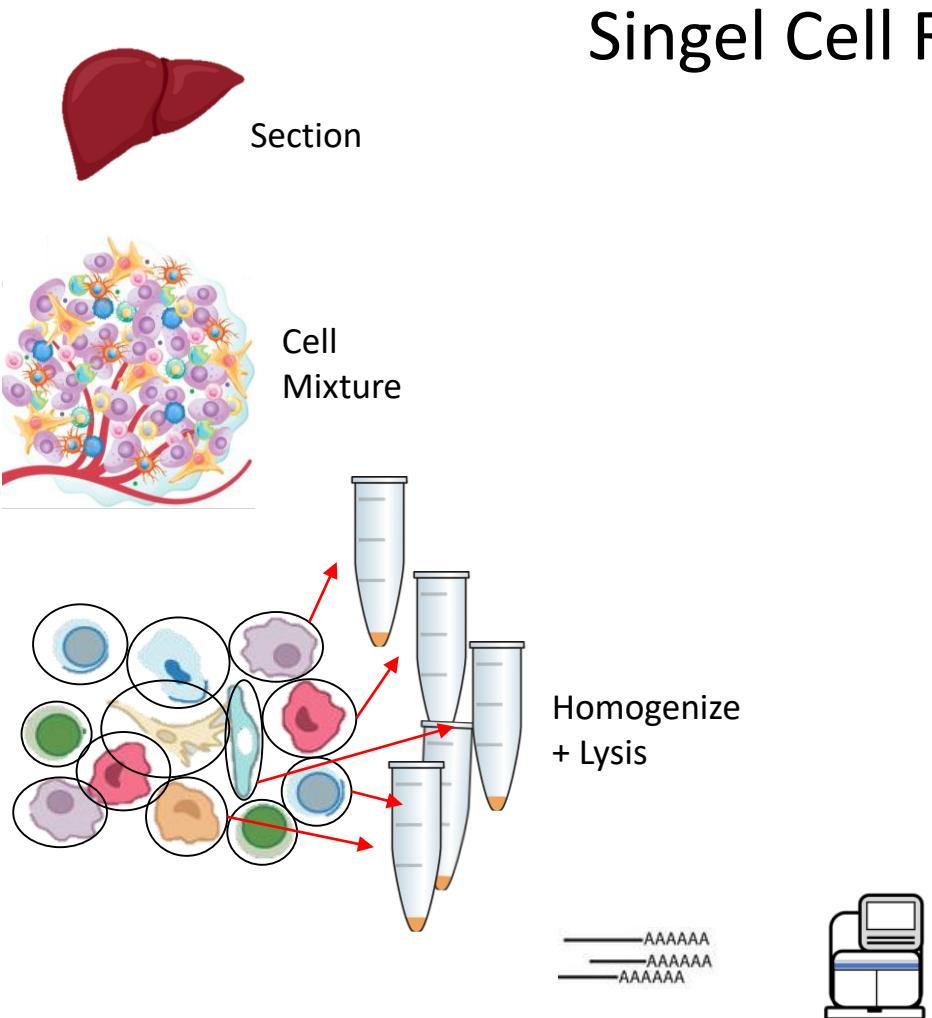
# From Bulk RNA-Seq to Single Cell RNA-Seq



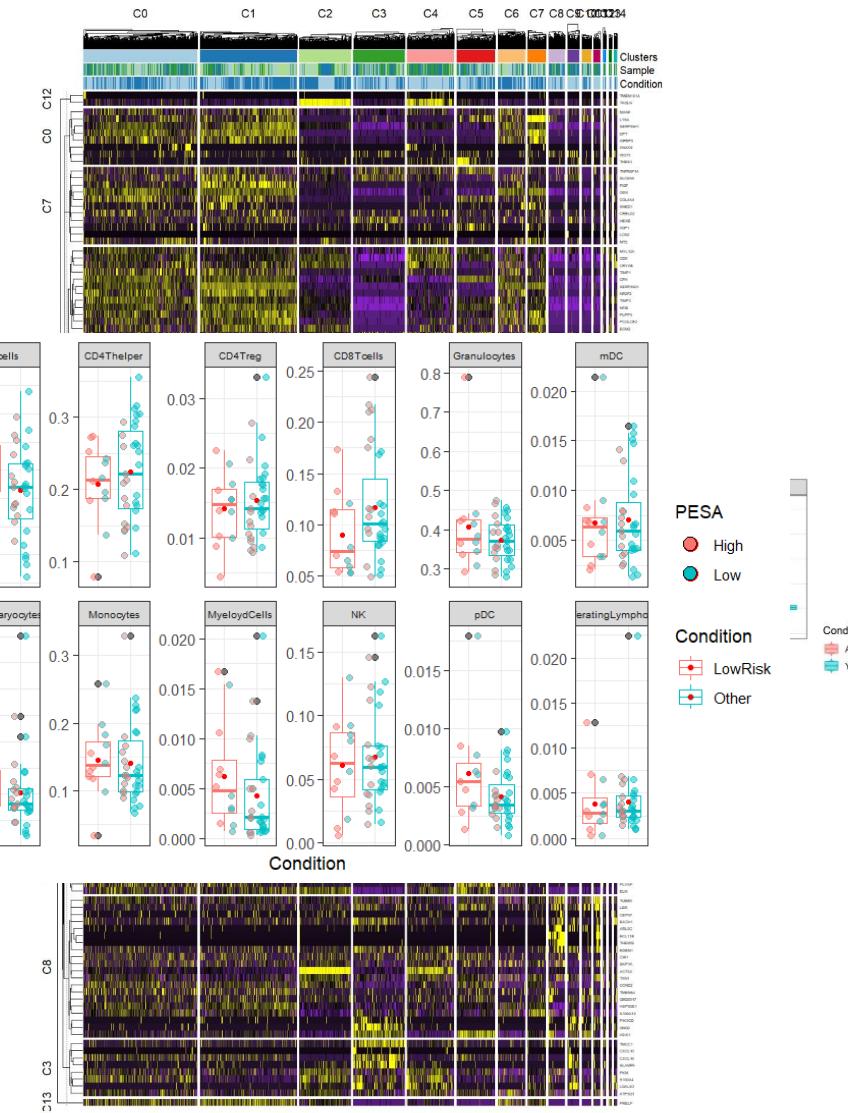
# From Bulk RNA-Seq to Single Cell RNA-Seq



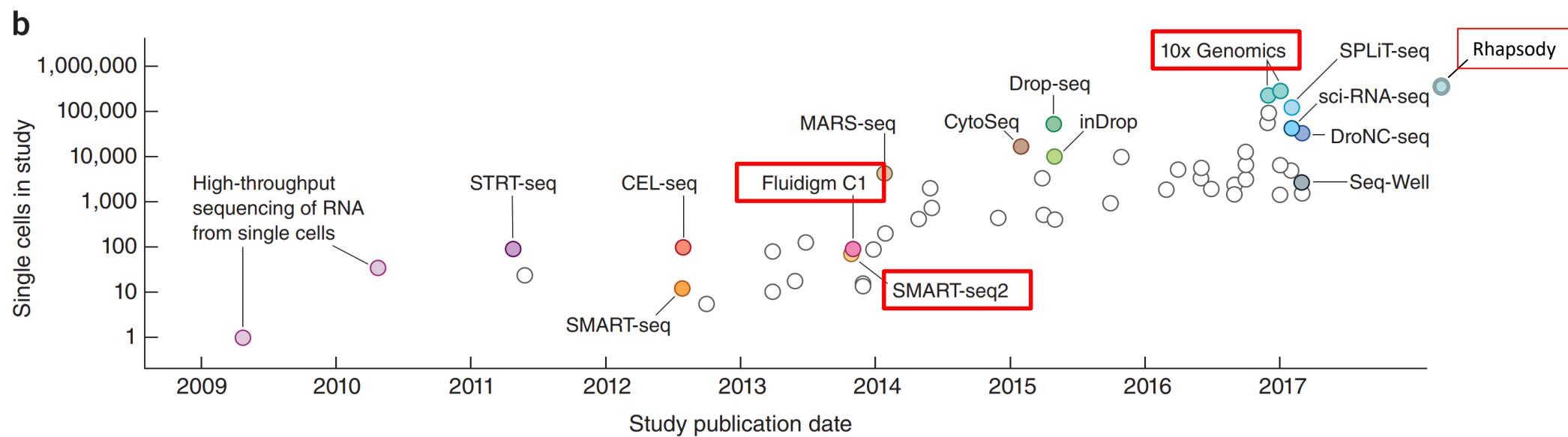
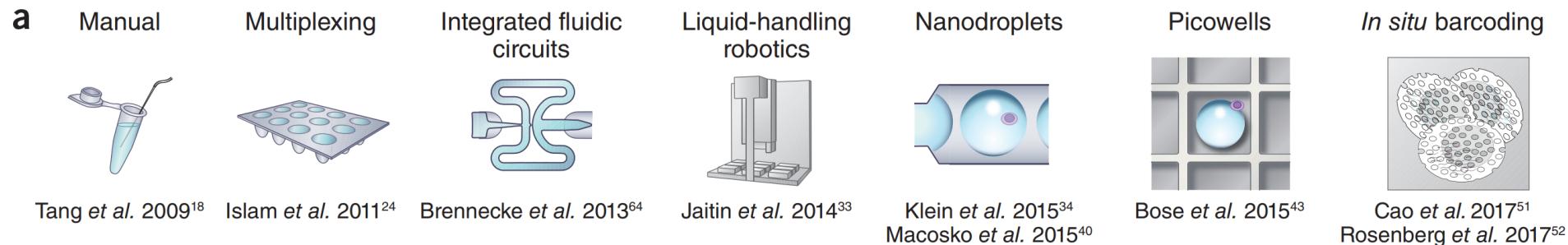
Are all changes the same in all cells?  
Are changes due to changes in cell proportions?  
Are changes exclusive of a particular cell type?



## Single Cell RNA-Seq experiment



# Popular single-cell RNA-seq protocols



# Popular single-cell RNA-seq protocols

	SMART-seq2	CEL-seq2	STRT-seq	Quartz-seq2	MARS-seq	Drop-seq	inDrop	Chromium	Seq-Well	sci-RNA-seq	SPLiT-seq
Single-cell isolation	FACS, microfluidics	FACS, microfluidics	FACS, microfluidics, nanowells	FACS	FACS	Droplet	Droplet	Droplet	Nanowells	Not needed	Not needed
Second strand synthesis	TSO	RNase H and DNA pol I	TSO	PolyA tailing and primer ligation	RNase H and DNA pol I	TSO	RNase H and DNA pol I	TSO	TSO	RNase H and DNA pol I	TSO
Full-length cDNA synthesis?	Yes	No	Yes	Yes	No	Yes	No	Yes	Yes	No	Yes
Barcode addition	Library PCR with barcoded primers	Barcoded RT primers	Barcoded TSOs	Barcoded RT primers	Barcoded RT primers	Barcoded RT primers	Barcoded RT primers	Barcoded RT primers	Barcoded RT primers	Barcoded RT primers and library PCR with barcoded primers	Ligation of barcoded RT primers
Pooling before library?	No	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Library amplification	PCR	In vitro transcription	PCR	PCR	In vitro transcription	PCR	In vitro transcription	PCR	PCR	PCR	PCR
Gene coverage	Full-length	3'	5'	3'	3'	3'	3'	3'	3'	3'	3'
Number of cells per assay											

[Chen et al. 2018. DOI:10.1146/annurev-biodatasci-080917-013452](https://doi.org/10.1146/annurev-biodatasci-080917-013452)

# Why would you do a SC experiment?

---

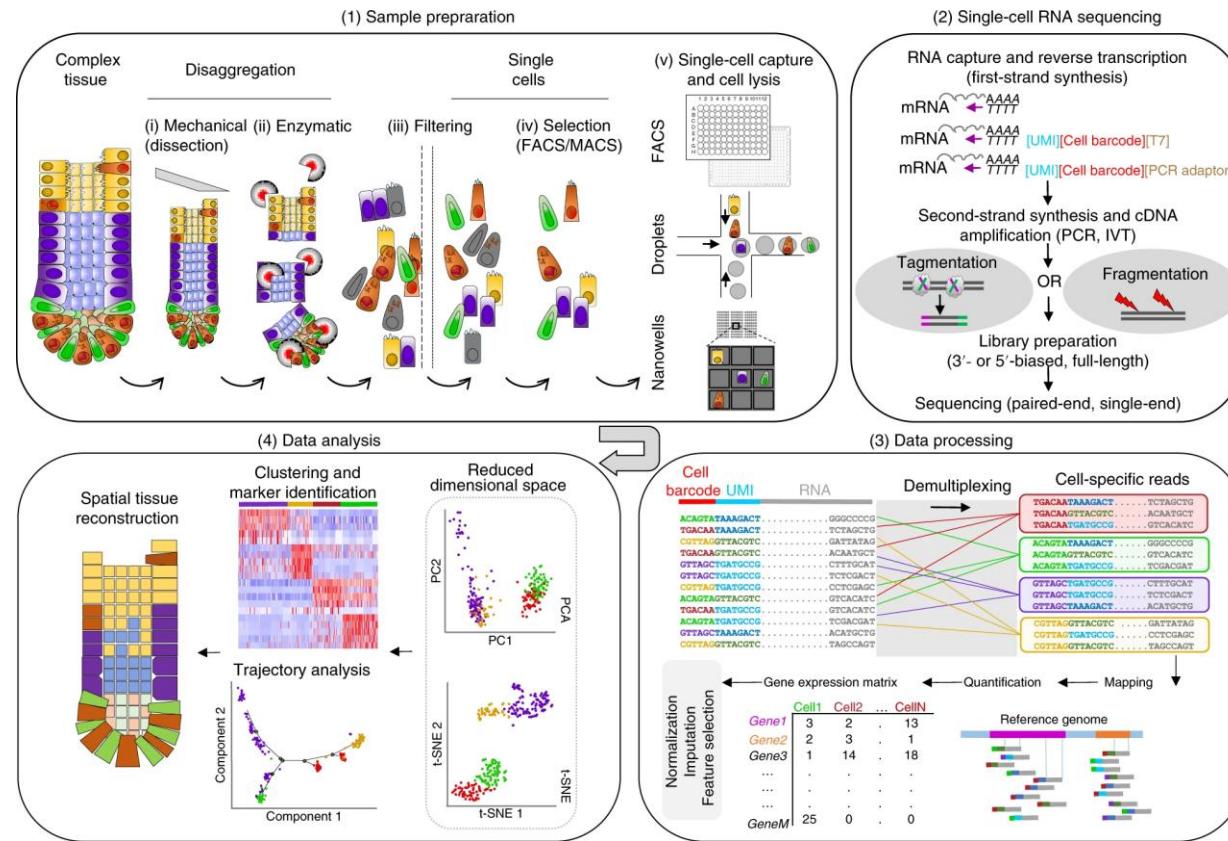
Identify populations in a experimental model

Increase the transcriptomic knowledge of already known populations

Compare populations across conditions / development

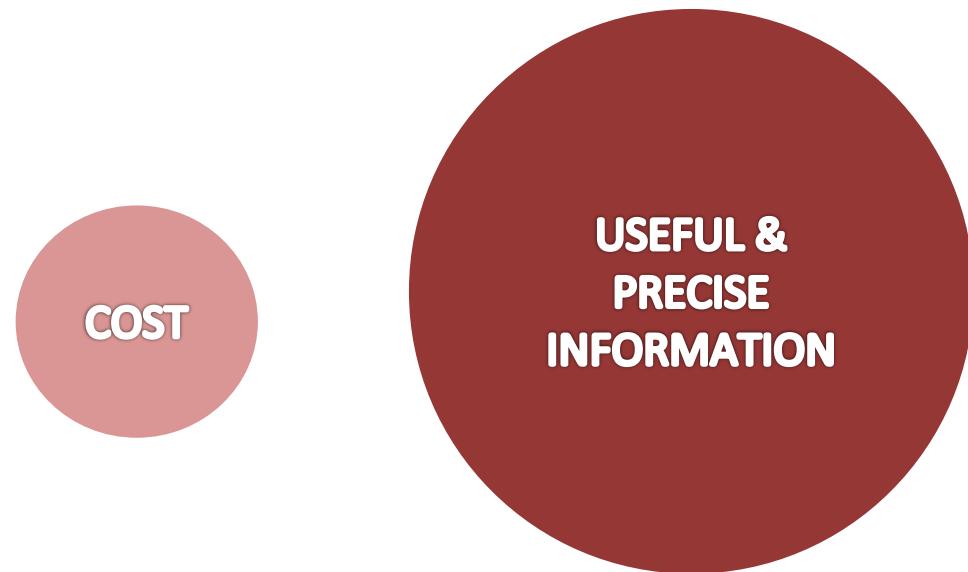
---

# Experimental Design on Single Cell



**Always use biological replicates. Not at single cell level but an individual level**

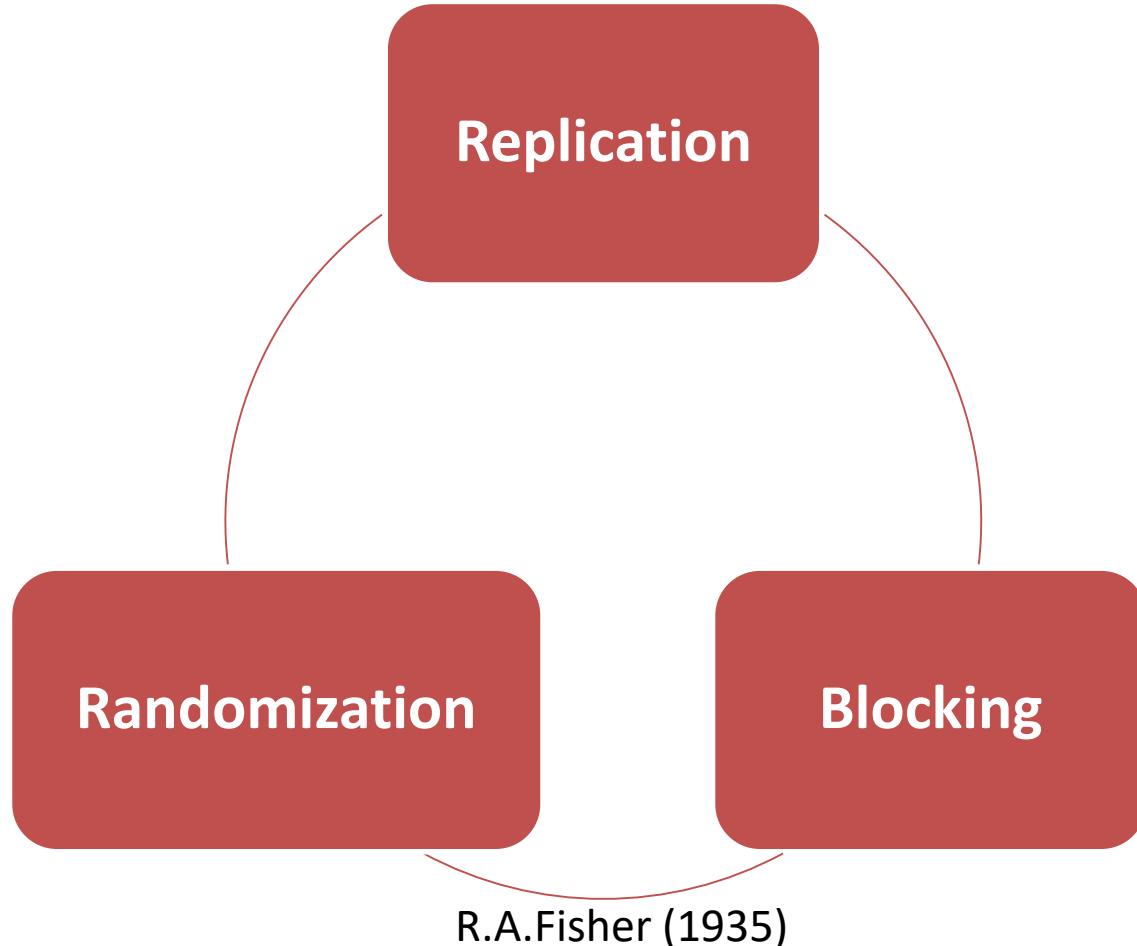
## AIMS OF EXPERIMENTAL DESIGN



## AIMS OF EXPERIMENTAL DESIGN

- Clear question we want to answer with the experiment:  
*Data-driven vs hypothesis-driven*  
but we know what we are looking for...
- Deep knowledge of the technology to:
  - Minimize random errors
  - Correct biases
  - Identify and control by nuisance variables

## PRINCIPLES OF EXPERIMENTAL DESIGN



## REPLICATION

- As Fisher (1935a) noted, without an estimate of variability (i.e., within treatment group), there is no basis for inference (between treatment groups).
- Although we can test for differential expression between treatment groups from unreplicated data, the results of the analysis only apply to the specific subjects included in the study (i.e., the results cannot be generalized)
- Add to this, the errors that can happen during data acquisition.

## TYPE OF ERRORS IN A EXPERIMENT

- Random errors
  - Not possible to calibrate
  - Minimize using repeated measurements
  - Account by nuisance variables
- Systematic errors
  - Possible to estimate and remove from the data
  - Normalization

## SAMPLING AND TECHNIQUE NOISE ARE HIERARCHICAL

- **Subject sampling:** individuals are ideally drawn from a larger population to which results of the study may be generalized. (biological)
- **Cells sampling:** from each individual we recover a sample of it's cells. (biological + technical)
- **RNA sampling:** occurs during the experimental procedure when RNA is isolated from the cell. (technical)
- **Fragment sampling:** only certain fragmented RNAs that are sampled from the cells are retained for amplification. Since the sequencing reads do not represent 100% of the fragments loaded into a flow cells, this is also at play. (technical)
- **Sequencing:** At the sequencing step there is a fragment sampling step and also the sequencing process that may introduce errors

## SAMPLING AND TECHNIQUE NOISE ARE HIERARCHICAL

Table 1: Replicate hierarchy in a hypothetical mouse single-cell gene expression RNA sequencing experiment

From: [Replication](#)

	Replicate type	Replicate category <sup>a</sup>
Animal study subjects	Colonies	B
	Strains	B
	Cohoused groups	B
	Gender	B
	Individuals	B
Sample preparation	Organs from sacrificed animals	B
	Methods for dissociating cells from tissue	T
	Dissociation runs from given tissue sample	T
	Individual cells	B
Sequencing	RNA-seq library construction	T
	Runs from the library of a given cell	T
	Reads from different transcript molecules	V <sup>b</sup>
	Reads with unique molecular identifier (UMI) from a given transcript molecule	T

<sup>a</sup>Replicates are categorized as biological (B), technical (T) or of variable type (V).

<sup>b</sup>Sequence reads serve diverse purposes depending on the application and how reads are used in analysis.

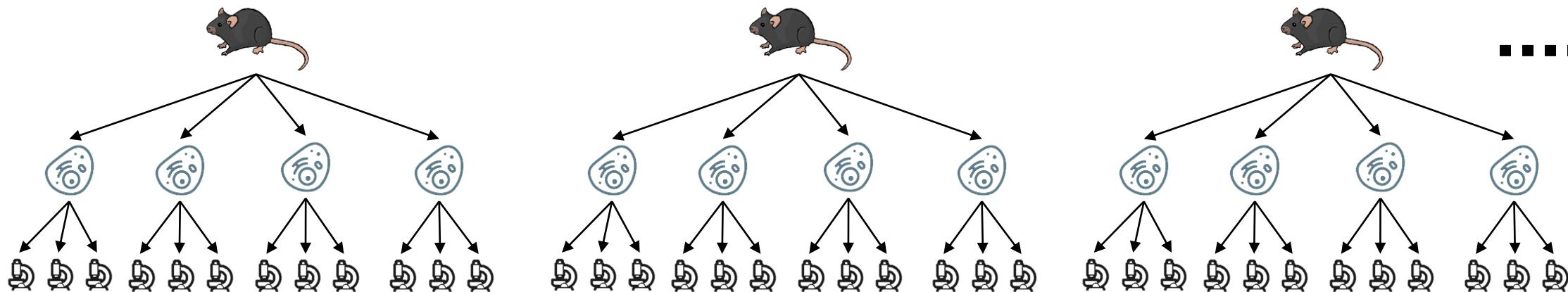
## REPLICATION HELPS TO CONTROL NOISE AND ERRORS

- Technical replicates
  - Minimize random errors through averaging
  - Test technology
- Biological replicates
  - Draw conclusions about the whole population, not about the individual
  - Control variability at different experimental steps

## REPLICATION

Suppose you perform an experiment in which you can measure the expression of a gene in a single cell and **you have money to perform 48 measurements**.

Blainey, P., Krzywinski, M. & Altman, N. Replication. *Nat Methods* **11**, 879–880 (2014).  
<https://doi.org/10.1038/nmeth.3091>



Which kind of replication is going on at every step?

Which sources of variation, error or noise can you think off?

## NUISANCE VARIABLES IN NGS

### Technical:

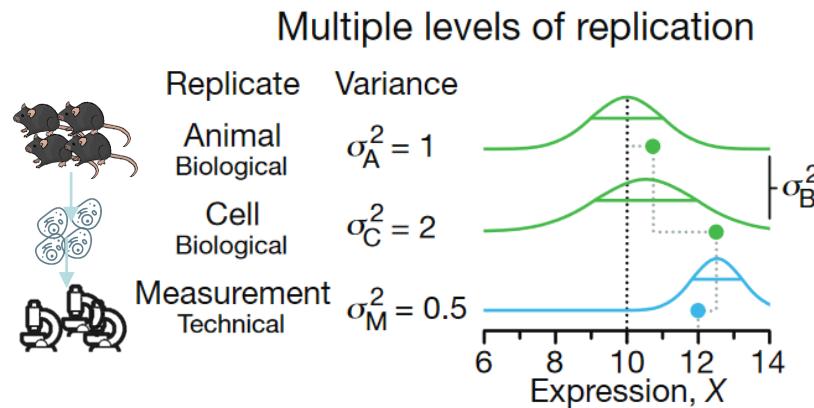
- RNA extraction
- Library preparation
- Flow-cell
- Lane
- Barcode
- Lab/technician

### Biological:

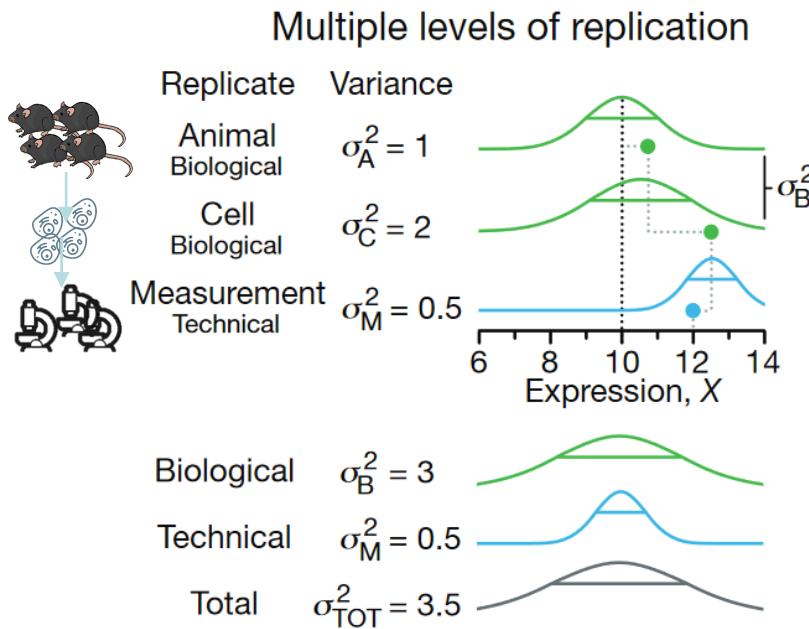
- Sex/gender
- Litter/family
- Age
- ...

Systematic Bias and Random Errors (Noise)

## REPLICATION



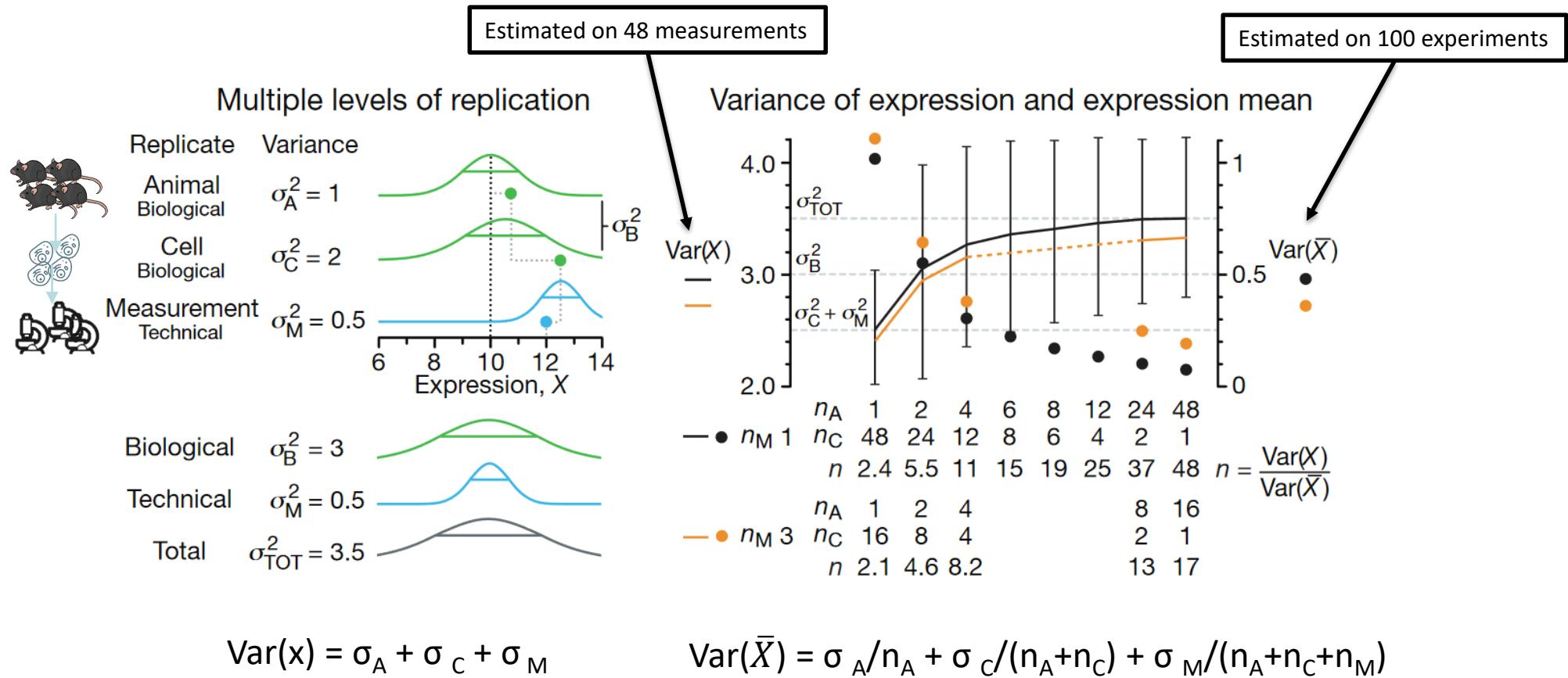
## REPLICATION



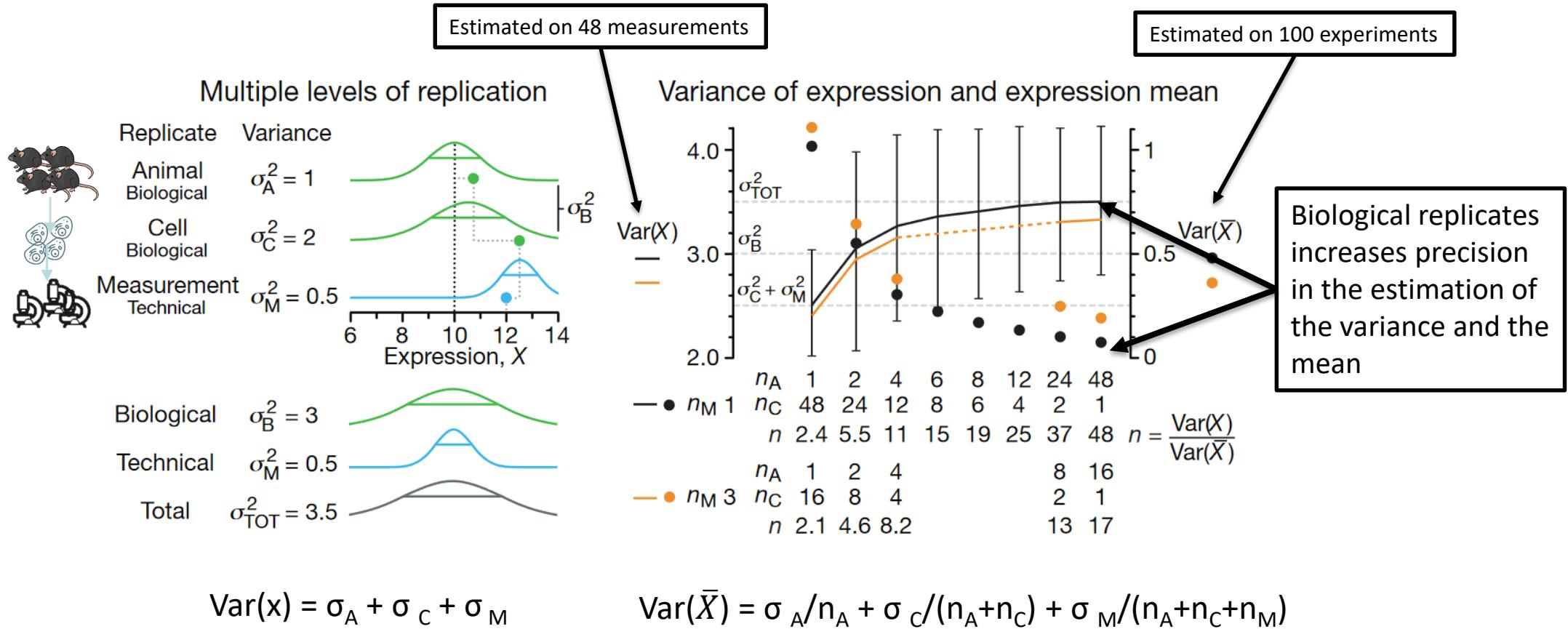
$$\text{Var}(x) = \sigma_A^2 + \sigma_C^2 + \sigma_M^2$$

Blainey, P., Krzywinski, M. & Altman, N. Replication. *Nat Methods* **11**, 879–880 (2014).  
<https://doi.org/10.1038/nmeth.3091>

## REPLICATION

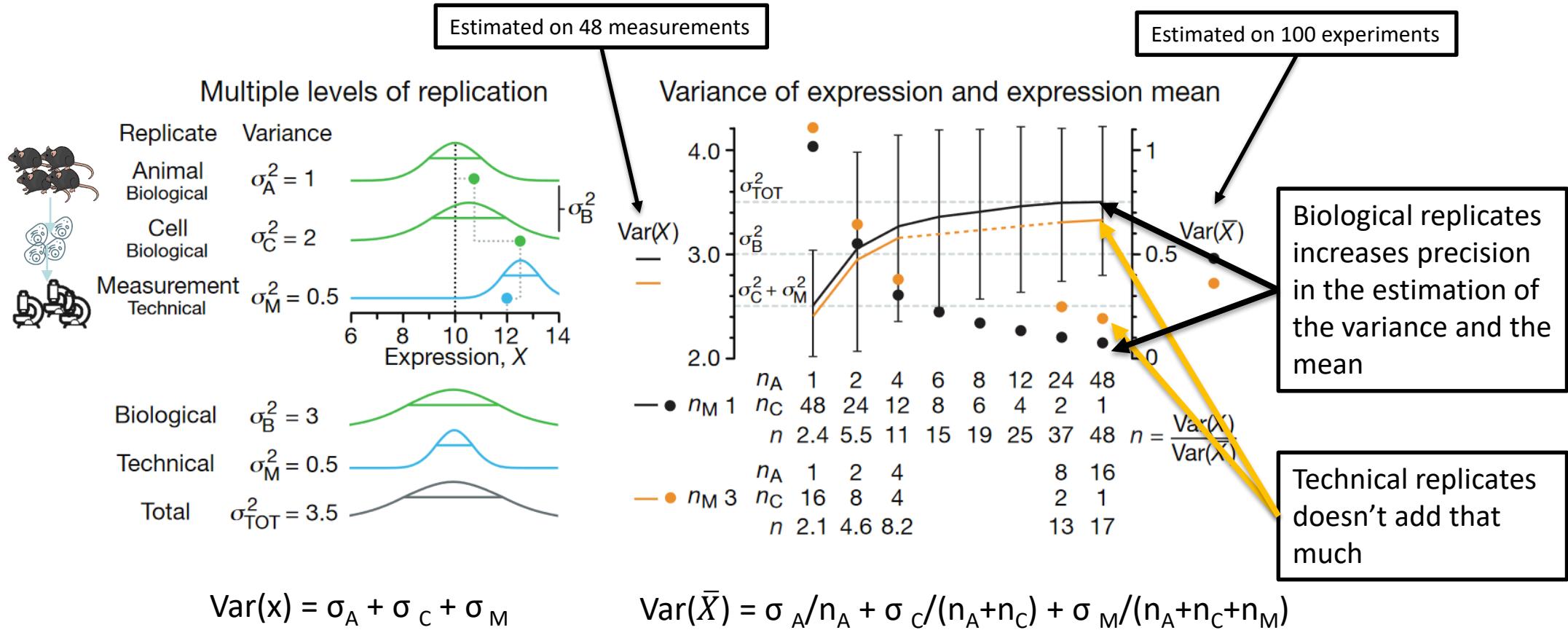


## REPLICATION



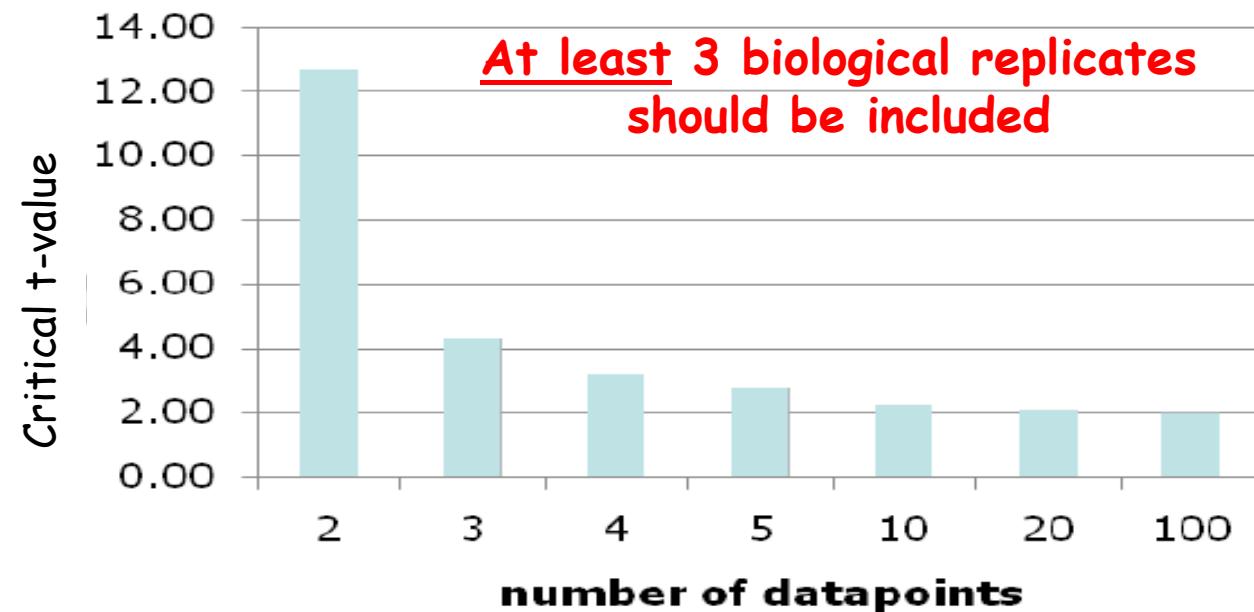
# Experimental Design

## REPLICATION



## REPLICATION

At least three biological replicates for omics technologies



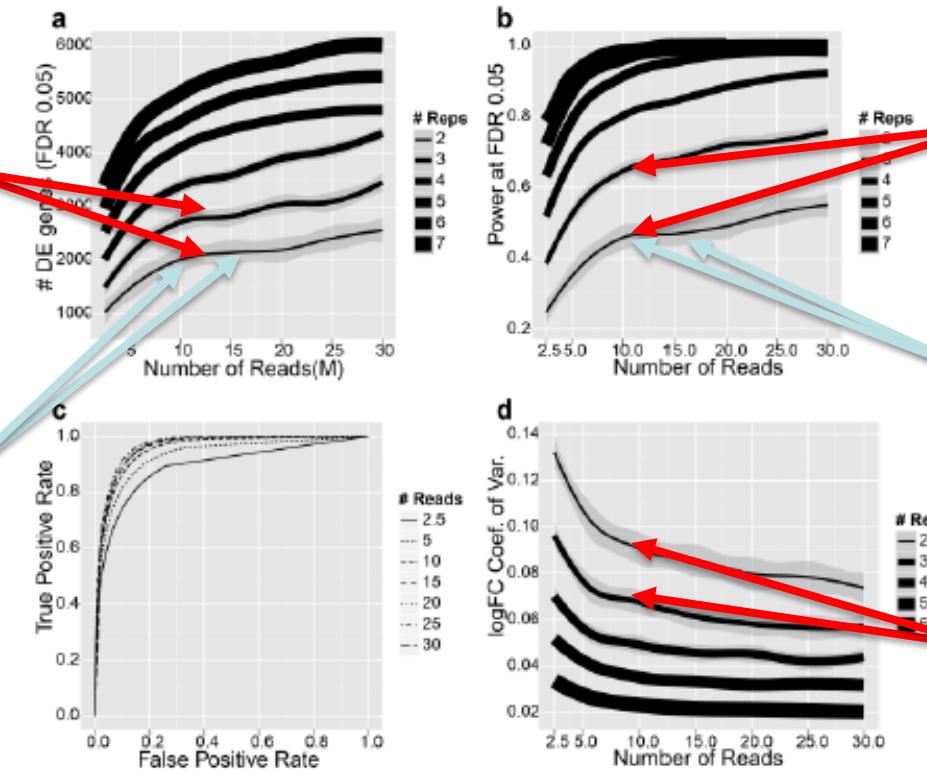
# Experimental Design

## REPLICATES OR DEPTH?

- RNA-Seq experiment from MCF7 tumor cell line (mammary gland) with 7 replicates per condition (17b Estradiol and Control).
- Sequenced up to 30M reads.

There is an increase of **35%** in the DE genes (2011 -> 2709) when increasing from 2 to 3 replicates of 10M (20M -> 30M) reads.

There is only an increase of **6%** (2011 -> 2139) in the DE genes when increasing by 50% the reads: from 2 replicates of 10M (20M) to 15M (30M) reads.



Power increases also **41%** (0.46 -> 0.65) when increasing from 2 to 3 replicates of 10M (20M -> 30M) reads.

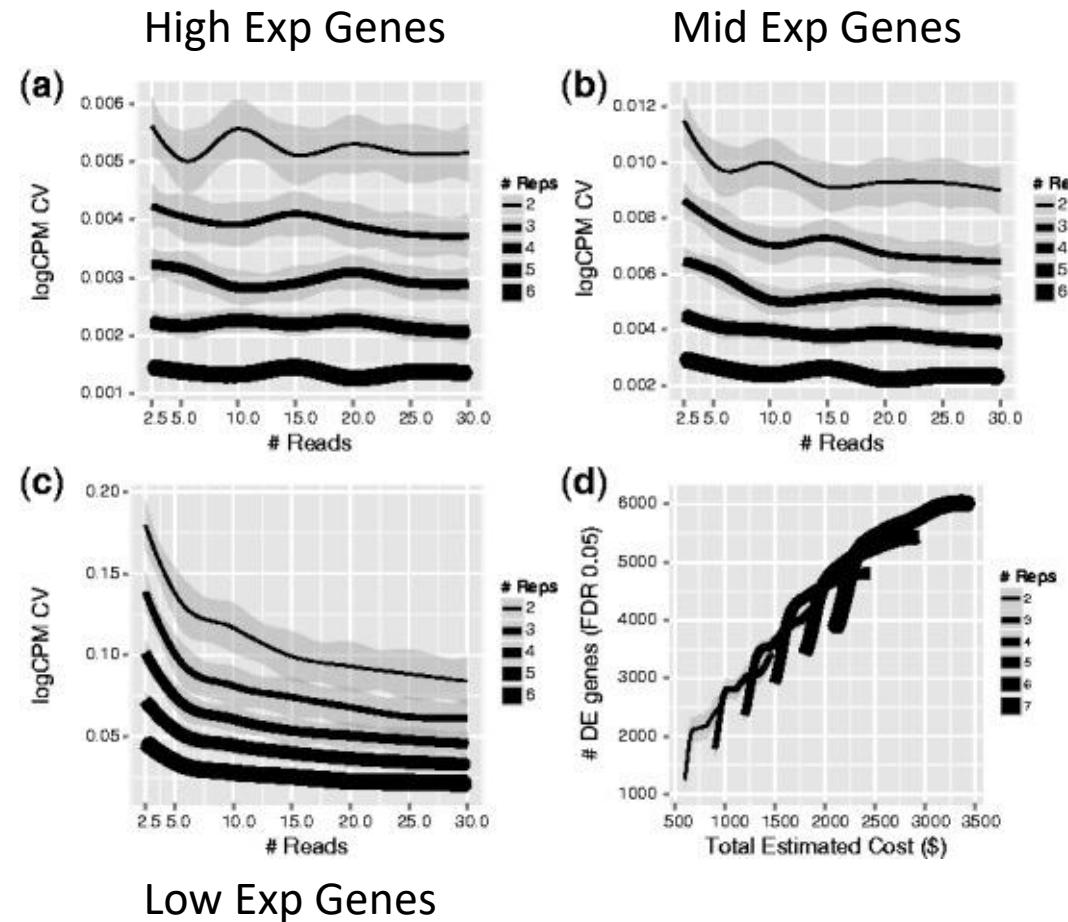
Power increase only **19%** (0.46 -> 0.55) when increasing by 50% the reads: from 2 replicates of 10M (20M) to 15M (30M) reads.

The logFC estimation precision increases with the replicates and not so much with the sequencing depth. (Even for low expressed genes.)

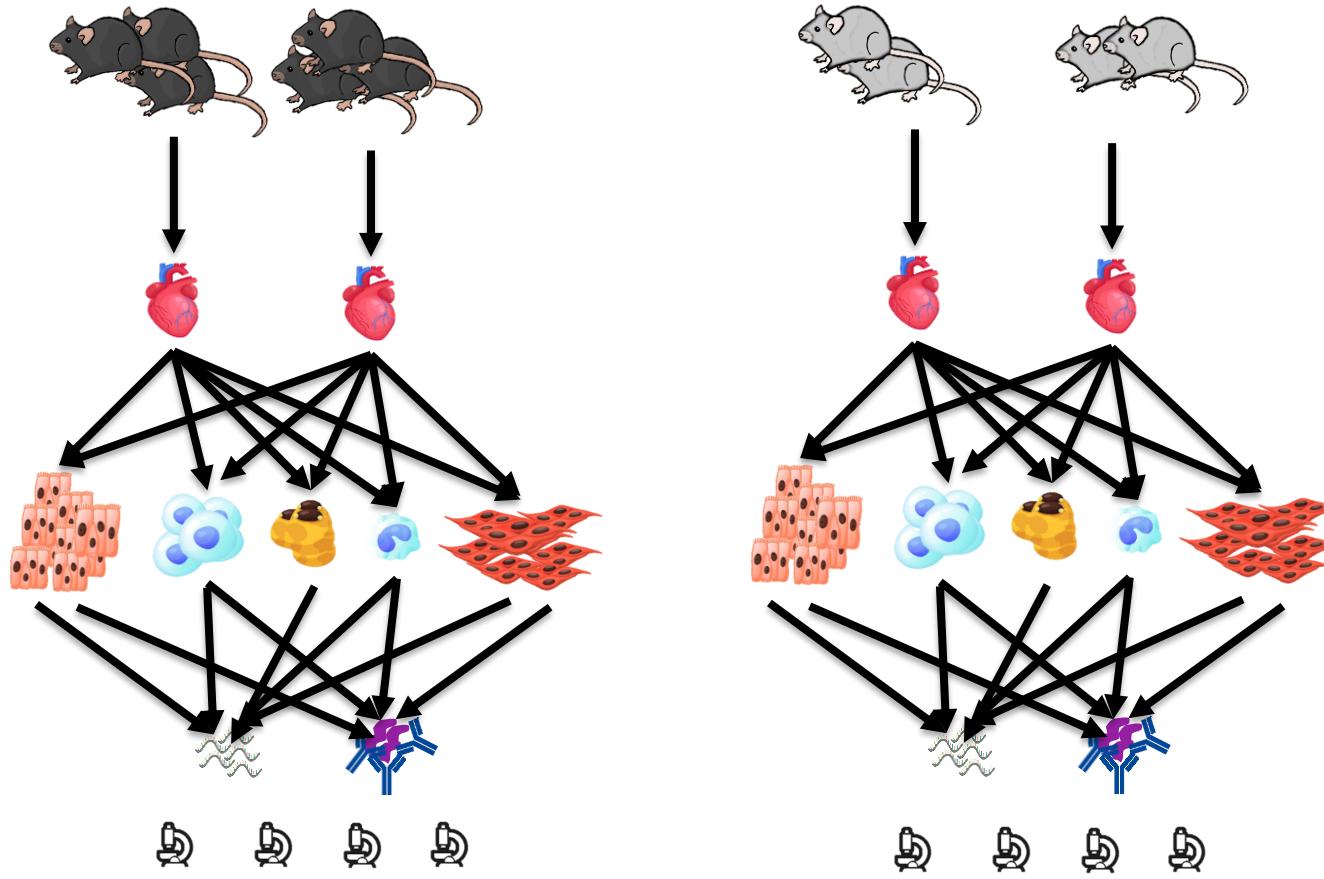
# Experimental Design

## REPLICATES OR DEPTH?

The gene expression estimation precision increases more with the replicates regardless of the level of expression.



# Experimental Design on Single Cell



Treatment / Genotype / Stage

Tissues / Condition

Cell Types / Condition

But how many cells should I sequence?

SCOPIT

[https://alexdavisscs.shinyapps.io/scs\\_power\\_multinomial/](https://alexdavisscs.shinyapps.io/scs_power_multinomial/)

<https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3167-9>

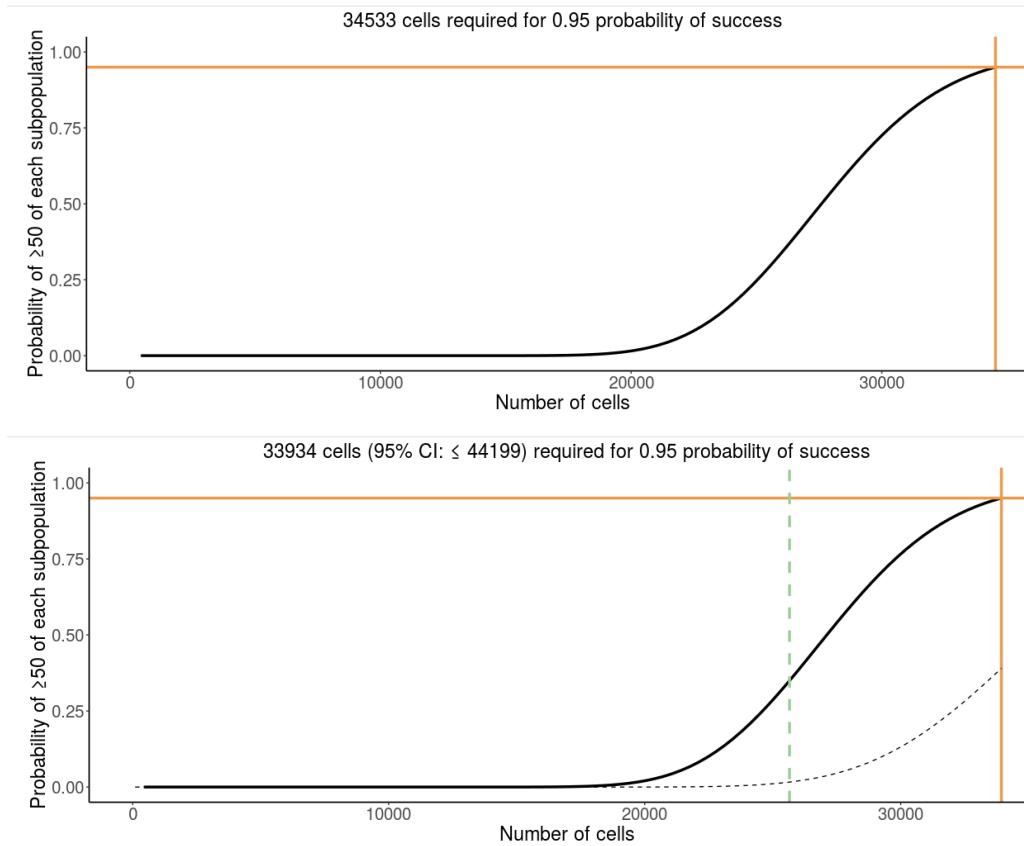
---

# Experimental Design on Single Cell

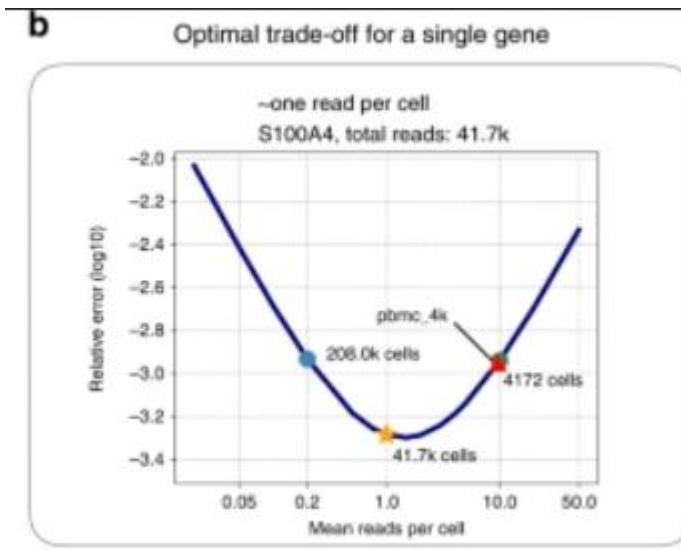
Example:

We have got an experiment with 25000 cells and 17 cell types:

Cluster	Cells	Freq
C0	5363	0.209019
C1	3327	0.129667
C2	3120	0.1216
C3	2662	0.103749
C4	2415	0.094123
C5	2393	0.093265
C6	1885	0.073466
C7	1177	0.045873
C8	1063	0.04143
C9	747	0.029114
C10	662	0.025801
C11	333	0.012978
C12	199	0.007756
C13	114	0.004443
C14	93	0.003625
C15	58	0.002261
C16	47	0.001832
Total	25658	1



## How many reads?



$$25000 \text{ cells} * 12000 \text{ genes} = 300M \text{ reads}$$

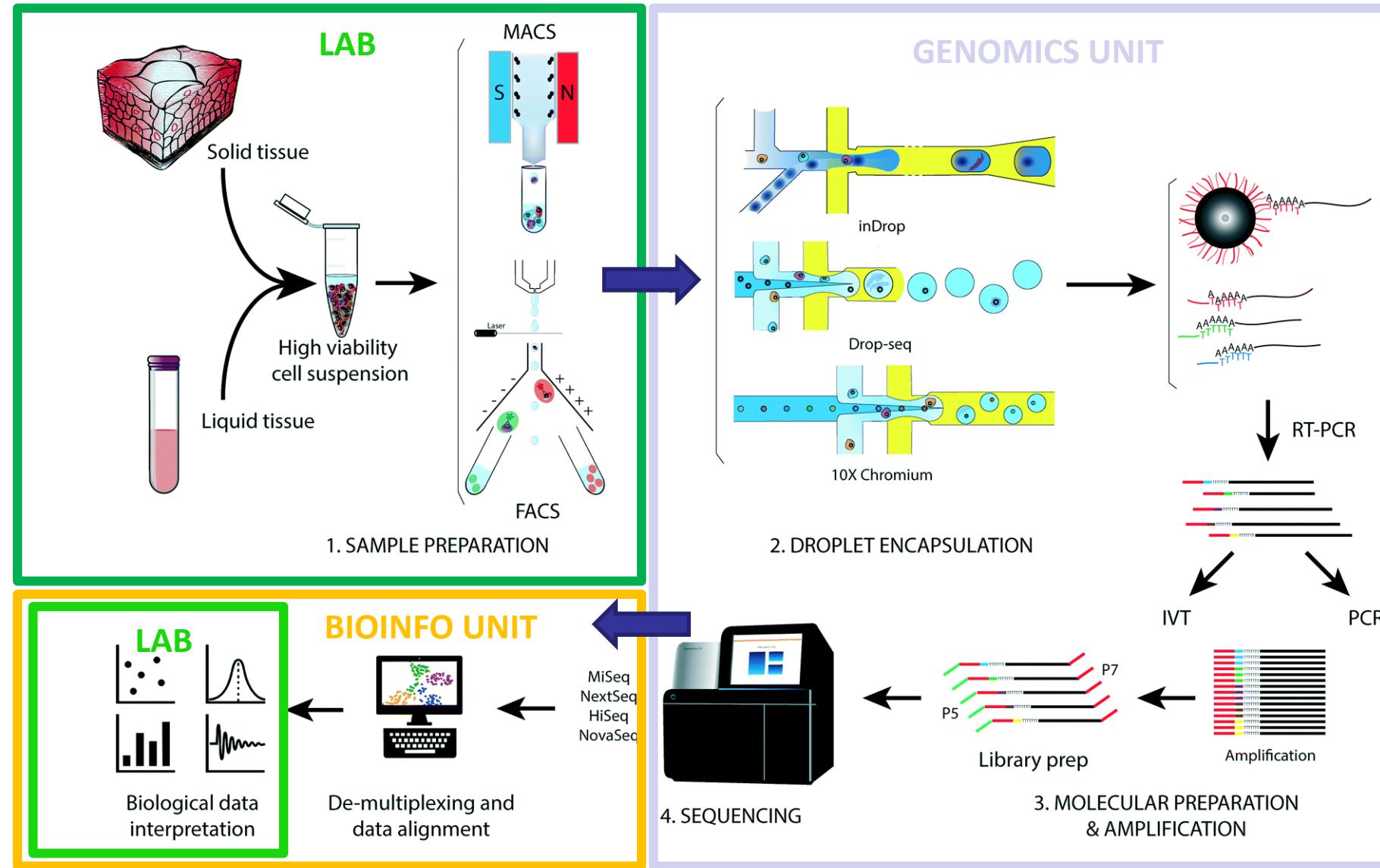
Zhang, M.J., Ntranos, V. & Tse, D. Determining sequencing depth in a single-cell RNA-seq experiment. *Nat Commun* **11**, 774 (2020). <https://doi.org/10.1038/s41467-020-14482-y>

In this paper, we introduced the sequencing budget allocation problem to provide a precise answer to this question; **given a fixed budget, sequencing as many cells as possible at approximately one read per cell per gene is optimal**, both theoretically and experimentally.

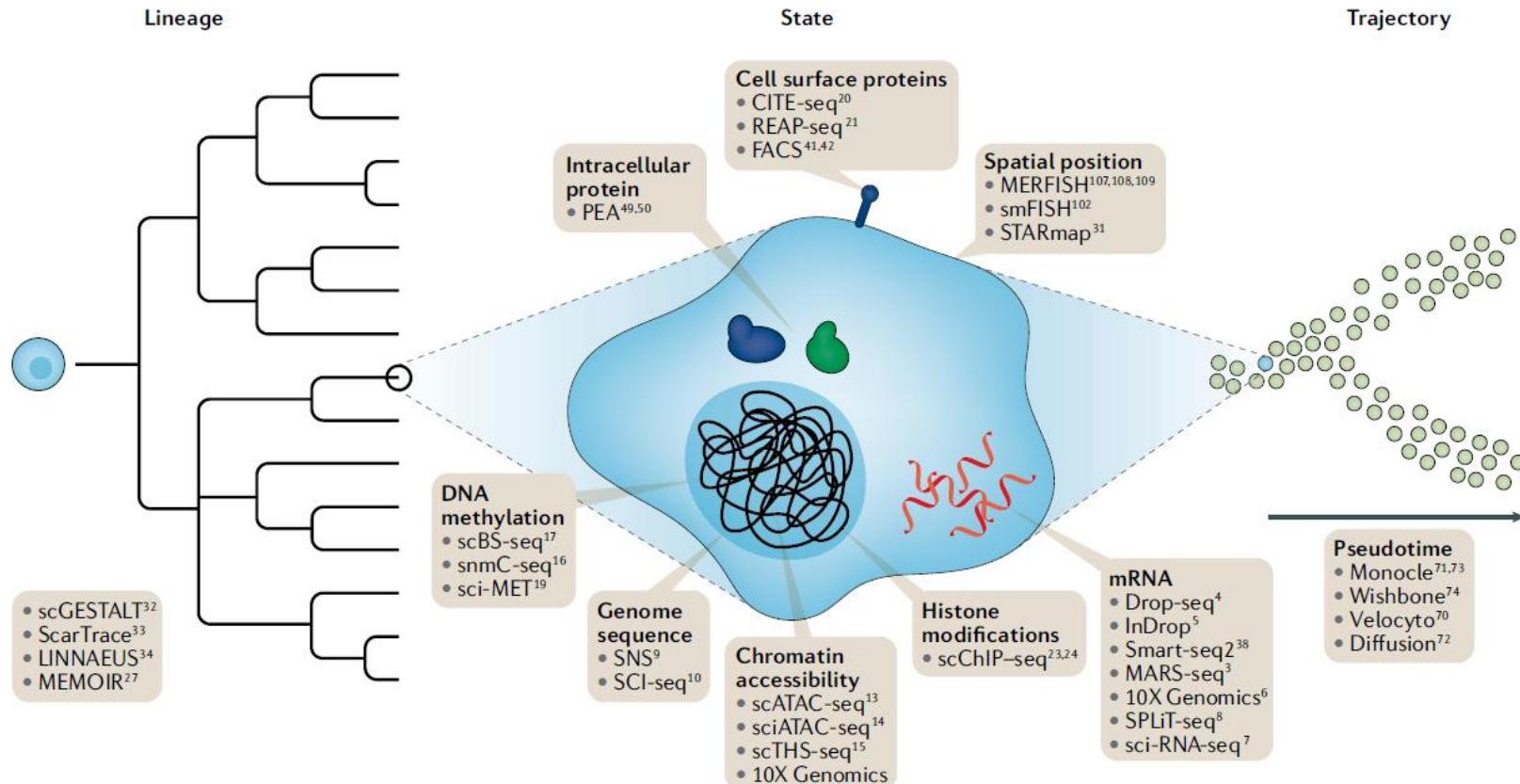
## **Always think about the experimental design before performing the experiment**

- Try to persuade the researcher to use replicates. (better replicates than more cells or sequencing depth)
- Try to reduce pooling as much as possible.
- Do not ever consider technical replicates (you will be lying to yourself).
- Estimate the minimum number of cells base on prior knowledge.

# Single Cell Experiment workflow



## Not just Transcriptome



# Starting Material

## Cell dissociation from solid tissues:

- Each tissue on each species has its own technique
- Protease digestion is required to destroy extracellular matrix (Pronase, Acutase, ...)
- Needs to be adjusted to avoid cell damage
- The procedure is quite dramatic.
- In cases where the cell integrity is compromised when dissociating tissues, single nuclei can be obtained.
- Nuclei are also used from fixed or frozen tissues.

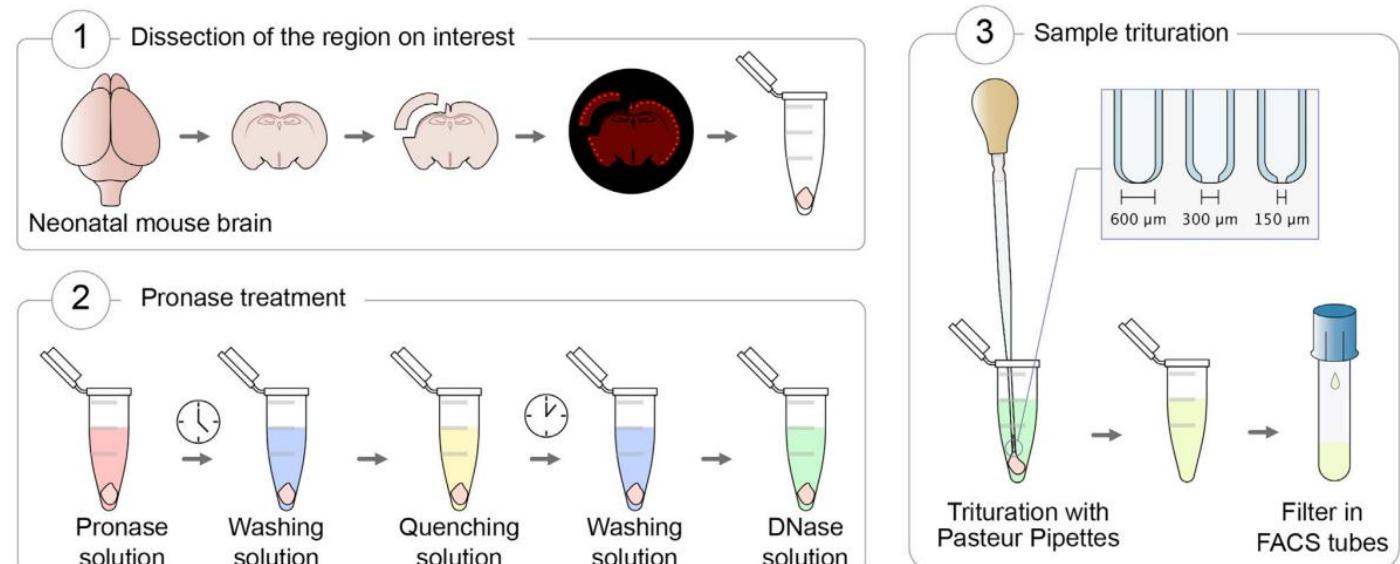
## Blood:

- Straight forward after erythrocyte digestion.

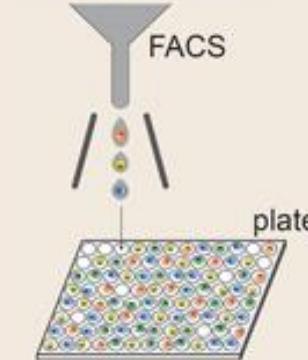
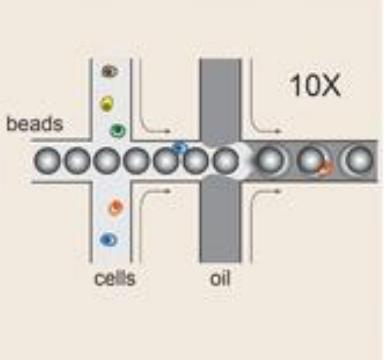
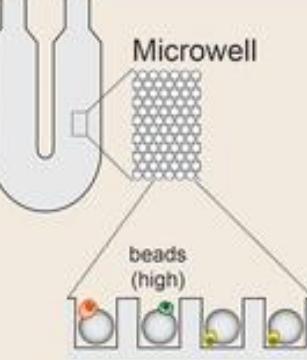
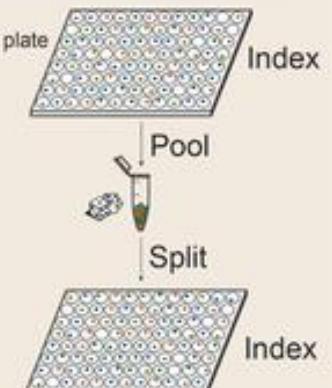
## Cell Sorting by FACS:

- Cells are a bit shocked, be careful

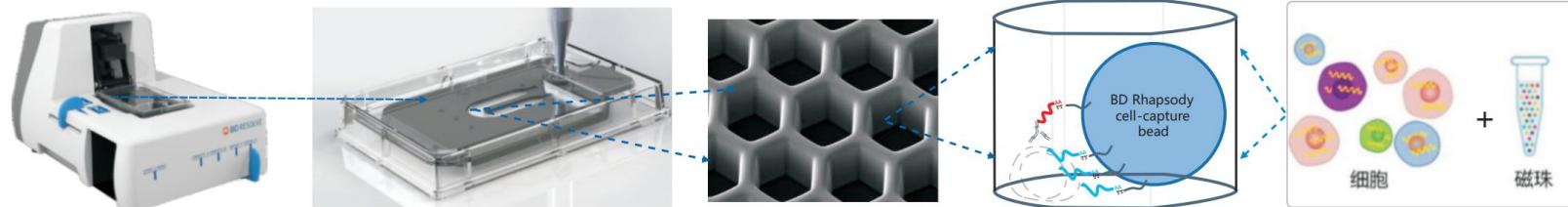
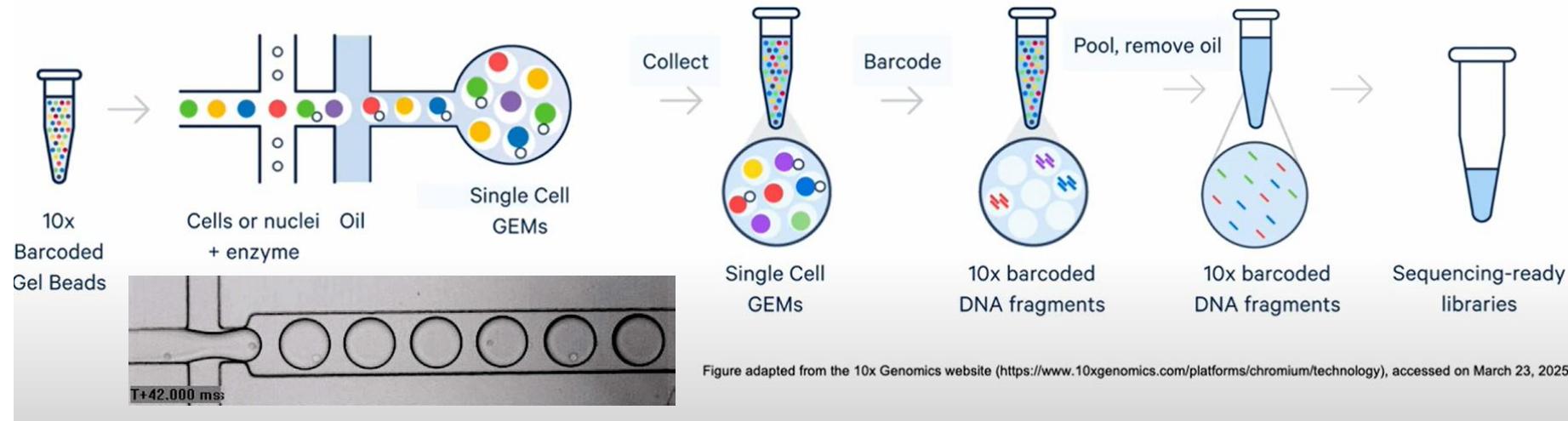
## Cells Enriched by primed beads

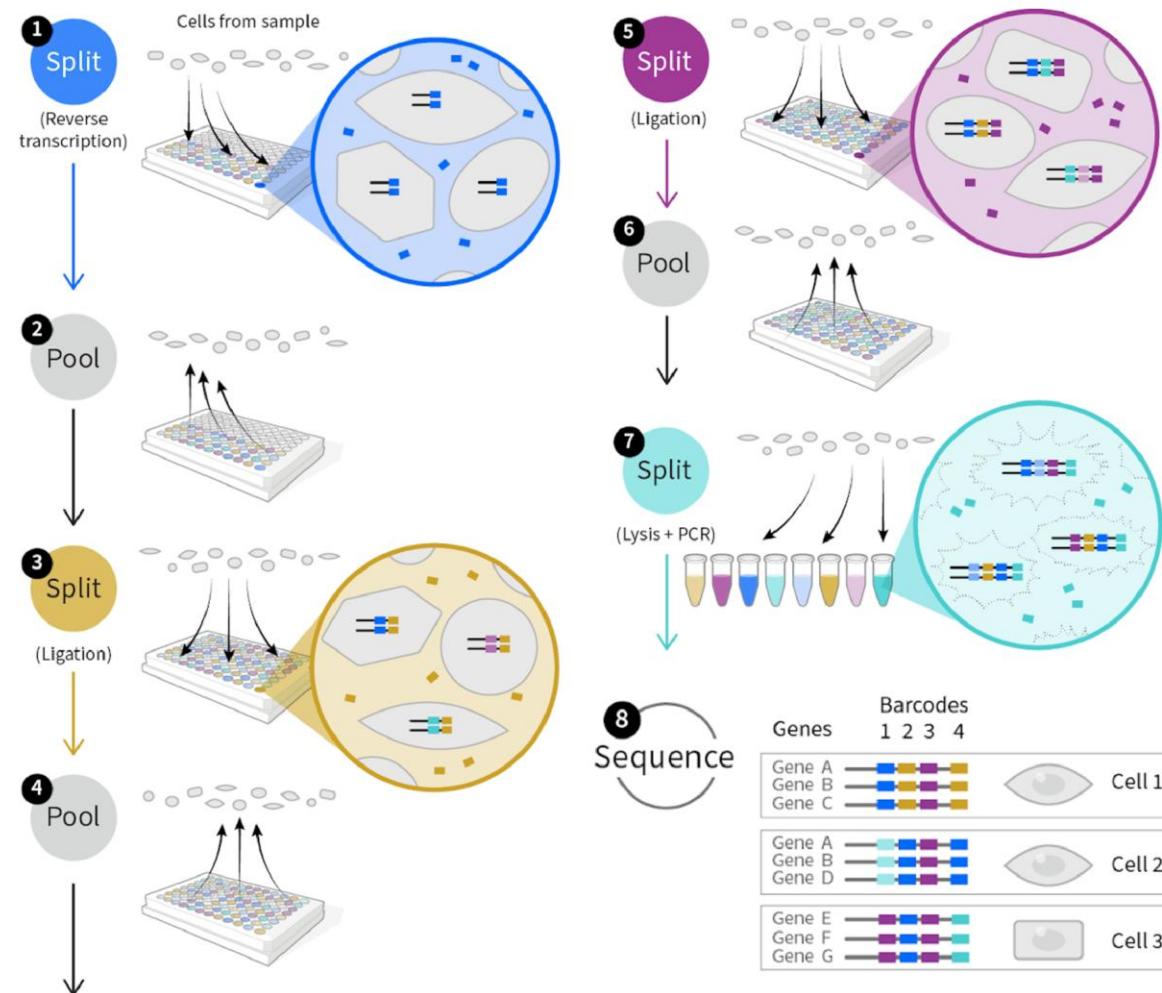


# How do I Isolate The Cell?

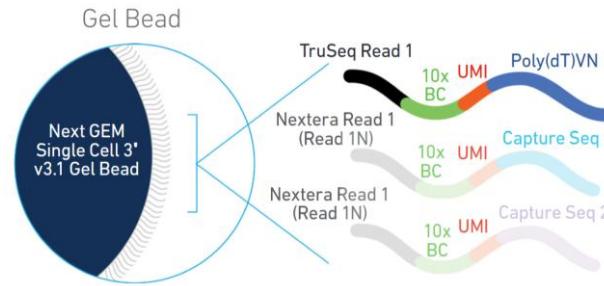
	SMART-seq	DropSeq/10X	BD Rhapsody	sci-RNA-seq
FACS-step	sc sorting into plate	bulk-sorting	bulk-sorting	sc sorting into plate
Principle	isolated cells in wells	microfluidics	microwell	split-pool indexing
Technology				

# 10X Chromium and Rhapsody Scheme





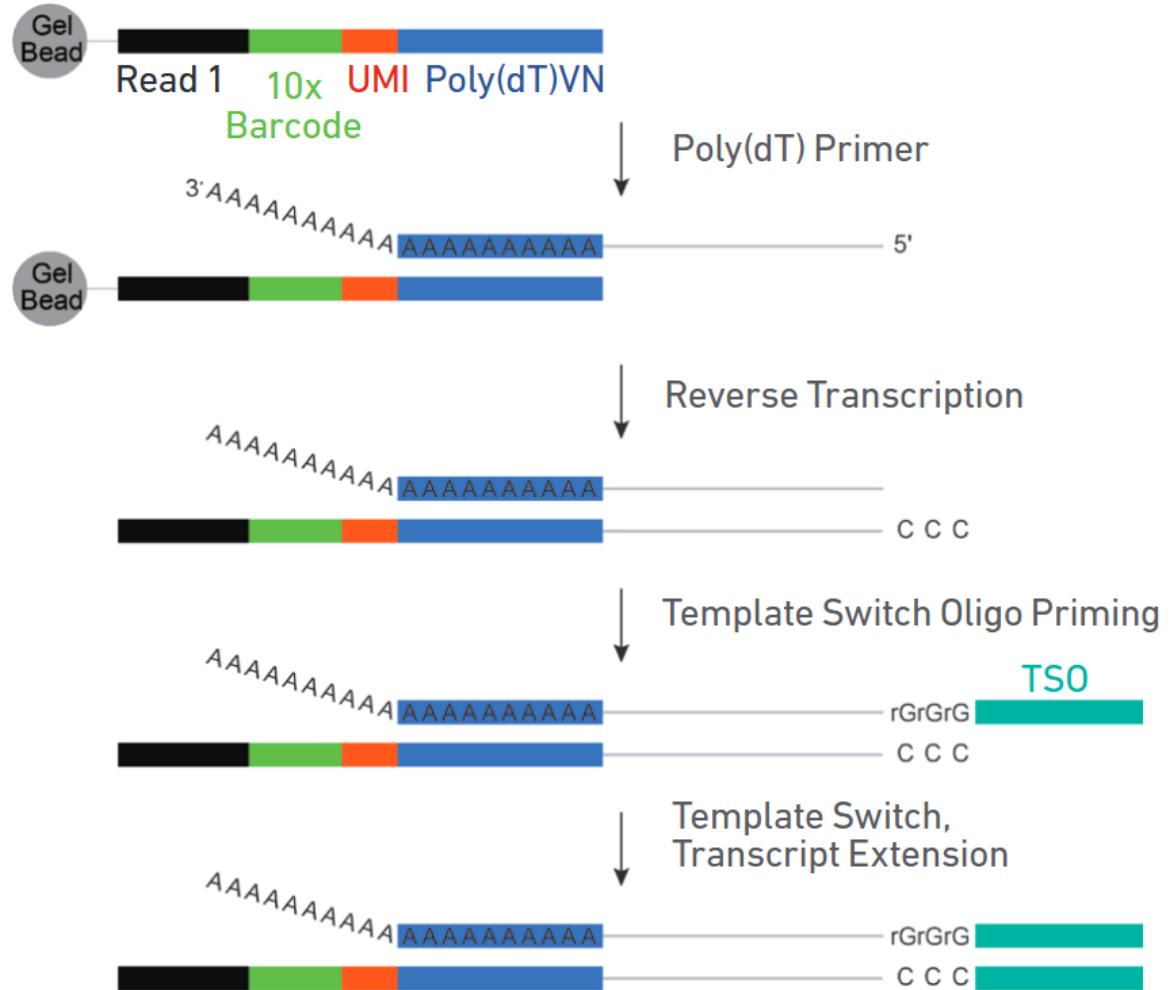
# Library Preparation



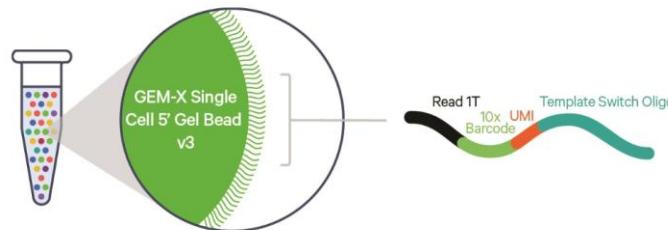
## Gene Expression Library: 3' Library

Within the drop:

- Cells lysate as soon the drop is formed
- Released cytoplasmic mRNA is captured by **PolyT** sequences attached to the bead.
- A cDNA is produced by the RT-Pol using the **TSO oligo** to complete a full cDNA.
- All attached adaptors of a bead carries the same **Cell Barcode**. Each bead has a different Cell Barcode.
- All attached adaptors of a bead carries a different **UMI** (Unique Molecular Identifier).



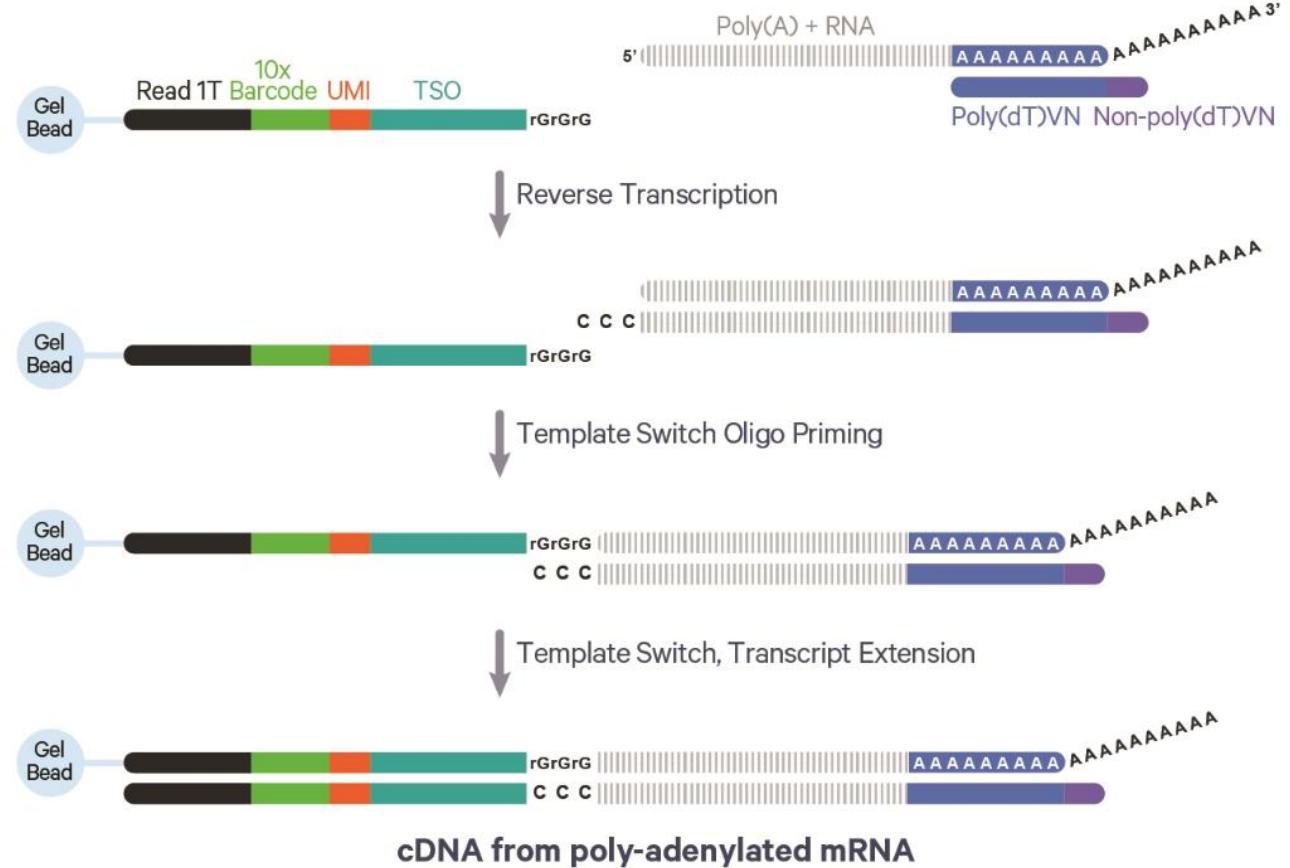
# Library Preparation



## Gene Expression Library: 5' Library

Within the drop:

- Cells lyse as soon the drop is formed
- Released cytoplasmic mRNA is reverse transcribed from a free polyT oligo by the RT-Pol.
- Hybrid RNA/DNA is captured by the TSO attached to the bead.
- A DNA polymerization is produced from the TSO oligo to complete a full cDNA.
- All attached adaptors of a bead carries the same **Cell Barcode**. Each bead has a different Cell Barcode.
- All attached adaptors of a bead carries a different **UMI** (Unique Molecular Identifier).

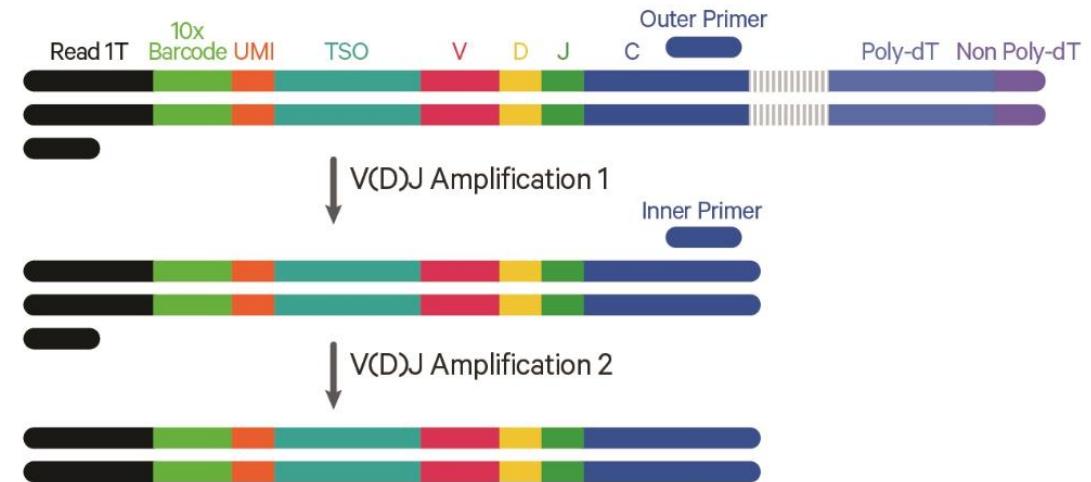


# Library Preparation

## VDJ Enrichment Library: 5' Library

In bulk:

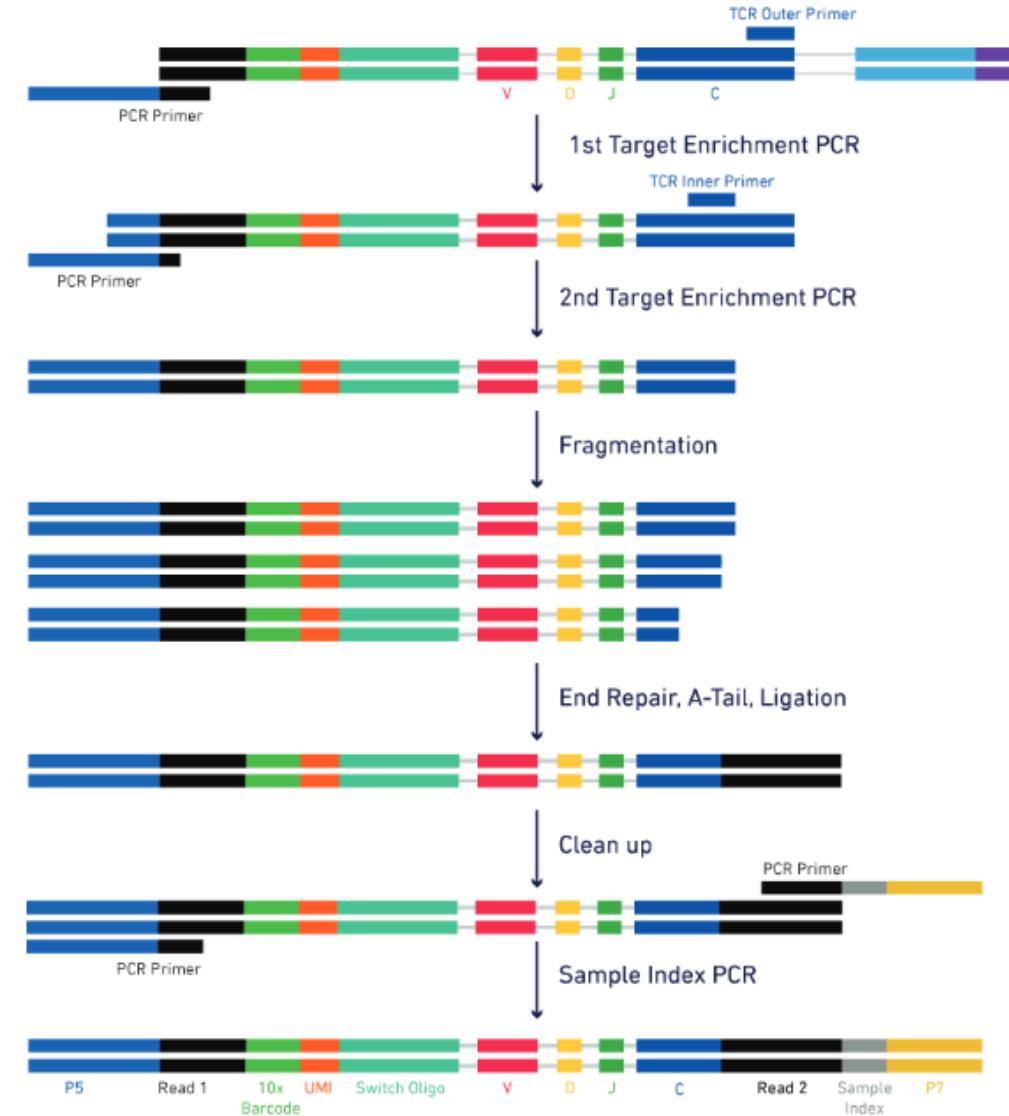
- Nested PCR amplification using same adaptor oligo and two consecutive constant region oligos



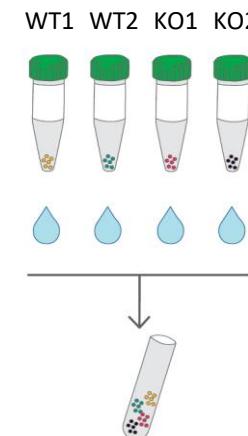
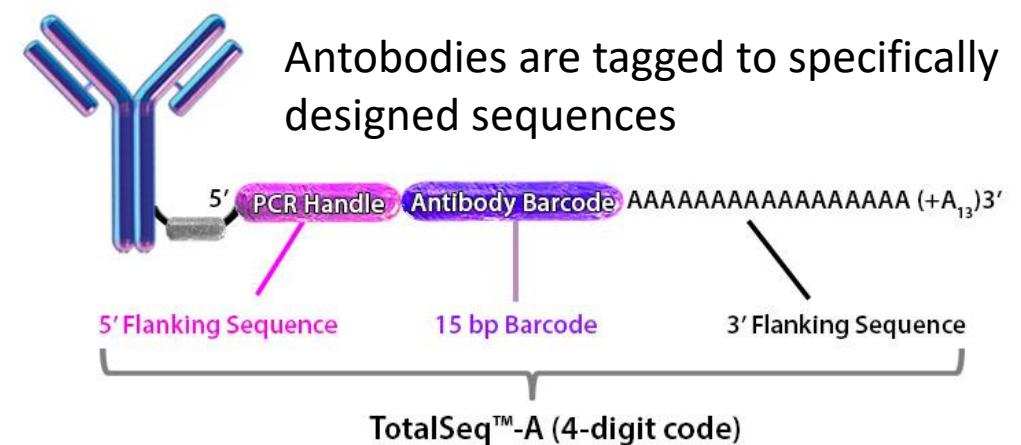
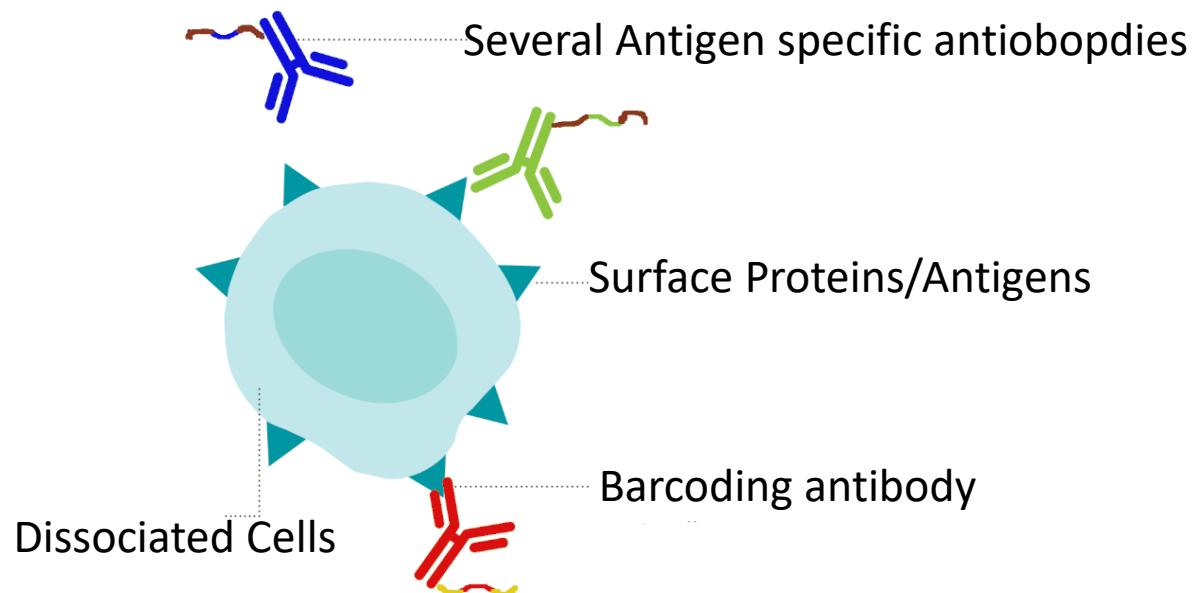
# Library Preparation

In bulk:

- cDNA amplification to increase the chances of sequence it.
- Enzymatic fragmentation to accommodate fragment size to the capacity of the sequencer.
- Further Illumina adaptors and library index ligation.



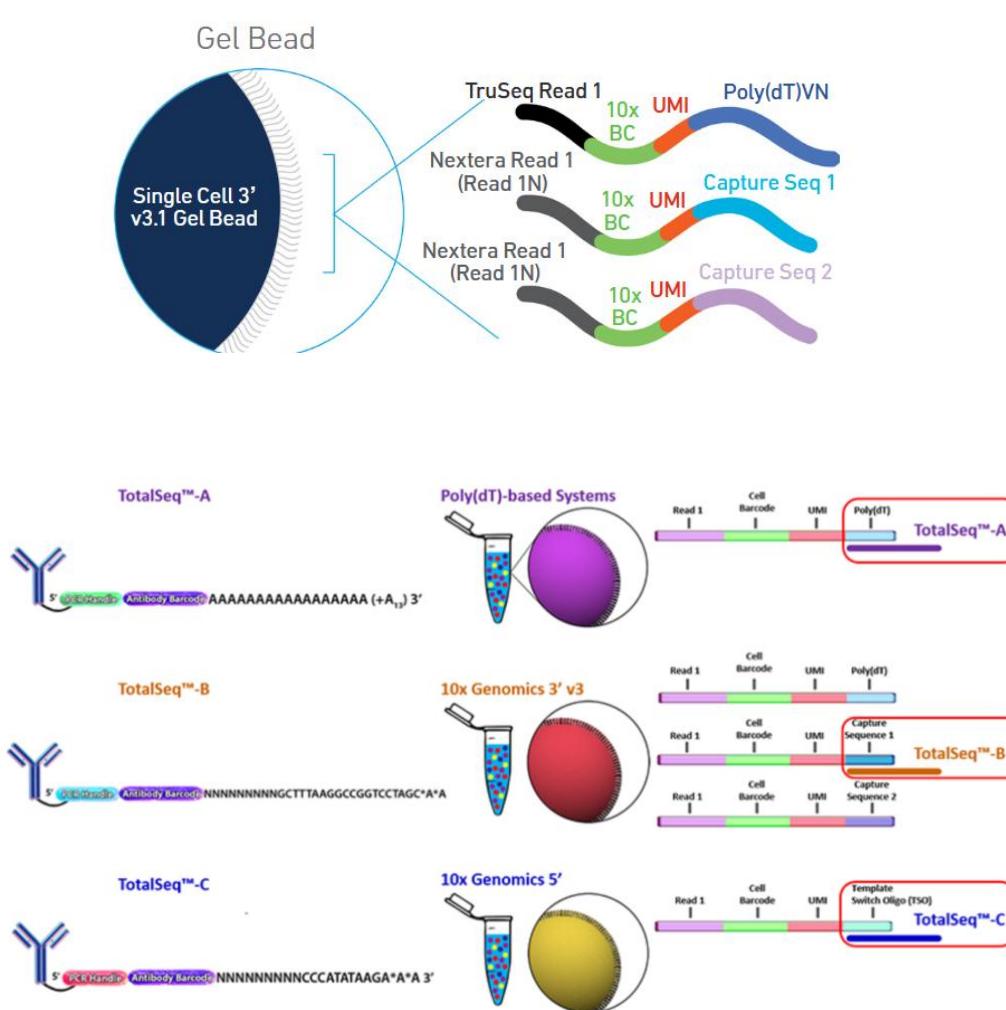
Dissociated Cells are incubated with surface antibodies to evaluate protein expression.



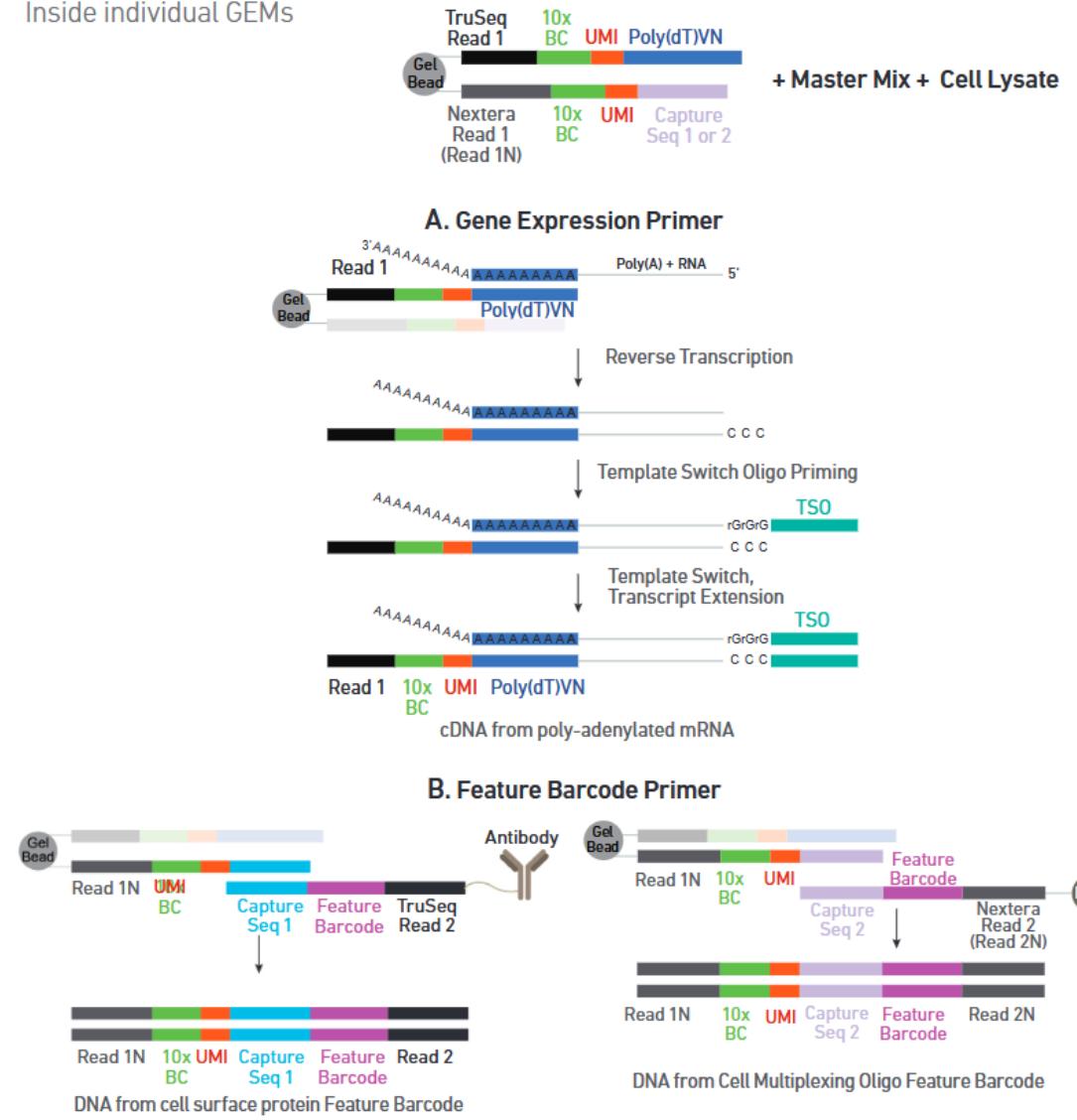
Antibodies against a general ubiquitous surface protein can be used to label samples. Same antibody but different associated barcodes for each tube.<

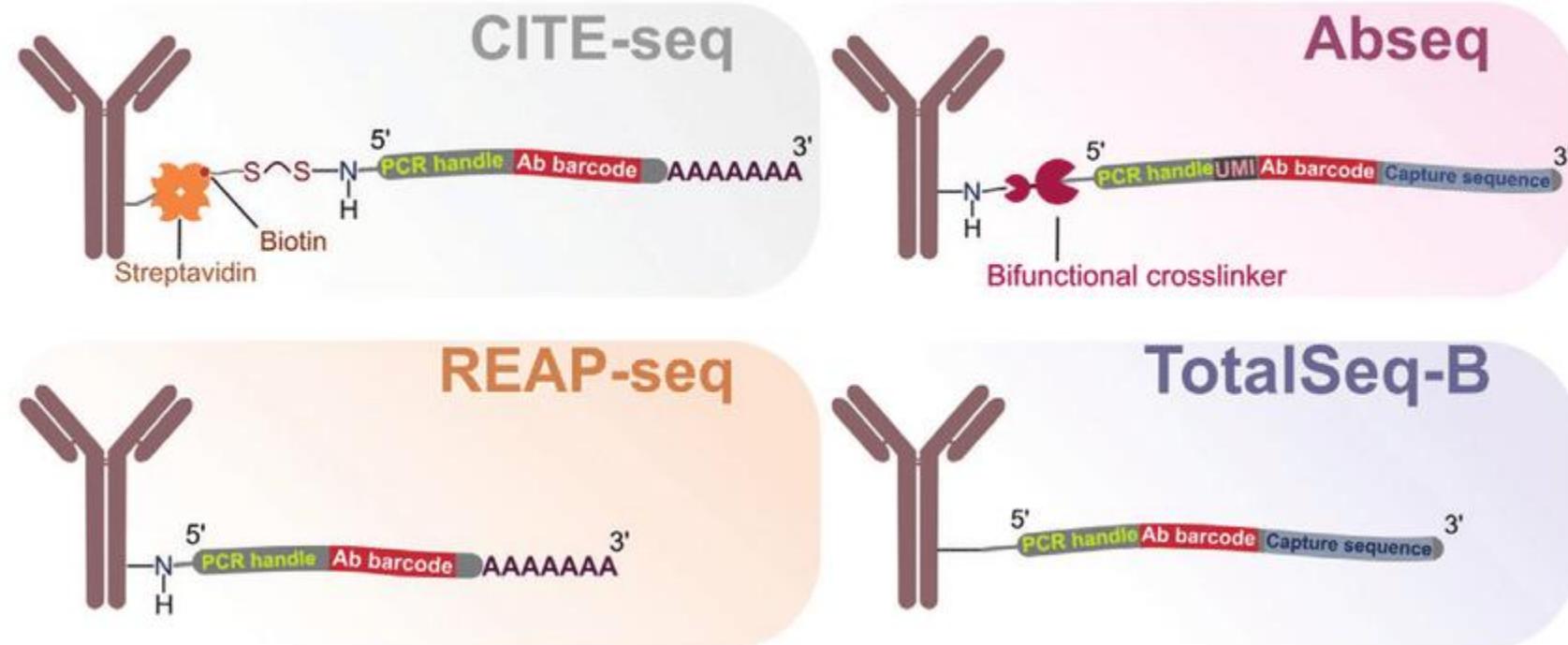
# Ab/Cite-Seq + GEX

*cnic*



## Inside individual GEMs

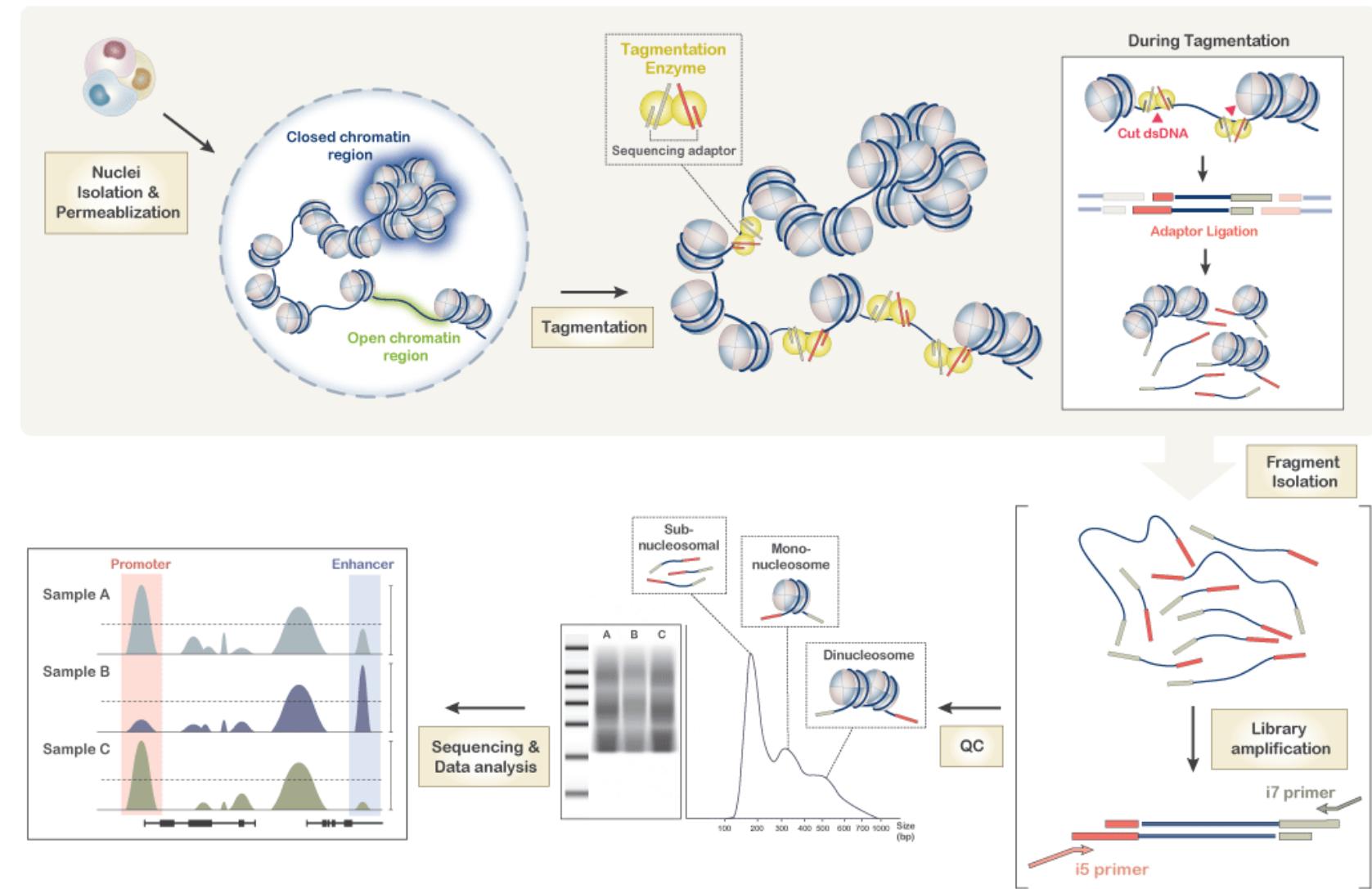




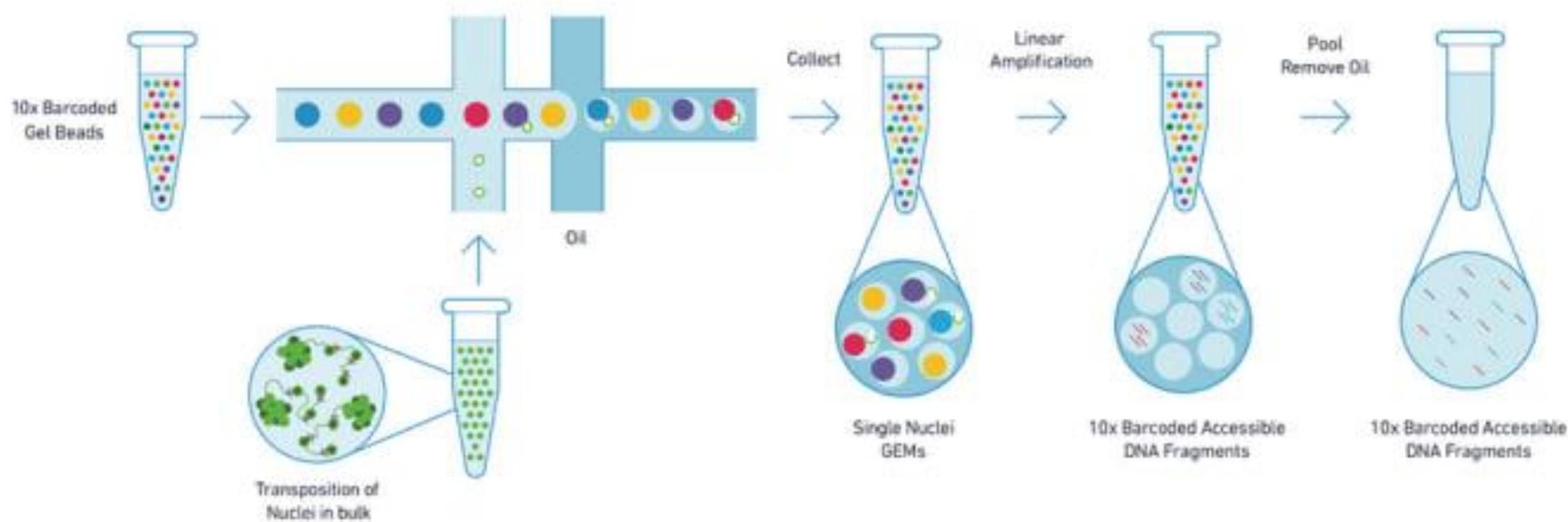
# ATAC-Seq

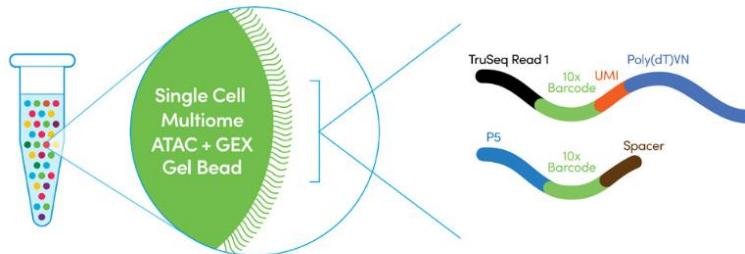
Allows to identify open chromatin regions which is a proxy of a transcriptionally active region.

Easier, more flexible, requires less starting material and cheaper than a ChIP-seq or 4C-Seq experiment



Starting material are isolated nuclei

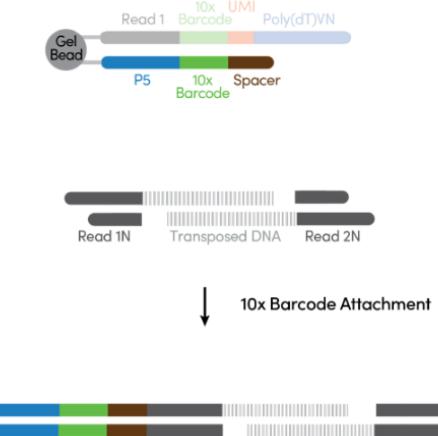




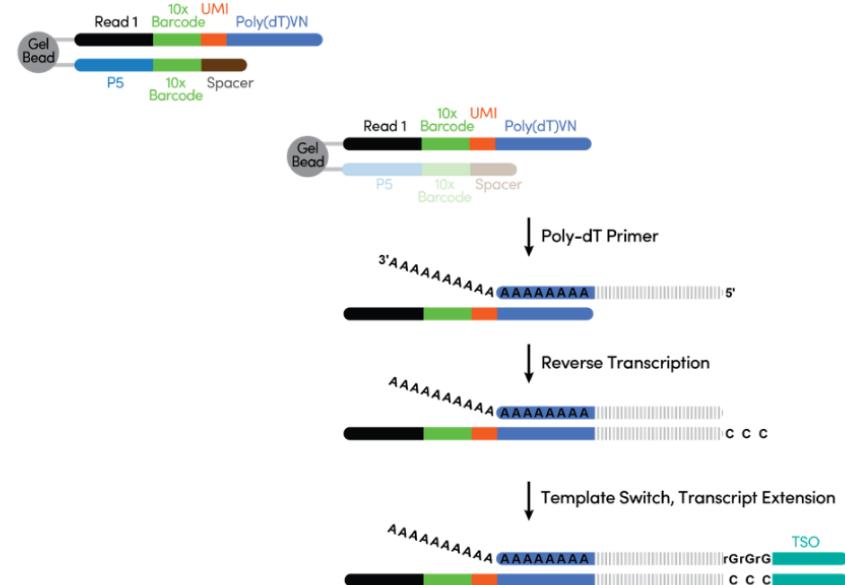
## ATAC + Gene Expression Library: 3' Library

Within the drop:

- Nuclei is lysate as soon the drop is formed.
- Released Tegmentase Fragments are released and captured by the Spacer (Tnf5 sequence complementary sequence)
- The bead barcode is attached to the fragment by PCR.
- Released nuclear mRNA is reverse transcribed from a free polyT oligo by the RT-Pol.
- A DNA polymerization is produced from the TSO oligo to complete a full cDNA.
- All attached adaptors of a bead carries the same **Cell Barcode**. Each bead has a different Cell Barcode.
- All attached adaptors of a bead carries a different **UMI** (Unique Molecular Identifier).

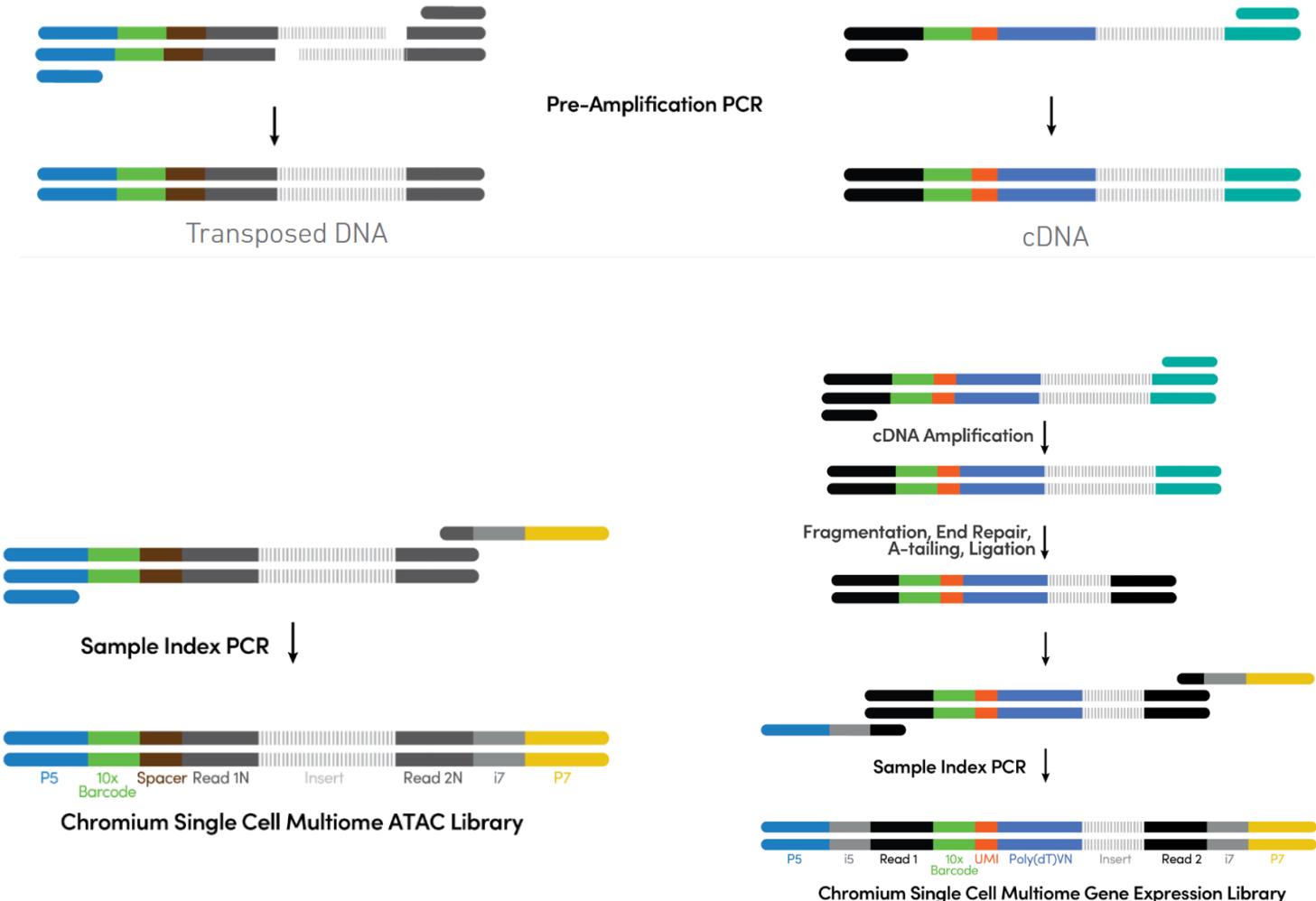


## GEM Generation & Barcoding



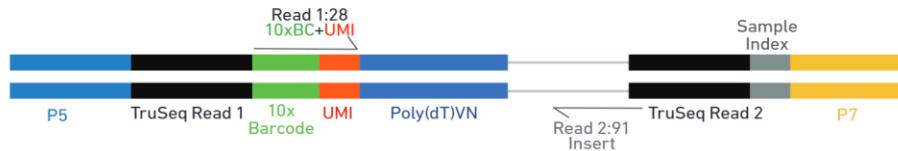
**In bulk:**

- cDNA amplification to increase the chances of sequencing it.
- Enzymatic fragmentation to accommodate fragment size to the capacity of the sequencer.
- Further Illumina adaptors and library index ligation and amplification.

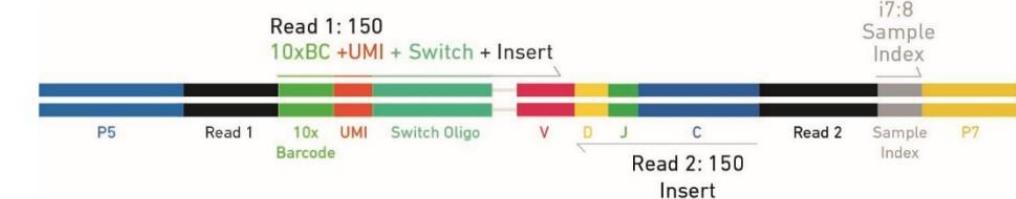


# Library Outputs

Chromium Single Cell 3' Gene Expression Library (Single Index)



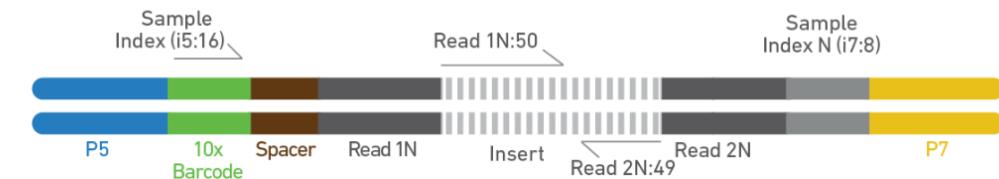
V(D)J Enriched Library Structure:



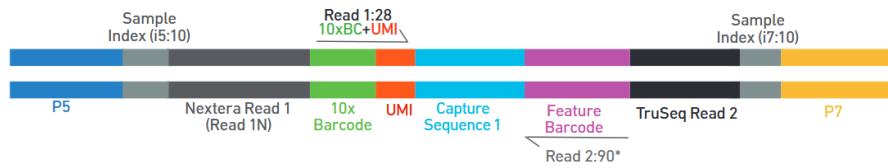
5' Gene Expression Library Structure:



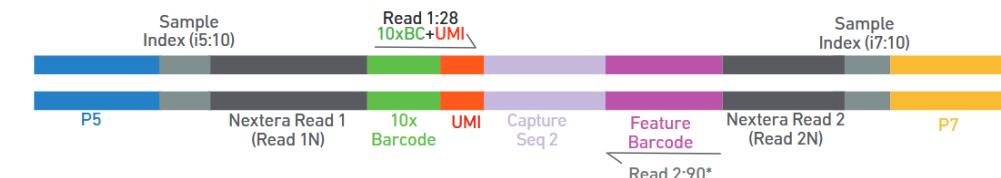
Chromium Single Cell Multiome ATAC Library

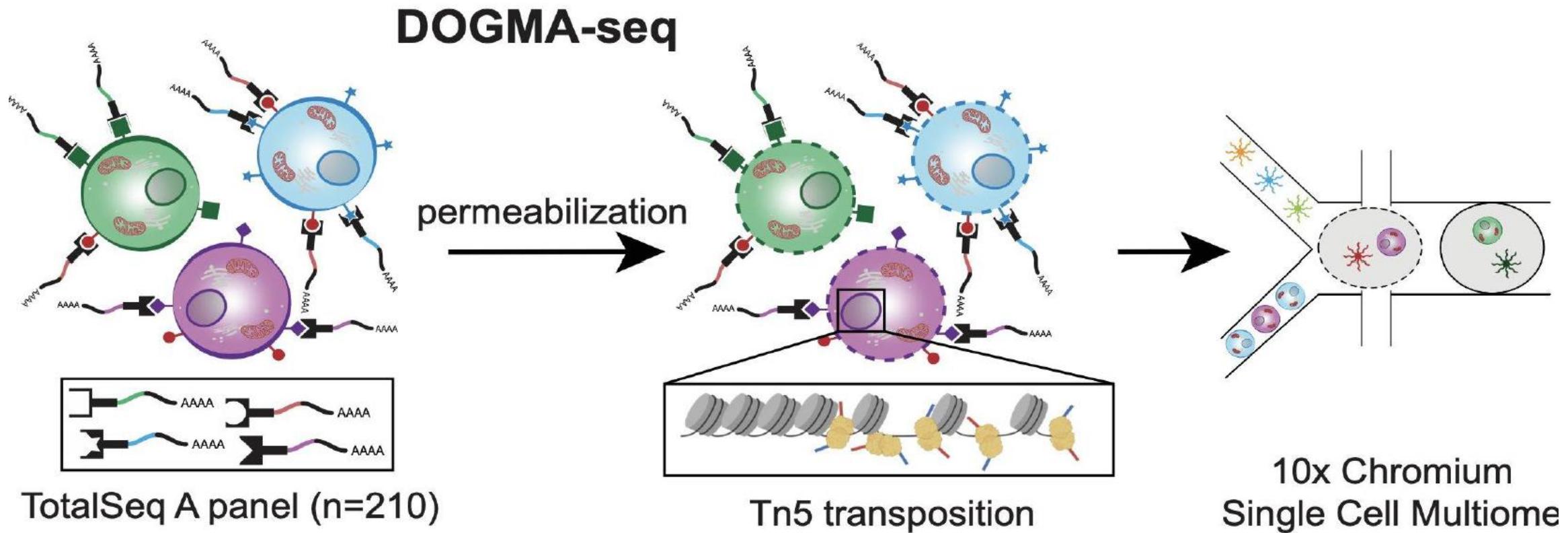


Chromium Single Cell 3' Cell Surface Protein Dual Index Library

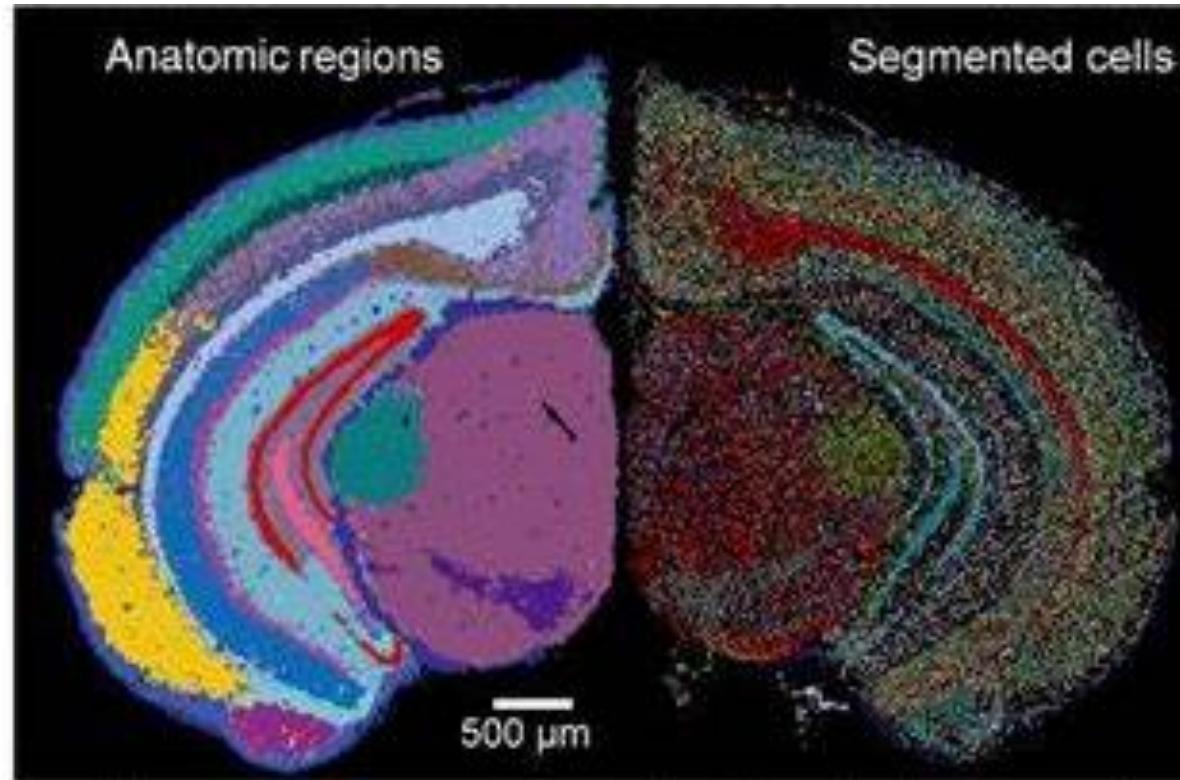


Chromium Single Cell 3' Cell Multiplexing Dual Index Library

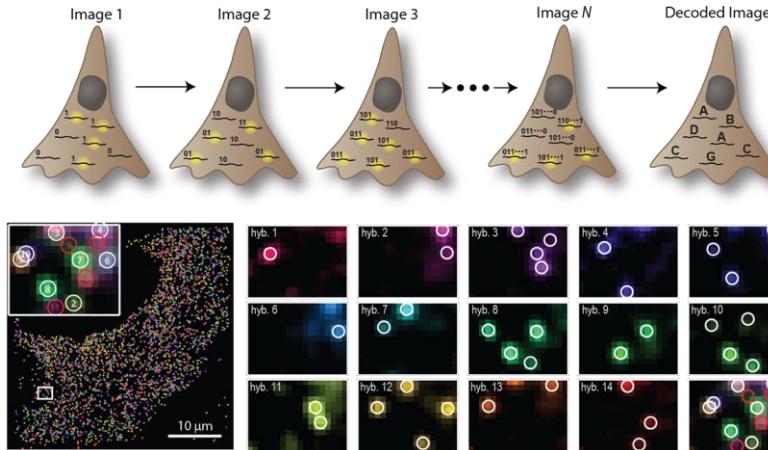




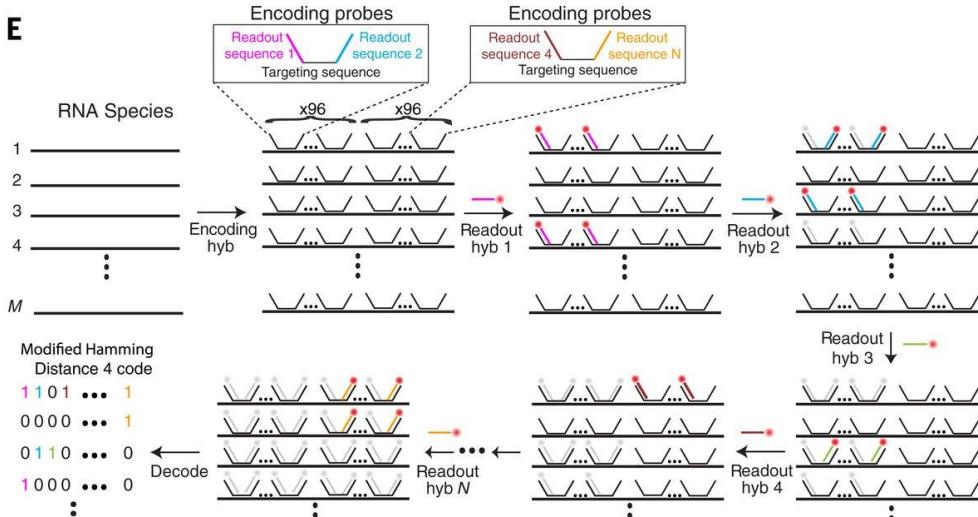
Spatial Single Cell Transcriptomics  
From Cell Function to Tissue Function



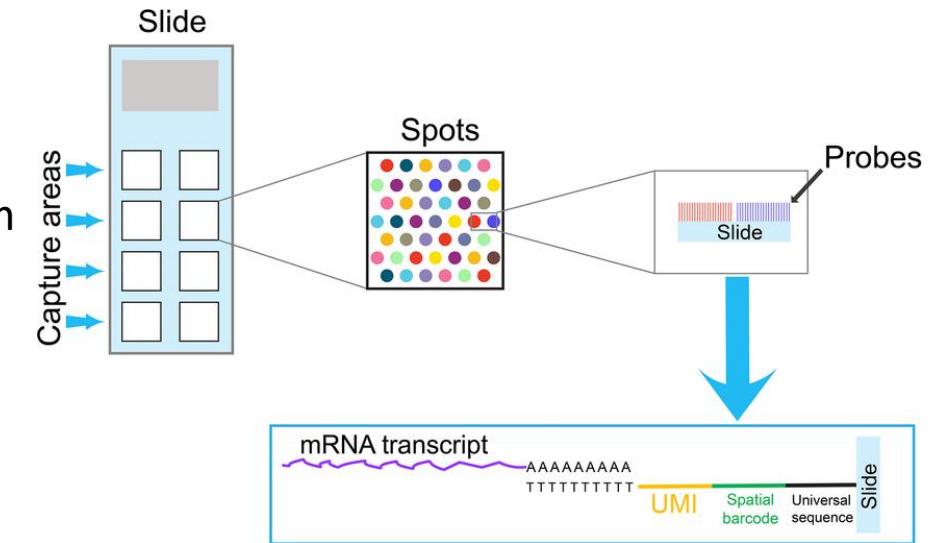
## In situ Image Based Detection MERFISH



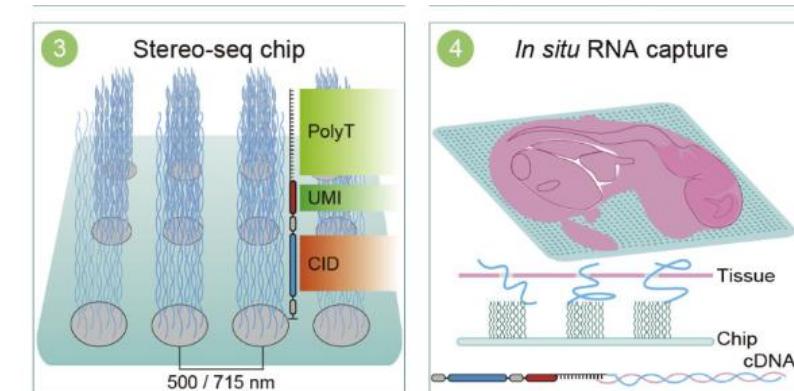
E



## Sequencing Array Based Detection

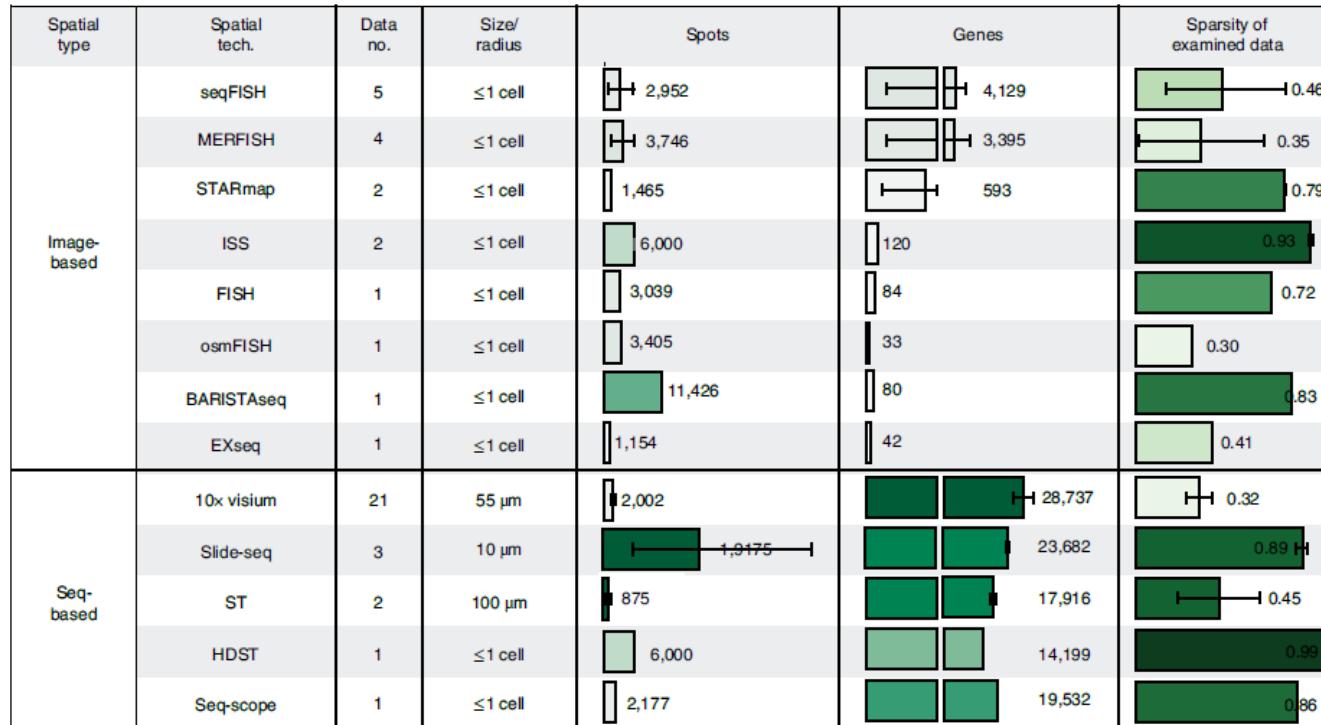


## Stero-Seq



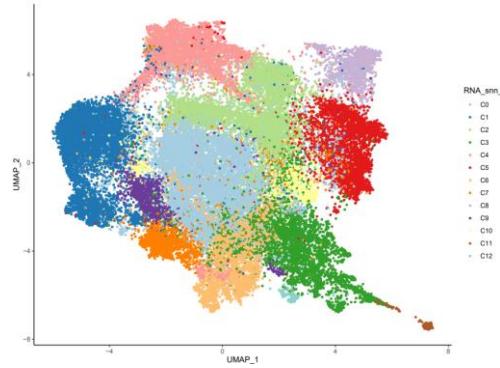
There are many technologies

Several Parameters are important when selecting the technology



## Cell Clustering based on Transcriptome Profile

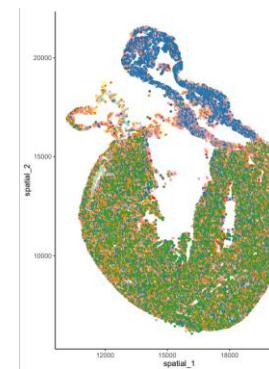
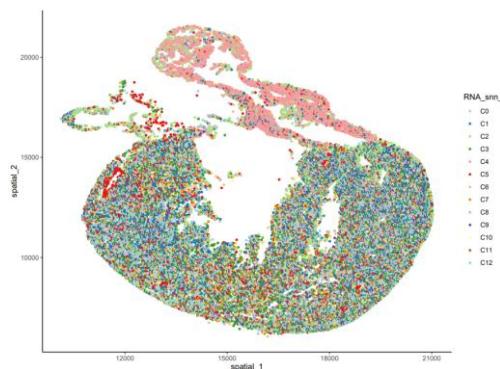
UMAP



## Manual Cell Type Identification

- scMCA\_annotations\_clean
- Atrial cardiomyocyte\_Acta2 high(Neonatal-Heart)
- Atrial cardiomyocyte(Neonatal-Heart)
- Asin2+ Myoblastogenic Progenitor cell\_Acta2 high(Lung-Mesenchyme)
- Cardiac muscle cell(Neonatal-Heart)
- Cardiac Muscle Progenitor cell(Lung-Mesenchyme)
- Endothelial cell\_Igfbp3 high(Neonatal-Heart)
- Epithelial Cell(Cardiac-Heart)
- Left ventricle cardiomyocyte\_Myf5 high(Neonatal-Heart)
- Mesenchyme(Neonatal-Heart)
- Mesenchymal Alveolar Niche Cell\_Don high(Lung-Mesenchyme)
- NA
- slow Muscle selected cell\_Tnct1 high(Muscle)
- Smooth muscle cell(Neonatal-Heart)
- Stromal cell\_Cxcr4 high(Neonatal-Heart)
- Stromal cell\_Syt high(Badder)
- Stromal cell\_Foxo1 high(Neonatal-Heart)
- Vascular endothelial cell(Neonatal-Heart)
- Ventricle cardiomyocyte\_Kcnq1 high(Neonatal-Heart)

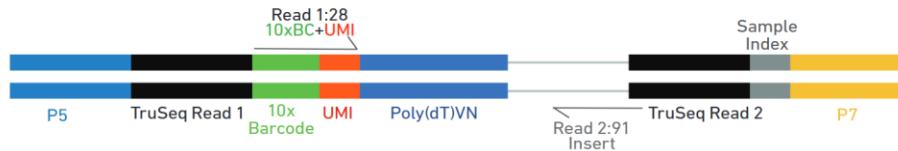
Spatial



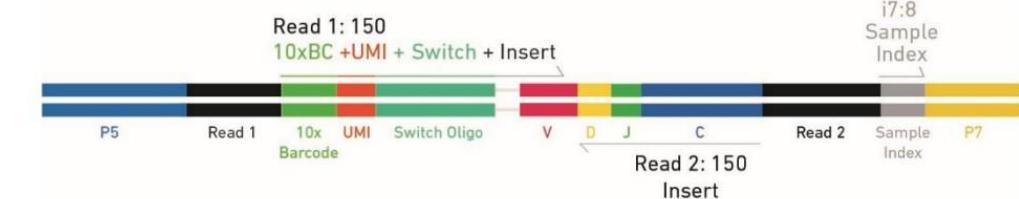
Chen, A., Liao, S., Cheng, M., Ma, K., Wu, L., Lai, Y., ... & Wang, J. (2022). Spatiotemporal transcriptomic atlas of mouse organogenesis using DNA nanoball-patterned arrays. *Cell*, 185(10), 1777-1792. DOI: 10.1016/j.cell.2022.04.003

# Library Outputs

Chromium Single Cell 3' Gene Expression Library (Single Index)



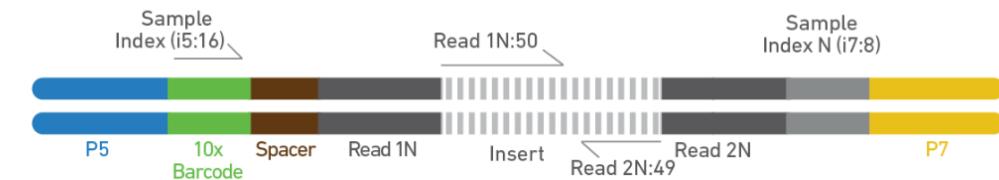
V(D)J Enriched Library Structure:



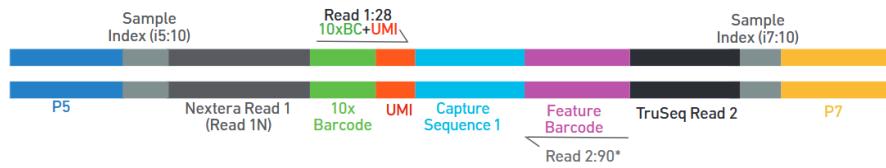
5' Gene Expression Library Structure:



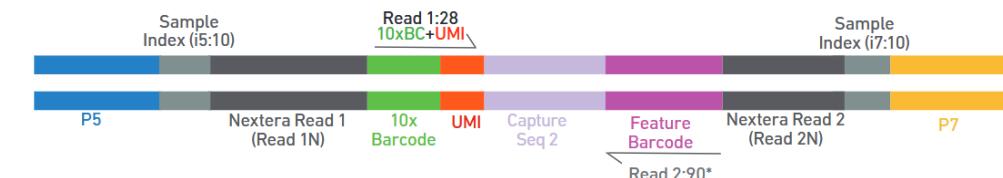
Chromium Single Cell Multiome ATAC Library



Chromium Single Cell 3' Cell Surface Protein Dual Index Library

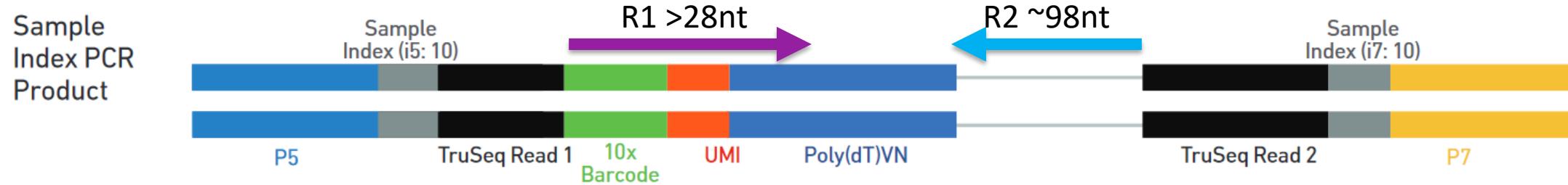


Chromium Single Cell 3' Cell Multiplexing Dual Index Library



# Sequencing

- Single Cell Libraries are sequenced in paired end mode.
- Different Single cell platforms and library protocols require different lengths for both reads



5'-AATGATAACGGCGACCACCGAGATCTACAC-N10-ACACTTTCCCTACACGACGCTTCCGATCT-N16-N12-TTTTTTTTTTTTTTTTTVN-cDNA\_Insert-AGATCGGAAGAGCACACGTCTGAACCTCCAGTCAC-N10-ATCTCGTATGCCGTCTCTGCTTG-3'  
3'-TTACTATGCCCTGGGGCTCTAGATGTG-N10-TGTGAGAAAGGGATGTGCTGGAGAAGGCTAGA-N16-N12-AAAAAAAAAAAAAAAABN-cDNA\_Insert-TCTAGCCTTCTCGTGTGAGACTTGAGGTCACTG-N10-TAGAGCATACGGCAGAACAGAAC-5'

- **R1** Contains Cell Barcode and UMI
- **R2** Contains sequence from the mRNA captured molecule

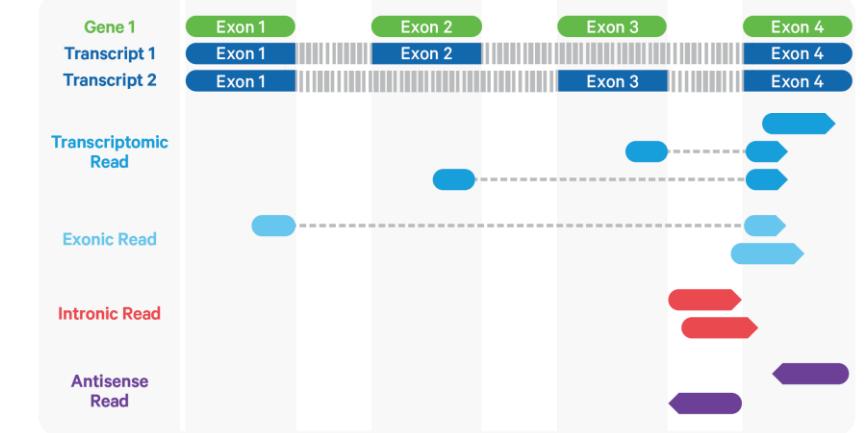
# Alignment and Quantification

## Pre-Processing:

- R2 is annotated with the Cell Barcode and UMI from R1
- From now only annotated R2 are used.

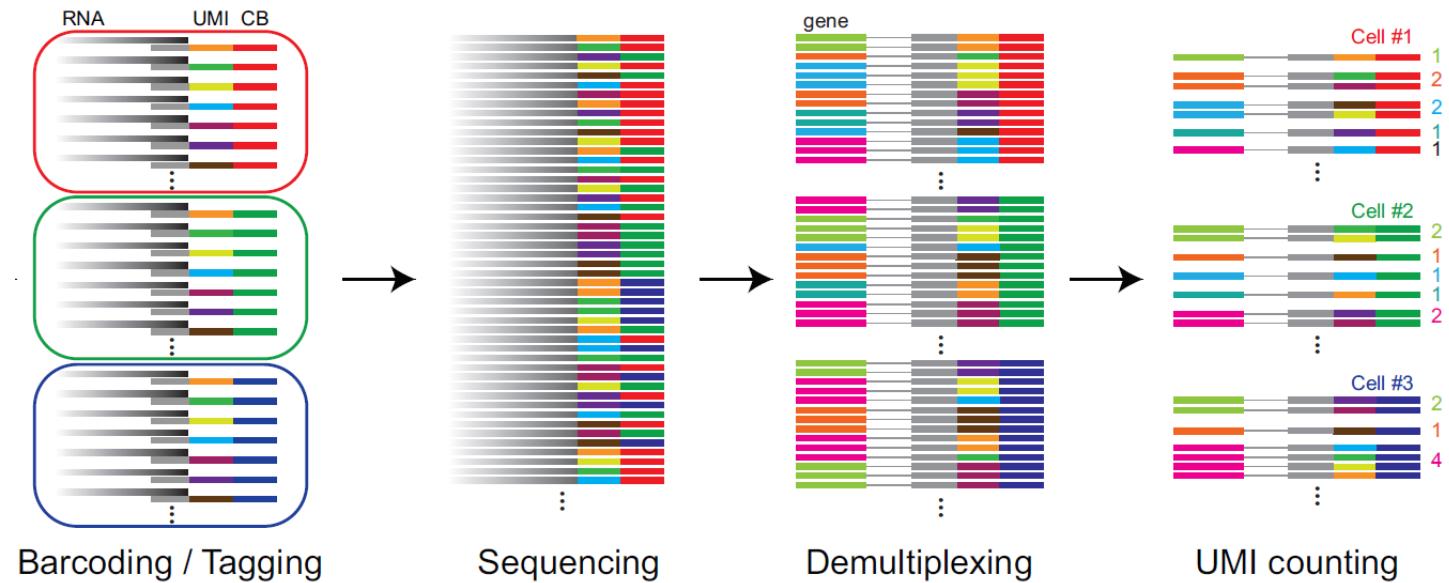
## Alignment:

- R2 is mapped to the genome using STAR.
- Only transcript compatible alignments associated to a unique gene are considered for next UMI counting steps.
- Each mapped R2 is annotated with the Gene it is compatible with.



## UMI quantification:

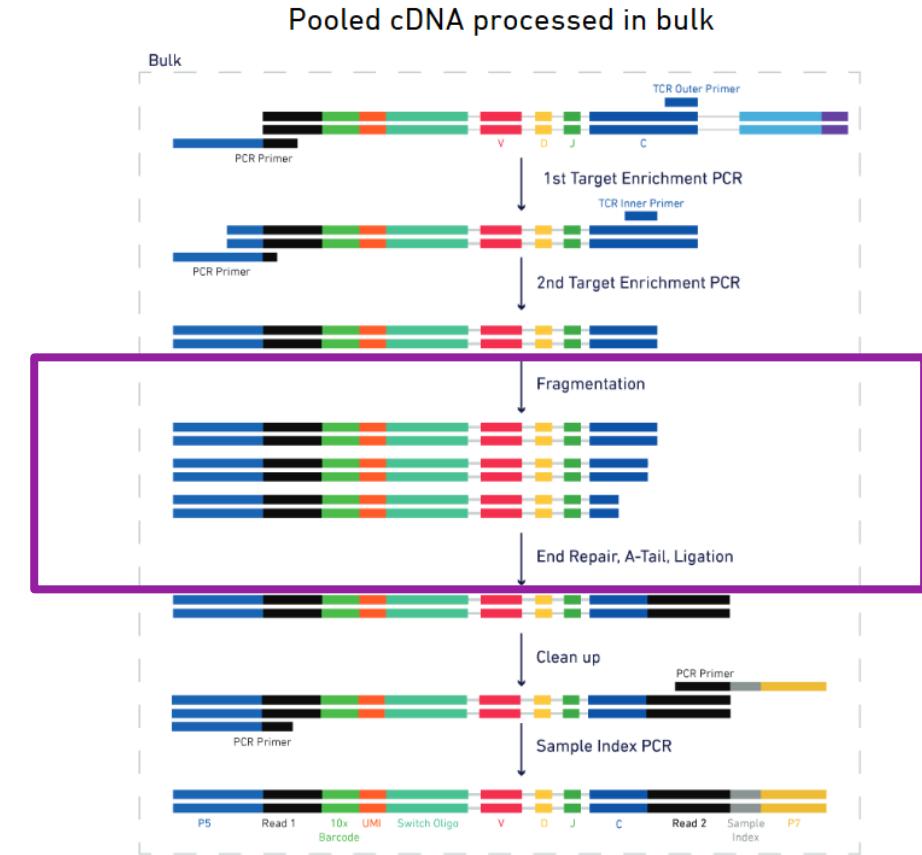
- Now we have a combination of Cell Barcode/UMI/Gene.
- All identical combinations are collapsed to one (they are PCR duplicates from one original mRNA on that cell).
- UMIs per gene per cell are computed.



# Library Preparation

## One UMI = One mRNA molecule

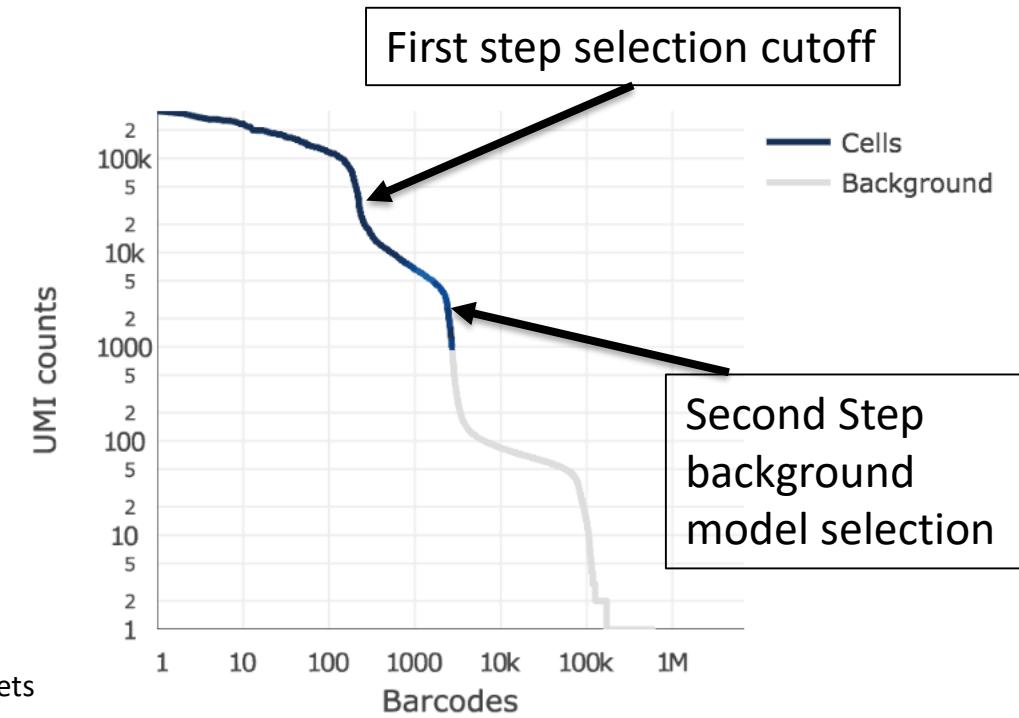
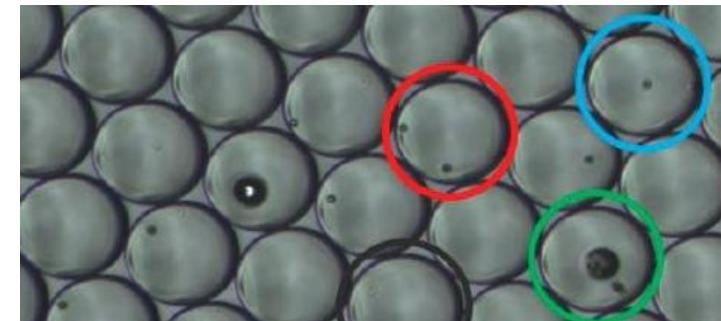
- We may observe the same UMI attached to different fragments of the same mRNA.
- Those different fragments are the product of the enzymatic fragmentation on a PCR amplified UMI-cDNA molecule.
- All identical UMIs have been originated from the same individual mRNA molecule



## Cell Detection:

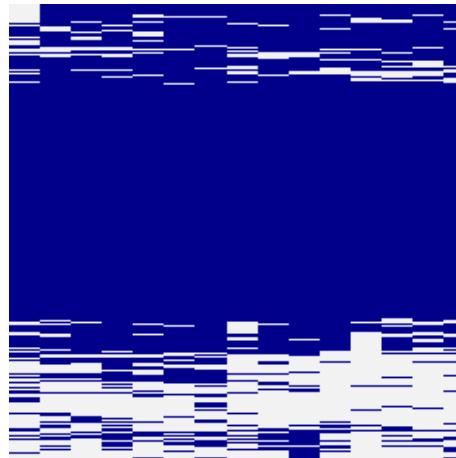
- Not all beads from a cell capture step have got a cell.
- Empty beads still may capture free background RNA and produce cDNAs.
- Cell calling step tries to differentiate true cells from background
- In general, true cells produce more RNA and accounting for more genes than empty beads.
- However there are certain cell types with low RNA content that might yield library product at similar levels to background.
- Cell calling is perform in two steps:
  - First a cut-off is placed at the point in which counts per cell drops dramatically.
  - On the lower yield cells a model of the background based on the data itself is produced and cells are called when they differ substantially from that model.
- Only transcript compatible alignments associated to a unique gene are considered for next UMI counting steps.

Lun, A., Riesenfeld, S., Andrews, T. et al. EmptyDrops: distinguishing cells from empty droplets in droplet-based single-cell RNA sequencing data. *Genome Biol* **20**, 63 (2019).  
<https://doi.org/10.1186/s13059-019-1662-y>

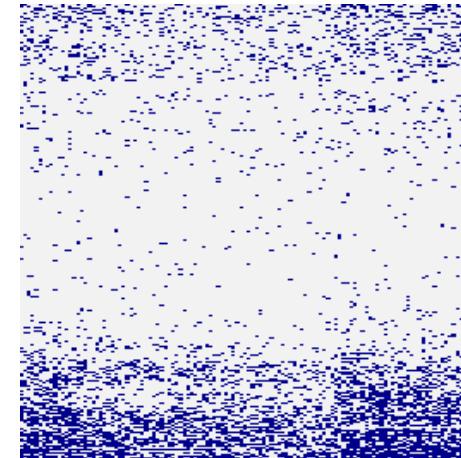


**Data is very sparse**

Bulk Samples



SingleCell Samples

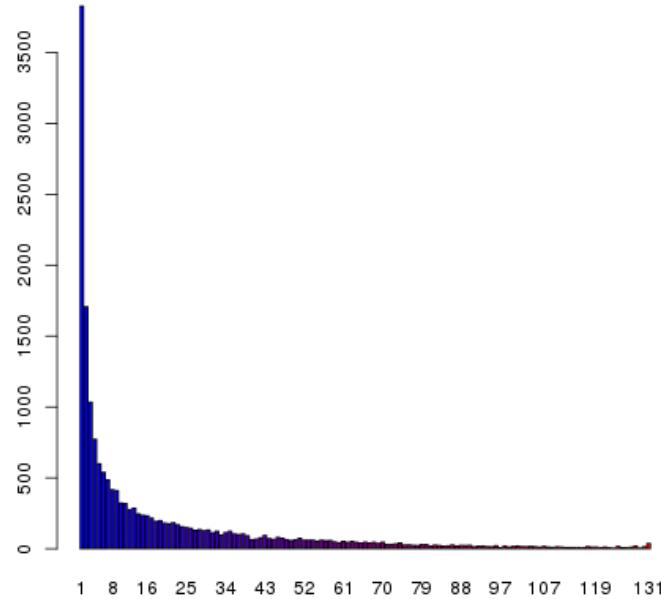


One Count in at least One Sample

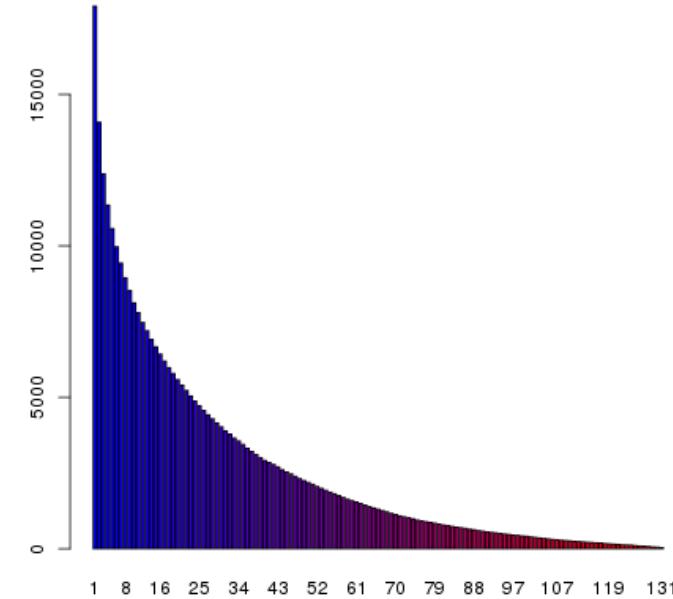
None of the assumptions for Bulk RNA-Seq applies to a scRNA-Seq experiment

# How the data looks like?

Genes detected in “n” Samples

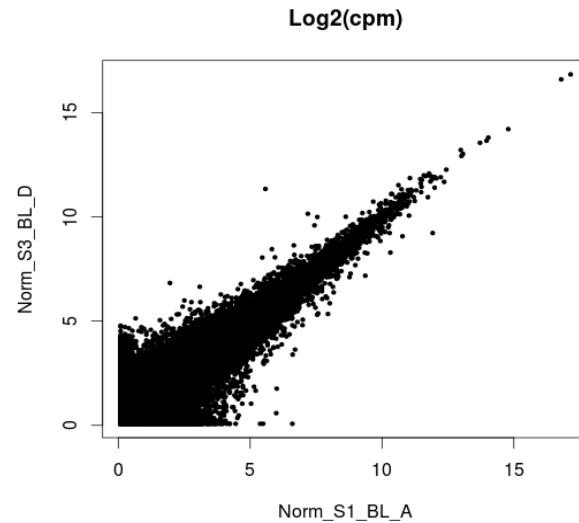


Genes detected in at least “n” Samples

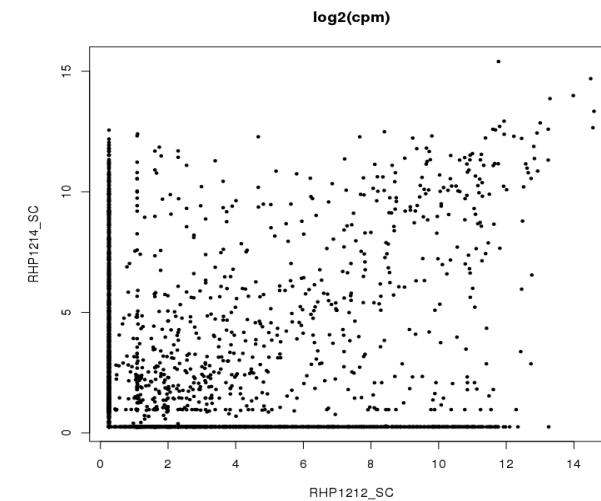


**Correlations between replicates is very poor**

Bulk Samples



SingleCell Samples



Seurat (R) : <https://satijalab.org/seurat/index.html>

Scater (R): <https://bioconductor.org/packages/release/bioc/html/scater.html>

Scran (R): <https://bioconductor.org/packages/release/bioc/html/scran.html>

Scanpy (Python): <https://scanpy.readthedocs.io/en/stable/>

---

*Cell.* 2019 June 13; 177(7): 1888–1902.e21. doi:10.1016/j.cell.2019.05.031.

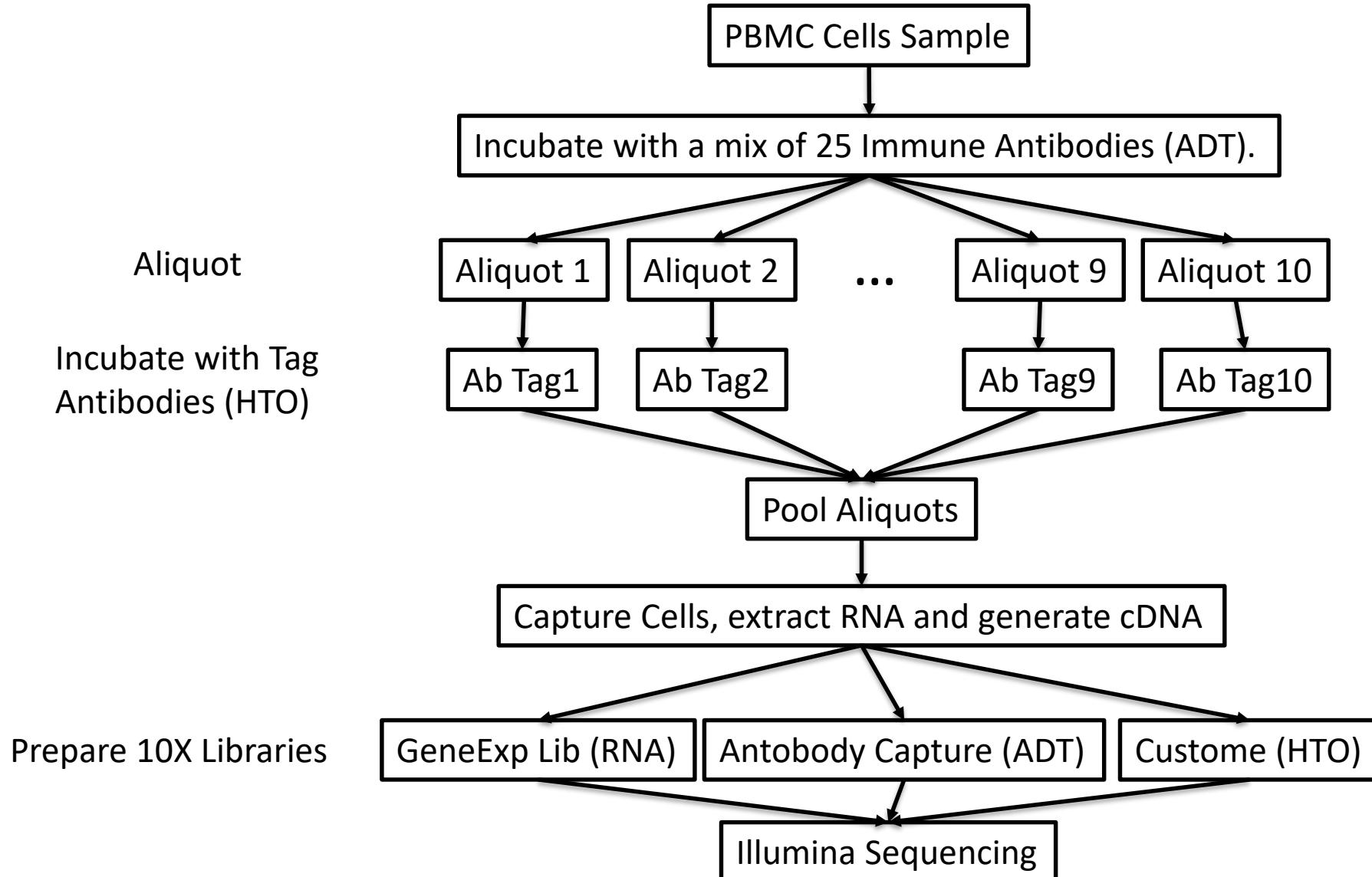
## Comprehensive integration of single-cell data

**Tim Stuart<sup>1,\*</sup>, Andrew Butler<sup>1,2,\*</sup>, Paul Hoffman<sup>1</sup>, Christoph Hafemeister<sup>1</sup>, Efthymia Papalexi<sup>1,2</sup>, William M. Mauck III<sup>1,2</sup>, Yuhan Hao<sup>1,2</sup>, Marlon Stoeckius<sup>3</sup>, Peter Smibert<sup>3</sup>, Rahul Satija<sup>1,2,\*\*</sup>**

<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE128639>

For the scRNA/CITE-seq experiment: a panel of 25 antibodies (ADTs) was used to stain the cells according to the CITE-seq protocol. To identify cell doublets, we split cells into 10 aliquots and stained each aliquot with a unique hashing (HTOs) antibody. Samples were washed and recombined and load on a 10genomics chip. 3prime 10x genomics Version 2 was ran and the manufacturers protocols was followed for cDNA library generation. To recover HTOs and ADTs we followed the CITE-seq protocol found here: <https://cite-seq.com/protocol/>

# Experiment Protocol



Using CellRanger pipeline versión 6. <https://www.10xgenomics.com> -> Support -> Software -> CellRanger

## Command

```
cellranger count --id=GSE128639_10XGenomics \
    --libraries=GSE128639_10XGenomics_Sample_fastqs.csv \
    --transcriptome=refdata-gex-GRCh38-2020-A \
    --feature-ref=GSE128639_10XGenomics_HTORefTable_feature_ref.csv \
    --chemistry=threeprime --localcores=40 --localmem=96 --disable-ui --nosecondary --no-bam
```

### Sample\_fastqs.csv

fastq	sample	library_type
HTOlib	GSE128639hto	Custom
ADTlib	GSE128639adt	Antibody Capture
GEXlib	GSE128639gex	Gene Expression

### HTORefTable\_feature\_ref.csv

id	name	read	pattern	sequence	feature_type
HTO1	LNH-94a	R2,5P(BC)	GTCAACTCTTAGCG		Custom
HTO2	LNH-94b	R2,5P(BC)	TGATGGCCTATTGGG		Custom
...					
HTO9	LNH-94i	R2,5P(BC)	CAGTAGTCACGGTCA		Custom
HTO10	LNH-94j	R2,5P(BC)	ATTGACCCGCGTTAG		Custom
CD3	CD3	R2,5P(BC)	CTCATTGTAACCT		Antibody Capture
CD56	CD56	R2,5P(BC)	TCCTTCCCTGATAGG		Antibody Capture
...					

Fastq Sample Prefix

Path to folder containing fastqs of the described library

# CellRanger Output folder

 antibody_analysis	15/03/2023 13:52	File folder
 filtered_feature_bc_matrix	15/03/2023 13:52	File folder
 raw_feature_bc_matrix	15/03/2023 13:52	File folder
 sample_feature_bc_matrix	15/03/2023 13:52	File folder
 feature_reference.csv	15/03/2023 04:40	Microsoft Excel Co... 2 KB
 filtered_feature_bc_matrix.h5	15/03/2023 04:12	H5 File 51,372 KB
 metrics_summary.csv	15/03/2023 04:40	Microsoft Excel Co... 2 KB
 molecule_info.h5	15/03/2023 04:18	H5 File 1,746,390 ...
 raw_feature_bc_matrix.h5	15/03/2023 03:58	H5 File 132,557 KB
 web_summary.html	15/03/2023 04:40	Firefox HTML Doc... 2,559 KB

# CellRanger Output: Matrix Market or h5 file

## barcodes.tsv.gz

```
AAACCTGAGACAAAGG-1  
AAACCTGAGACGCACA-1  
AAACCTGAGAGGGCTT-1  
AAACCTGAGAGTAAGG-1  
AAACCTGAGAGTACAT-1  
AAACCTGAGAGTGAGA-1  
AAACCTGAGAGTTGGC-1  
AAACCTGAGCTAACTC-1  
AAACCTGAGCTGAAAT-1  
AAACCTGAGCTTATCG-1  
AAACCTGAGGTCGGAT-1
```

## features.tsv.gz

```
ENSG00000243485 MIR1302-2HG Gene Expression  
ENSG00000237613 FAM138A Gene Expression  
ENSG00000186092 OR4F5 Gene Expression  
ENSG00000238009 AL627309.1 Gene Expression  
ENSG00000239945 AL627309.3 Gene Expression  
ENSG00000239906 AL627309.2 Gene Expression  
ENSG00000241860 AL627309.5 Gene Expression  
ENSG00000241599 AL627309.4 Gene Expression  
...  
HT08 LNH-94h Custom  
HT09 LNH-94i Custom  
HT010 LNH-94j Custom  
CD3 CD3 Antibody Capture  
CD56 CD56 Antibody Capture  
CD19 CD19 Antibody Capture  
CD11c CD11c Antibody Capture
```

## matrix.mtx.gz:

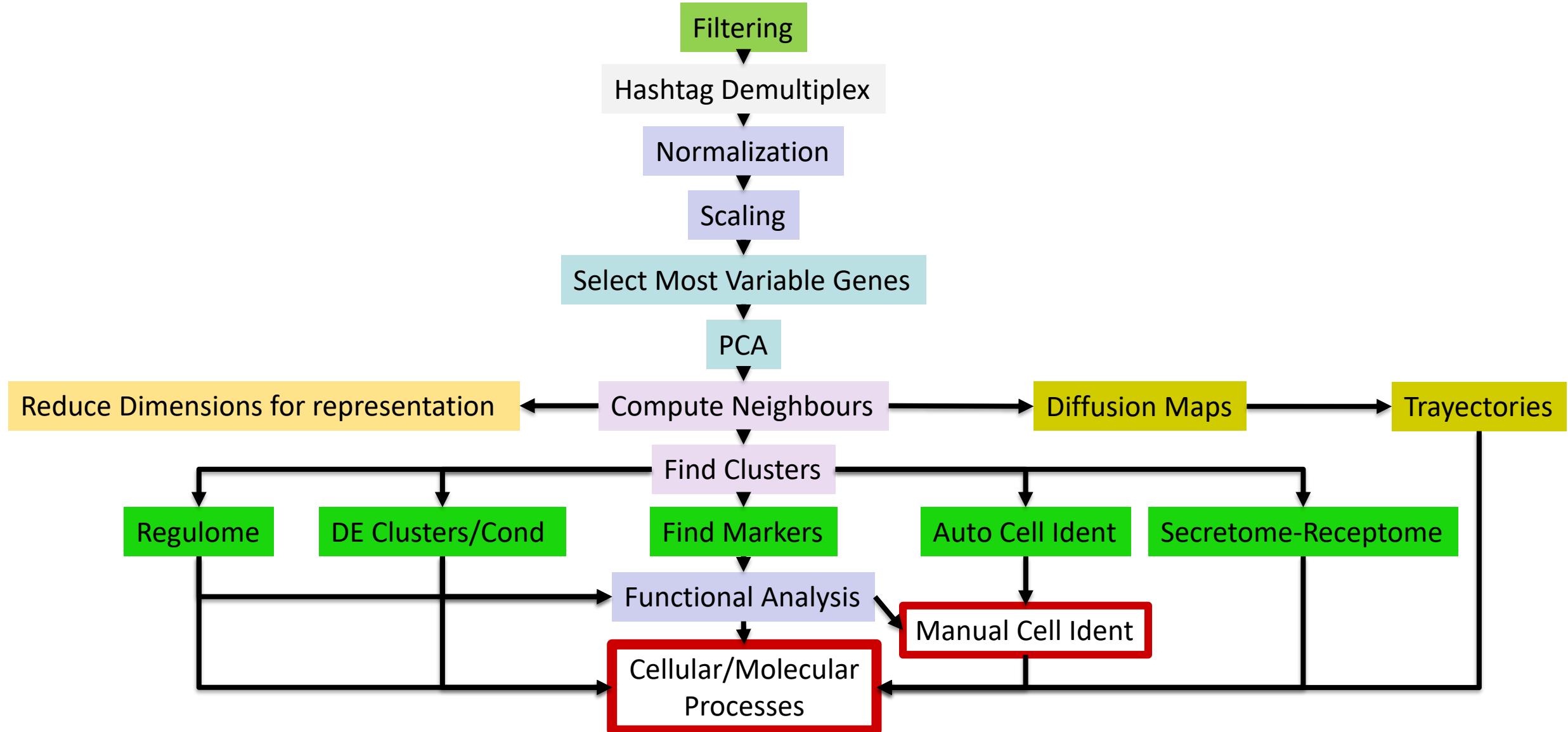
- C1 -> features (rows),
- C2 -> barcodes/cells (columns),
- C3 -> Value (value)

```
%%MatrixMarket matrix coordinate integer general  
%metadata_json: {"software_version": "cellranger-6.1.1", "format_version": 2}  
36636 30674 34324997  
76 1 1  
171 1 4  
203 1 1  
437 1 1  
440 1 1  
493 1 1  
525 1 6  
566 1 1
```

## SPARSE MATRIX

```
ENSG00000224051 . . . . . . 1 . .  
ENSG00000169962 . . . . . . . . .  
ENSG00000107404 . . . . . . . . .  
ENSG00000162576 . . . . . . . . .  
ENSG00000175756 . . . 1 3 4 . . 1 .  
ENSG00000221978 . . . . . . . . .  
ENSG00000224870 1 . . . . . . . . .  
ENSG00000242485 1 . . . 6 . . 3 . .  
ENSG00000272455 . . . . . . . . .  
ENSG00000235098 . . . . . . . . .  
ENSG00000225905 . . . . . . . . .
```

# scRNA-Seq Workflow



## Filtering Cells:

Filtering of low quality cells helps to improve clustering by removing unwanted sources of variation

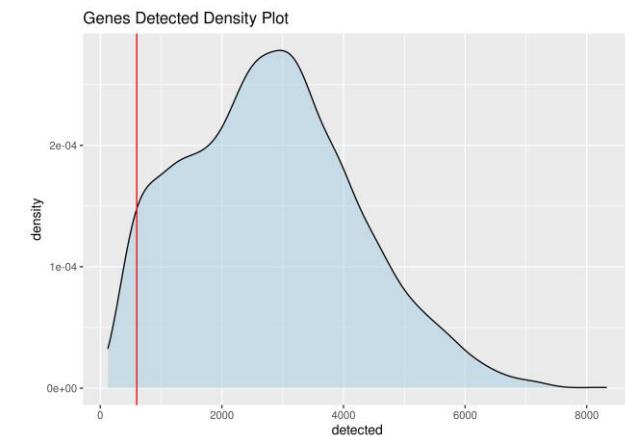
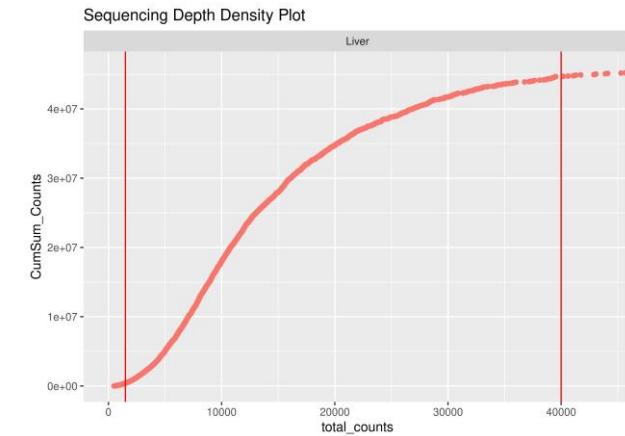
Typical filters used:

### Counts:

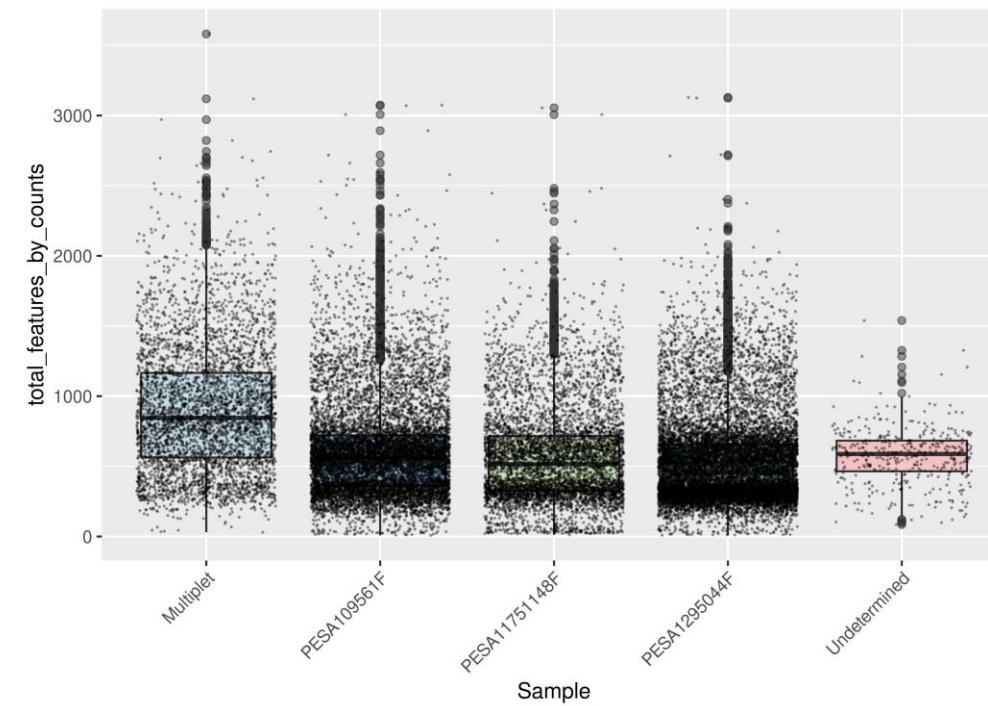
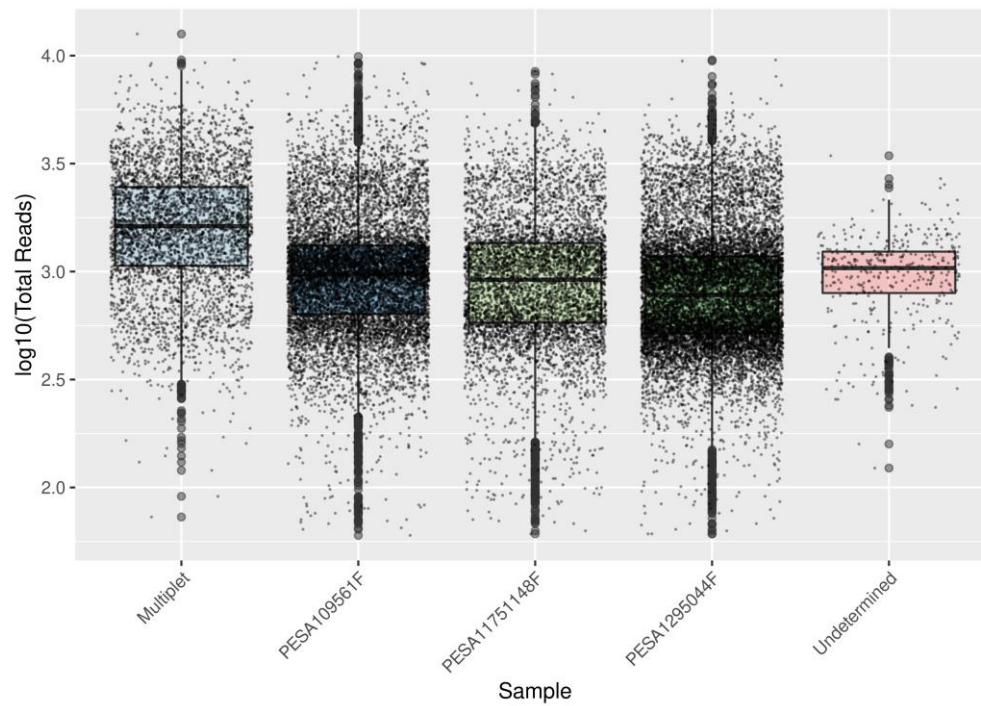
- Low counts levels reduce the ability to profile the gene expression
- High counts may indicate multiplets (several cells captured in the same drop/well)

### Genes:

- Low level of genes detected may indicate problems in the cell viability or particular chemistry problems.
- High levels may indicate also multiplets (several cells captured in the same drop/well).
- Caution: Several cell types have particularly low levels of genes, like the neutrophils
- Caution: Several cell types have particularly high levels of genes like Undifferentiated or Stem Cells.

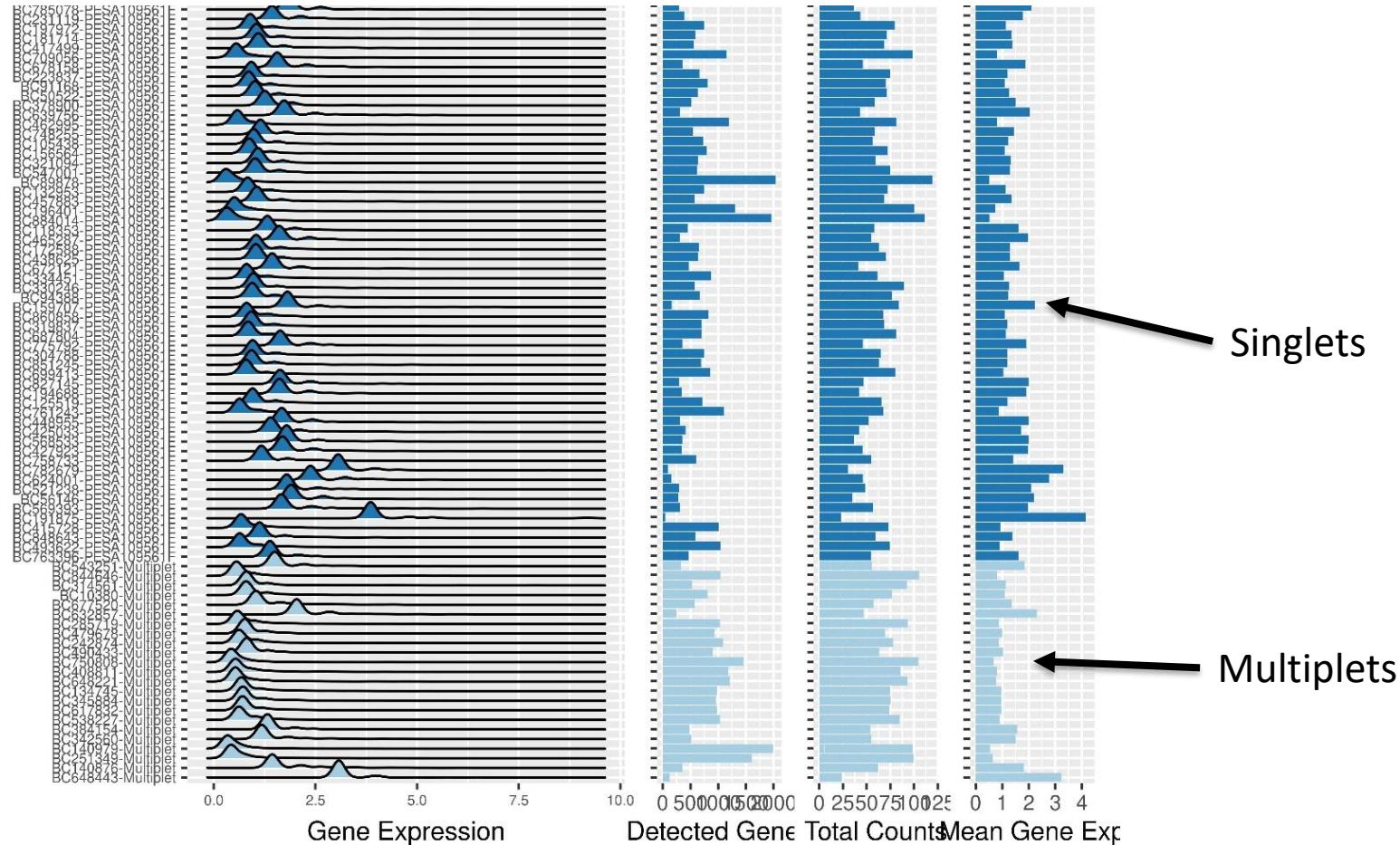


# scRNA-Seq Workflow



# scRNA-Seq Workflow

Counts, Genes Detection and Expression are related

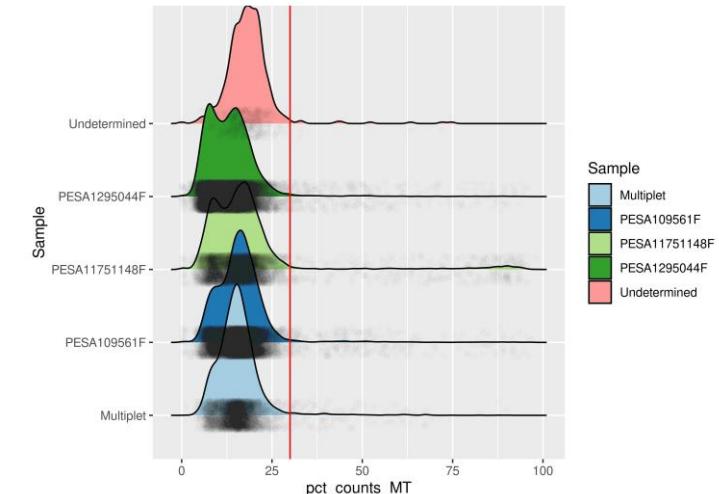


The more counts obtained from a cell, the more genes are usually captured and average gene expression decreases

## Filtering Cells:

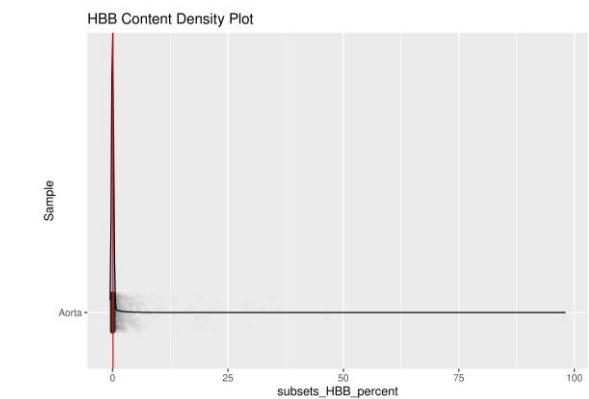
### Mitochondrial Genes:

- In situations when a cell is dying or there is a problem with the chemistry, the mRNA is degraded or more sensitive than the Mitochondrial genes RNA and these becomes dominant in the profile.
- High levels of MT RNA indicates low quality cells (dying cells or low quality capture).
- Caution: Filters have to be adjusted according to the biology in study. Muscle cells have high levels of MT.



### HBB Genes:

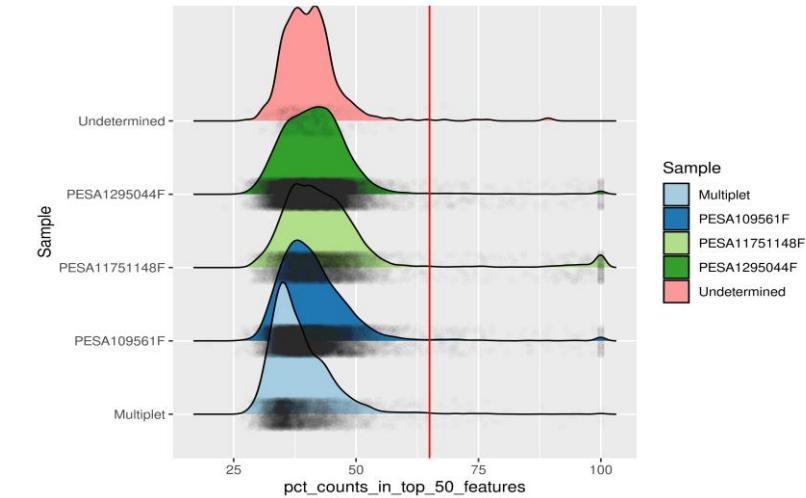
- There are experiments that involves vessels with erythrocytes that could contaminate the sample.
- Erythrocytes have high levels of HBB RNA and a few other RNAs and are very sensitive and easy to lyse.
- They may release their content and increase the background RNA to a point that an empty drop/well might be considered a cell.
- High content HBB cells are usually removed unless biologically relevant



## Filtering Cells:

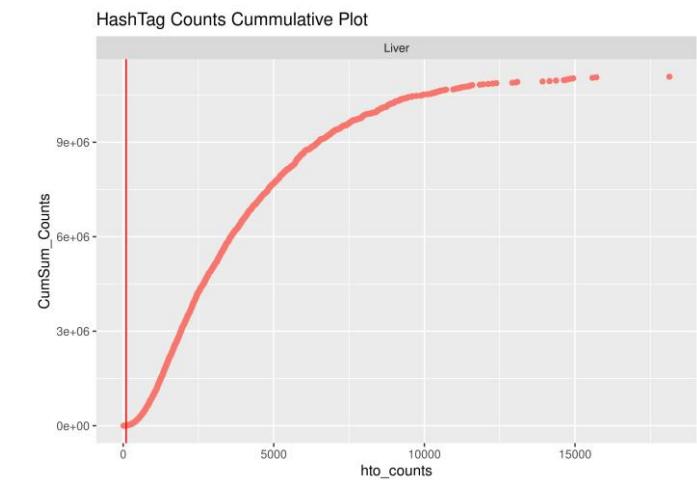
### Gene Complexity:

- The idea is to detect those cells that having many reads, and sufficient genes detected, still they have very few genes with most of the reads.
- Those cells are also very low quality and usually reflects a RNA capture problem or a dying cell.

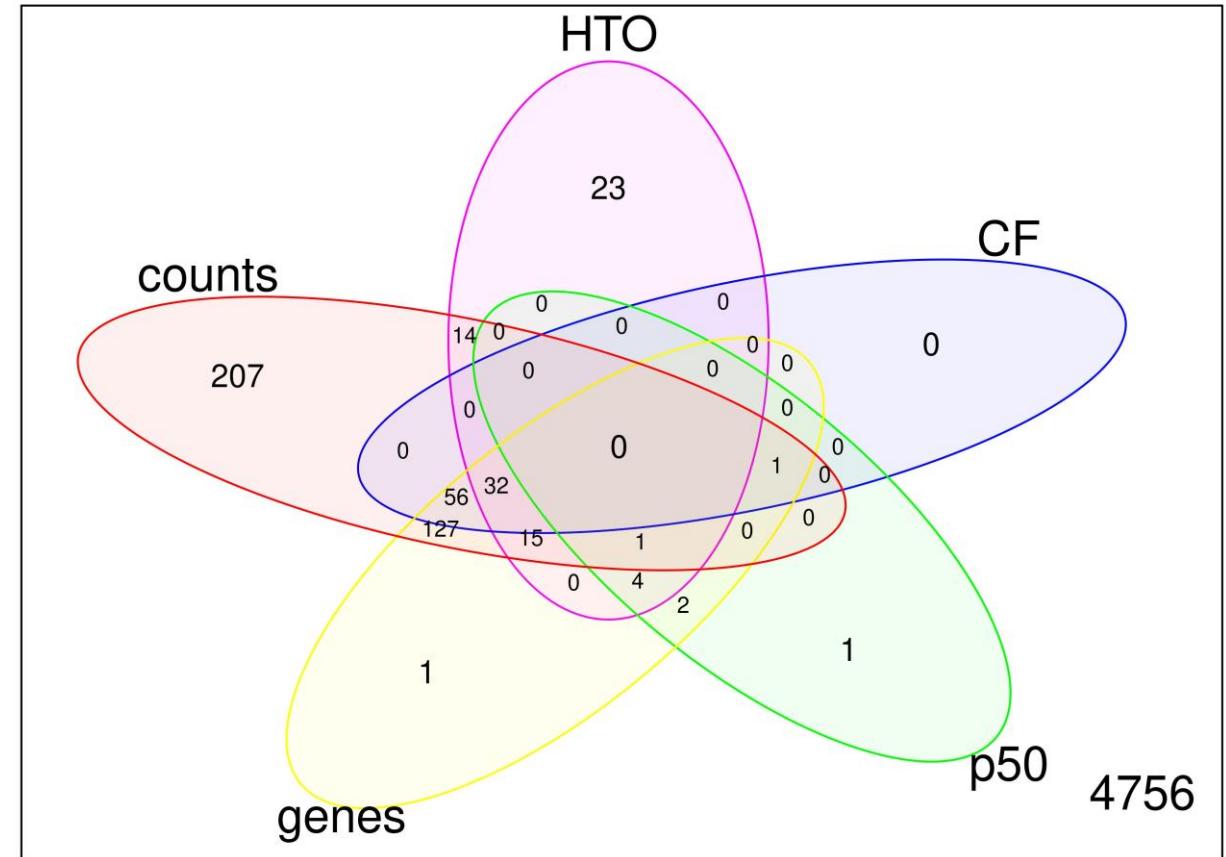
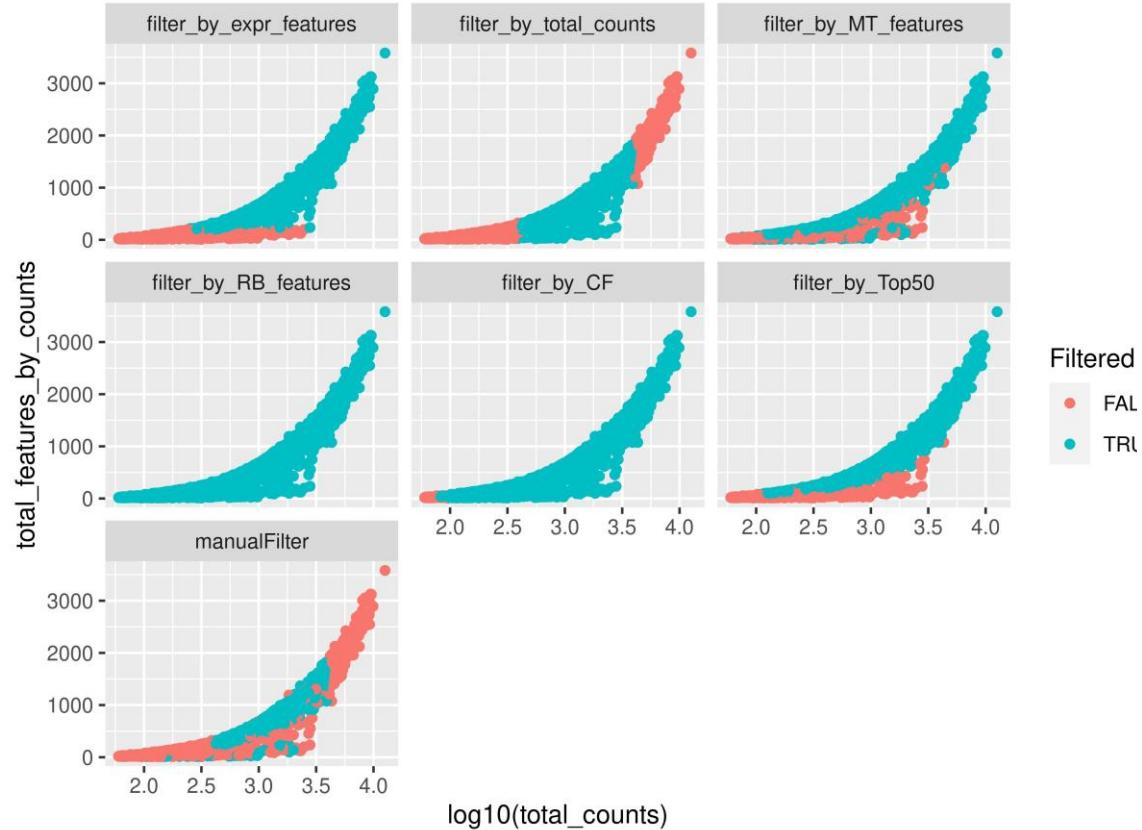


### Hashtag Content:

- When there are several samples pooled in the experiment and they are labelled with antibodies/hashtags, a hashtag filter is also recommended.



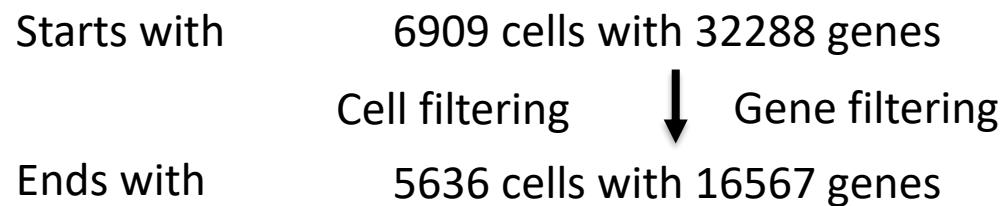
## Summary Of Cell Filtering Process



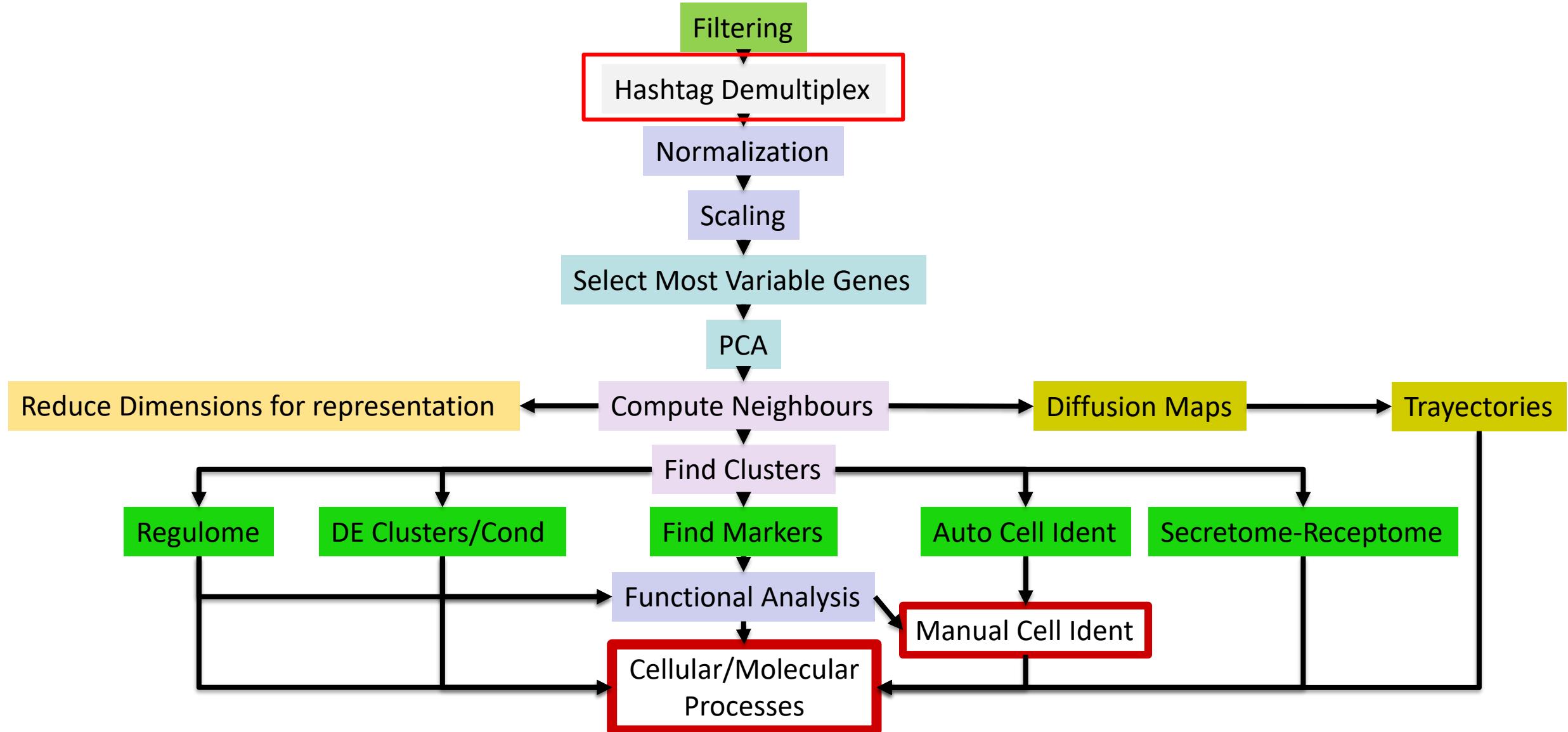
## Filtering Genes

- The idea is to reduce the number of genes to those more relevant to increase efficiency.
- Obviously we will remove all the genes with 0 counts across the experiment, but sometimes it is also good to remove some that are only detected in a few cells.
- This will depend on the total number of cells in the experiment and minimum amount of cells expected to be part of a cluster.
- In many cases we don't know whether we are going to have small relevant clusters.
- So, to start with, just remove those present in no more than 10 cells. Later the filter can be increased or reduced

A normal situation:



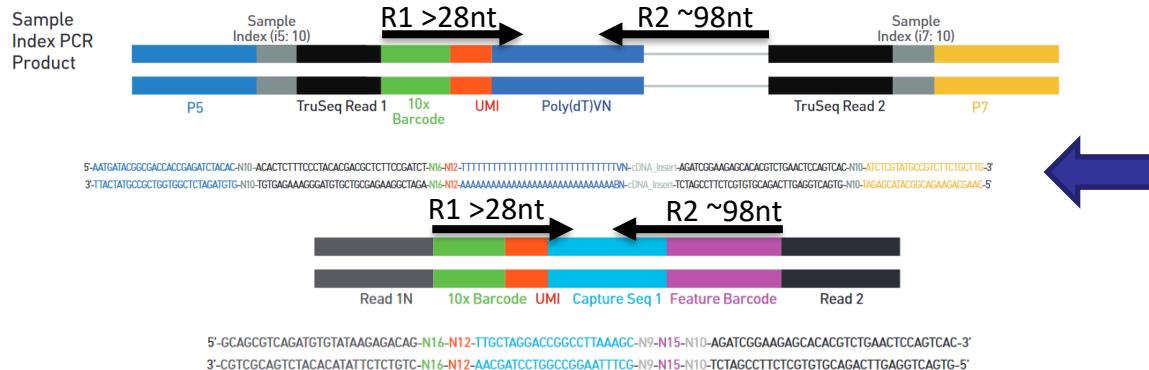
# scRNA-Seq Workflow



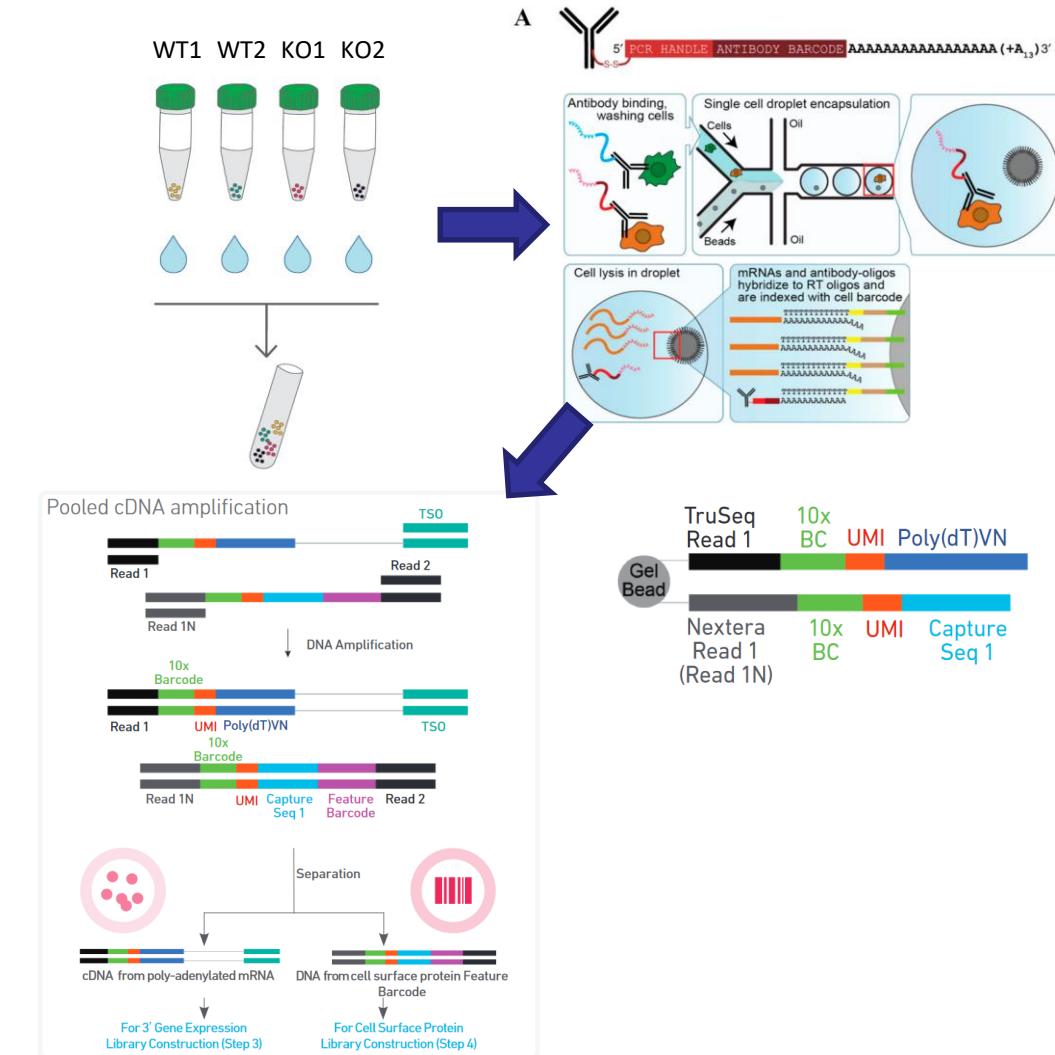
# scRNA-Seq Workflow

## Sample Hashtag Multiplexing:

- Several samples can be processed in one experiment by labelling the cells of each sample with an antibody against an ubiquitous membrane antigen.
- This reduces the capture and library preparation costs.
- Sample Hashtag (Sample Ab-barcode) libraries are prepared alongside but independently of mRNA libraries once cDNA is produced and amplified.
- Hashtag libraries are sequenced together with the mRNA derived libraries and sequences from both libraries are demultiplexed based on Illumina Sample barcoding system.



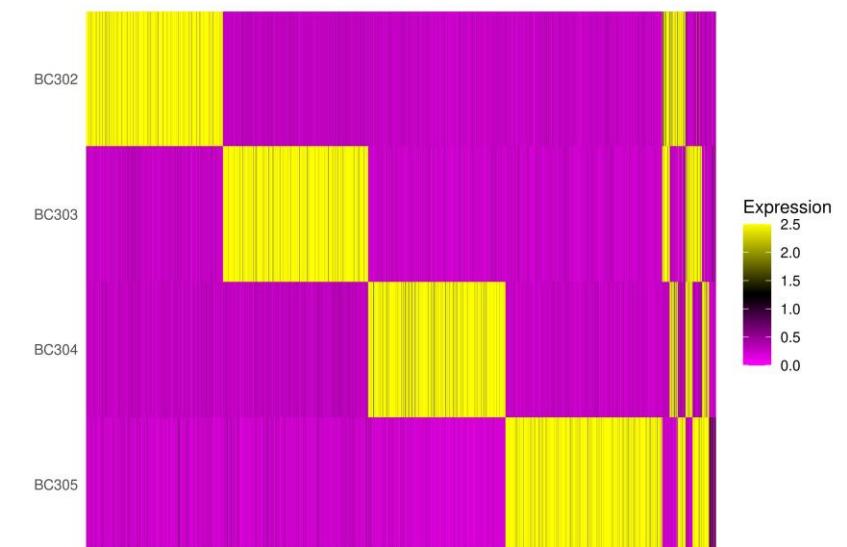
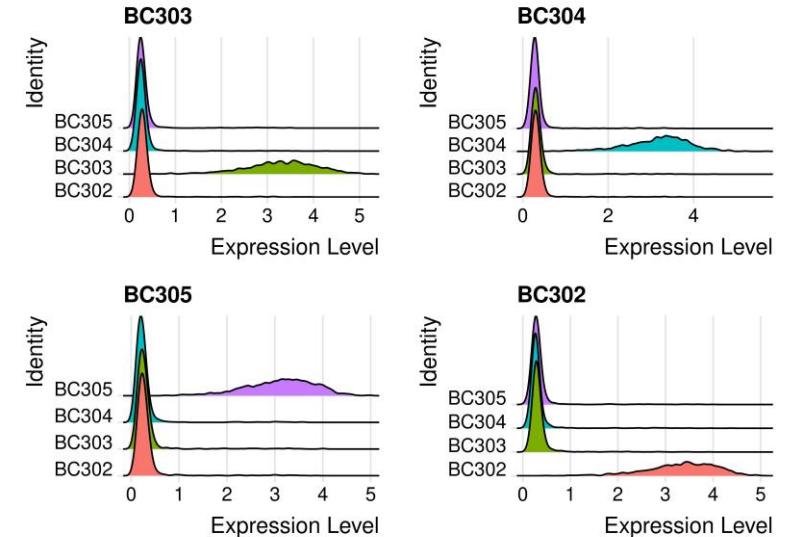
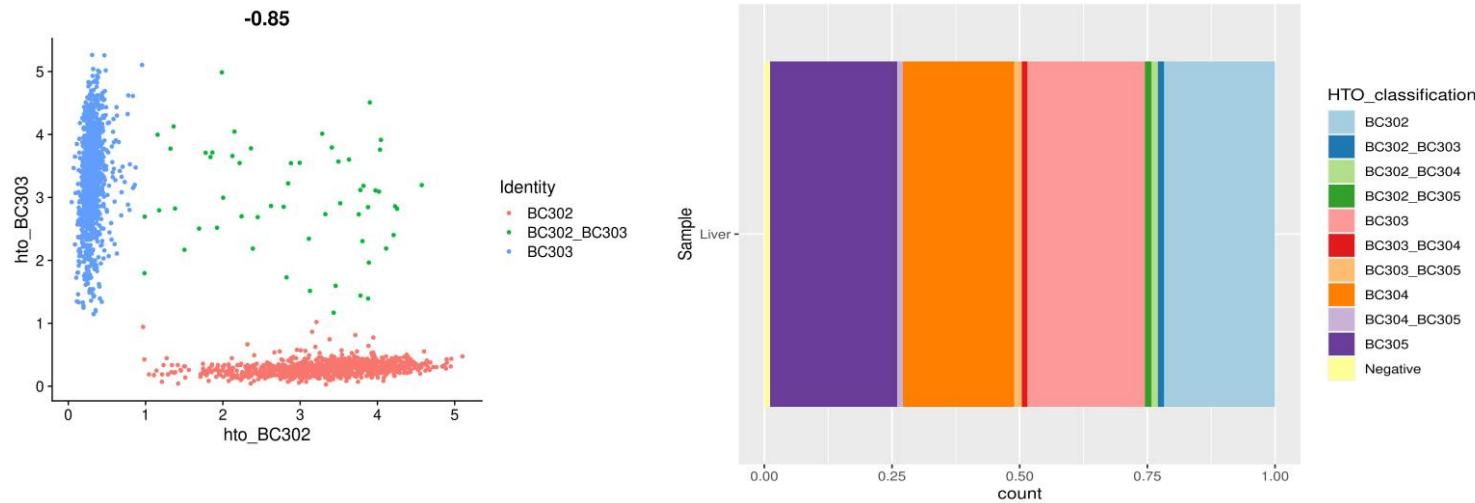
## Single Cell Samples:



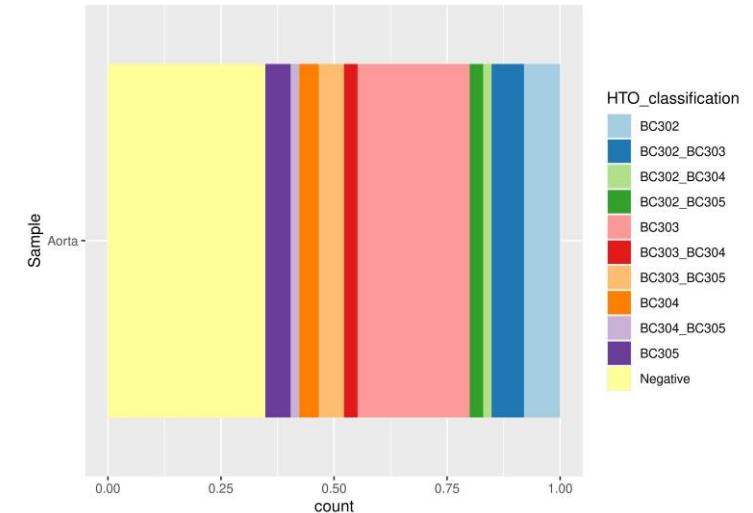
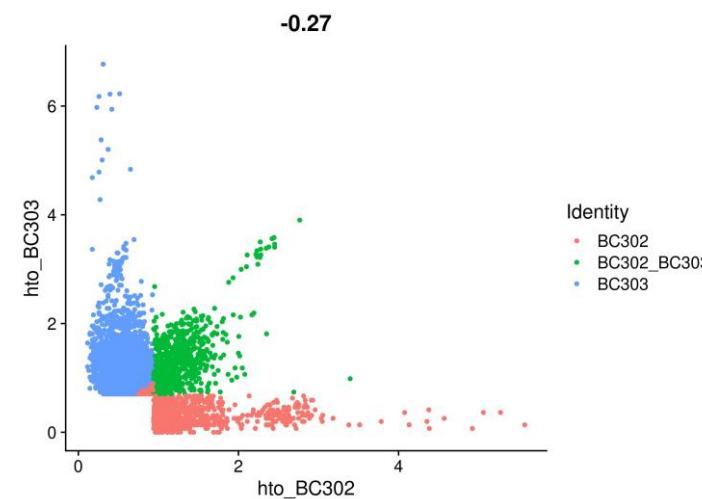
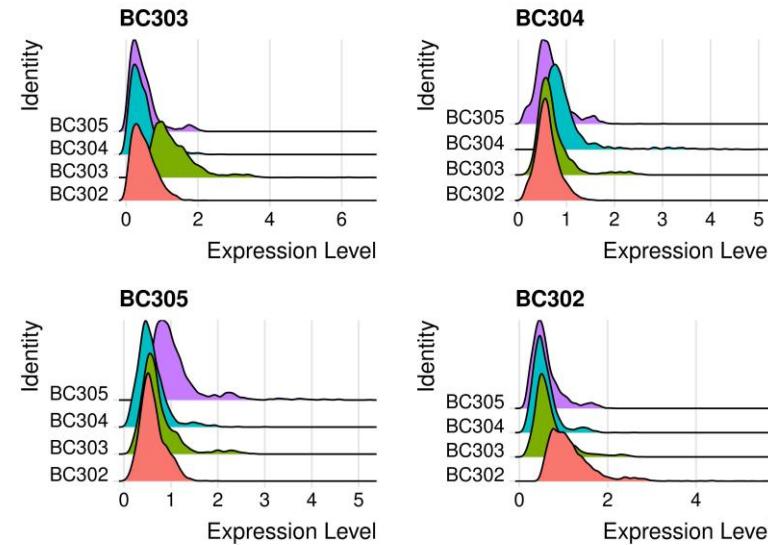
# scRNA-Seq Workflow

## Cell Sample Assignment

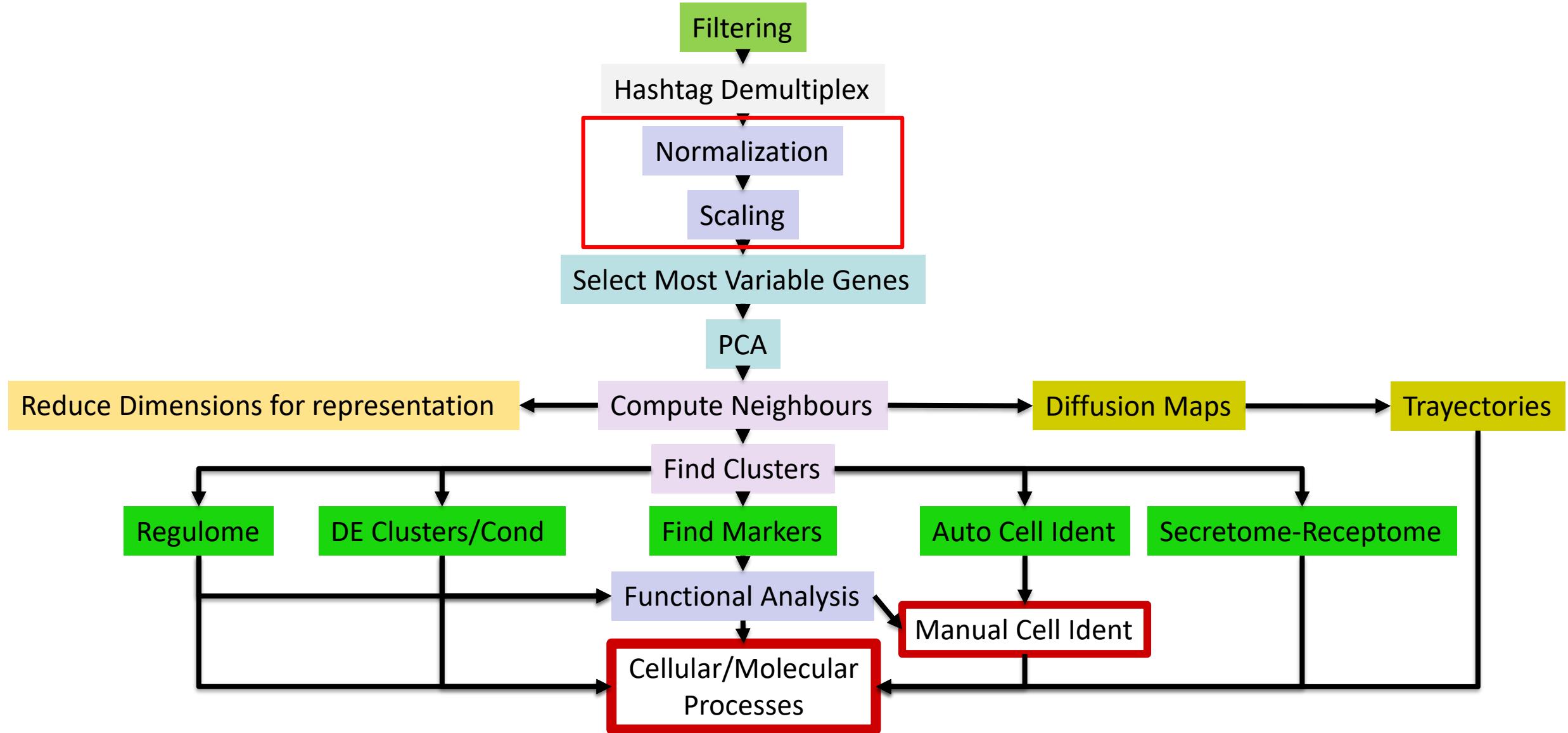
- R2 Sequence is annotated with Cell Barcode and UMIs from R1
- R2 are mapped to the Hashtag sequences used to label each sample
- Counts are generated in the same way as the genes: unique UMI/Barcode/Hashtag combination
- Only detected cell barcodes based on gene expression are kept.
- A background model is generated for each hashtag based on the expression of all the hashtags minus the highest expressed on each cell.
- Hashtags significantly far from this background model are considered as detected.
- In this way cells are labelled according to the hashtags detected



# scRNA-Seq Workflow



# scRNA-Seq Workflow



→ ERCC Control Sequences: Set of 96 mammal like sequences not present in model genomes used to control efficiency of the protocols

→ Normalization algorithms are able to use the information from these control sequences to estimate technical variability between replicates and batches and correct by that.

Characterizing noise structure in single-cell RNA-seq distinguishes genuine from technical stochastic allelic expression

Jong Kyoung Kim<sup>1</sup>, Aleksandra A. Kolodziejczyk<sup>1,2</sup>, Tomislav Ilicic<sup>1,2</sup>, Sarah A.

Normalization of RNA-seq data using factor analysis of control genes or samples

Davide Risso<sup>1</sup>, John Ngai<sup>2–4</sup>, Terence P Speed<sup>1,5,6</sup> & Sandrine Dudoit<sup>1,7</sup>

→ New normalization algorithms are designed to reduce the effect of the vast amounts 0 values like those based on pooling and deconvolution techniques.

Pooling across cells to normalize single-cell RNA sequencing data with many zero counts

Aaron T. L. Lun<sup>1\*</sup>, Karsten Bach<sup>2</sup> and John C. Marioni<sup>1,2,3\*</sup>

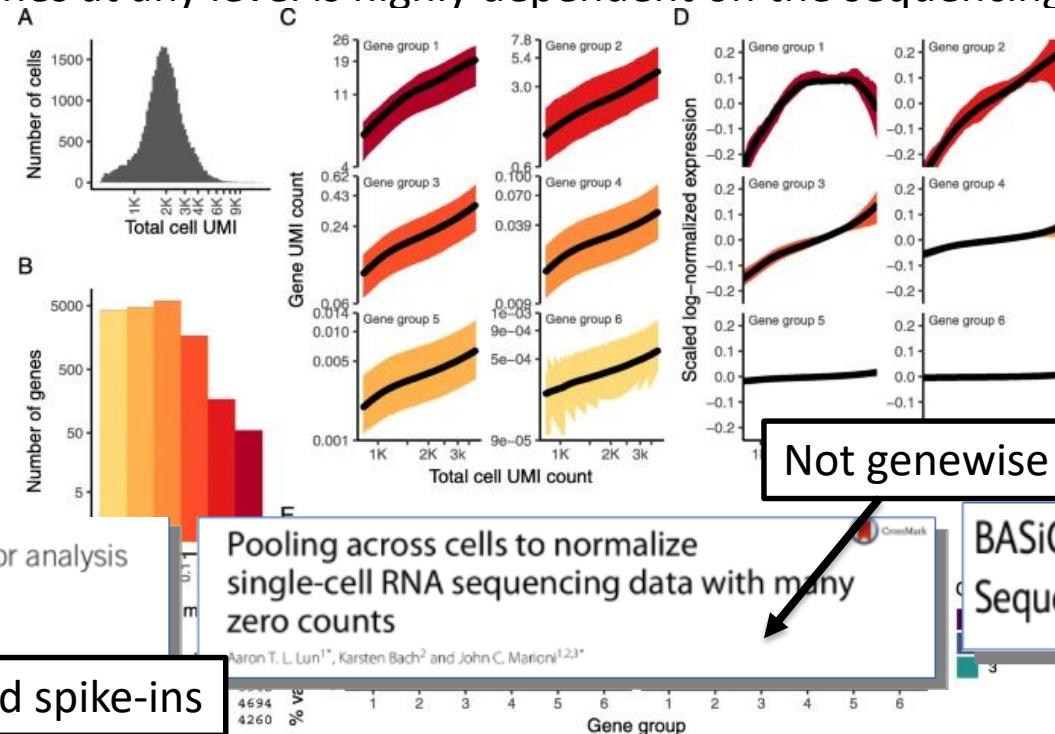
BASiCS: Bayesian Analysis of Single-Cell Sequencing Data

## Normalization & Scaling

Normalization on bulk samples is based on the fact that most of the genes do not change across replicates and conditions.

In SC experiments is close to the opposite. Most of the genes detected are different between cells.

Also, the expression of genes at any level is highly dependent on the sequencing depth of the cell.



Not genewise

Overfit

Normalization of RNA-seq data using factor analysis of control genes or samples

Davide Risso<sup>1</sup>, John Ngai<sup>2-4</sup>, Terence P Speed<sup>1,5,6</sup> & Sandrine Dudoit<sup>1,2</sup>

Need spike-ins

Pooling across cells to normalize single-cell RNA sequencing data with many zero counts

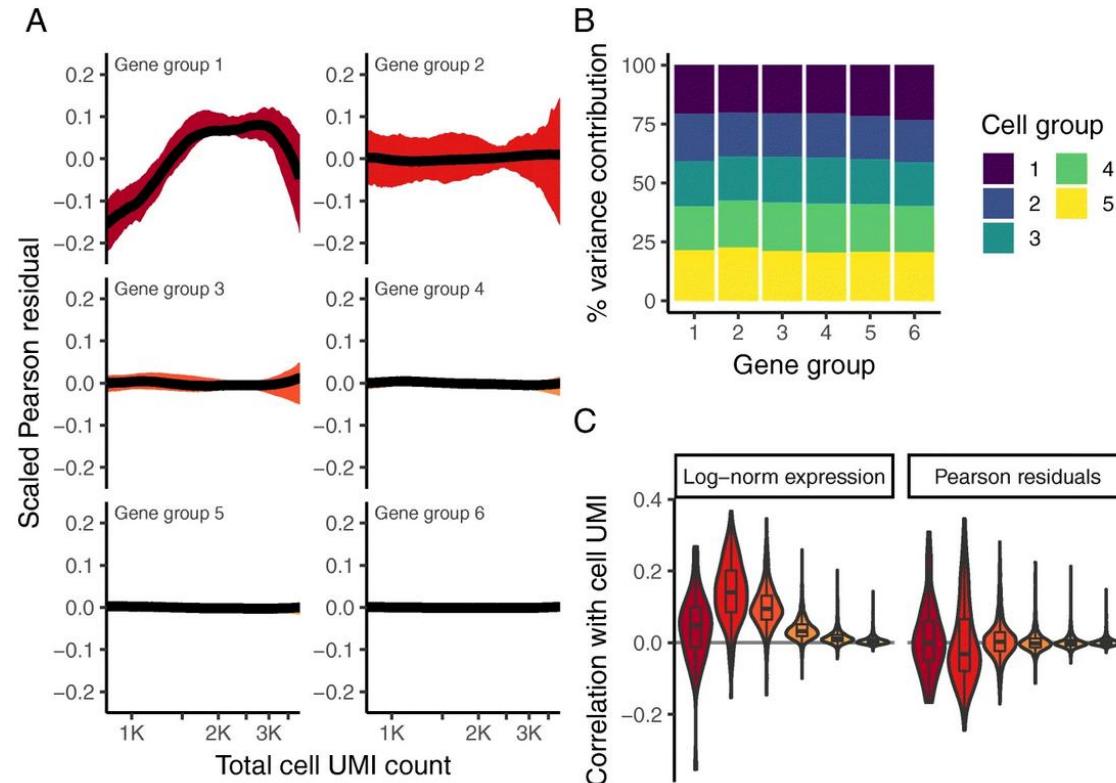
Naron T. L. Lun<sup>1\*</sup>, Karsten Bach<sup>2</sup> and John C. Marioni<sup>1,2,3\*</sup>

BASiCS: Bayesian Analysis of Single-Cell Sequencing Data

3

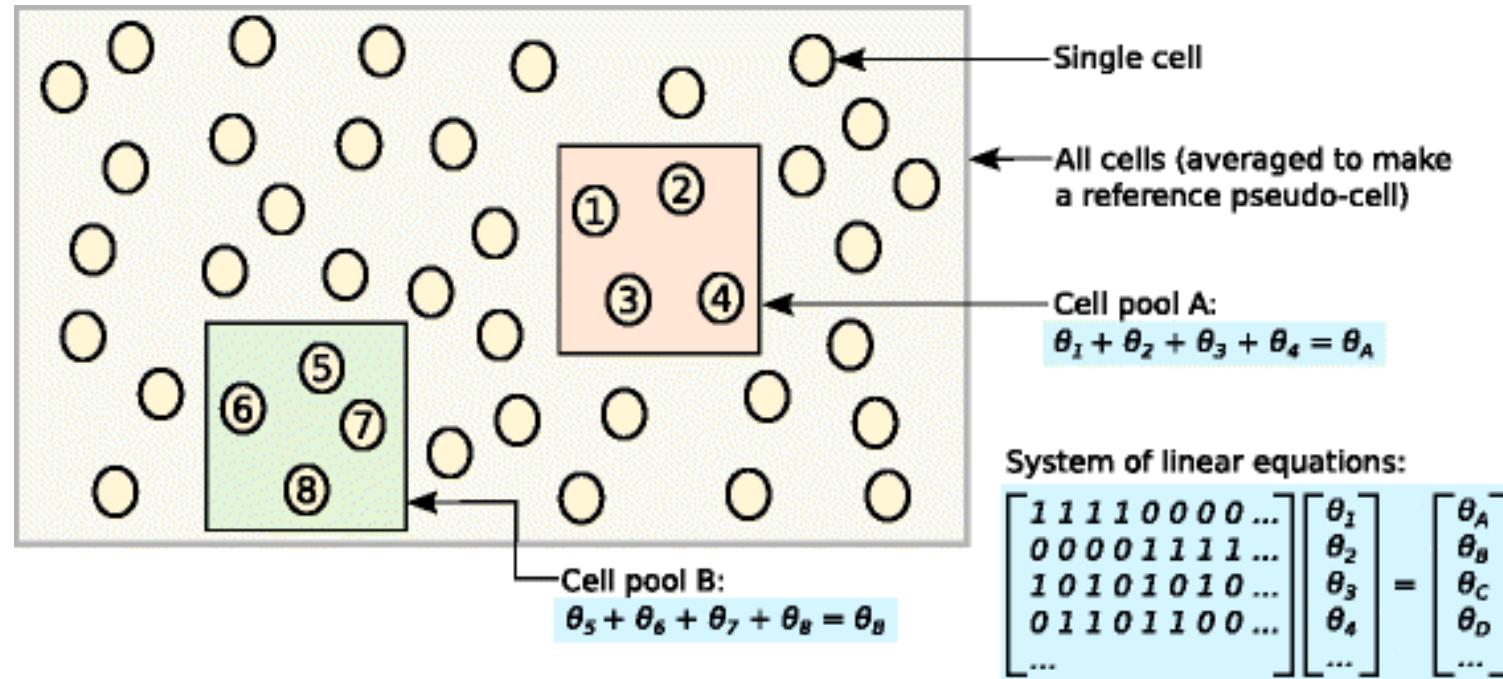
## Normalization & Scaling

### Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression



Hafemeister, C., Satija, R. Normalization and variance stabilization of single-cell RNA-seq data using regularized negative binomial regression. *Genome Biol* **20**, 296 (2019). <https://doi.org/10.1186/s13059-019-1874-1>

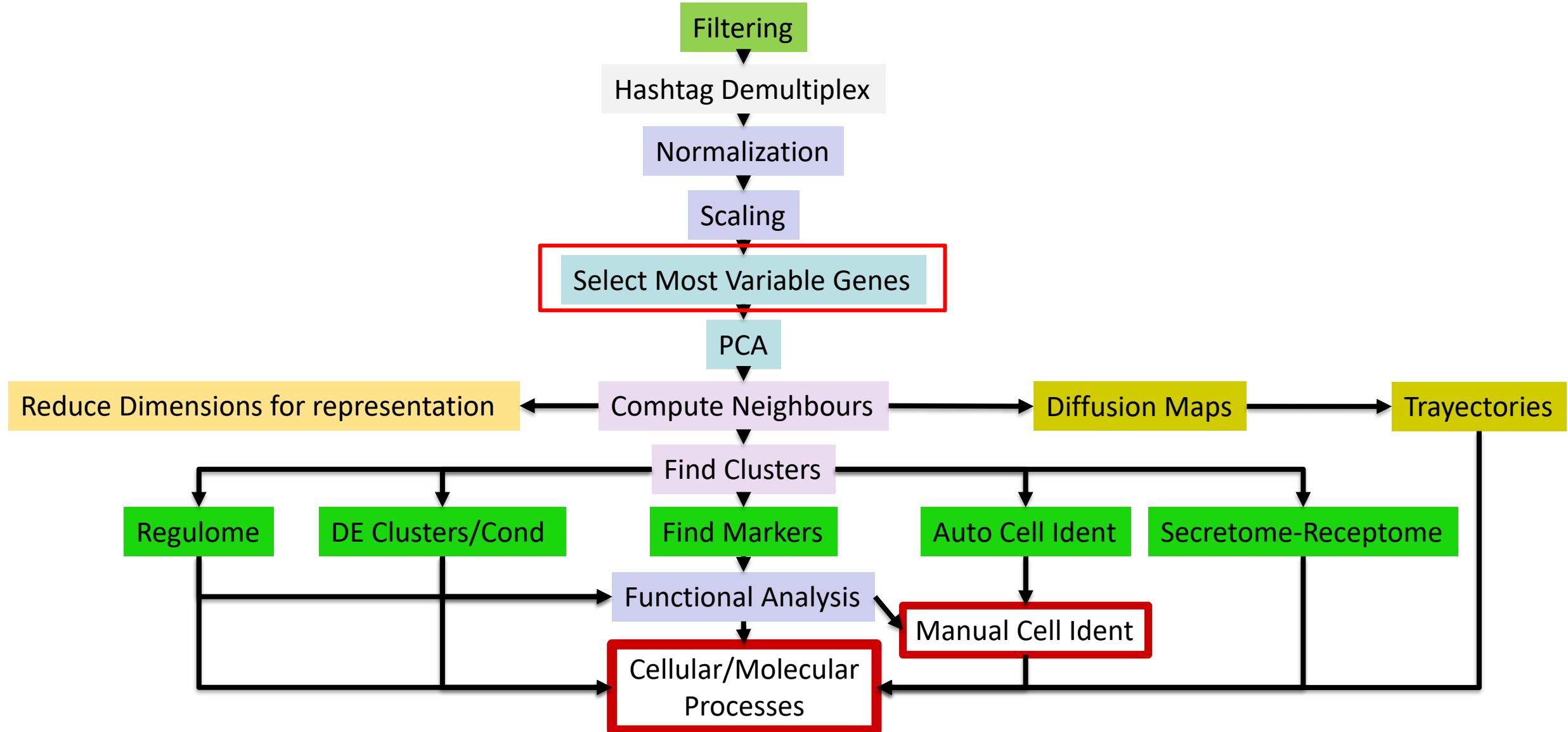
# Pooled Normalization



- Each cell is pooled many times with other different cells from different experimental groups giving a set of linear equations from which deconvolute a normalization factor for each cell.
- Incorporates intrinsically a batch correction as each cell will be combined with cells from other batches in each pool.

L. Lun, A.T., Bach, K. & Marioni, J.C. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol* 17, 75 (2016). <https://doi.org/10.1186/s13059-016-0947-7>

# scRNA-Seq Workflow

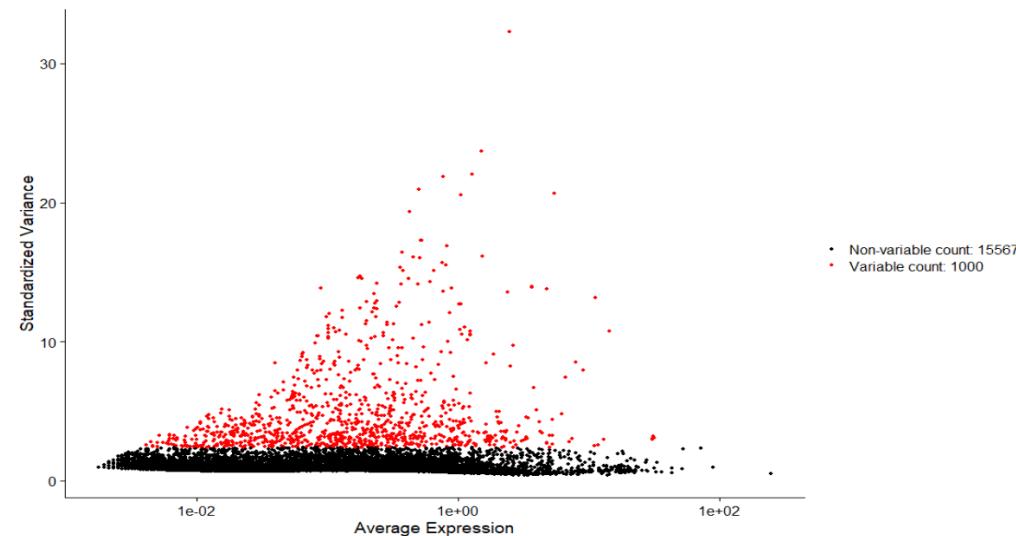


## Most Variable Genes:

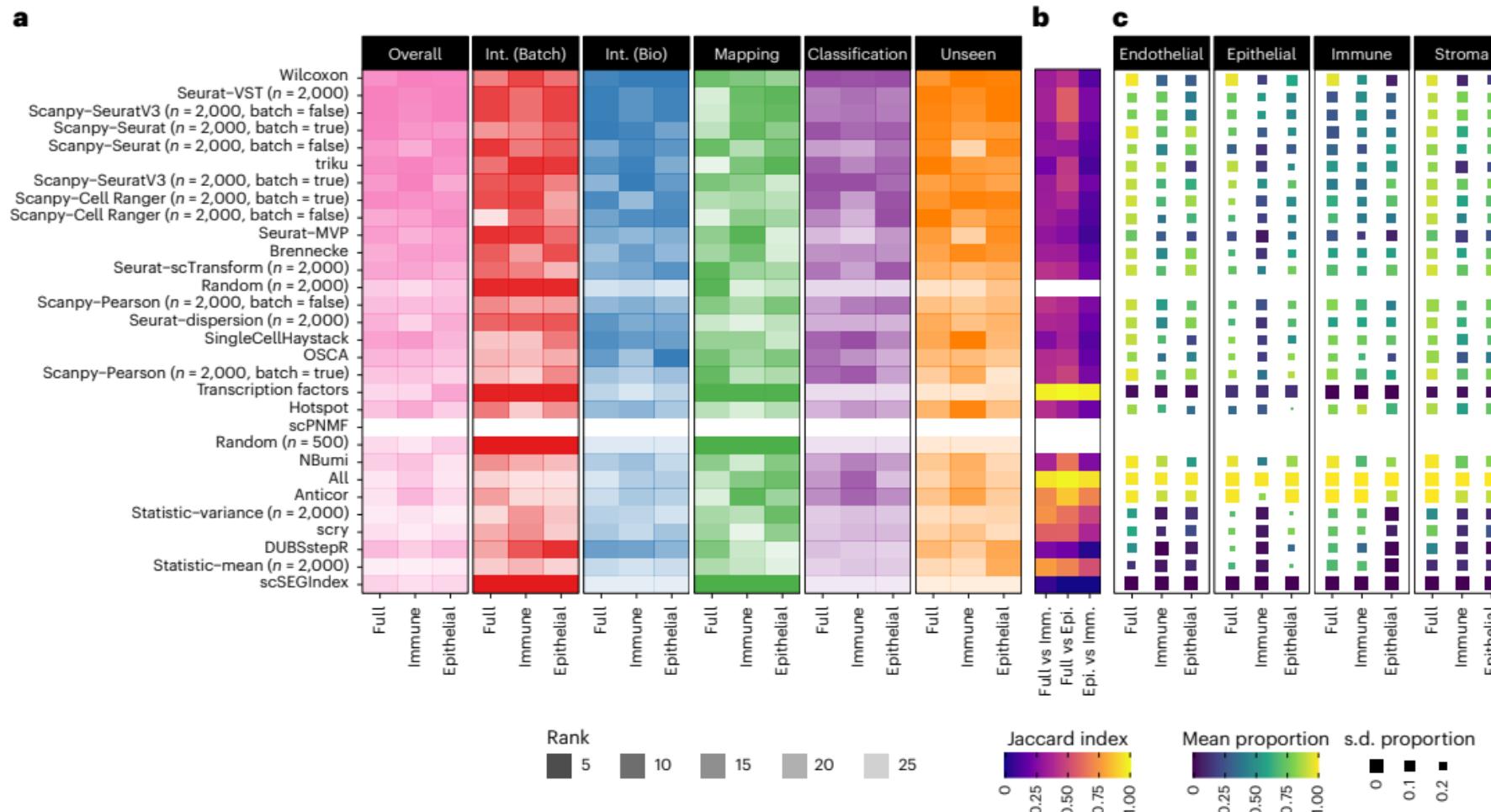
To reduce the amount of noise and increase computational efficiency there is a feature selection step to reduce the number of dimensions/features to use in the clustering analysis.

Several approaches are possible:

- Using some cut-offs in the mean expression and standard deviation
- Selecting the n top genes with the highest stdev (i.e. 2000 genes).

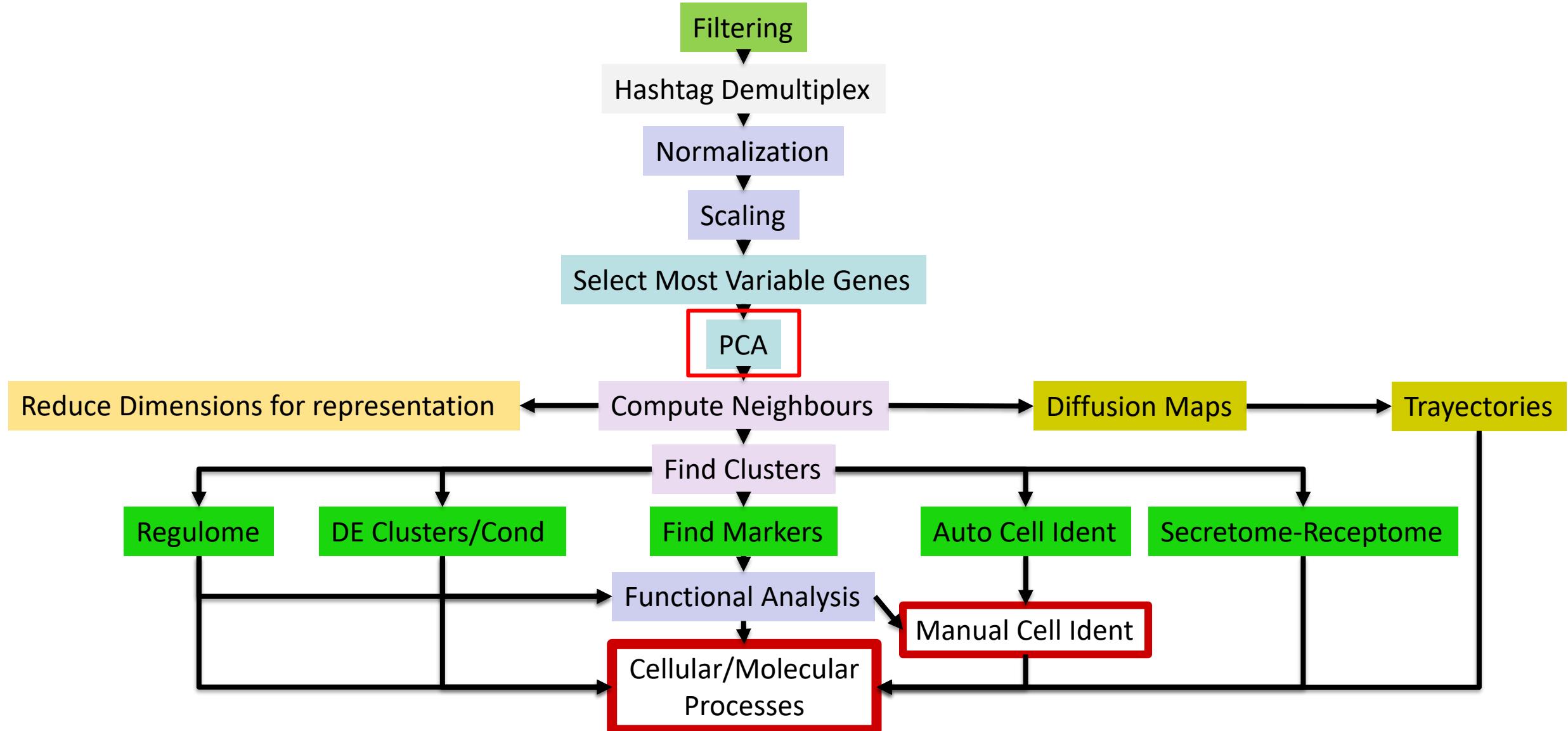


# Select Most Variable Genes



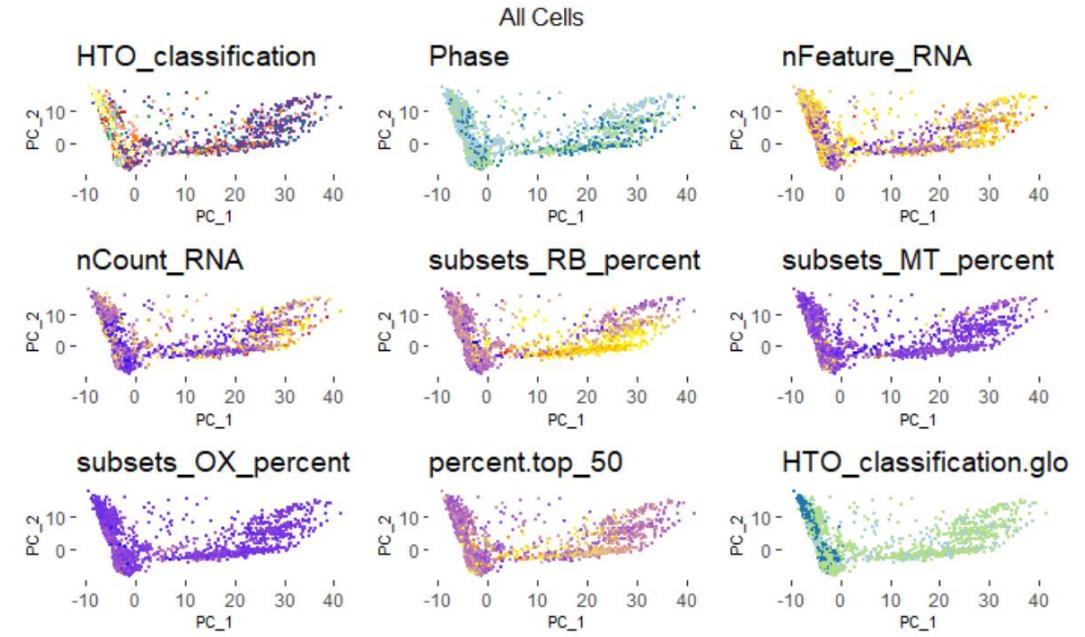
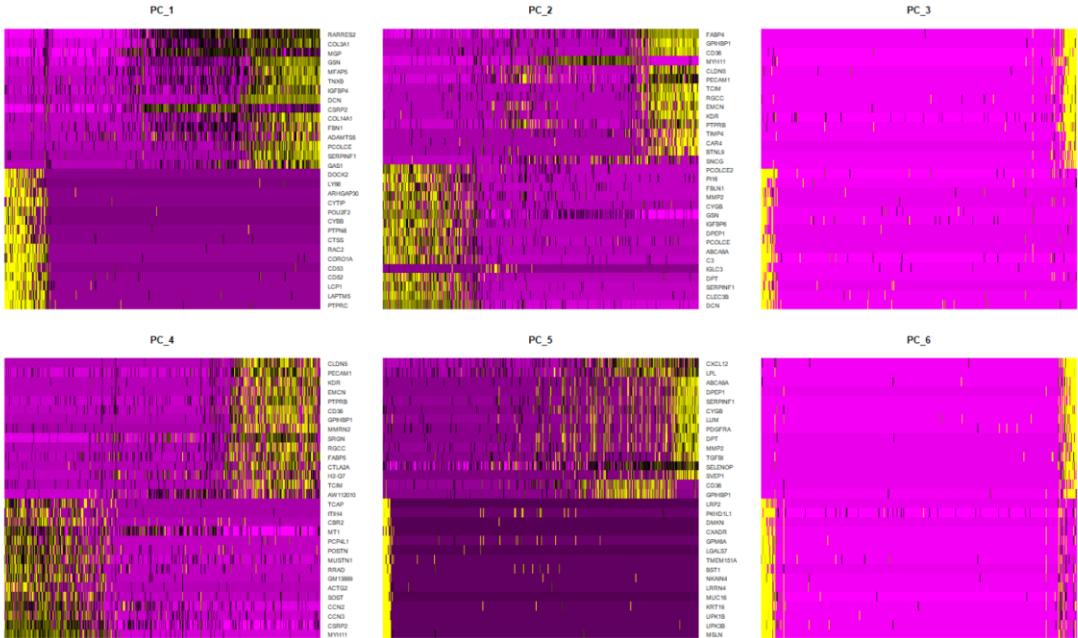
Zappia L, Richter S, et. al., and Theis FJ. Feature selection methods affect the performance of scRNA-seq data integration and querying. Nat Methods. 2025 Apr;22(4):834-844. doi: 10.1038/s41592-025-02624-3. Epub 2025 Mar 13. PMID: 40082610; PMCID: PMC11978513.

# scRNA-Seq Workflow



## PCA

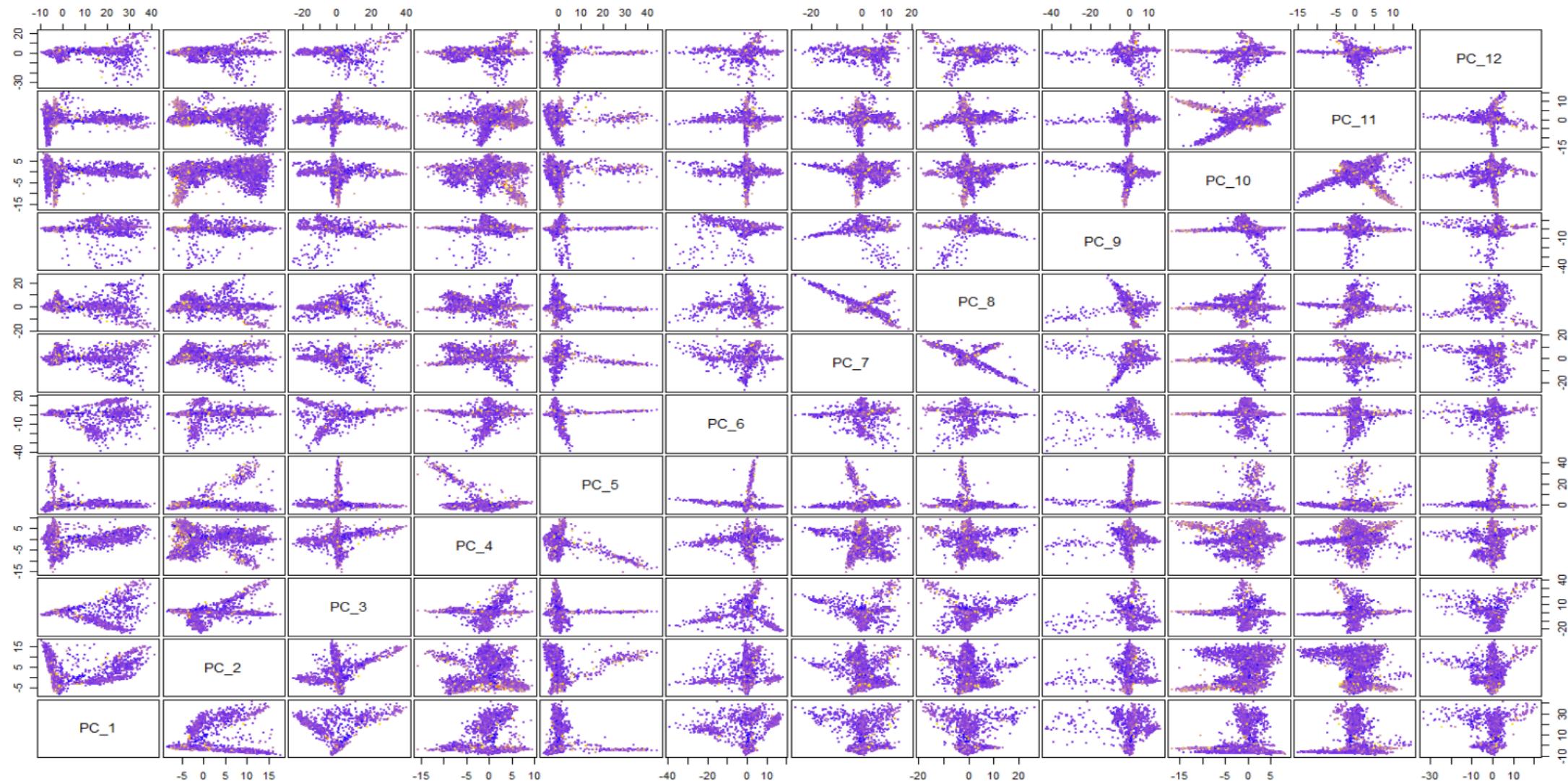
- A Principal Component Analysis is at the basis of all following steps.
- The goal is to reduce even further the dimensions before all posterior clustering and further dimensionality reductions.
- Now the dimensions over which clustering will work are the different components. They summarise the dimensions in which the experiment express its higher variability.



We can start exploring which genes are relevant in the experiment or maybe some technical variability

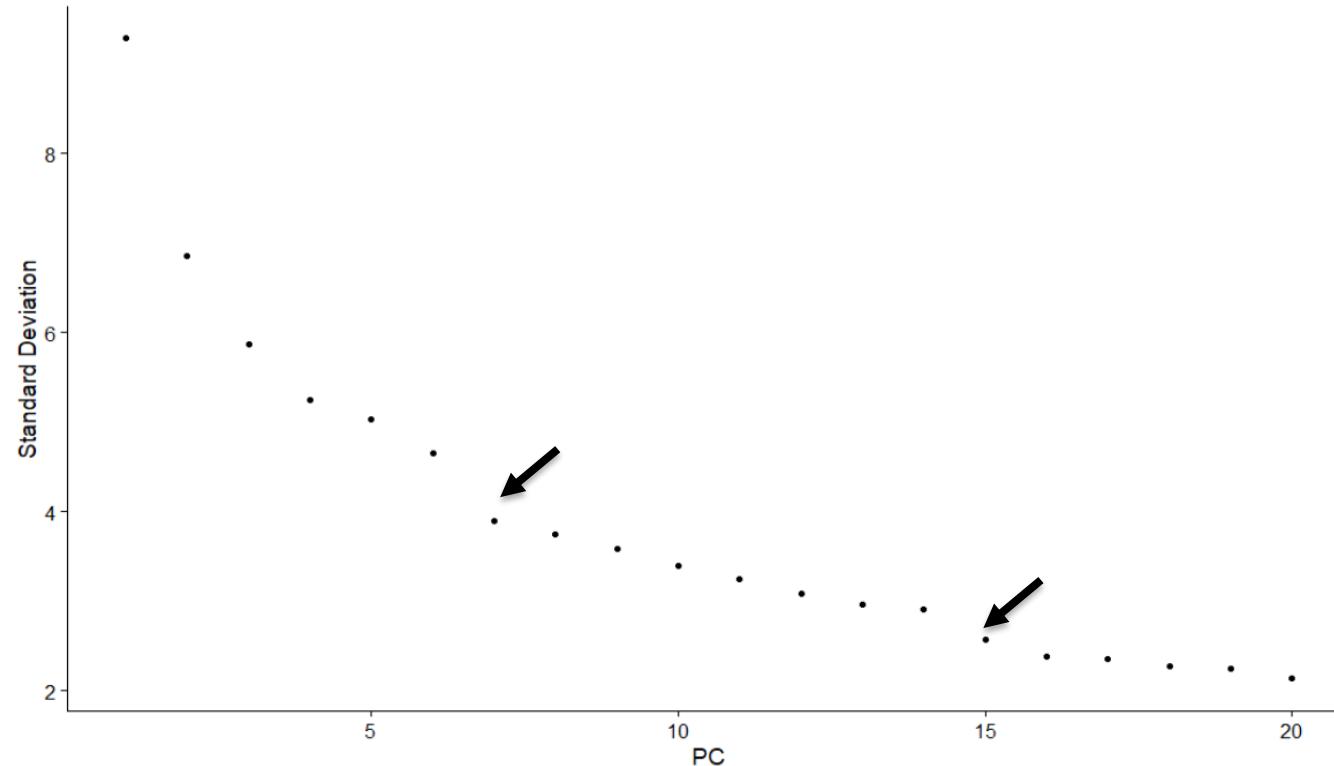
# scRNA-Seq Workflow

Cell MT Content across PCAs

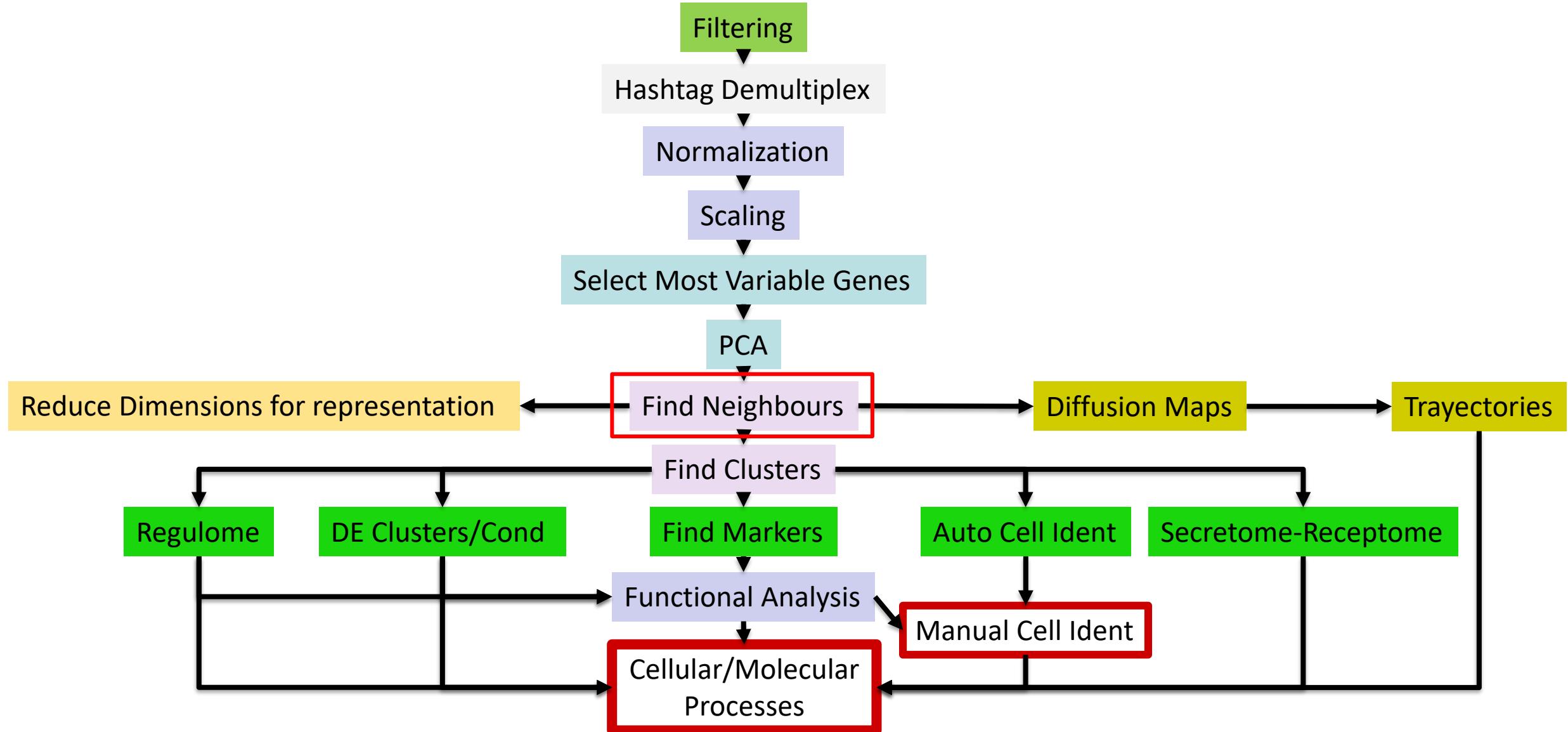


## PCA

Again, to reduce spurious noise, not all PCs from the PCA are used and a PC selection is carried out again. To select the most relevant PCs we usually evaluate the variability explained by each PC and select those PCs that accounts for most of the variability looking at places where high decreases are observed.



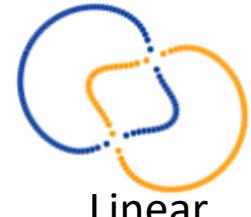
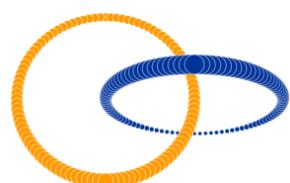
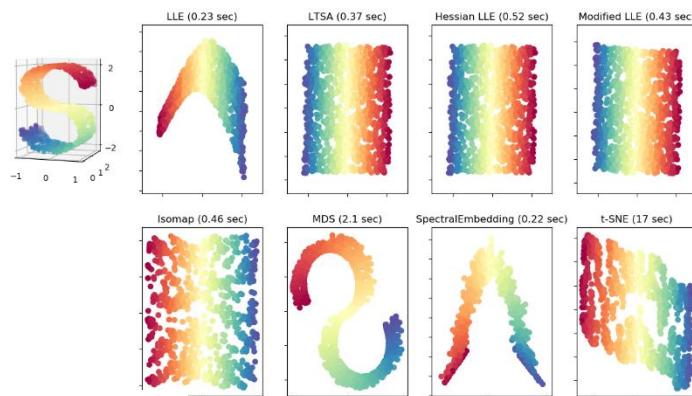
# scRNA-Seq Workflow



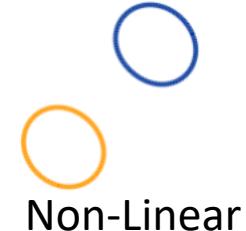
## □ Visualization

### Dimensionality reduction techniques

- Linear → PCA, LDA,...
- Non-linear (manifolds) → tSNE, UMAP, Isomap, LLE, MLLE, Spectral embedding, MDS...



Linear

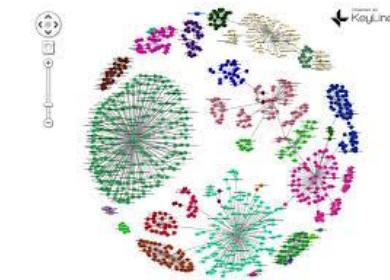
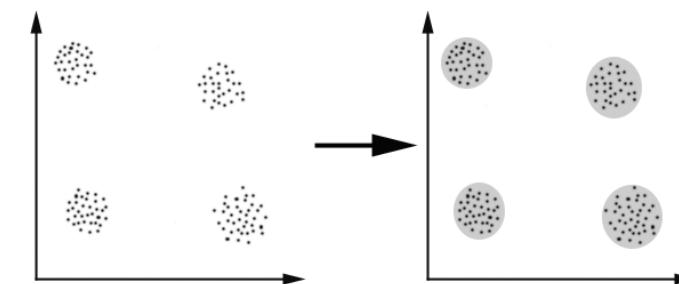


Non-Linear

## □ Detection of Populations

### Clustering Algorithms

- Supervised/Informed → Kmeans, Hierarchical, GMM,...
- Unsupervised/Uninformed → Mean-shift, X-shift, Phenograph, ClusterX



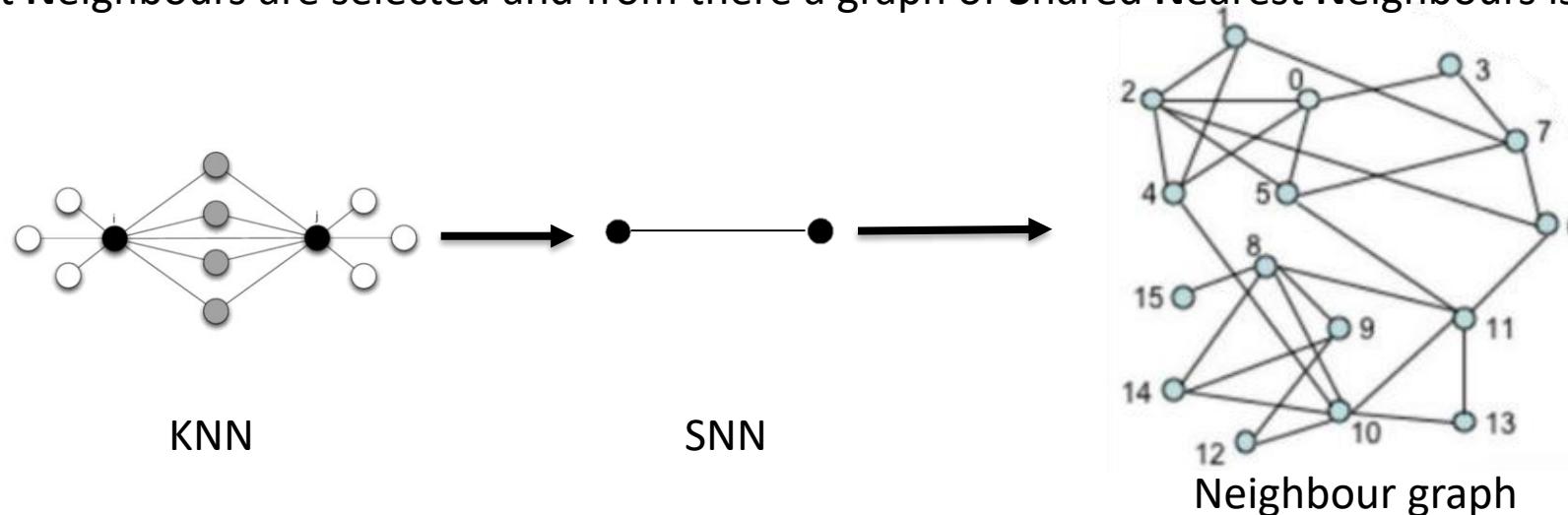
## Compute Neighbours

The goal is to define in the multidimensional space which cells are next to each other. Find the neighbours of each cell.

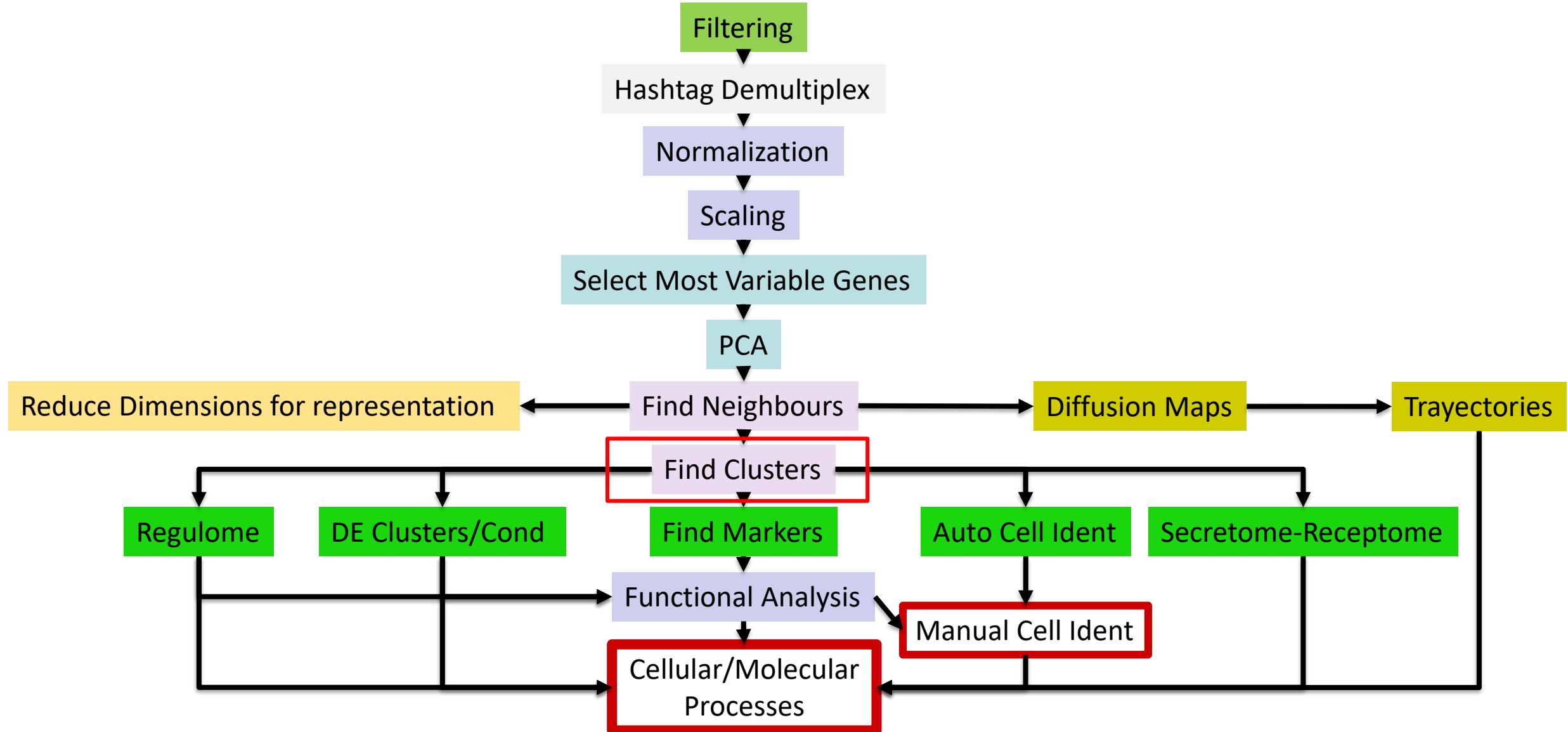
All posterior analysis, clustering, dimensionality reduction, trajectories,... depends on this set of neighbours found for each cell.

In the first steps the distance between cells is calculated. To improve efficiency some a kernel density is applied to reduce the amount of distances to be calculated.

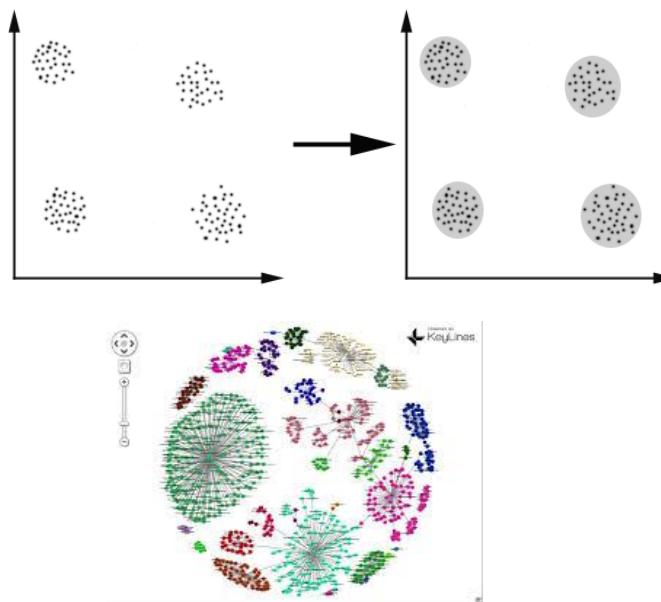
The **K Nearest Neighbours** are selected and from there a graph of **Shared Nearest Neighbours** is constructed



# scRNA-Seq Workflow

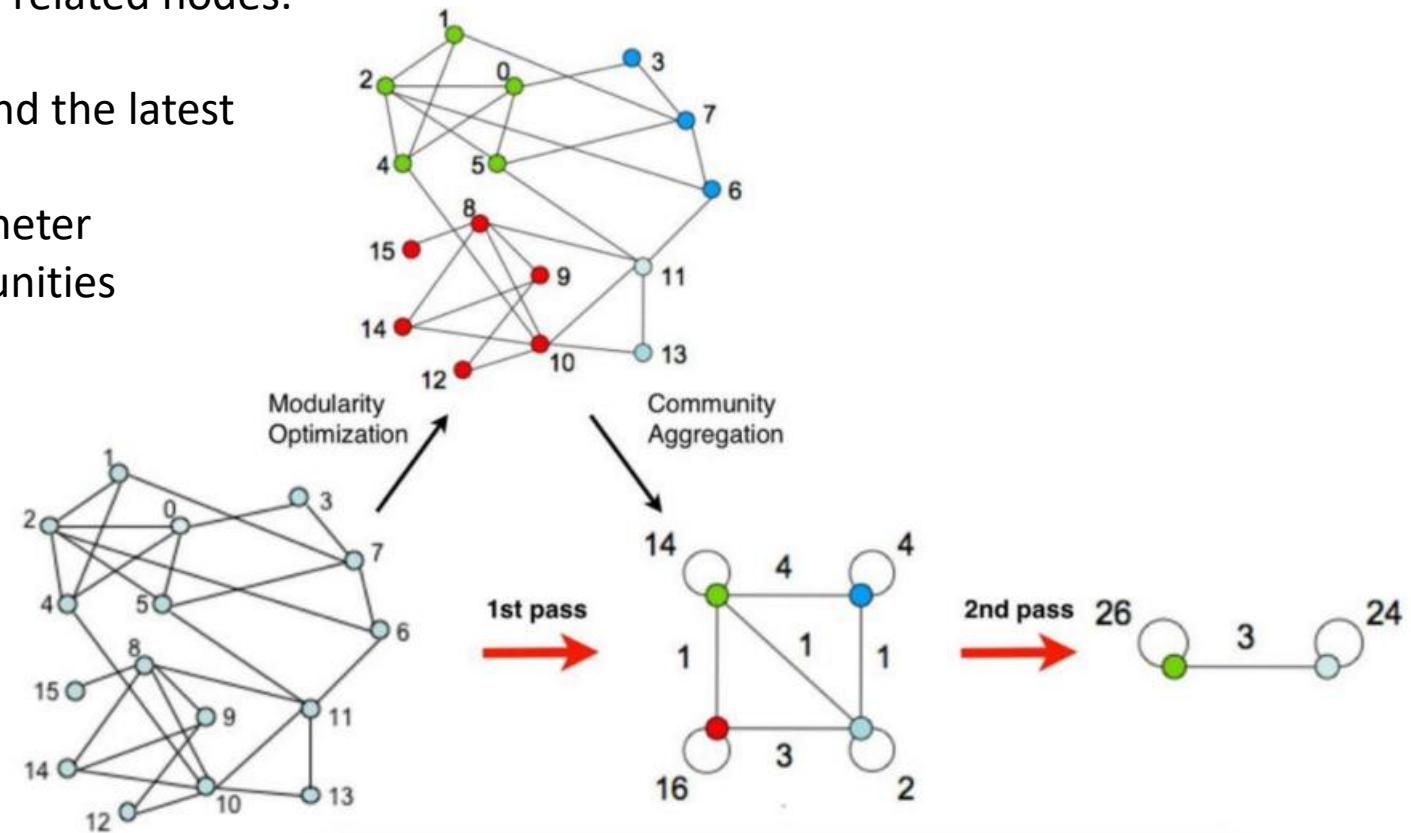


- Detection of Populations (**Clustering** Algorithms)
  - Supervised/Informed → Kmeans, Hierarchical, GMM,...
  - Unsupervised/Uninformed → Mean-shift, X-shift, Phenograph, ClusterX, Louvain



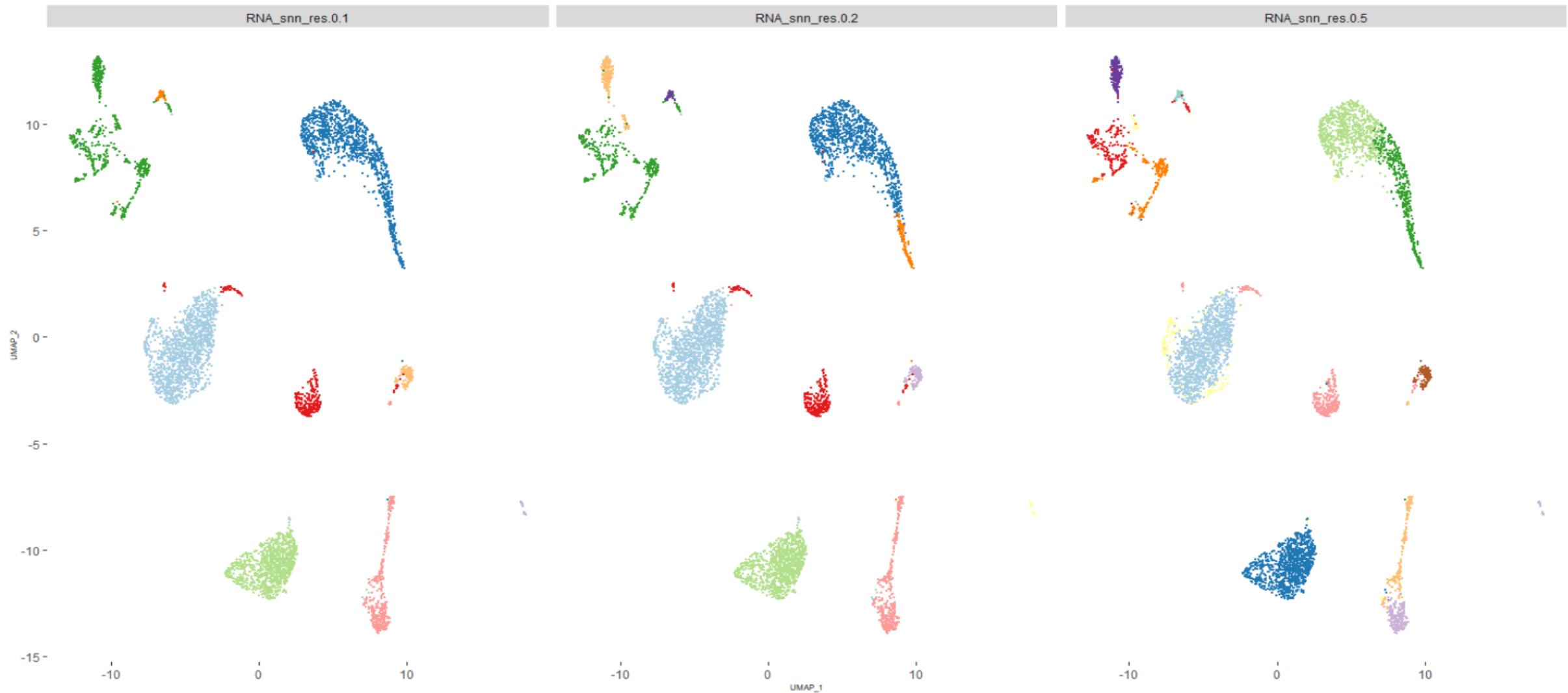
## Find Clusters

- This step tries to find local communities of very related nodes.
- There are several algorithms available.
- Most widely used is Louvain based clustering and the latest Leiden improvement.
- Communities are defined by a resolution parameter
- The higher resolution, the more are the communities subdivided.
- The resolution parameter is arbitrary and here a bit of biological interpretation is needed to refine the automated clusters.

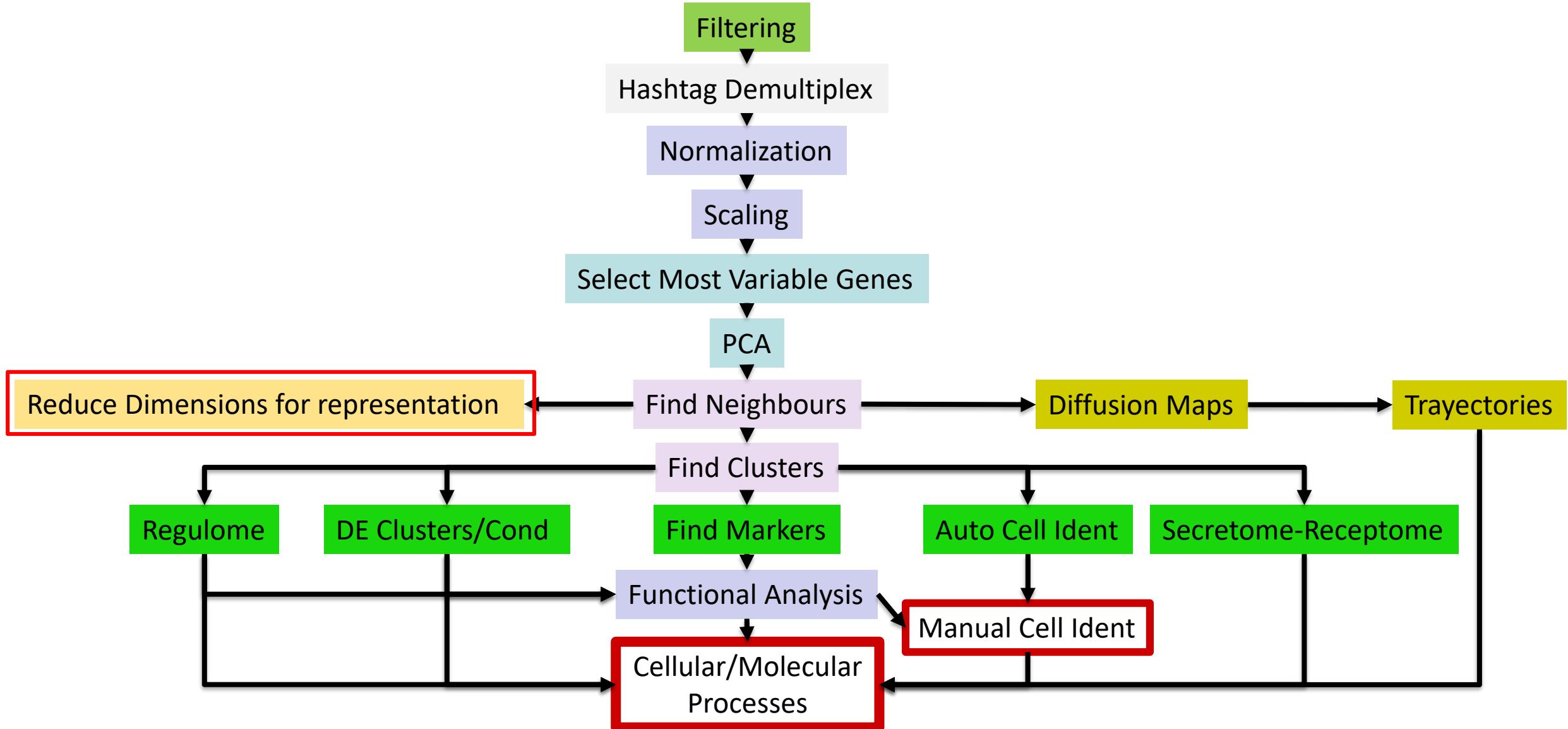


# scRNA-Seq Workflow

*cnii*c

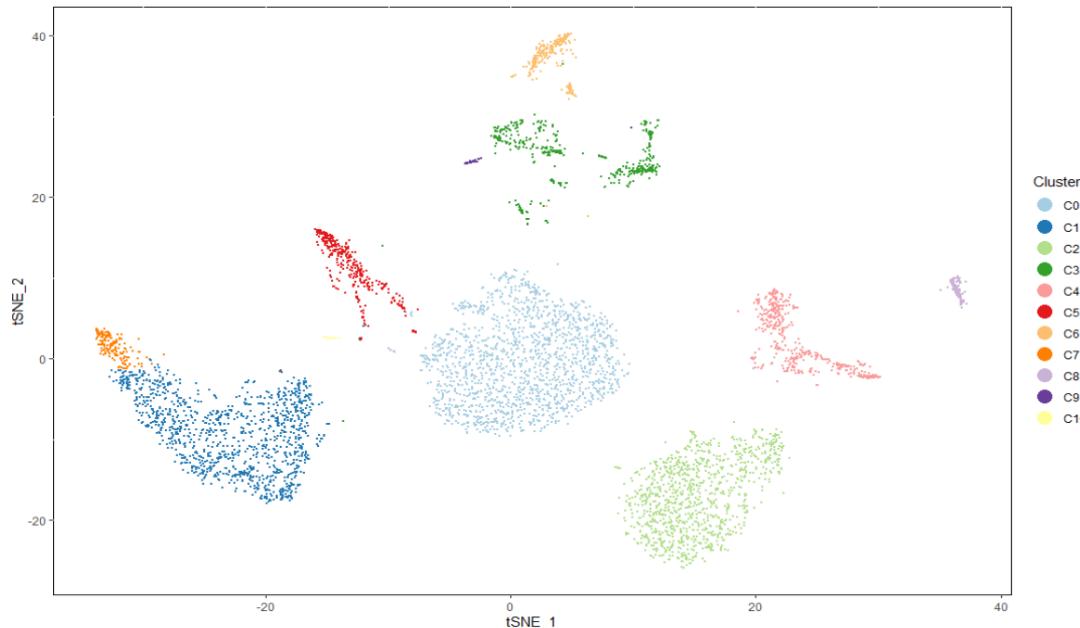


# scRNA-Seq Workflow

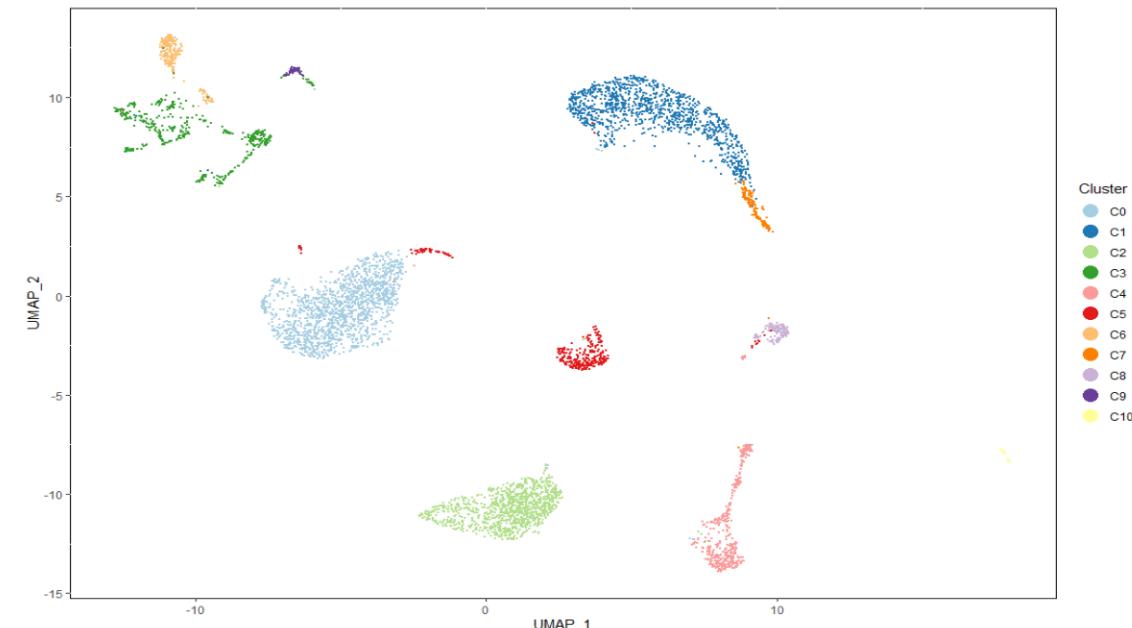


## Reduce Dimensions for representation

tSNE

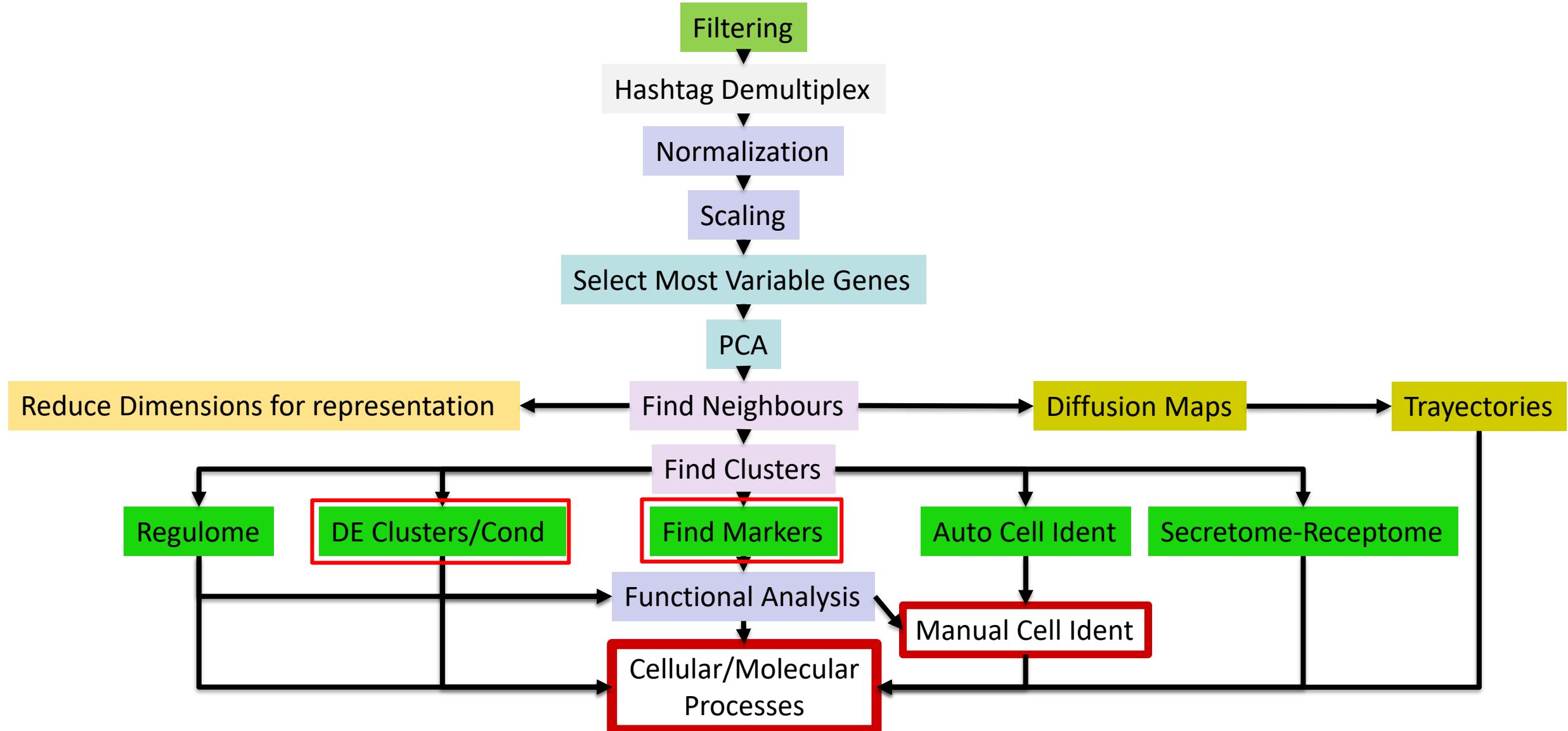


UMAP



Both are perfectly valid for representation purposes. UMAP better represents global structure although still, distances between separated clusters are not indicative.

# scRNA-Seq Workflow



## Find Markers

- The goal is to define the most relevant genes that define a cluster.
- The test compares the average expression of all the genes in a particular cluster to the average expression in the rest.
- Several clusters may share marker genes.
- Several methods can be applied.

However the most robust methods are:

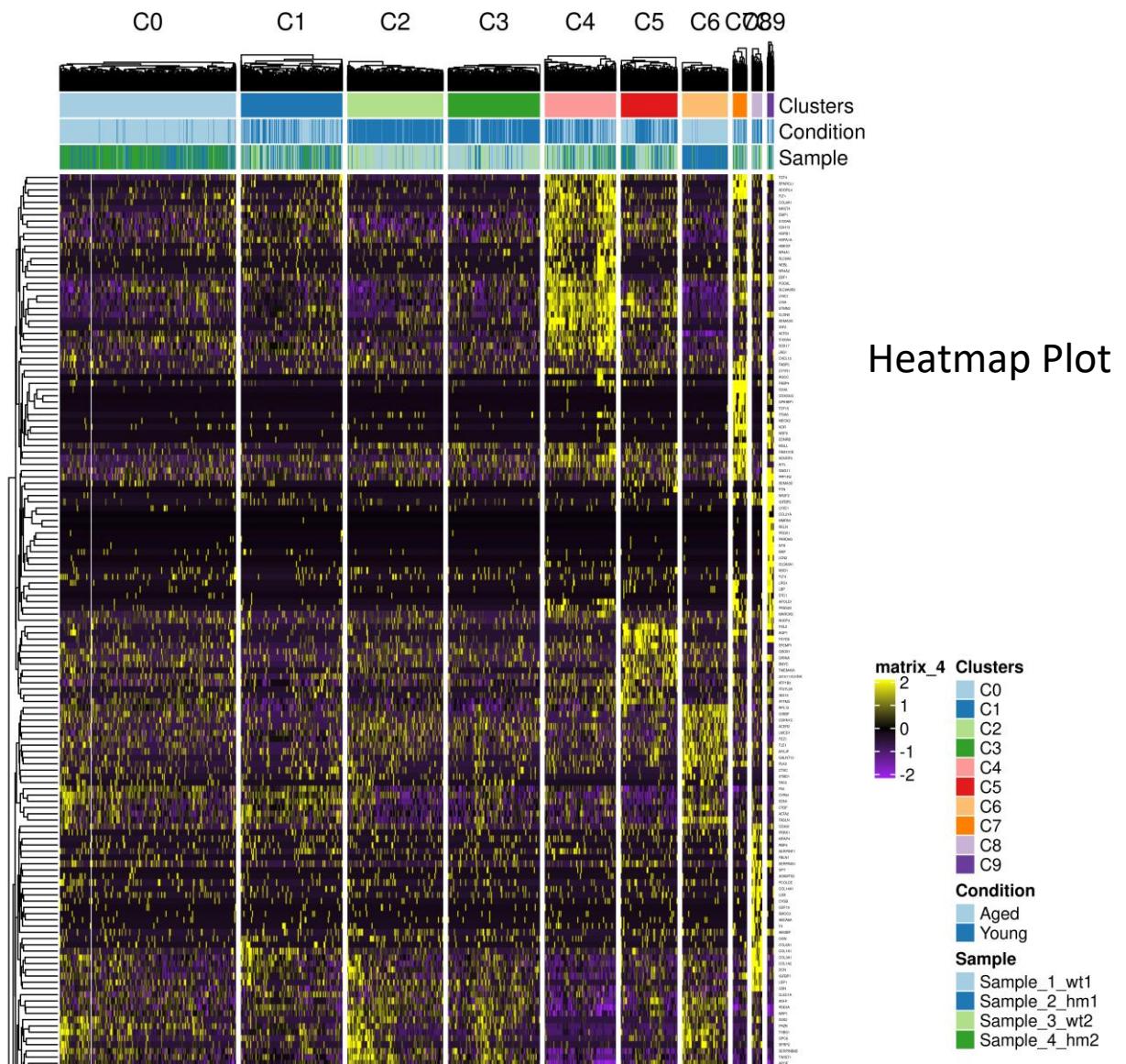
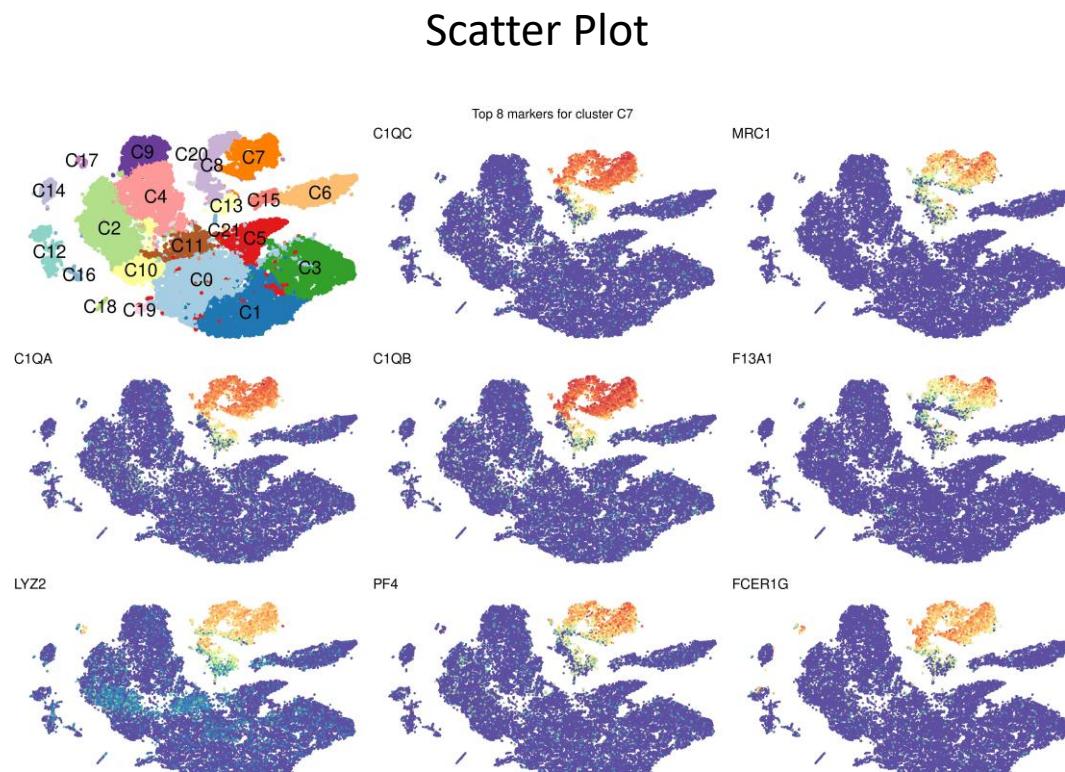
- Wilcoxon: A non-parametric test which is robust to the underlying distribution of the data although not very powerful.
- MAST: a two-part generalized linear model that simultaneously models the rate of expression over the background of various transcripts, and the positive expression mean. This model can accommodate complex designs.

Finak, G., McDavid, A., Yajima, M. *et al.* MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data. *Genome Biol* **16**, 278 (2015). <https://doi.org/10.1186/s13059-015-0844-5>

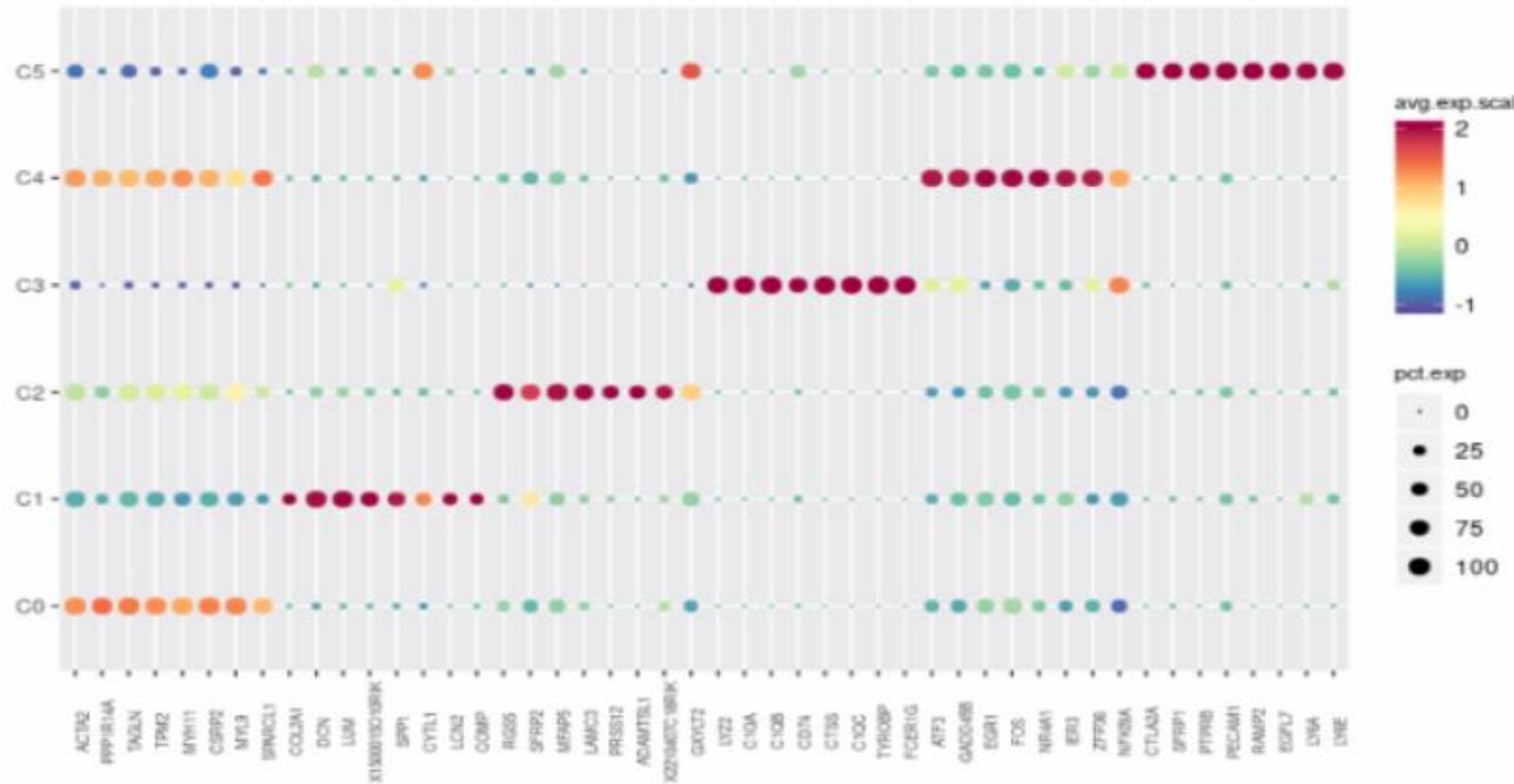
DE between cluster pairs and DE between conditions can also been computed

# scRNA-Seq Workflow

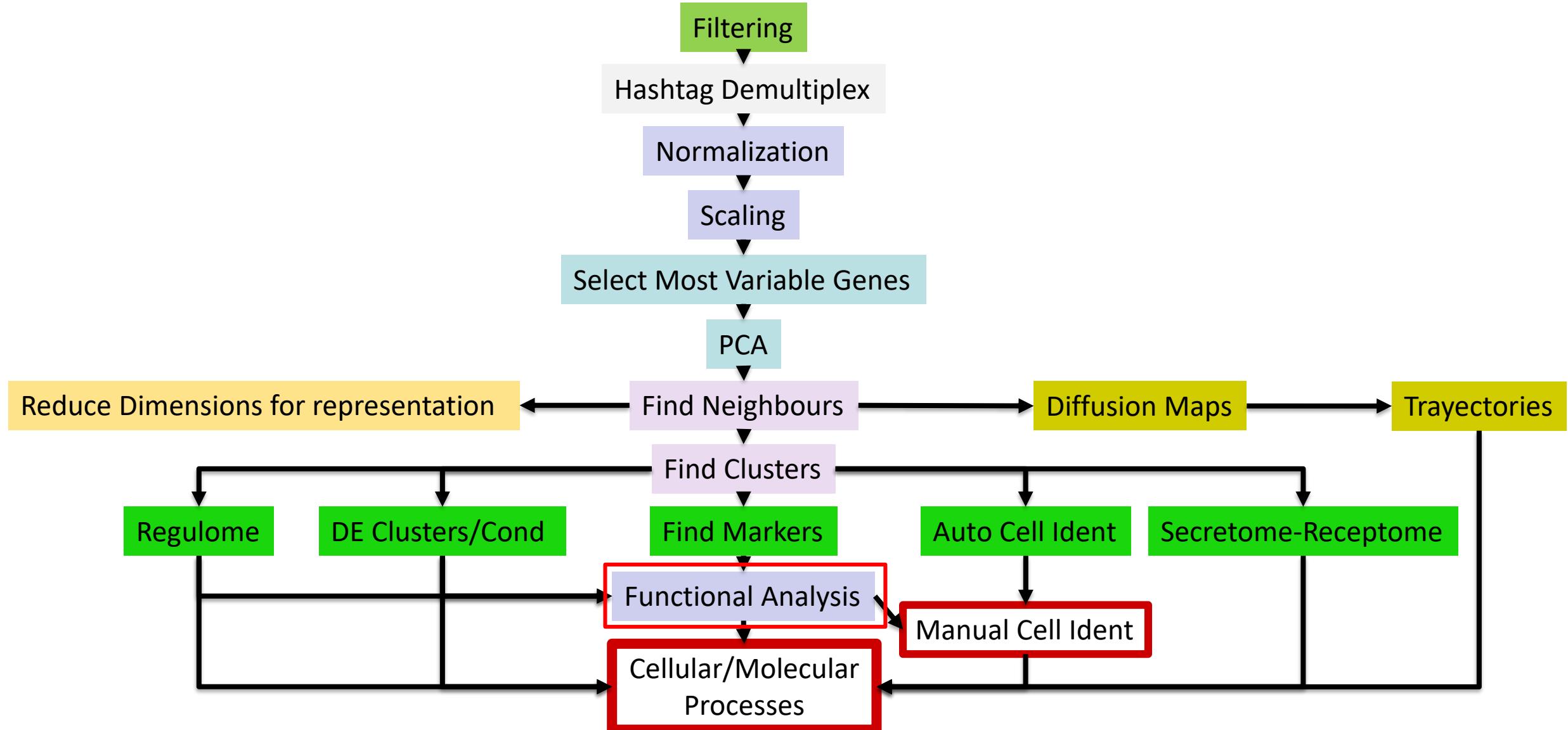
Find Markers



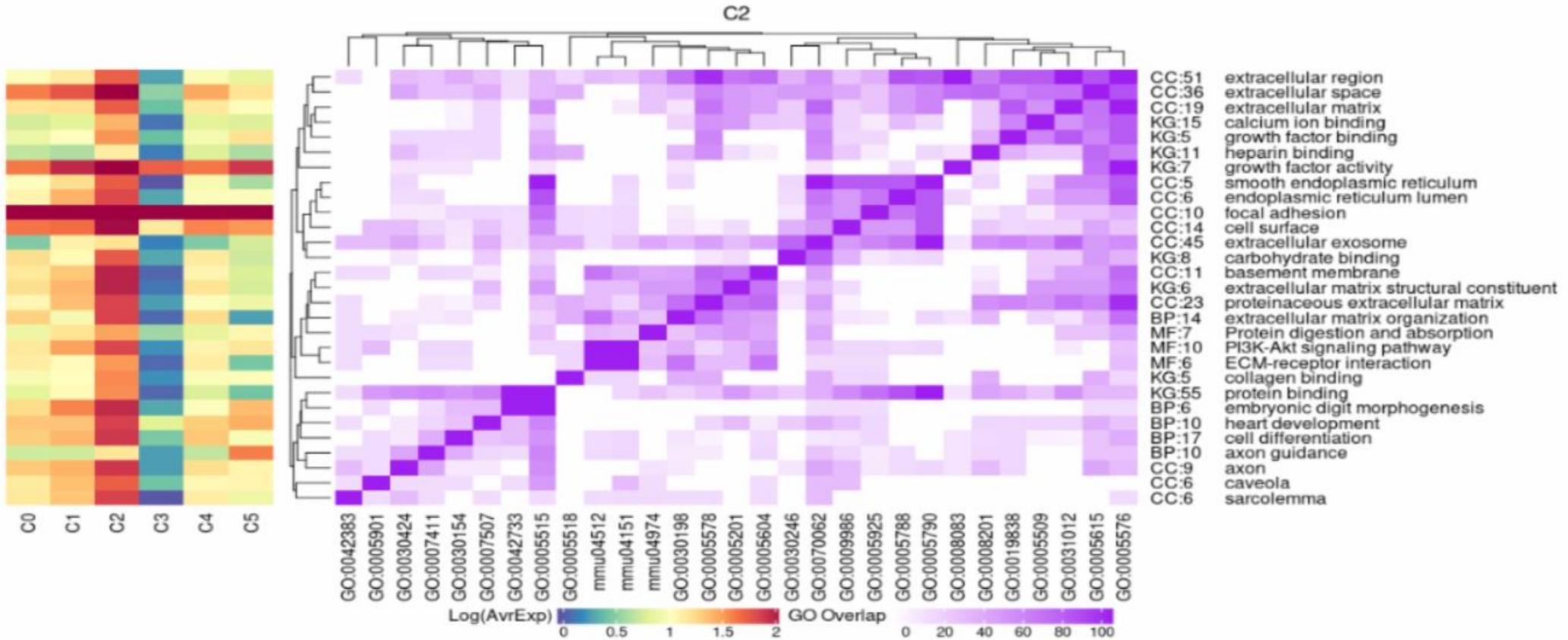
Top Marker Genes Dot plot



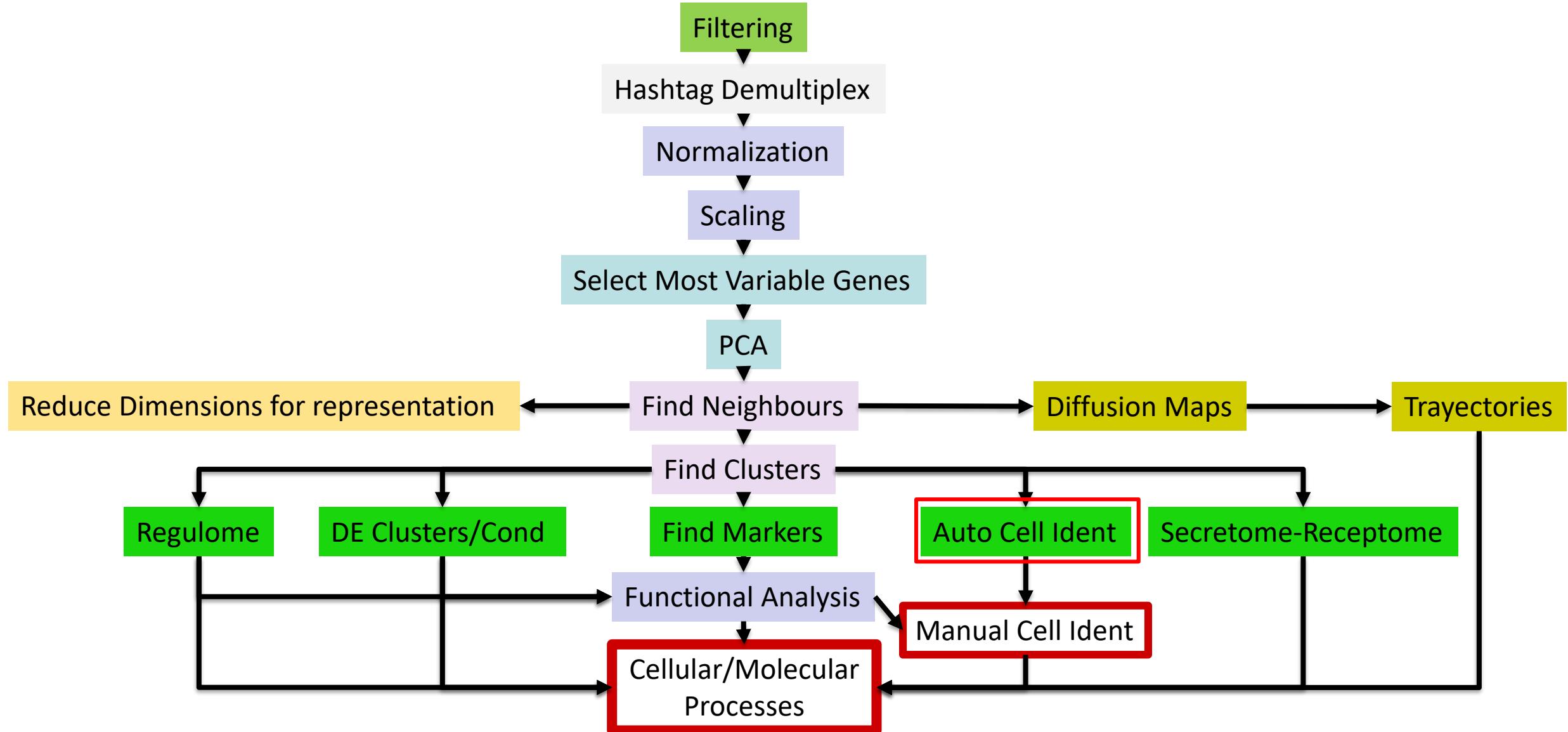
# scRNA-Seq Workflow



## Functional Analysis



# scRNA-Seq Workflow



## Auto Cell Type Identification

SingleR: A reference correlation based cell type identification using robust variant of nearest-neighbors classification, with some tweaks to improve resolution for closely related labels

### **Human primary cell atlas (HPCA)**

The HPCA reference consists of publicly available microarray datasets derived from human primary cells (Mabbott et al. 2013). Most of the labels refer to blood subpopulations but cell types from other tissues are also available.

### **Blueprint/ENCODE**

The Blueprint/ENCODE reference consists of bulk RNA-seq data for pure stroma and immune cells generated by Blueprint (Martens and Stunnenberg 2013) and ENCODE projects (The ENCODE Project Consortium 2012).

### **Mouse RNA-seq**

This reference consists of a collection of mouse bulk RNA-seq data sets downloaded from the gene expression omnibus (Benayoun et al. 2019). A variety of cell types are available, again mostly from blood but also covering several other tissues.

### **Novershtern hematopoietic data**

The Novershtern reference (previously known as Differentiation Map) consists of microarray datasets for sorted hematopoietic cell populations from [GSE24759](#) (Novershtern et al. 2011).

### **Immunological Genome Project (ImmGen)**

The ImmGen reference consists of microarray profiles of pure mouse immune cells from the [project of the same name](#) (Heng et al. 2008). This is currently the most highly resolved immune reference - possibly overwhelmingly so, given the granularity of the fine labels.

### **Database of Immune Cell Expression/eQTLs/Epigenomics (DICE)**

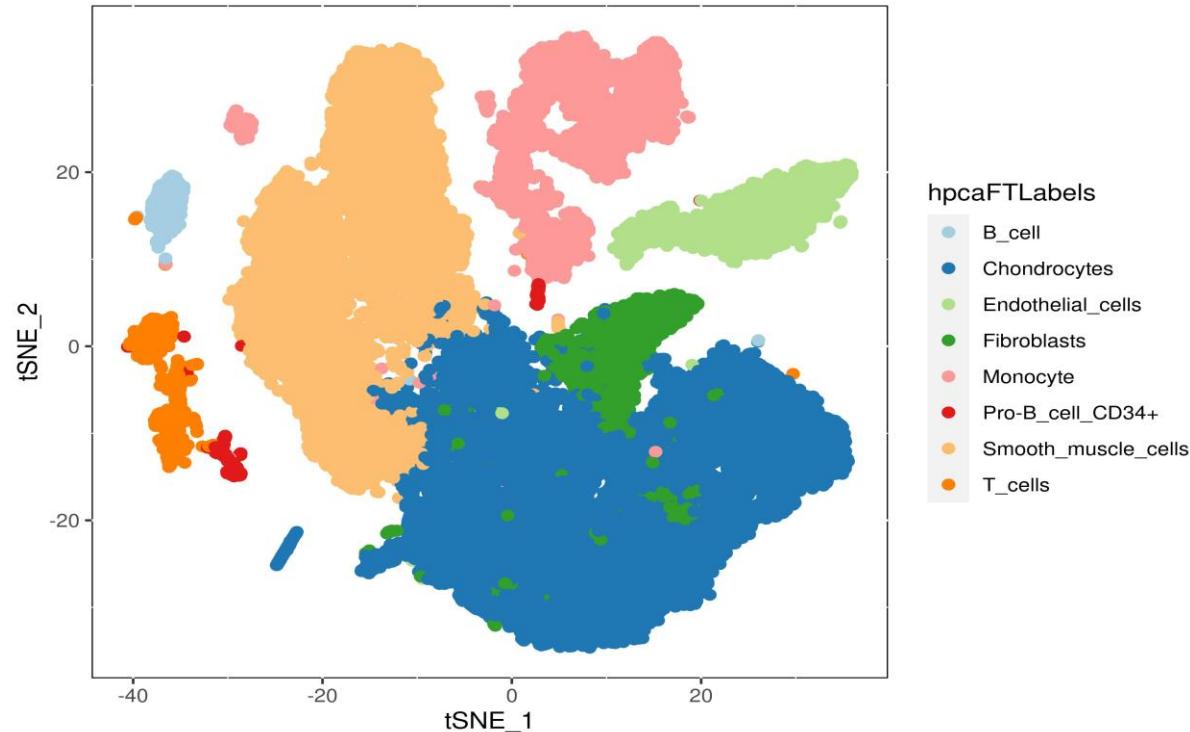
The DICE reference consists of bulk RNA-seq samples of sorted cell populations from the [project of the same name](#) (Schmiedel et al. 2018).

### **Monaco immune data**

The Monaco reference consists of bulk RNA-seq samples of sorted immune cell populations from [GSE107011](#) (Monaco et al. 2019).

## Auto Cell Type Identification

SingleR Cell Type Classification based on Human Cell Atlas Reference

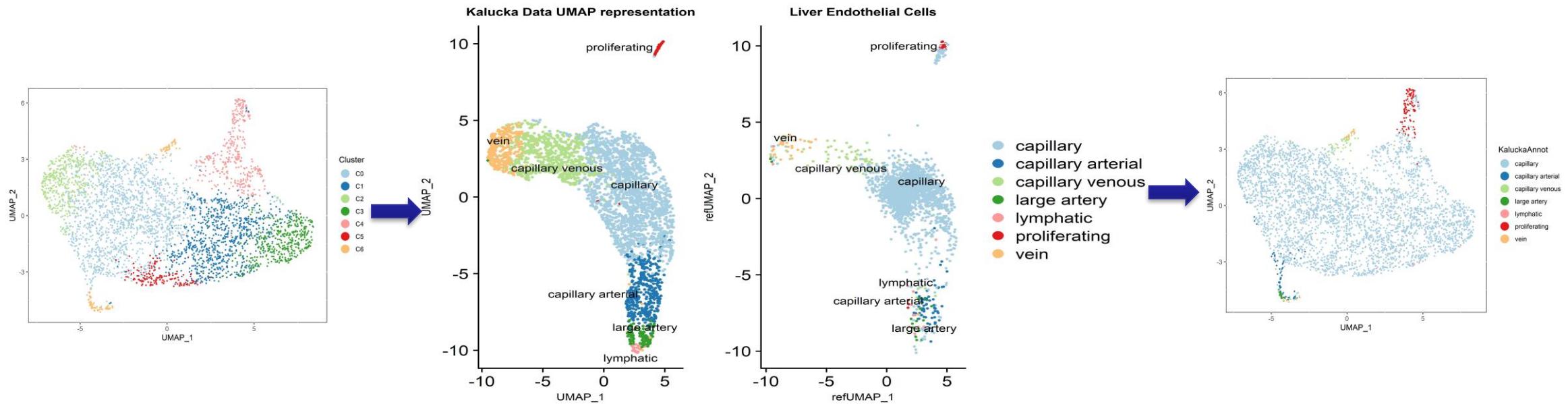


Classifications are far from perfect but they are a way to start and incorporate more knowledge using gene expression and functional analysis results to refine the annotations

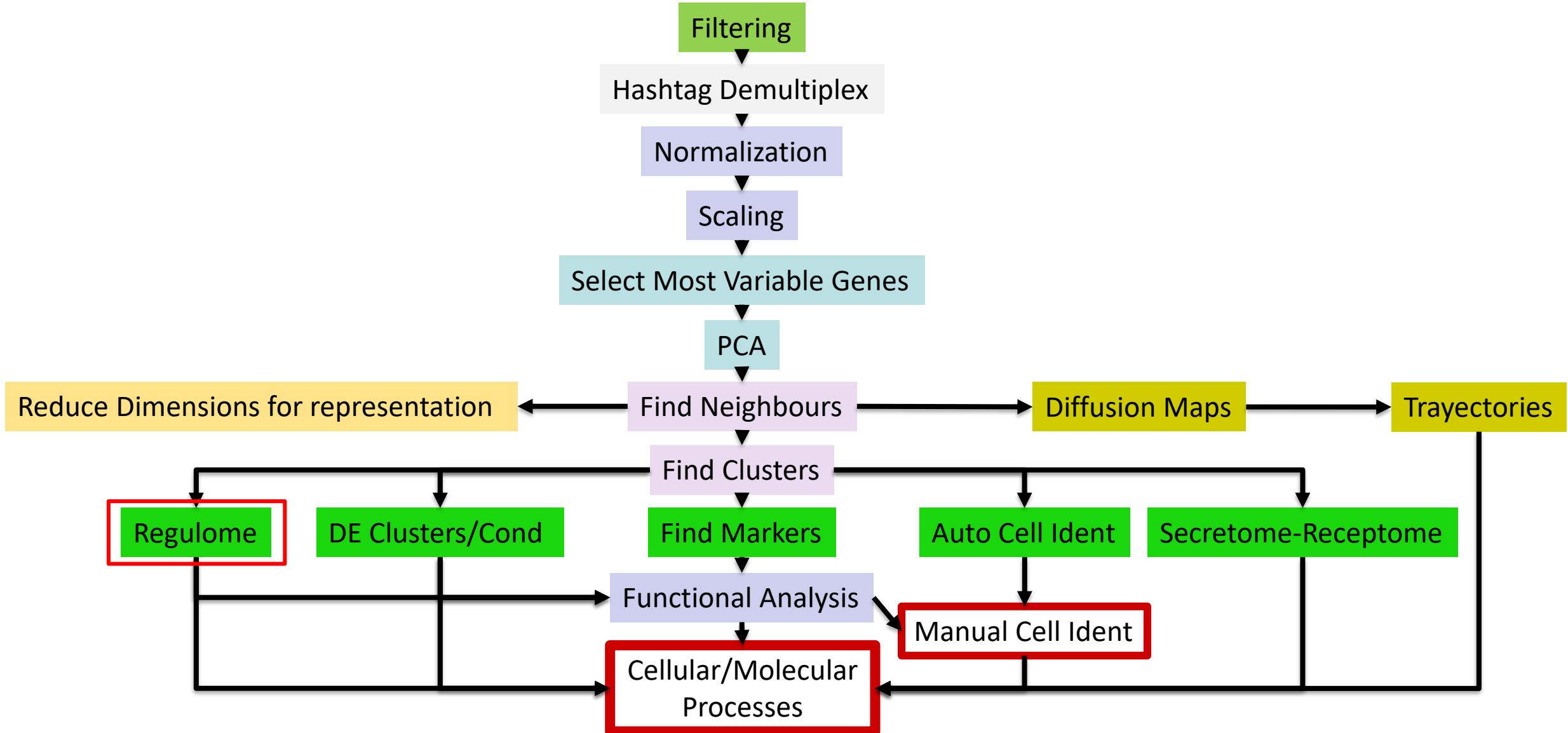
## Auto Cell Type Identification

Map single cell experiment to an already annotated single cell experiment:

Using **Canonical Croscorelation Analysis** and **Mutual Nearest Neighbours**.

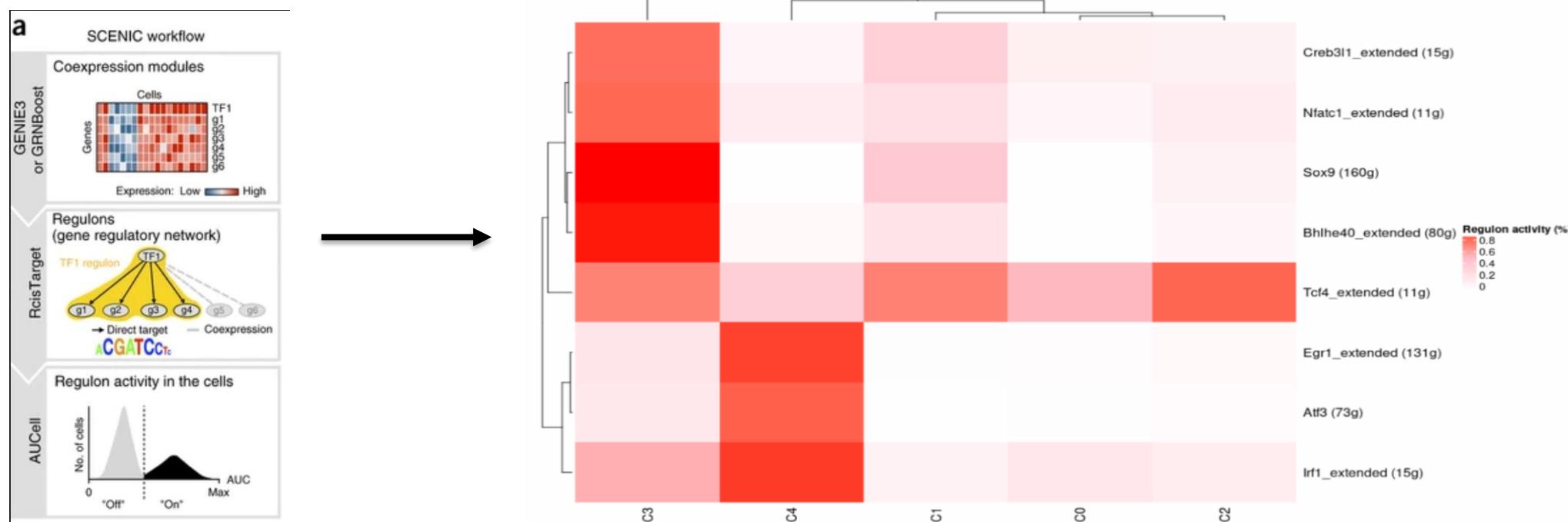


# scRNA-Seq Workflow



## Regulome

This analysis tries to find transcription factors that might explain the gene expression variability found between clusters

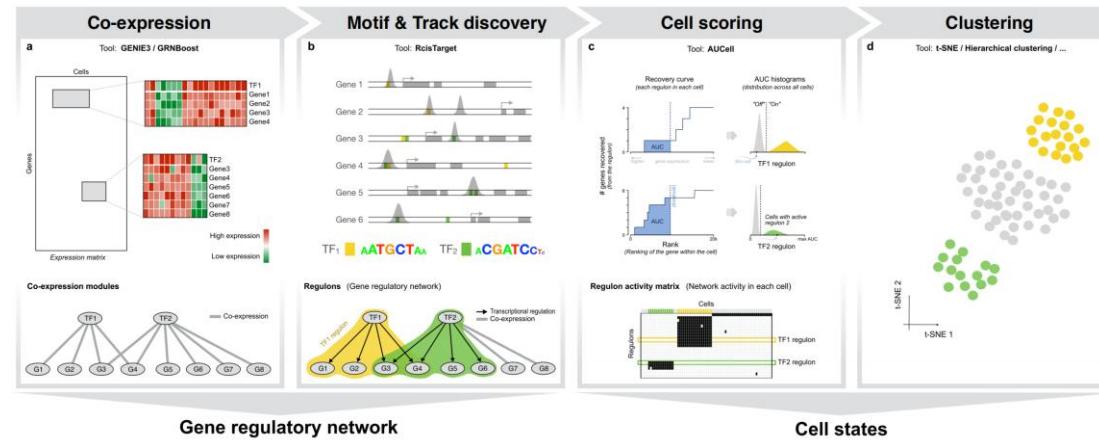


Aibar, S., González-Blas, C., Moerman, T. et al. SCENIC: single-cell regulatory network inference and clustering. *Nat Methods* **14**, 1083–1086 (2017).  
<https://doi.org/10.1038/nmeth.4463>

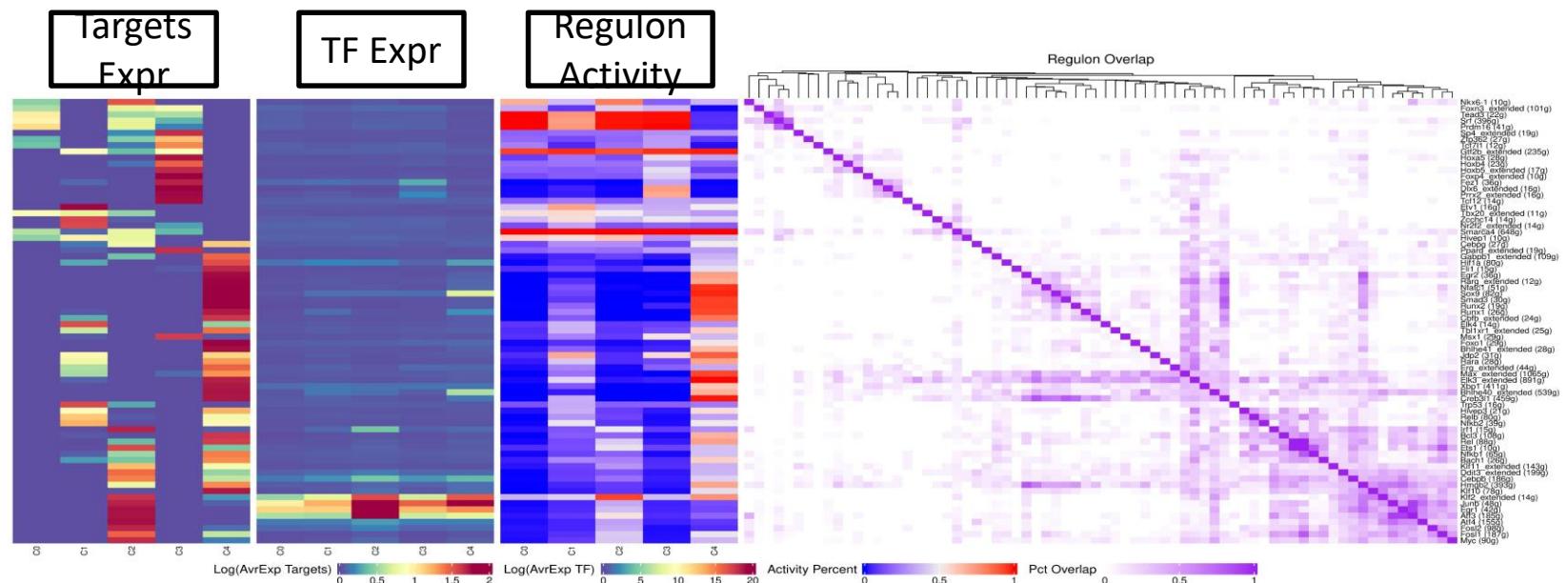
# scRNA-Seq Extended

cnic

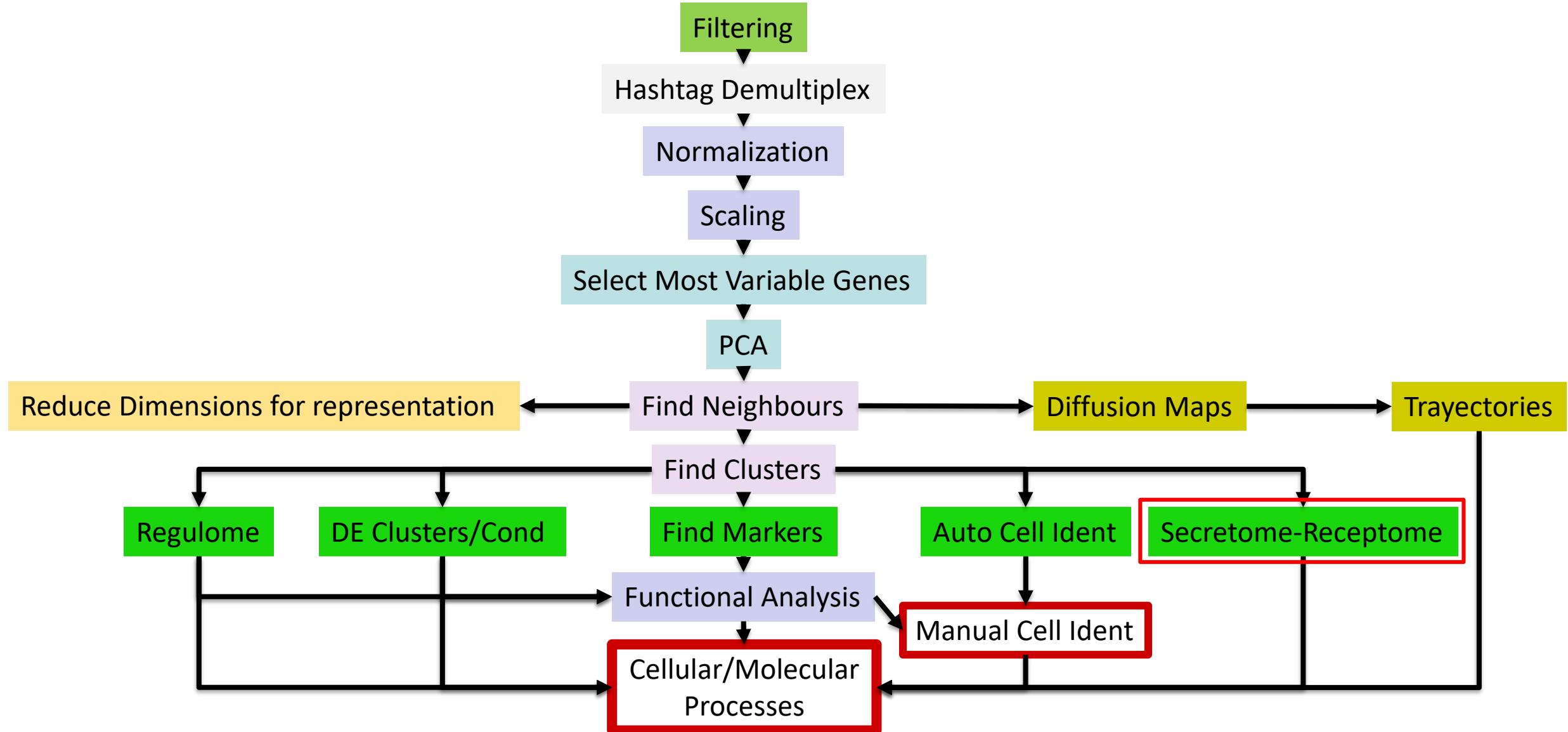
## Regulome



<https://scenic.aertslab.org/>



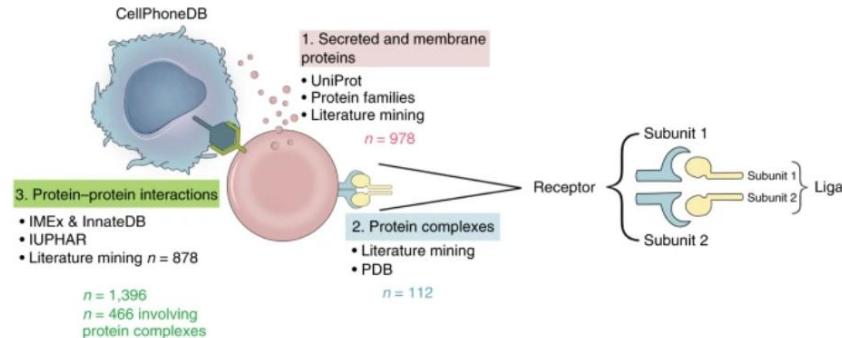
# scRNA-Seq Workflow



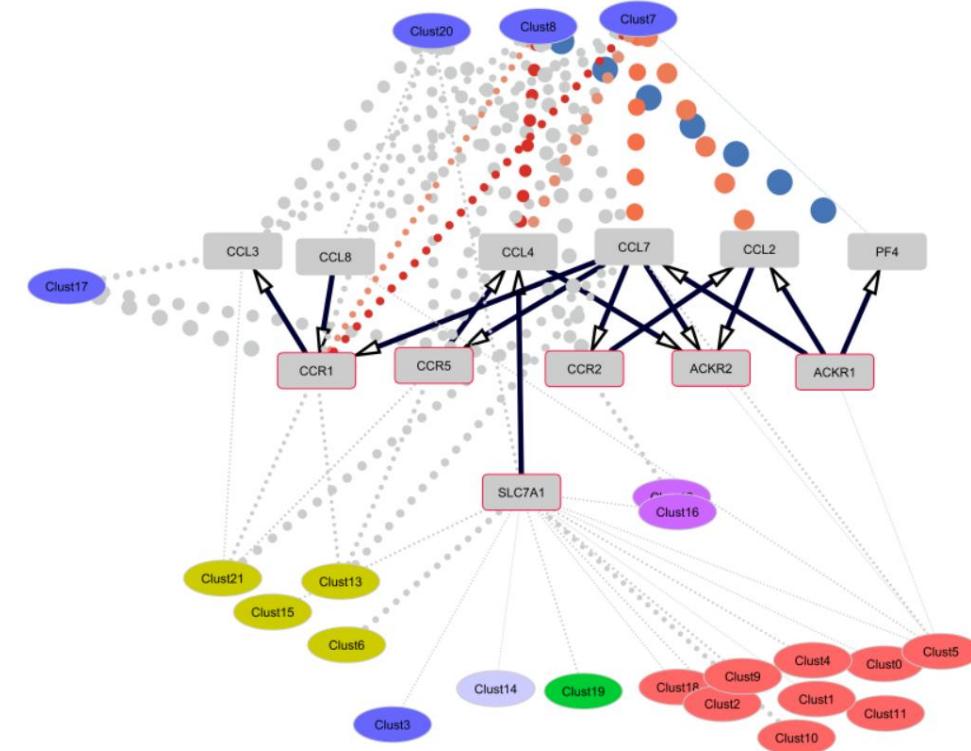
## Secretome-Receptome

### CellPhoneDB, CellChat, ...

- Public repository of ligands, receptors and their interactions to enable a comprehensive, systematic analysis of cell–cell communication molecules
- derive enriched ligand–receptor interactions between two cell states on the basis of expression of a receptor by one cell state and a ligand by another cell state and use empirical shuffling to calculate which ligand–receptor pairs display significant cell-state specificity

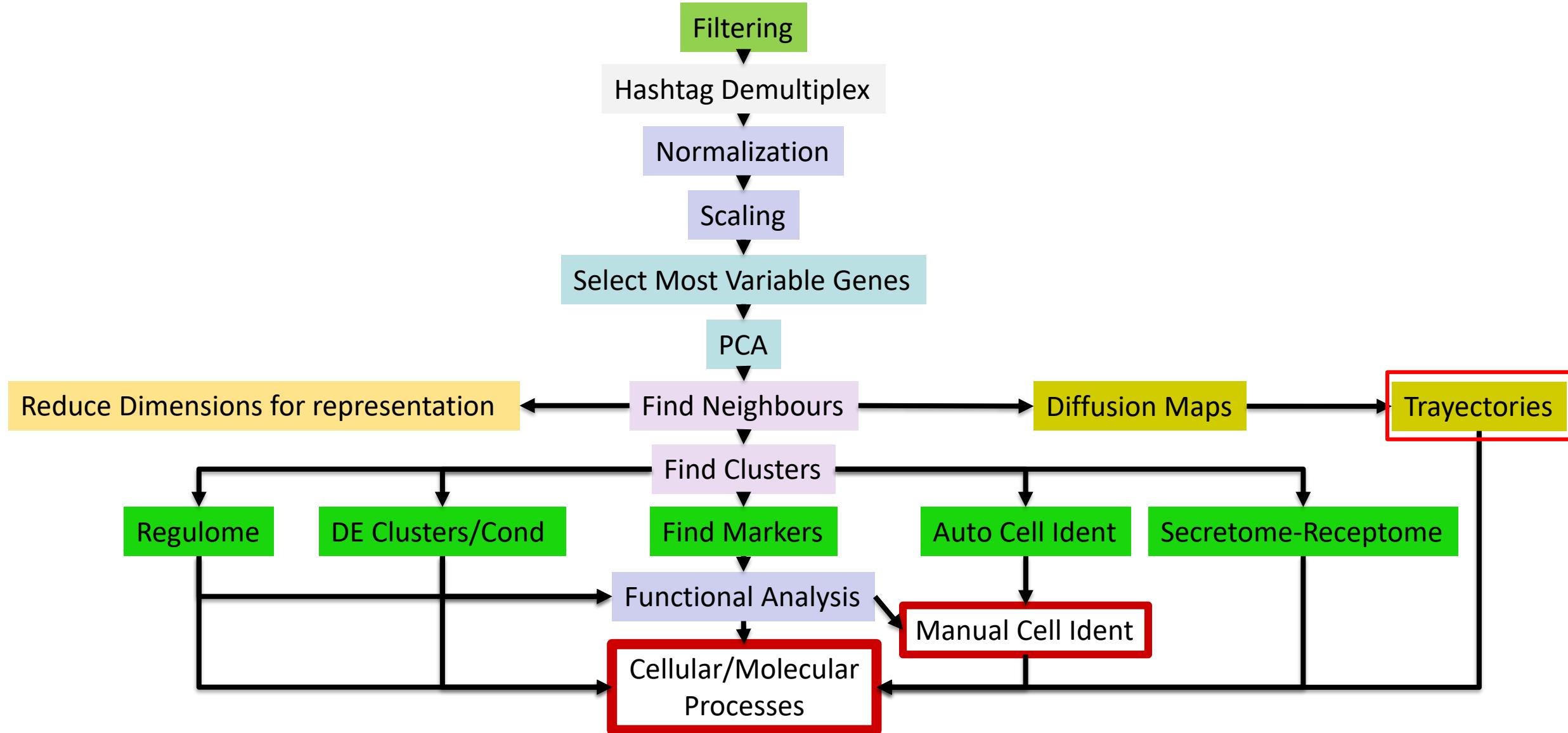


## Interaction Module from Aged Aorta Samples

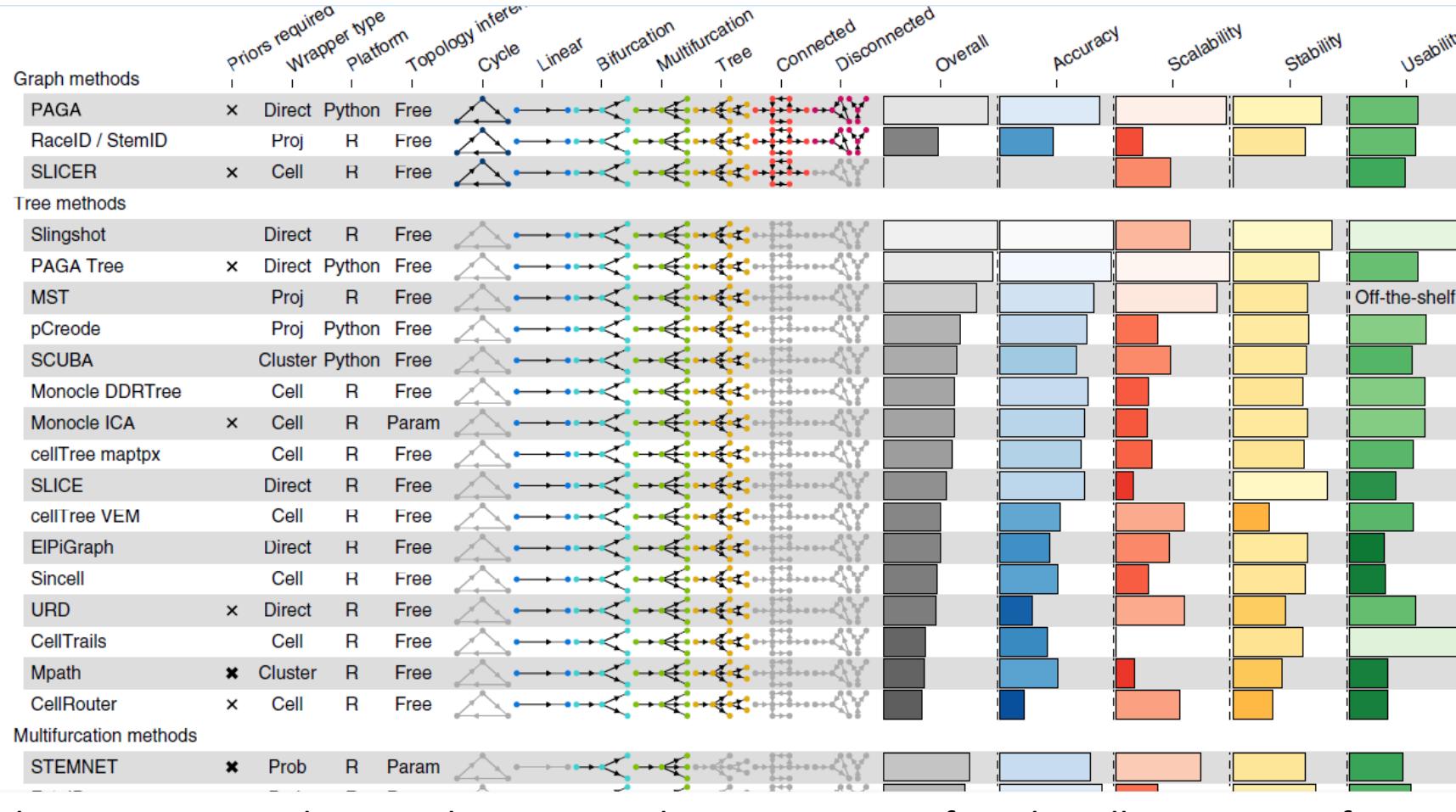


Efremova, M., Vento-Tormo, M., Teichmann, S.A. et al. CellPhoneDB: inferring cell–cell communication from combined expression of multi-subunit ligand–receptor complexes. *Nat Protoc* **15**, 1484–1506 (2020). <https://doi.org/10.1038/s41596-020-0292-x>

# scRNA-Seq Workflow

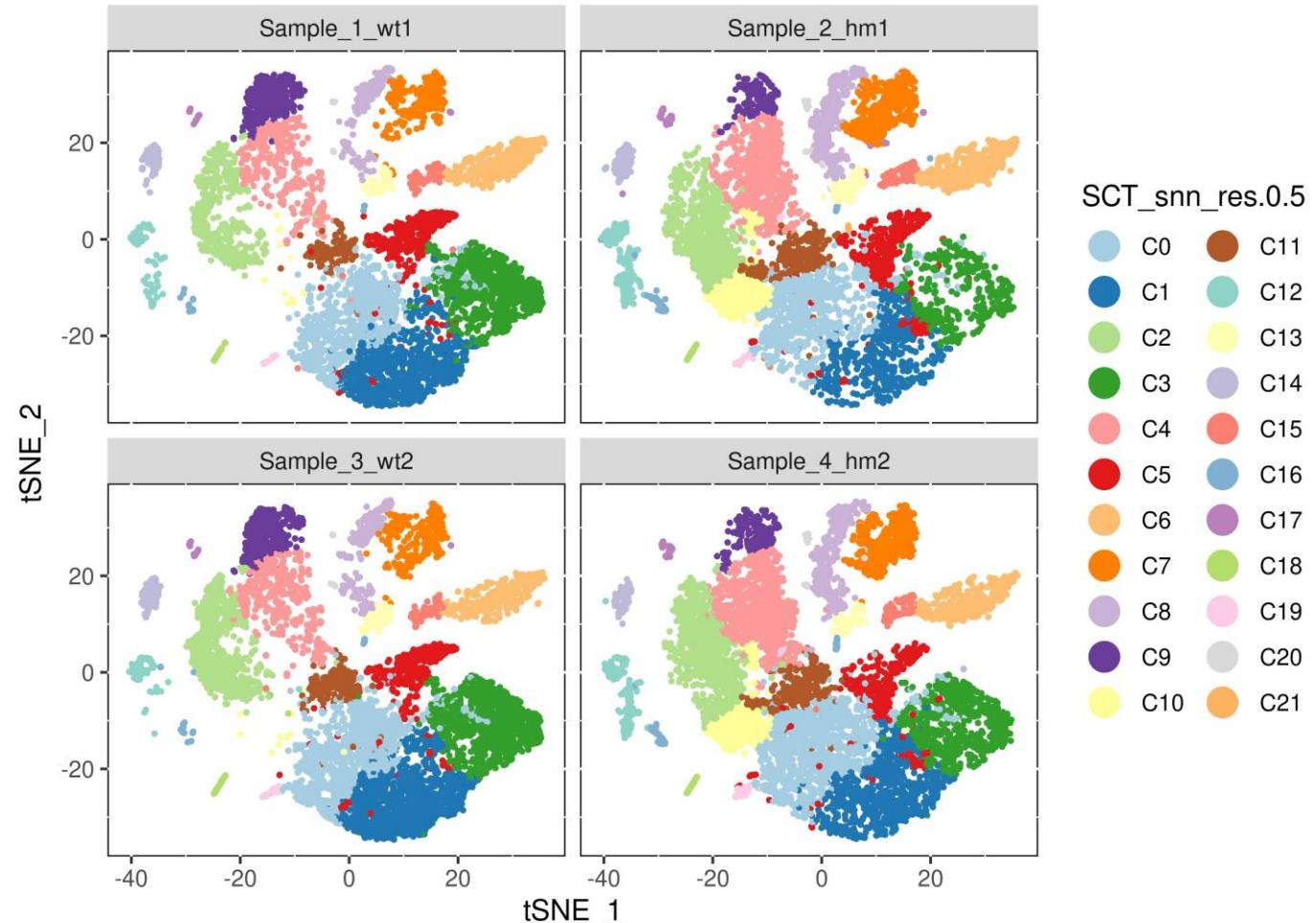


# scRNA-Seq Workflow

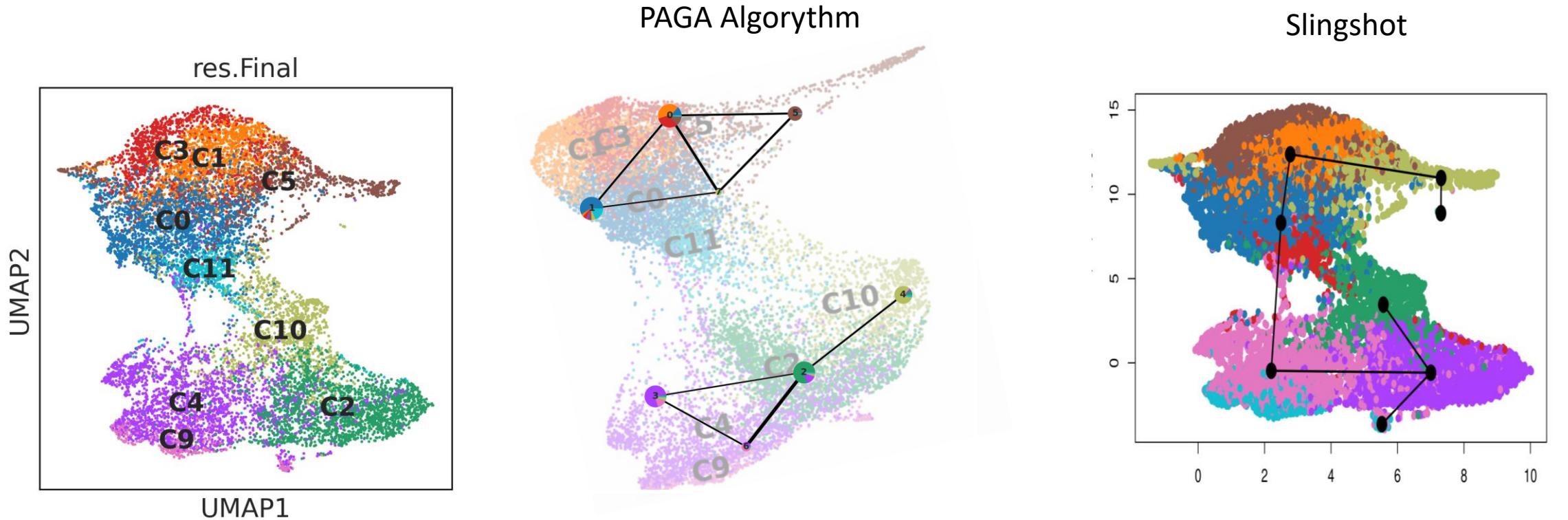


Saelens, W., Cannoodt, R., Todorov, H. et al. A comparison of single-cell trajectory inference methods. *Nat Biotechnol* **37**, 547–554 (2019). <https://doi.org/10.1038/s41587-019-0071-9>

Aorta Aged Samples from Ana Barrentino. LAB\_VA



Trajectory Analysis helped to identify the origin of the new cell population





# scDAVIS

<https://bioinfo.cnic.es/scdavis/>

- Explore Published Datasets
- Upload your own experiment
- Interrogate the data:
  - Manually define clusters
  - Find markers for the cluster
  - Plot results in different ways
- Download your analysis



**END**

---