



# Single cell RNAseq basics

*Anna Lorenc (Ania)*  
*al16@sanger.ac.uk*



## Aims for this lecture

- Recap of the technology
- Understand challenges & ways to solve them
- Steps from sequenced reads to cell-level data
- From cell-level data to better cell level data

Principles

From FASTQ  
to expression  
matrices

Better  
expression  
data

3 parts – questions after each

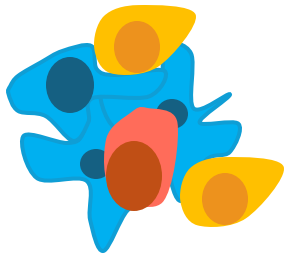
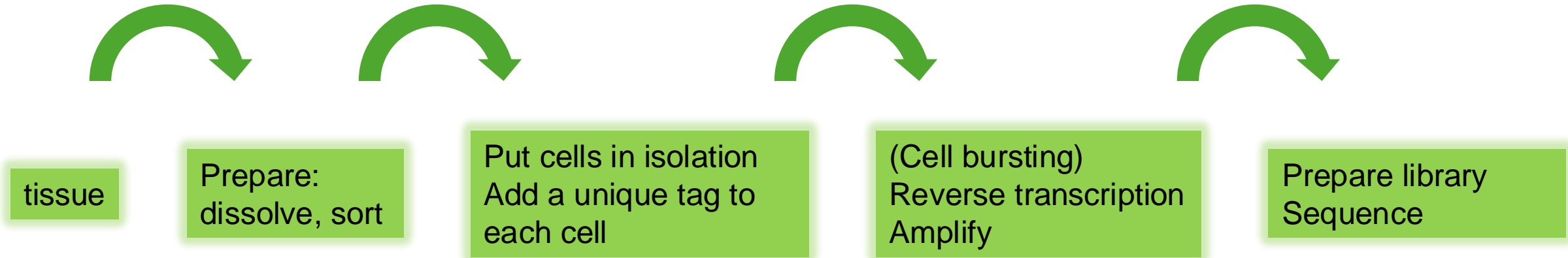
We want to know expression of ~~all~~-many genes in-~~all~~-many cells

...to understand how cells and tissues differ between each other

...to understand how cells and tissues change in a perturbation (experiment, disease, development)

# principles

We want to know expression of many genes in many cells

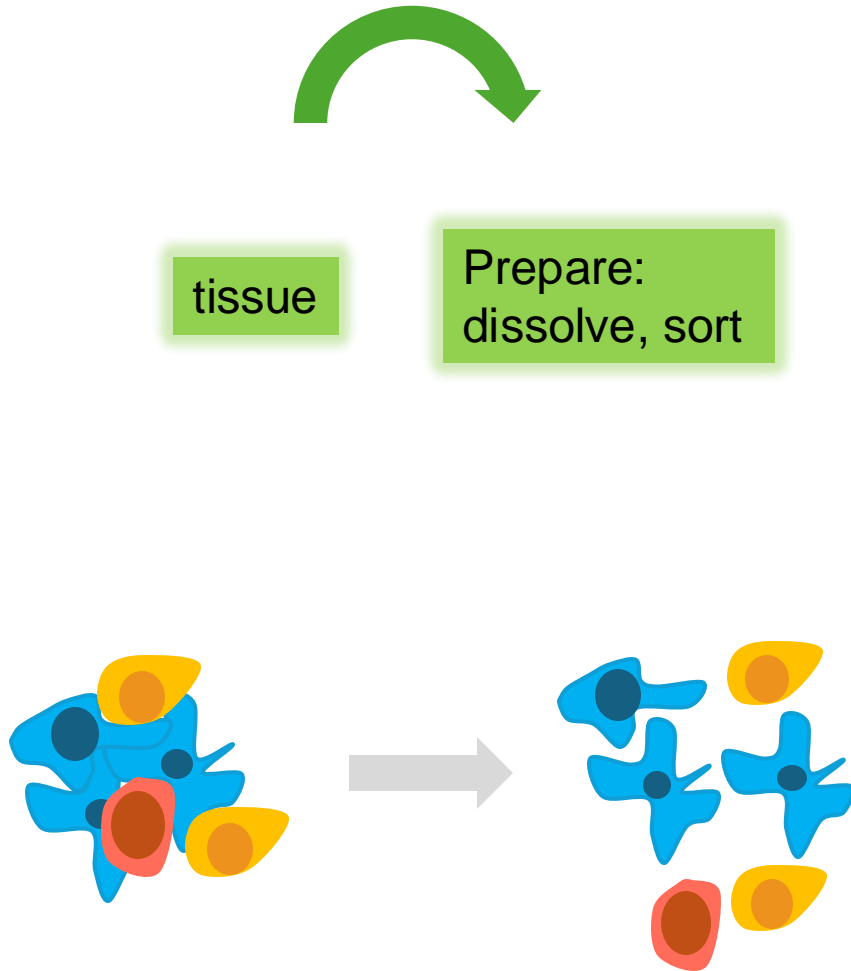


```
@A00815:607:H5V32DRX3:1:2101:1108:1000 2:N:0:ACAATGTGAA+TAACGGTACG
GAACCACTGAGCACAAAGTTTCTTCATCGTTCCTCAGATTCAGTAACATTATTAATTTAGACAATCCCGTGAAGGCCAATTCATCAGTG
+
FF:FFFF:F,:FFFF,F,:F,FFF,FFFFFFFF,,,,:F::,FFFFFF,FF,FF,FF::,:FFF,,F,:FFF:F,FF,F:F:FF,F,
@A00815:607:H5V32DRX3:1:2101:1181:1000 2:N:0:AAAATGTGAA+TAACGGTACG
CGTCAAGGCAGGGCCTCGTTCCTTGCTGGGCACCAAGAGCAGATGACATATATAGCACAGTGCCTCCCCAGGACAGGAAGATGAGGCTAG
+
F,FFFF,FF,:FF,FF,,F:F:,FFF,F:FF:FFF,F,F,,FFF,,FFF:,F:FFF,FF,FFF::F:F:,FF,,FF,,FFFF:,
@A00815:607:H5V32DRX3:1:2101:1597:1000 2:N:0:ACAATGTGAA+TAACGGTACG
GAAGTGAGGGATGCTGAGGGCCGGGACAAGCTATCGGACTGTCTGCTGCCATCGGTAATGAGTCTCAGTAGACCTGGAACGTCACCTCGC
+
FFFFFFFFFFFFFFFF,FFFF:FFFFFFFFFFFFFFFF,FFFF,:FFF,:FFFF,FFFFFFFFFFFFFFFFFFFFFFFF:FFFFFF:
@A00815:607:H5V32DRX3:1:2101:1615:1000 2:N:0:ACAATGTGAA+TAACGGTACG
ATGTGACTATAGGCTCATAGCCATCTCATTATGCAAAATGTATTCACTACTGCTTTGTATGTCTCAATAGTCTCCAGATATACGGCGGT
+
,,, :FFFF,,FFFF,FFFF,F:F,F,FF,FF,:FFF:,F,:F,FFF,:FF::,:F,F,FF,,F,F,FFFFFF,FF,,F,F,F,:,,F:
```

## Challenges

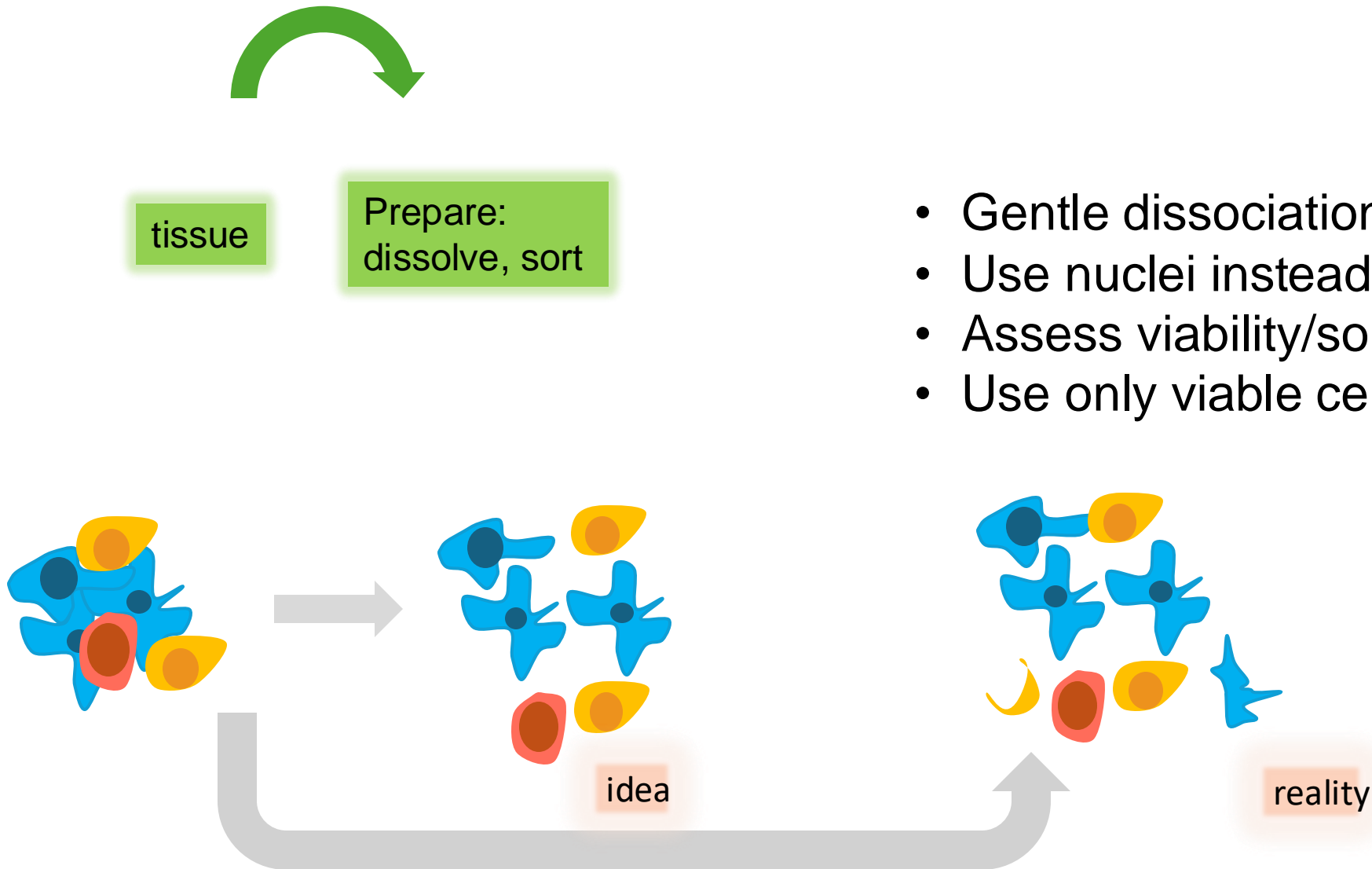
- very little RNA in each cell
- PCR steps – amplification bias
- separate cells so we can measure each cell separately
  - cells difficult to separate, fragile
  - empty droplets/multiplets
- cells dying/bursting during the procedure
- very wide dynamic range of expression of genes
- not enough of cells of a specific type

## principles



- Gentle dissociation protocols
- Use nuclei instead of cells
- Assess viability/sort
- Use only viable cells

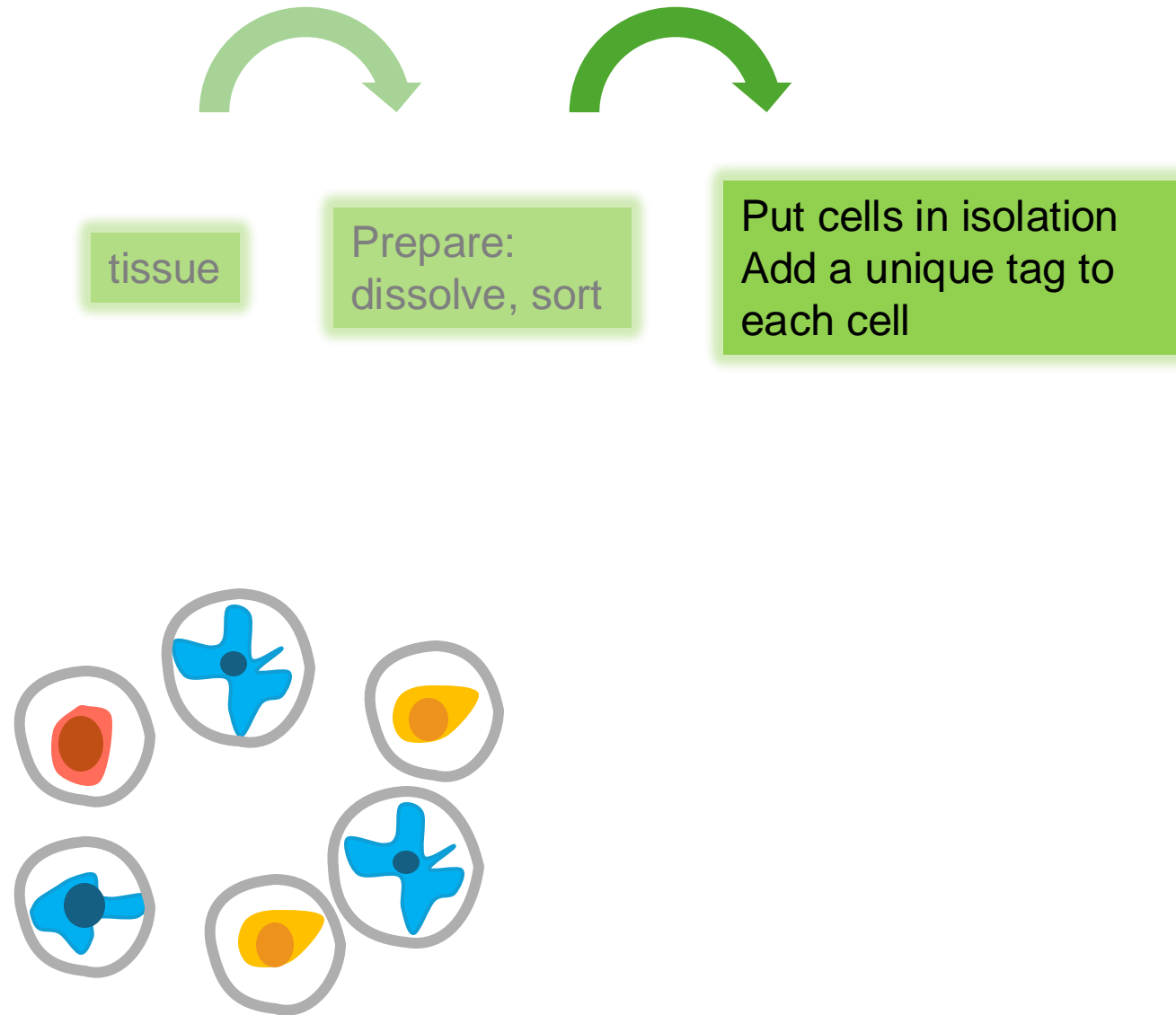
## principles



- Gentle dissociation protocols
- Use nuclei instead of cells
- Assess viability/sort
- Use only viable cells

And presence of debris is a challenge because...

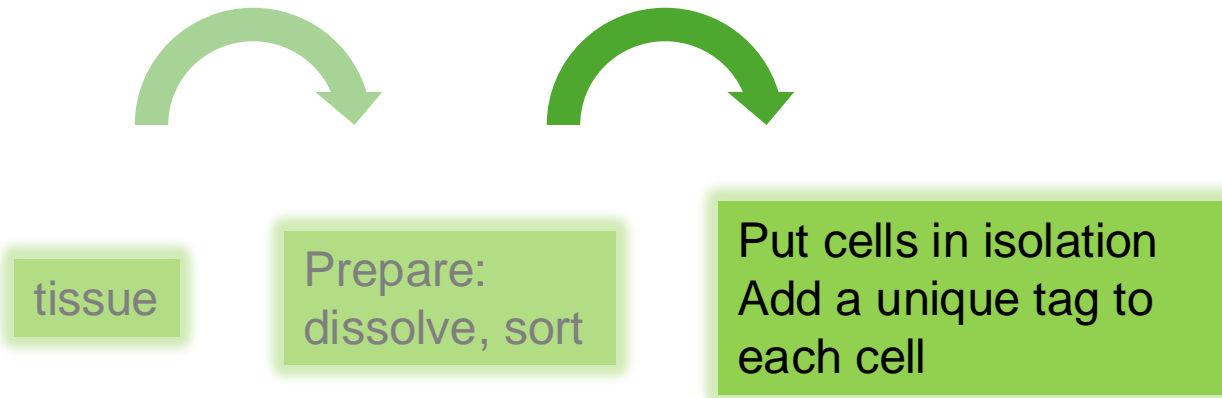
## principles



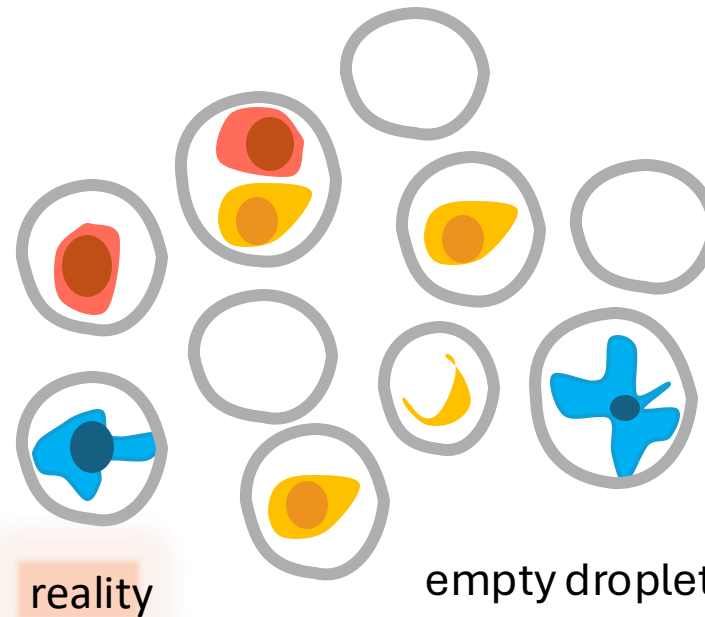
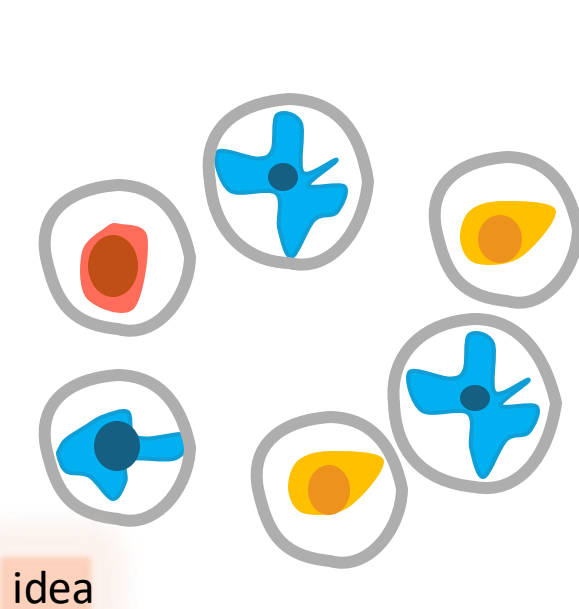
- wells
- microwells (BDRhapsody, Hive)
- droplets (10x)
- combinatorial tags (PARSE, Scala)



# principles

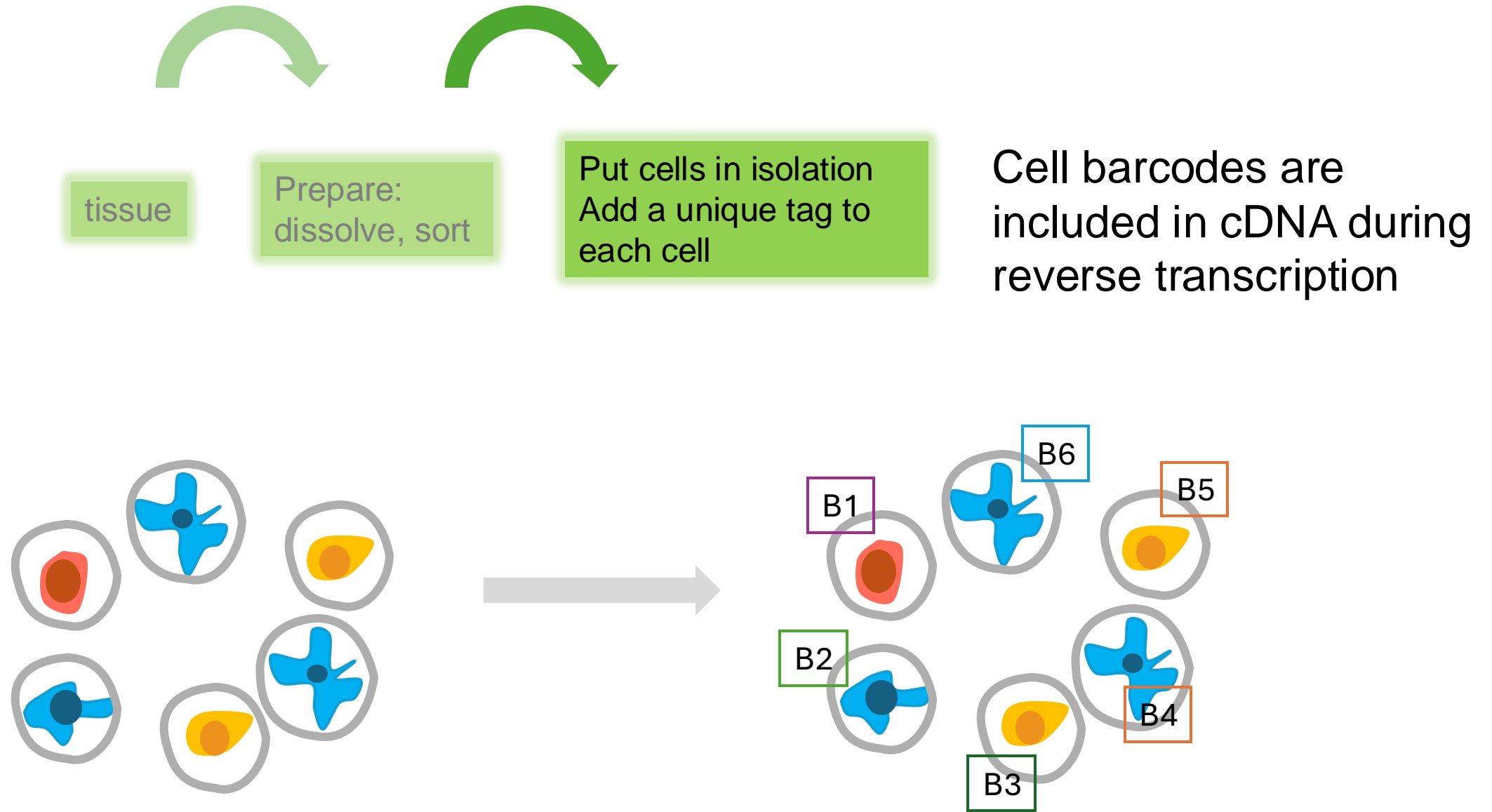


- wells
- microwells (BDRhapsody, Hive)
- droplets (10x)
- combinatorial tags (PARSE, Scala)

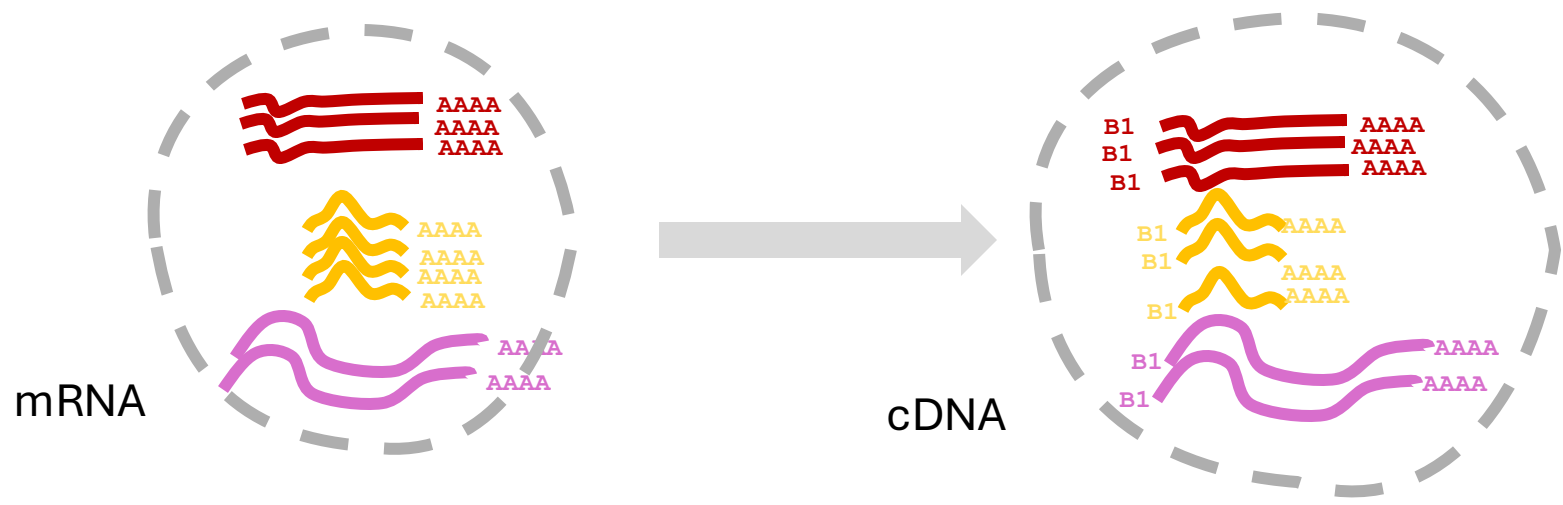
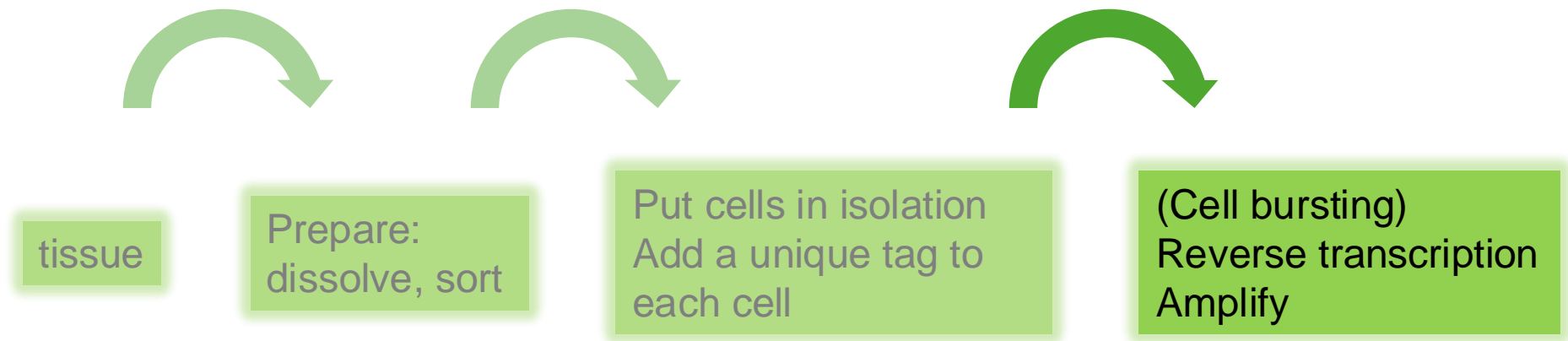


empty droplets/wells ← **single cell** → multiplets

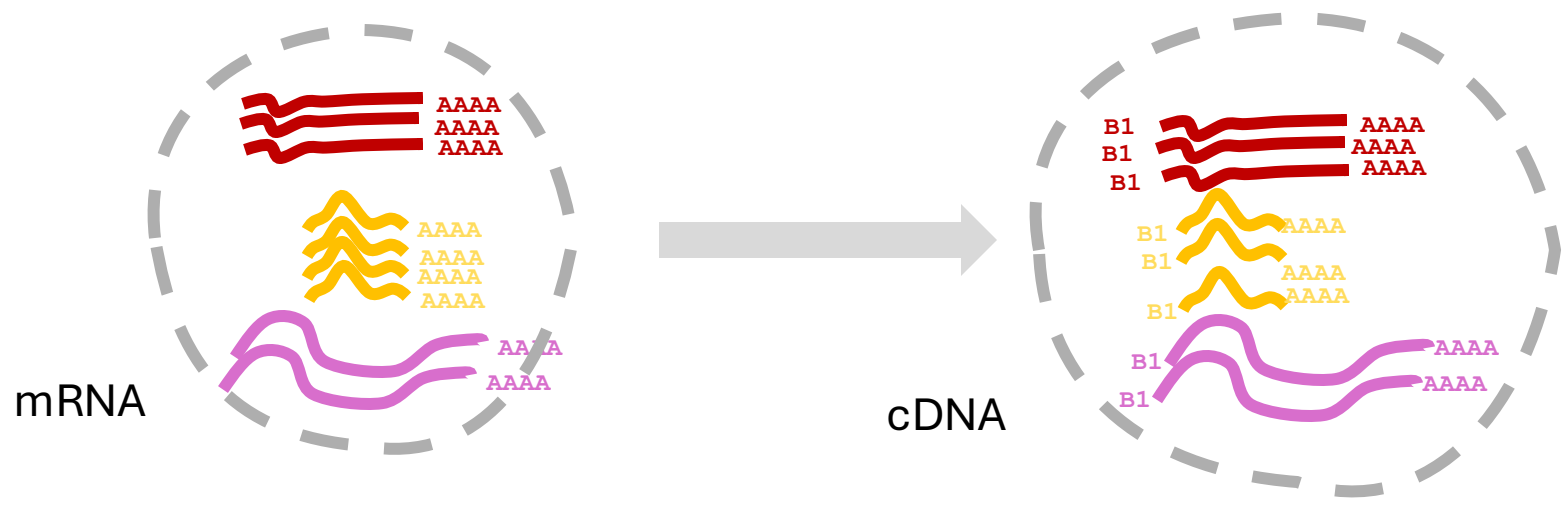
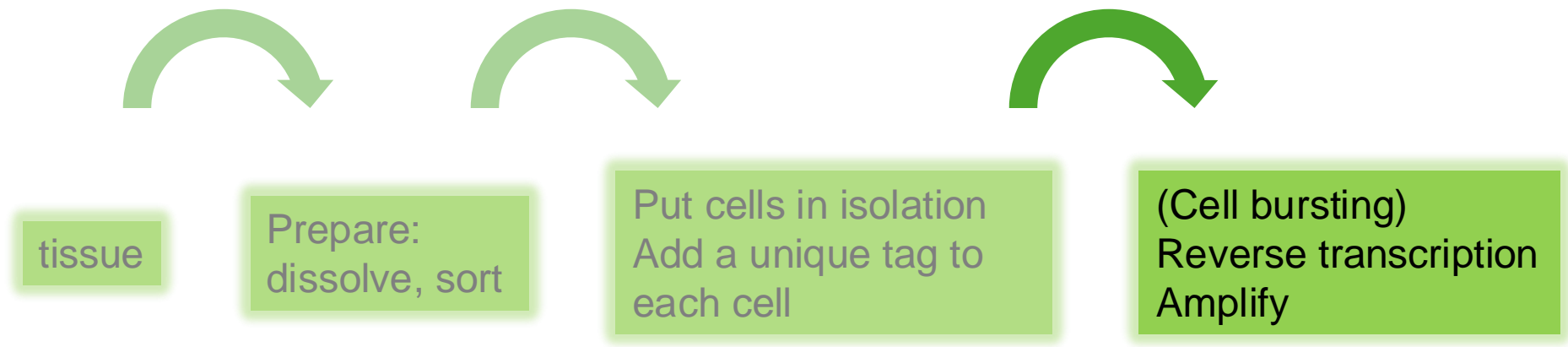
## principles



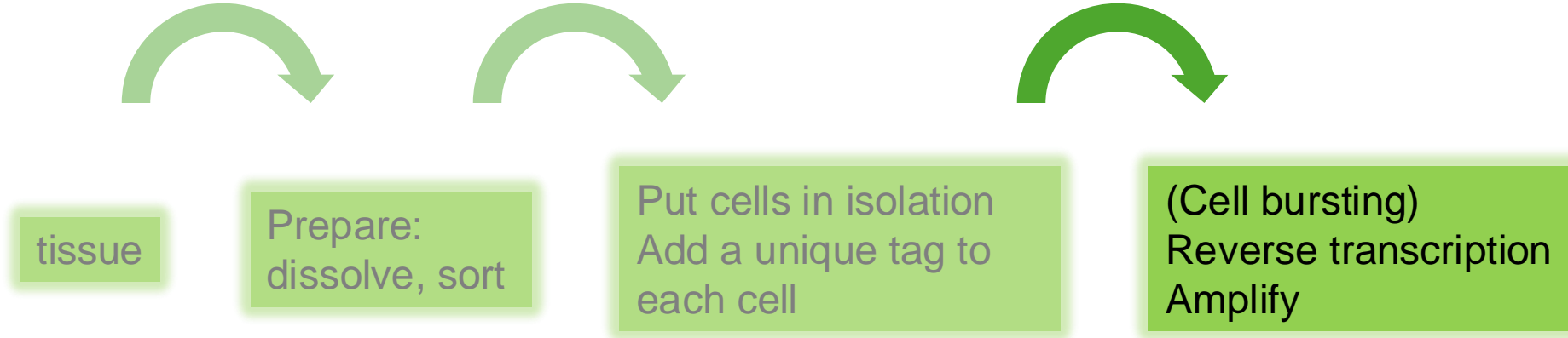
# principles



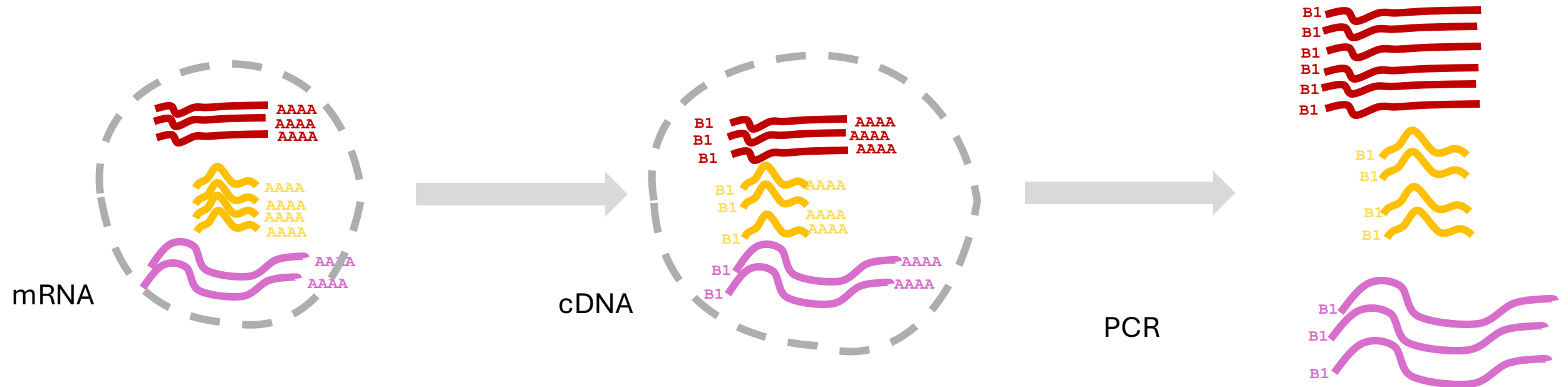
# principles



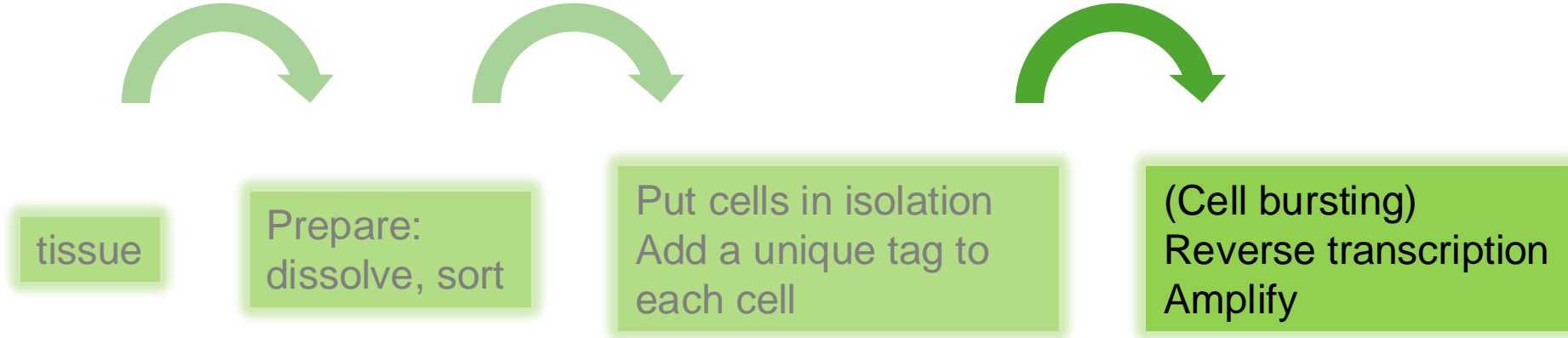
## principles



- Each cell has very little mRNA → amplification
- PCR bias → UMI included when cDNA constructed



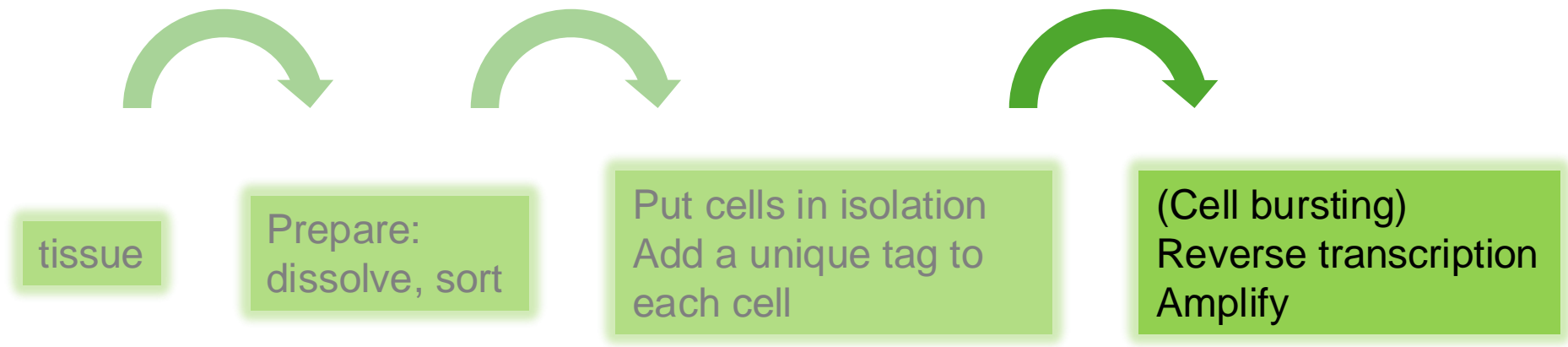
# principles



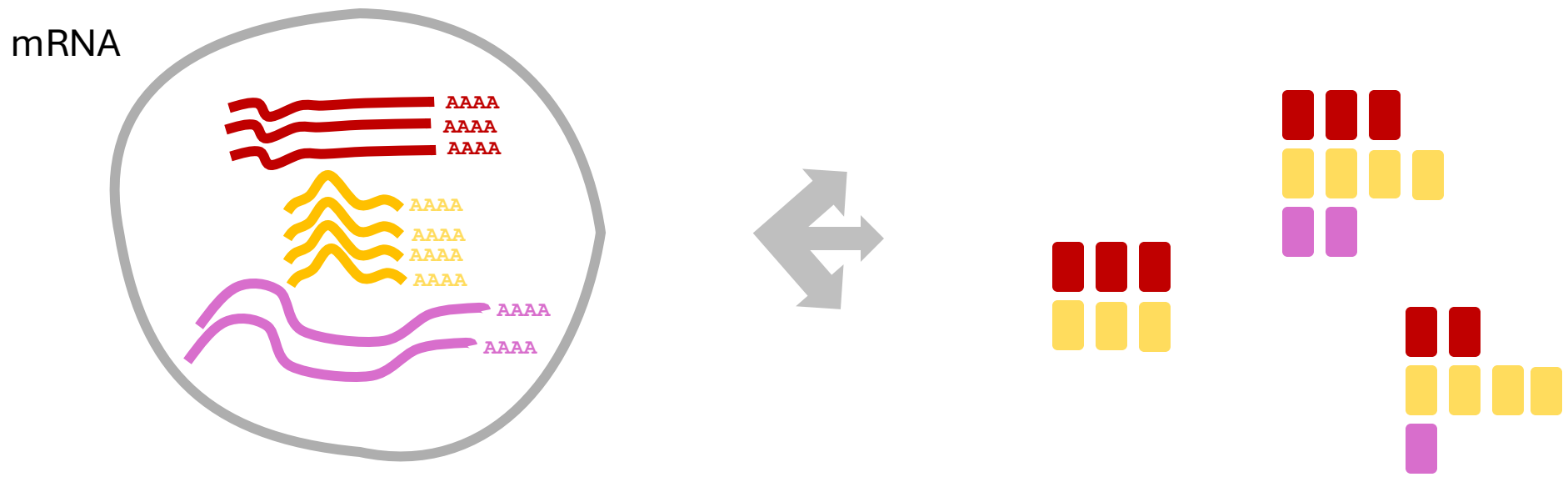
- Each cell has very little mRNA → amplification
- PCR bias → UMI included when cDNA constructed



# principles



- Drop outs → sparsity of the data



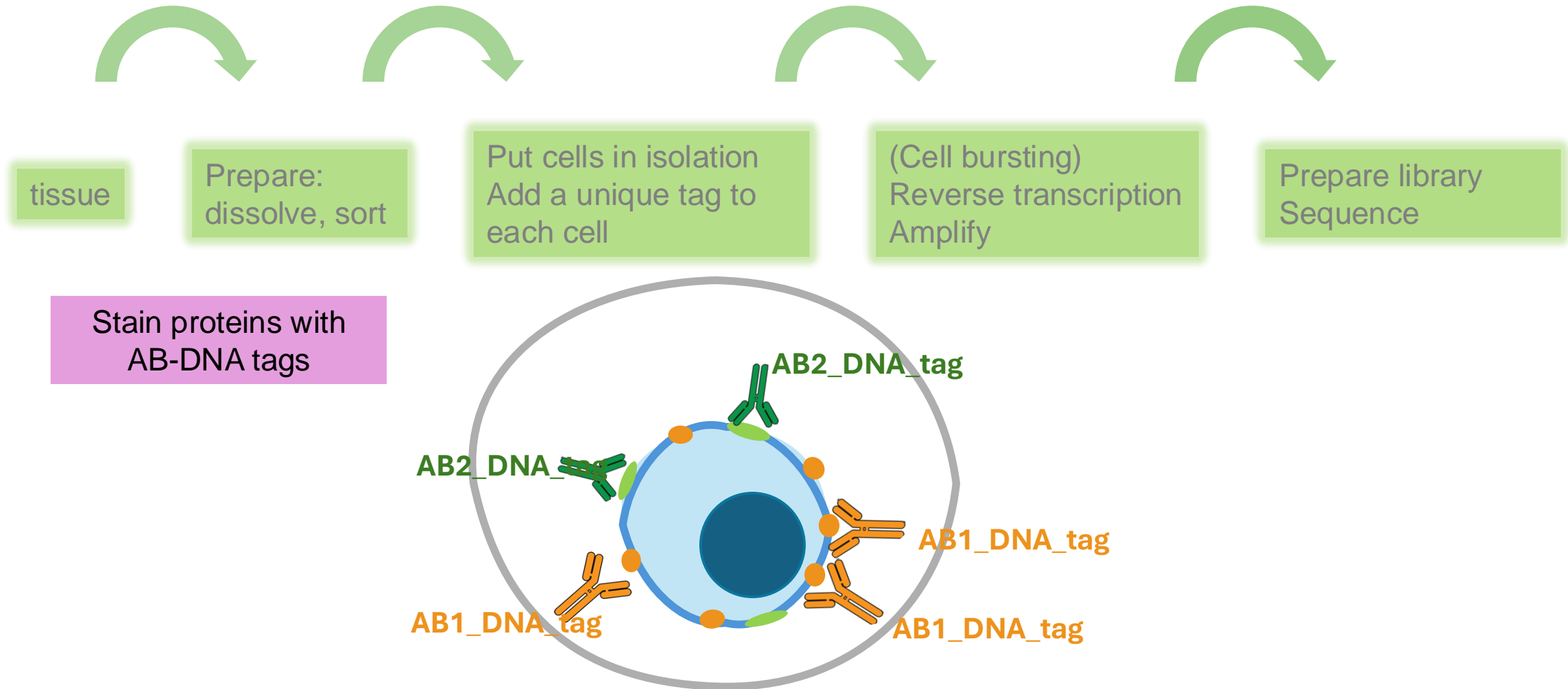
## principles



Many platforms



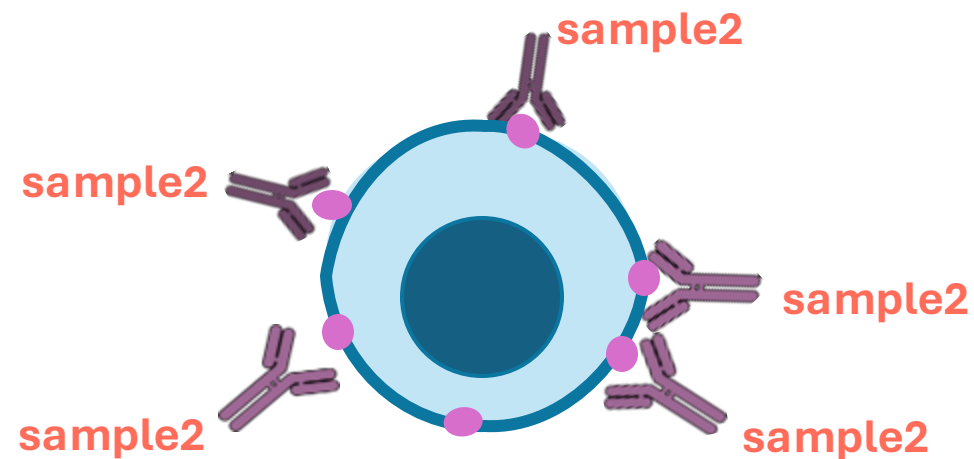
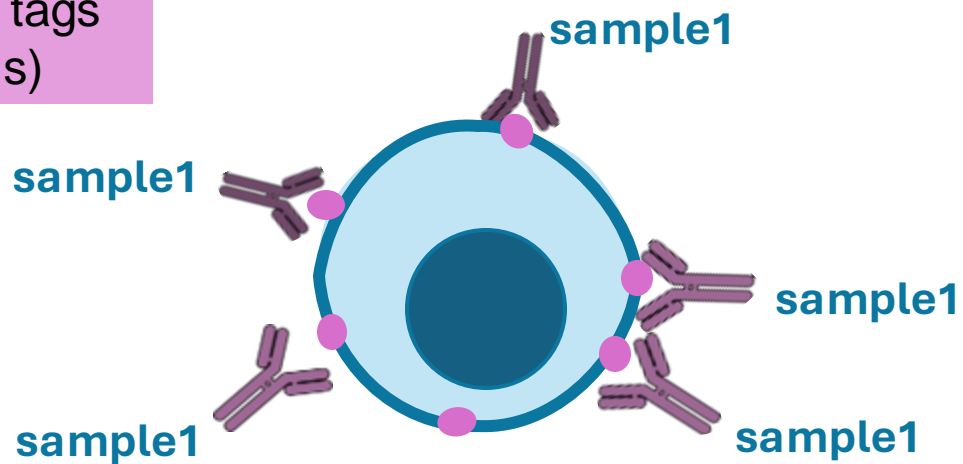
# principles



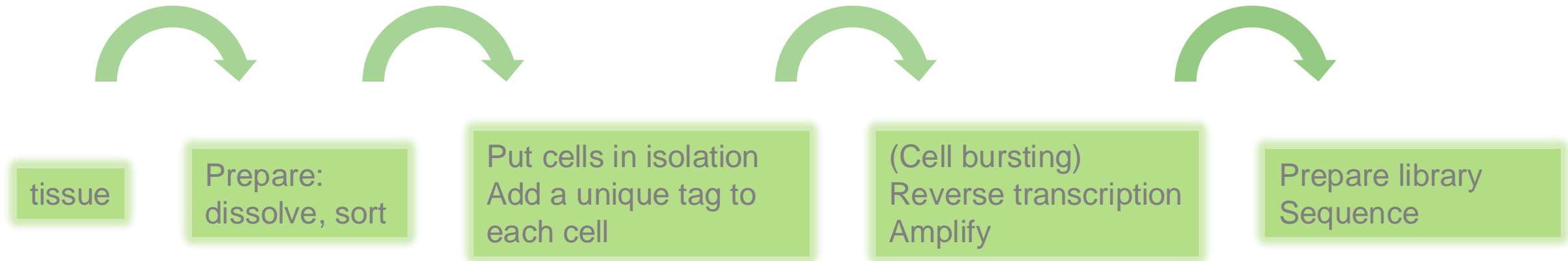
# principles



Sample tags  
(#tags)



## principles

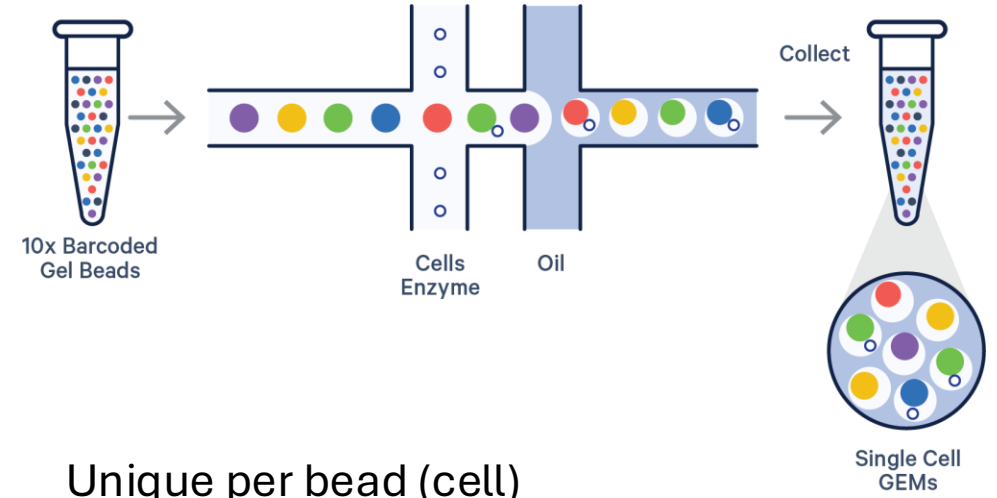


Example of specific platform



# principles 10x

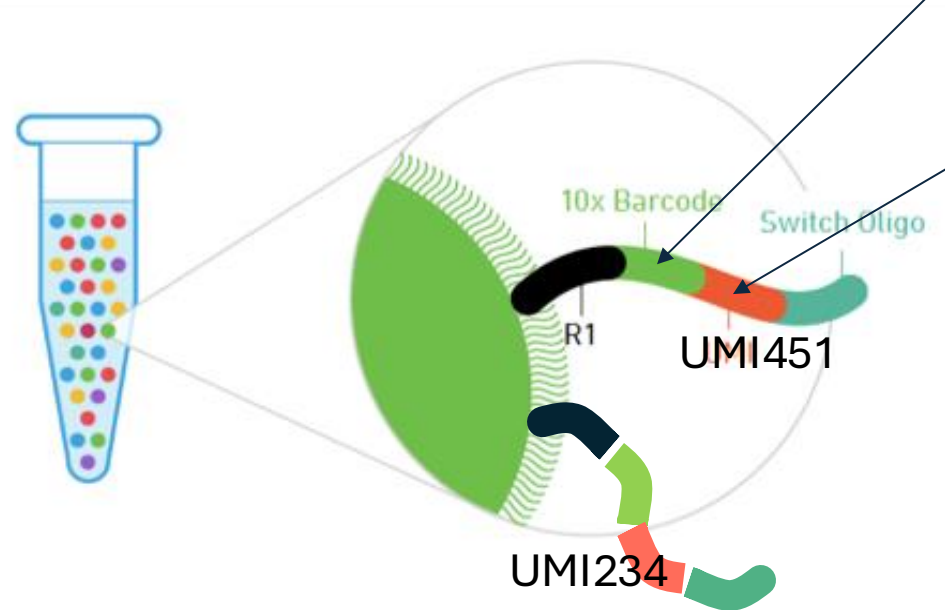
## 10x Next GEM Technology for Single Cell Partitioning



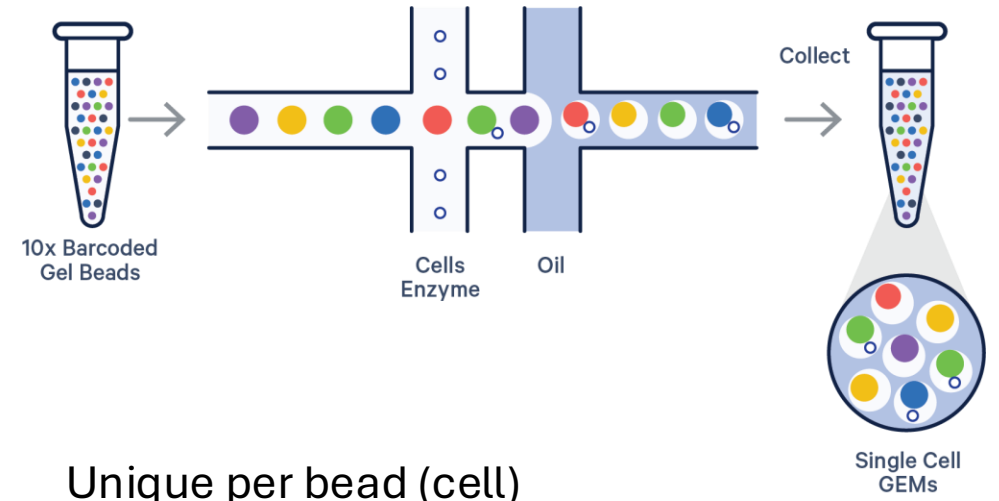
Unique per bead (cell)

Unique per oligotail: UMI

5' GEX +  
CITEseq



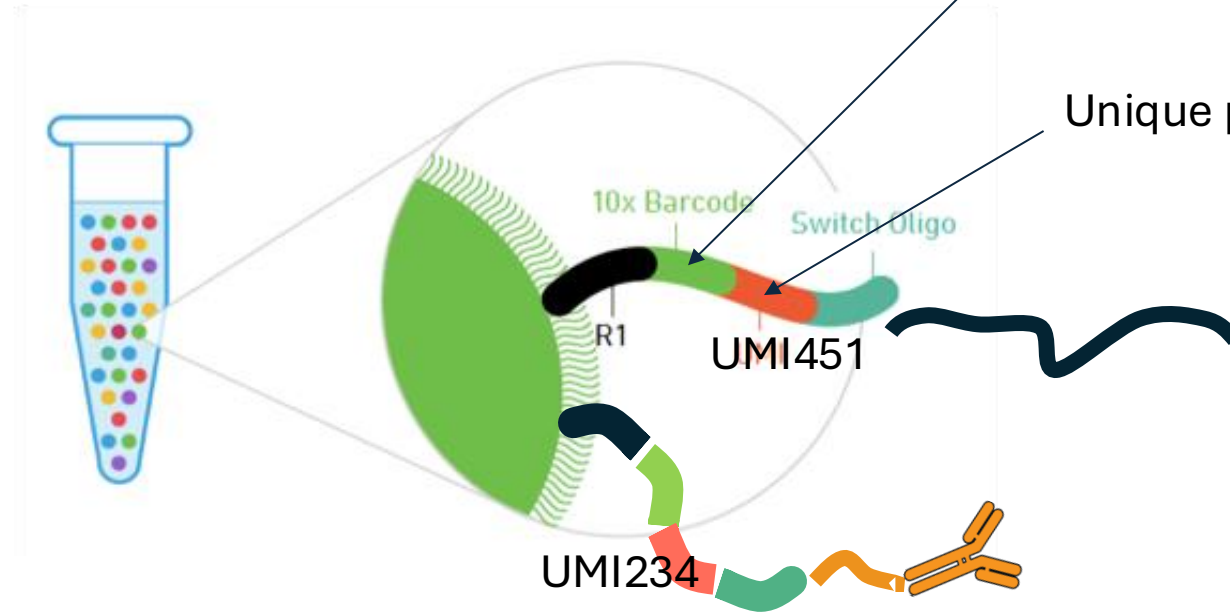
# principles 10x



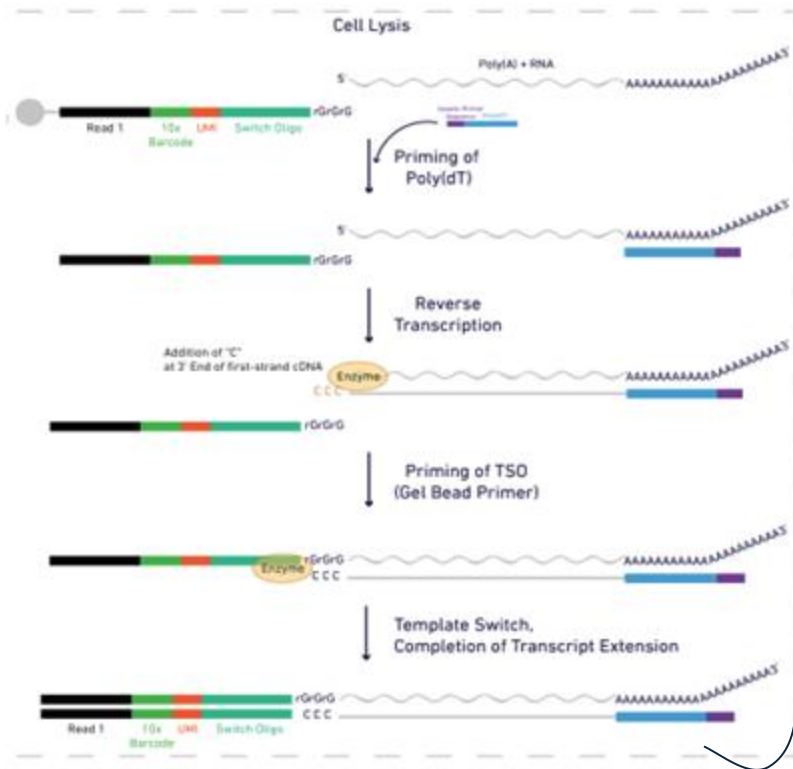
Unique per bead (cell)

Unique per oligotail: UMI

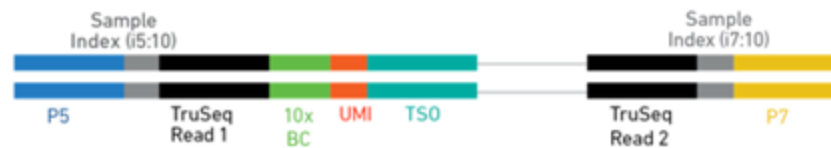
5' GEX



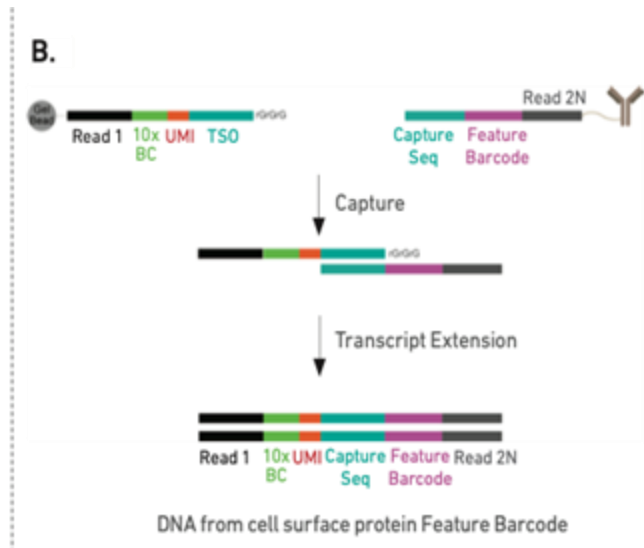
## Inside individual GEMs



## Pooled cDNA processed in bulk

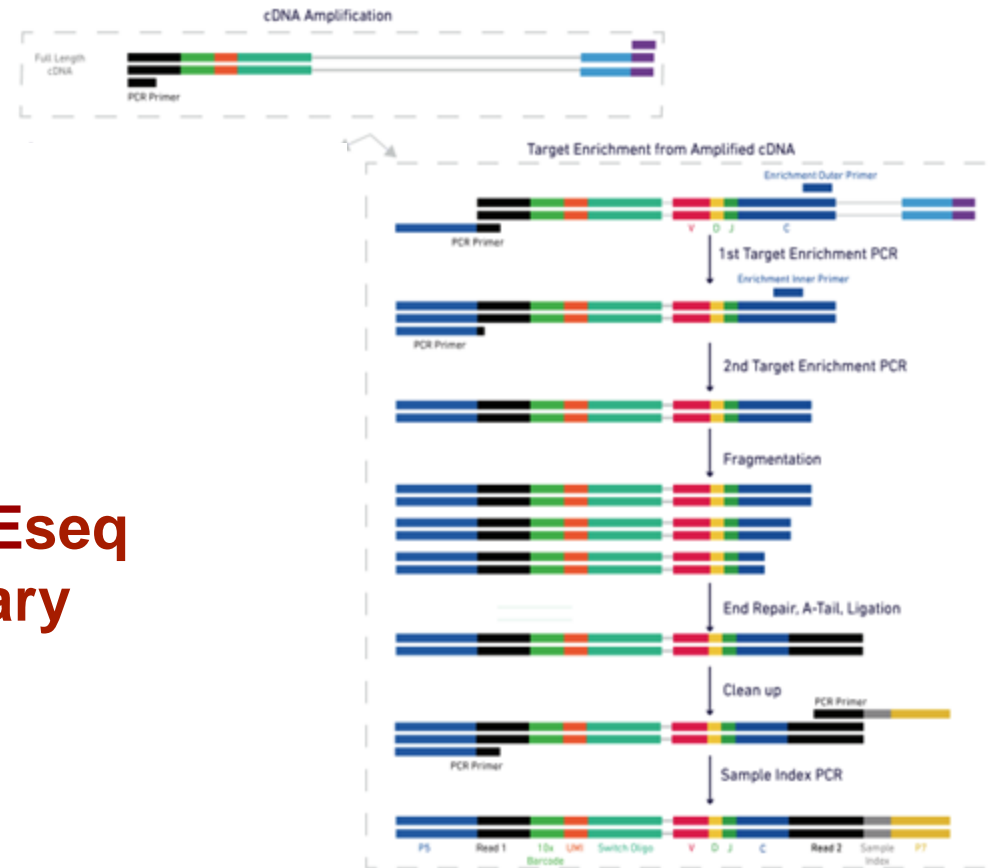


## Inside individual GEMs



**CITEseq  
library**

## Pooled cDNA processed in bulk



**VDJ  
library**

How much to sequence? ← Depends on the technology & desired resolution

### scRNAseq

- 10-15K reads per cell
- 50K reads per cell
- Saturation: the bottleneck is the amount of starting RNA
- Optimise: more transcripts, more cells, more biological replicates

### CITEseq cell hashtags

- Less complexity
- Small pool of possible sequences
- Few types of reads per cell (depends on the number of Abs)

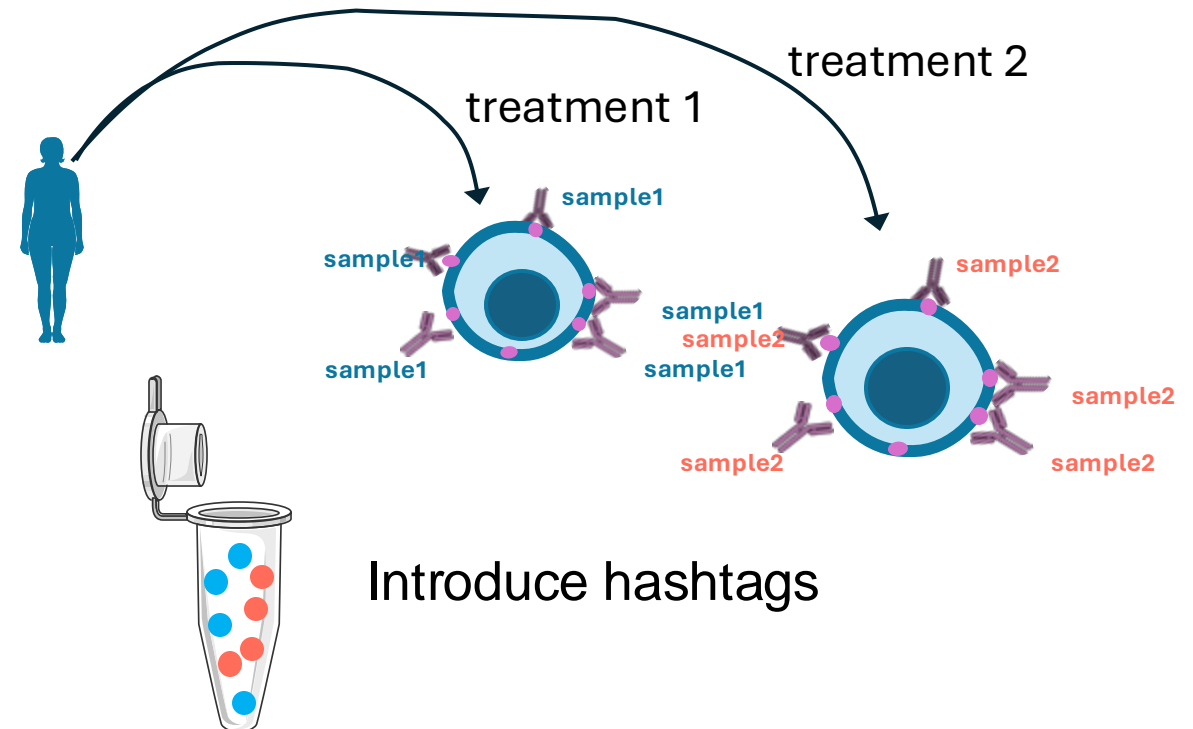
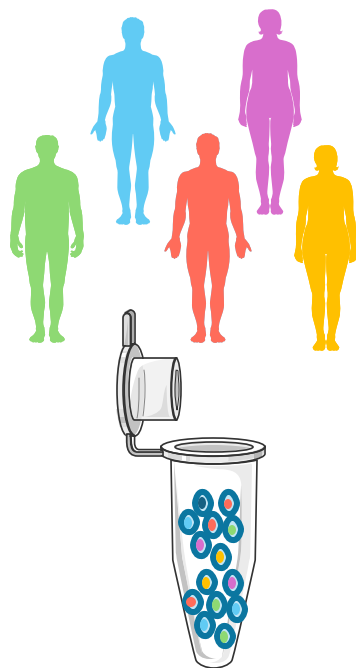


questions

Pool as much as possible, as soon as possible

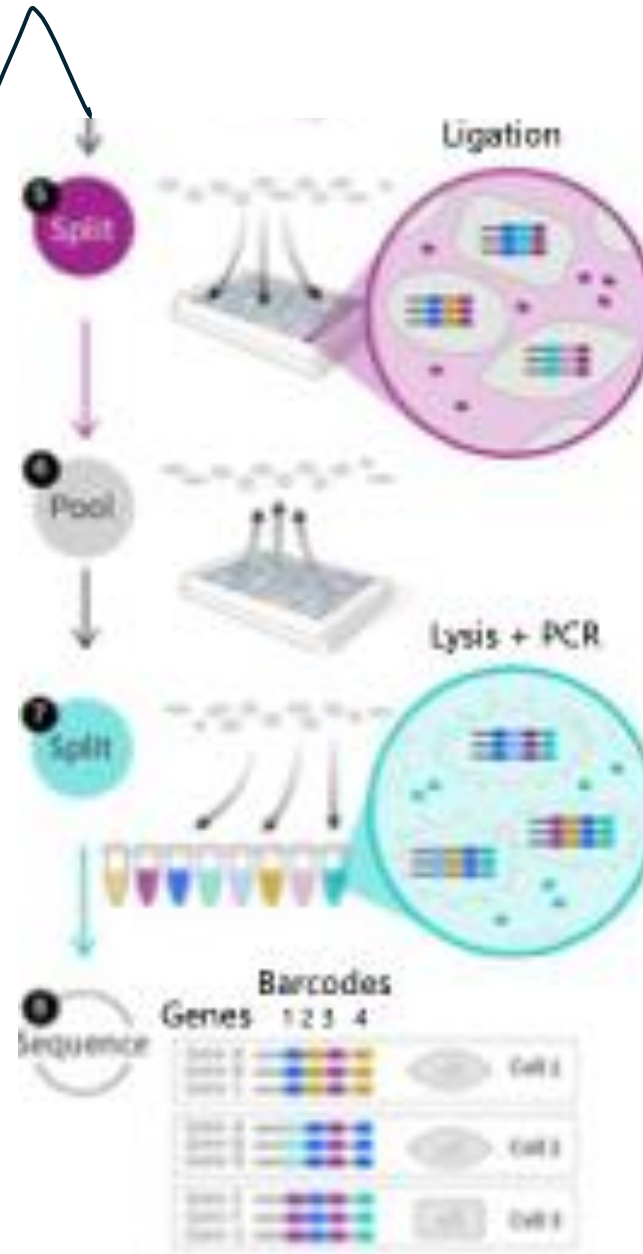
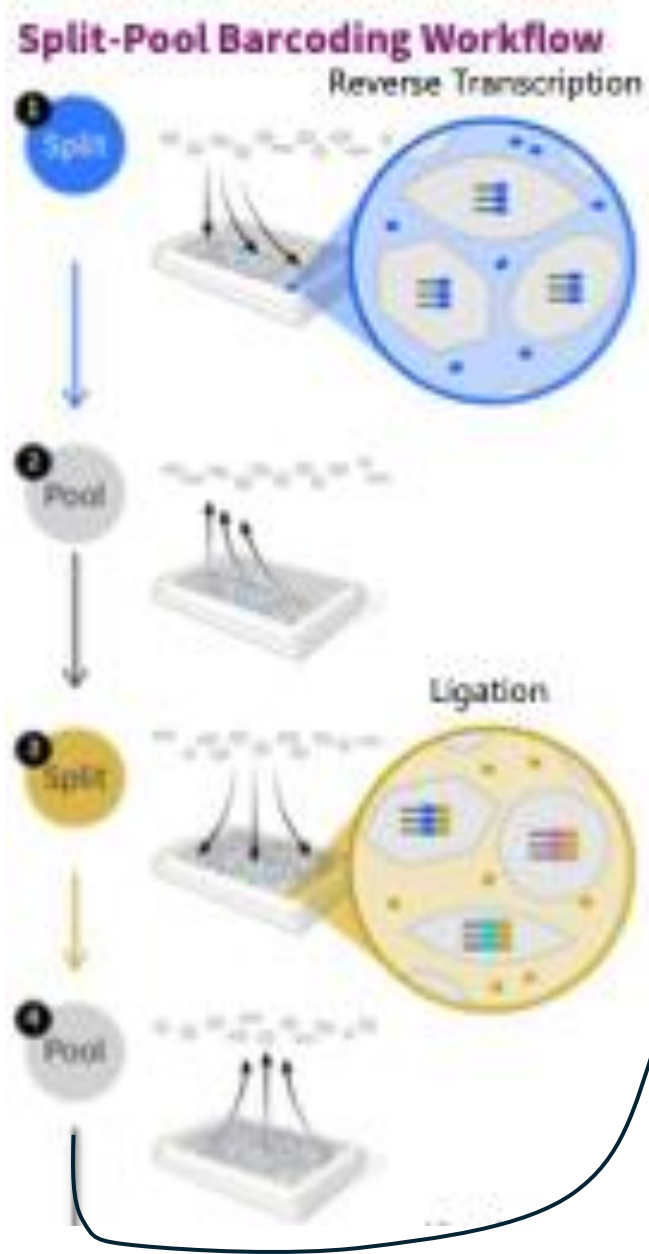
- Reduce batch effect
- Identify multiplets
- Pool samples to reach minimal cell numbers

Use  
samples'  
genotypes



# principles

Fixed cells



PARSE  
SCALE Bio

# From FASTQ to expression matrices

## FASTQ files

```
@A00815:607:H5V32DRX3:1:2101:1108:1000 2:N:0:ACAATGTGAA+TAACGGTACG
GAACCACTGAGCACAAGTTTCTTCATCGTTCTCAGATTAGTAACATTATTAATTTTAGACAATCCCGTGAAGGCCAATTCATCAGTG
+
FF:FFFF:F,:FFFF,F,:F,FFF,FFFFFFFFF,,, :F::,FFFFFFFF,FF,FF,FF::,:FFF,,F,:FFF:F,FF,F:F:FF,F,
@A00815:607:H5V32DRX3:1:2101:1181:1000 2:N:0:AAAATGTGAA+TAACGGTACG
CGTCAAGGCAGGGCCTCGTTCTTGCTGGGCACCAAGAGCAGATGACATATATAGCACAGTGCCTCCCCCAGGACAGGAAGATGAGGCTAG
+
F,FFFF,FF,:FF,FF,,,F:F::,FFF,F:FF:FFF,F,F,,,FFF,,FFF:::F:FFF,FF,FFF:::F:F:,FF,,FF,,,FFFF:,
@A00815:607:H5V32DRX3:1:2101:1597:1000 2:N:0:ACAATGTGAA+TAACGGTACG
GAAGTGAGGGATGCTGAGGGCCGGGACAAGCTATCGGACTGTCTGCTGCCATCGGTAATGAGTCTCAGTAGACCTGGAACGTCACCTCGC
+
FFFFFFFFFFFFFFFFFFFF,FFFF:FFFFFFFFFFFFFFFFFFFF,FFFF,:FFF,:FFFF,FFFFFFFFFFFFFFFFFFFF:FFFFFF:
@A00815:607:H5V32DRX3:1:2101:1615:1000 2:N:0:ACAATGTGAA+TAACGGTACG
ATGTGACTATAGGCTCATAGCCATCTCATTATGCAAAATGTATTCACTGTCTTTGTATGTCTCAATAGTCTCCAGATATACGGCGGT
+
,,, :FFFF,,FFFF,FFFF,F:F,F,FF,FF,:FFF:,F,:F,FFF,:FF,::::F,F,FF,,F,F,FFFFFF,FF,,F,F,F,:,,F:
```

- Done by proprietary software/free alternatives available
- Basis for any downstream analysis
- Some steps might be improved with additional community-developed software

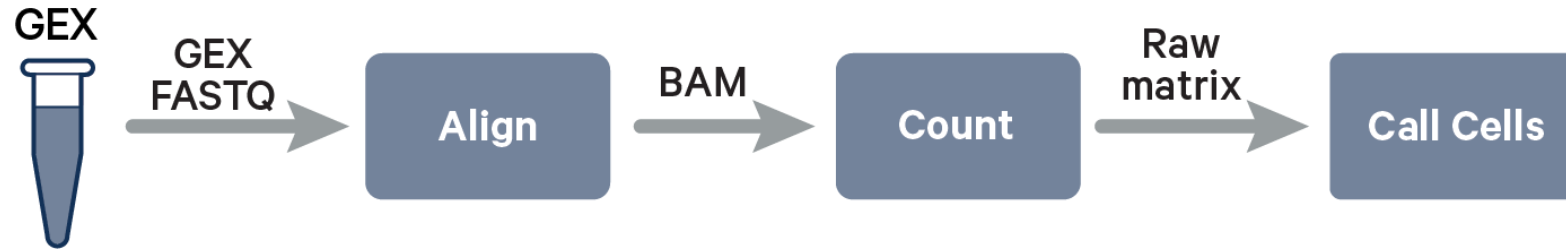
	cell1	cell2	cell3
gene1			
gene2			
gene3			

	cell1	cell2	cell3
abs1			
abs2			
abs3			

	cell1	cell2	cell3
#1			
#2			
#3			

**10 X Cell Ranger** A set of analysis pipelines that perform sample demultiplexing, barcode processing, single cell 3' and 5' gene counting, V(D)J transcript sequence assembly and annotation, and Feature Barcode analysis from single cell data.

## From FASTQ to expression matrices



- Reads with valid barcodes
- Map to a genome
- Exclude/annotate intergenic reads
- Align to a transcriptome
- Transcriptomic reads
- UMI correction
- exclude low support molecules/resolve duplicated molecules (barcode+UMI+gene)

	barc1	barc2	barc3
gene1			
gene2			
gene3			

Cell Ranger

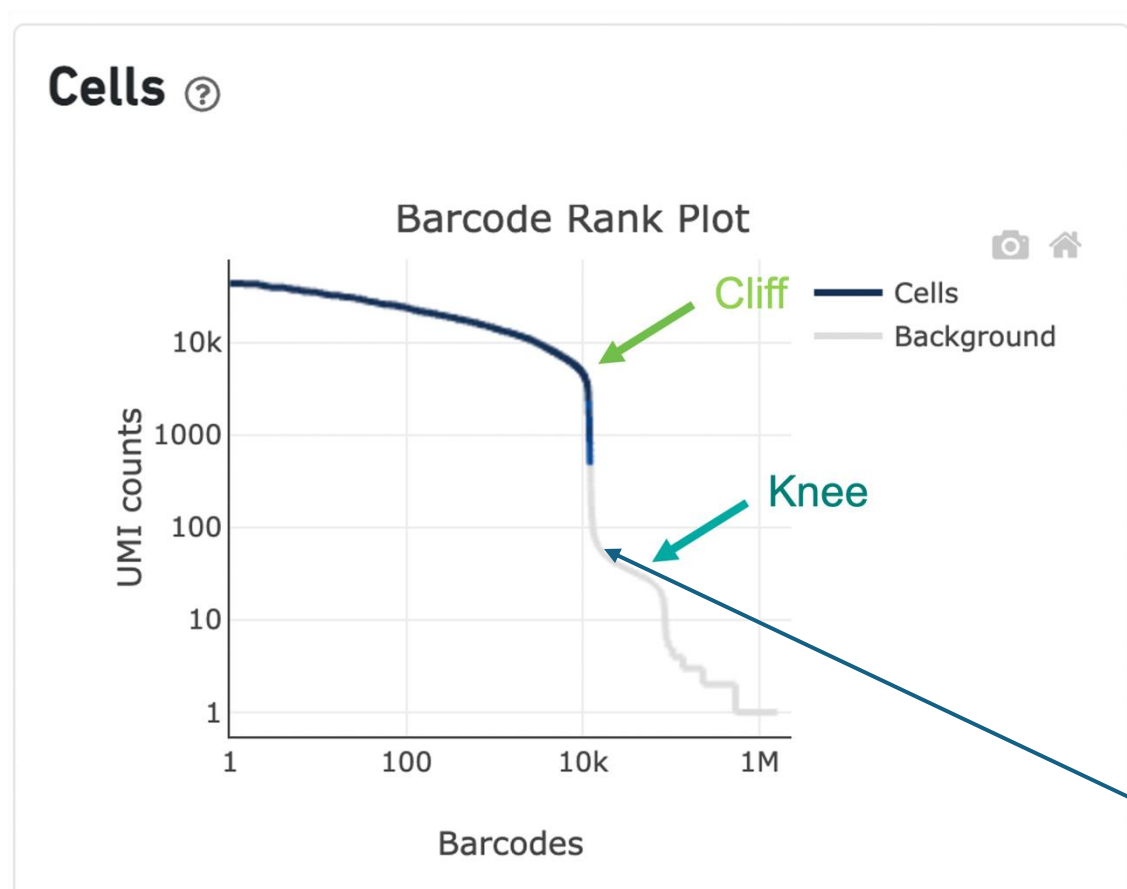
## From FASTQ to expression matrices

Is barcode a cell or an empty droplet?

→ Order barcodes by UMI counts

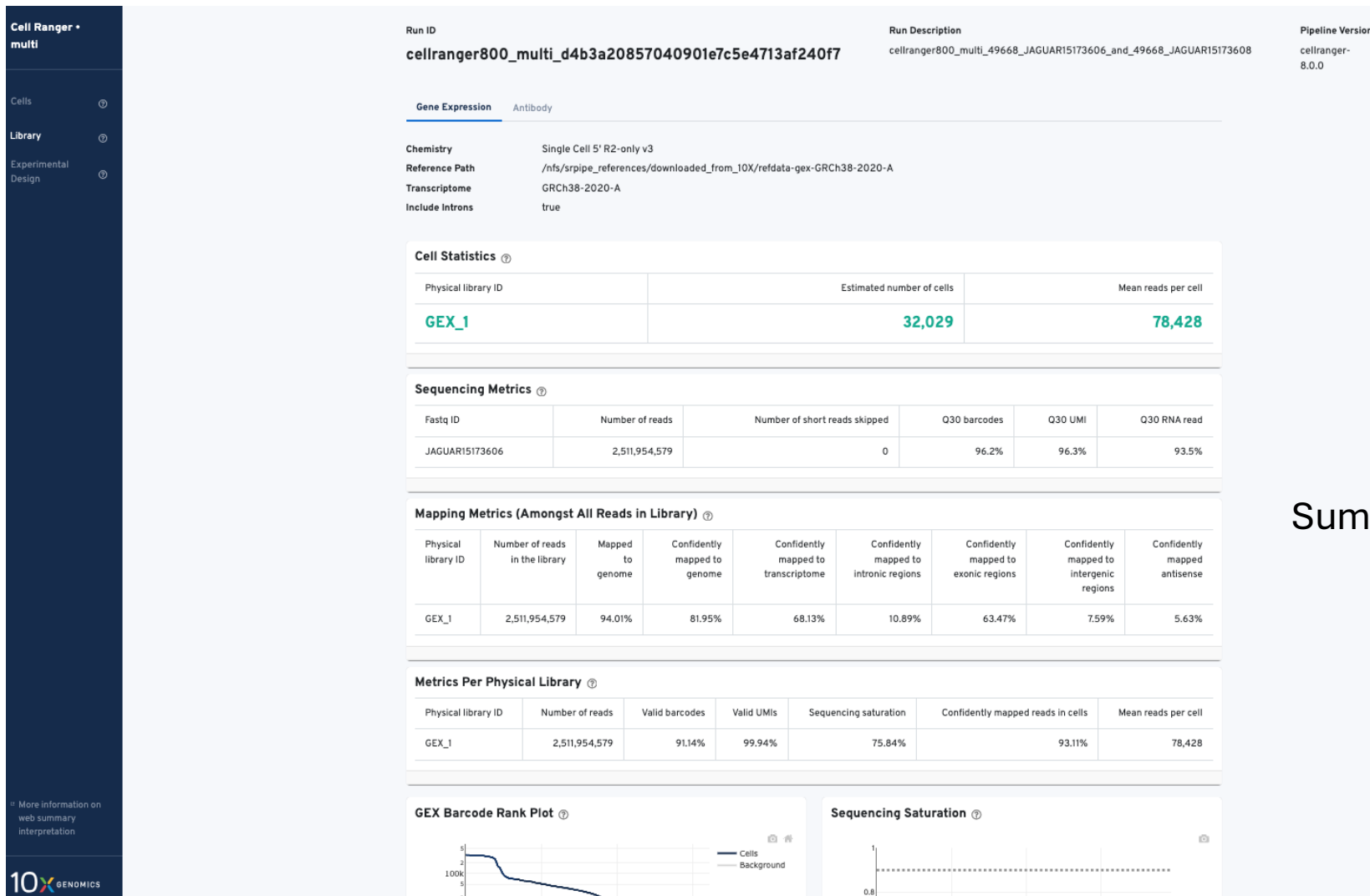
→ Barcode rank plot – identify the cutoff

	barc1	barc2	barc3
gene1			
gene2			
gene3			
Total			



Cell Ranger

# From FASTQ to expression matrices



Summary of reads/transcripts/cells...

## Expression matrices

	cell1	cell2	cell3
gene1			
gene2			
gene3			

	cell1	cell2	cell3
#1			
#2			
#3			

	cell1	cell2	cell3
abs1			
abs2			
abs3			

Not ready for comparing gene expression etc. yet!

Demultiplexing – if cell hashtags  
Demultiplexing by genotypes

Cell-level QC and removal of bad quality cells

Additional cleanup – removal of ambient RNA  
Alternative cell calling  
Additional cleanup of abs-based data



questions

Collab notebooks

### Demultiplexing of cell hashing

- Within Cell Ranger
- Might be done also from Seurat/scanpy environment
- Per barcode signal from all possible hashtags
- Winning hashtag per cell/if too mixed → a multiplet

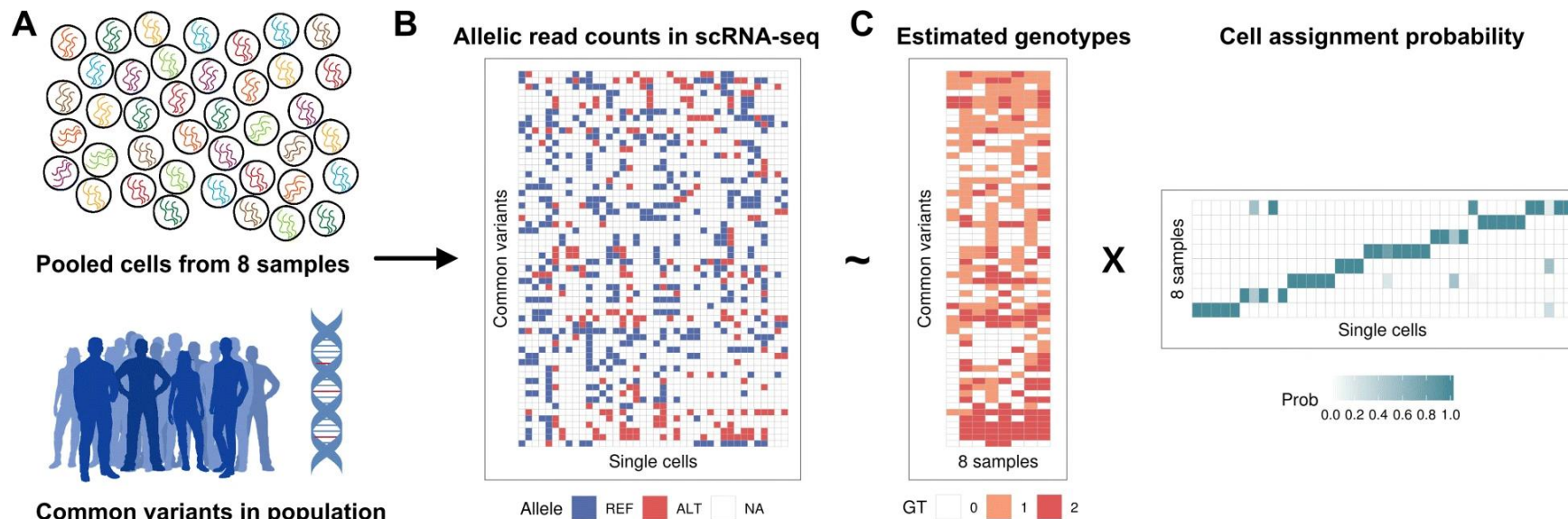
	#tag1	#tag2	#tag3	#tag4
Cell1	<b>4071</b>	4	5	1
Cell2	3	<b>2380</b>	16	2
Cell3	<b>1341</b>	5	21	7
Cell4	5	4	<b>5434</b>	8
Cell5	4	2	6	<b>1474</b>
Cell6	4	<b>1203</b>	3	<b>1020</b>

Seurat

Scanpy

CellRanger

## Better expression data



Demultiplexing by genotypes

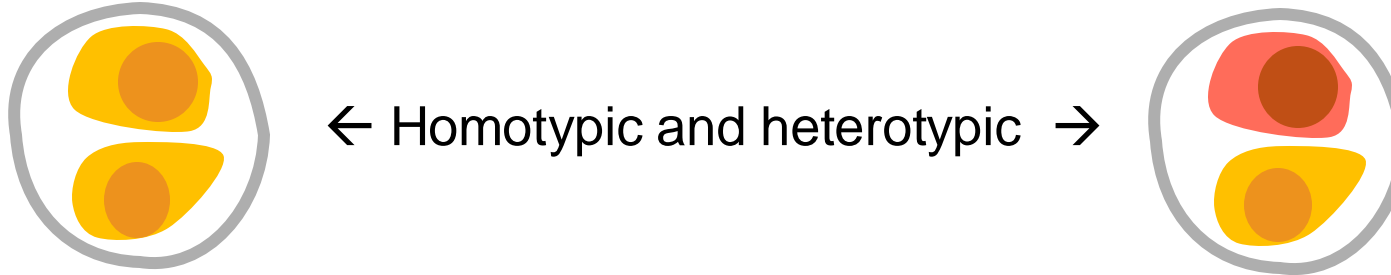
If genotypes available +++

**Demuxafy**

For each barcode: SNP calling from BAM files  
Minimal input: a reference list of SNPs (eg 1KG data),  
expected number of individuals

Neavin, D., Senabouth, A., Arora, H. *et al.* *Demuxafy*:  
improvement in droplet assignment by integrating multiple  
single-cell demultiplexing and doublet detection  
methods. *Genome Biol* **25**, 94 (2024)

## Multiplets



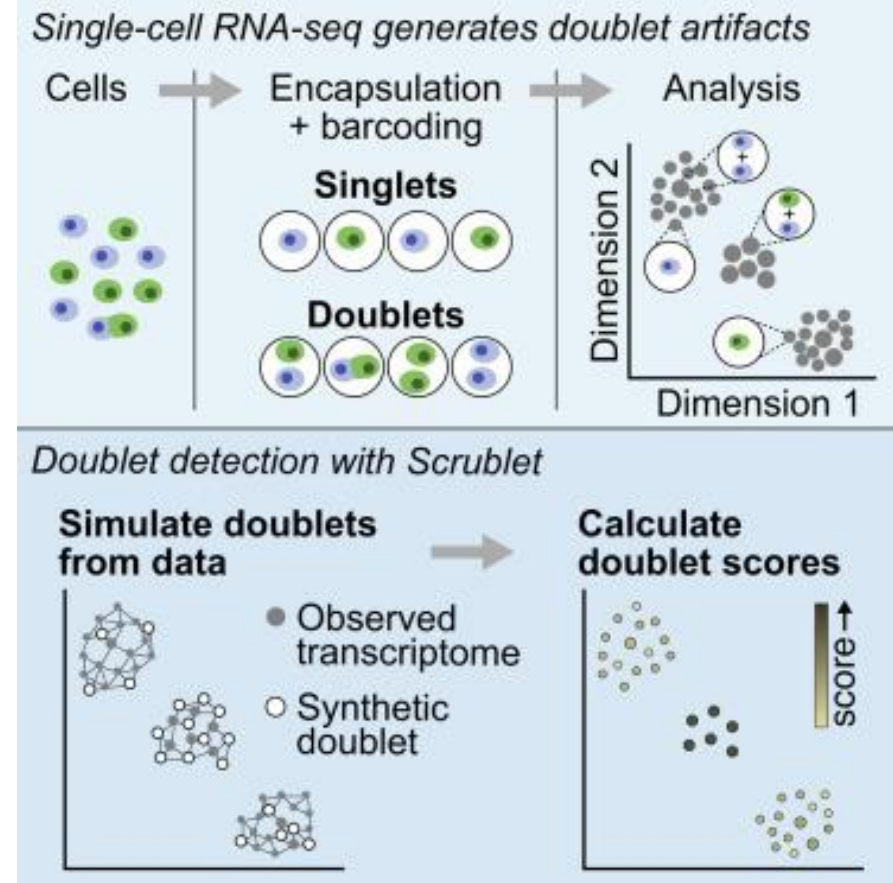
## How to identify

- Genotypes
- Hashtags
- Extremely high number of expressed genes/UMI count
- Impossible biology
  - TCRs
  - Mutually exclusive proteins/genes CD3, CD19, CD14, CD16

Better expression data

Modelling artificial doublets from the actual experiment data, using this to define doublets

Demuxafy

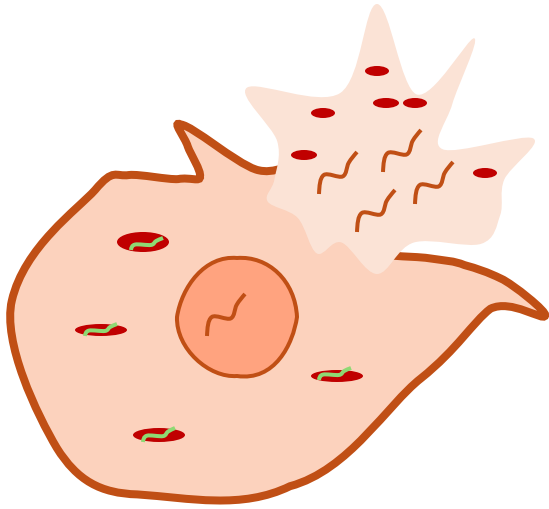


Seurat

Scanpy

Standalone software

Damaged, early burst cells

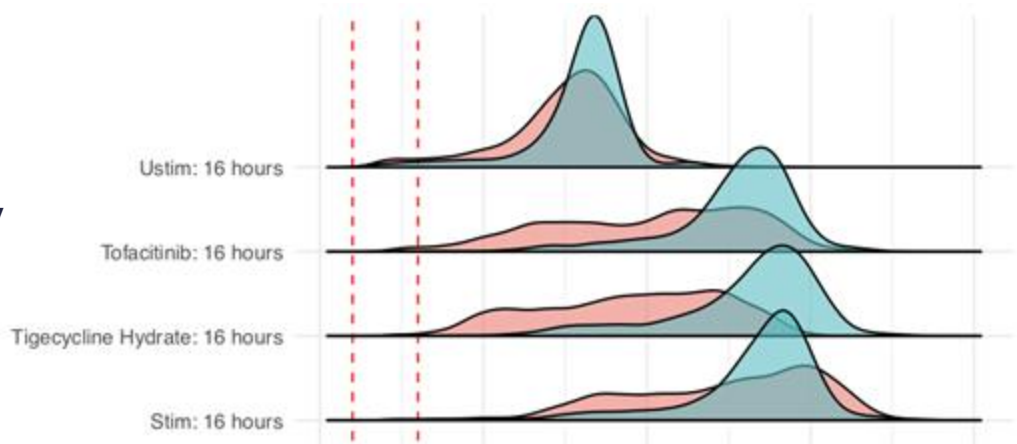


Relatively

- More of mt transcripts
- Less of cytoplasmic transcripts (ribosomal RNA)
- More of nuclear transcripts (PBMCs: lncRNA MALAT)
- Lower number of detected genes
- Lower total number of transcripts

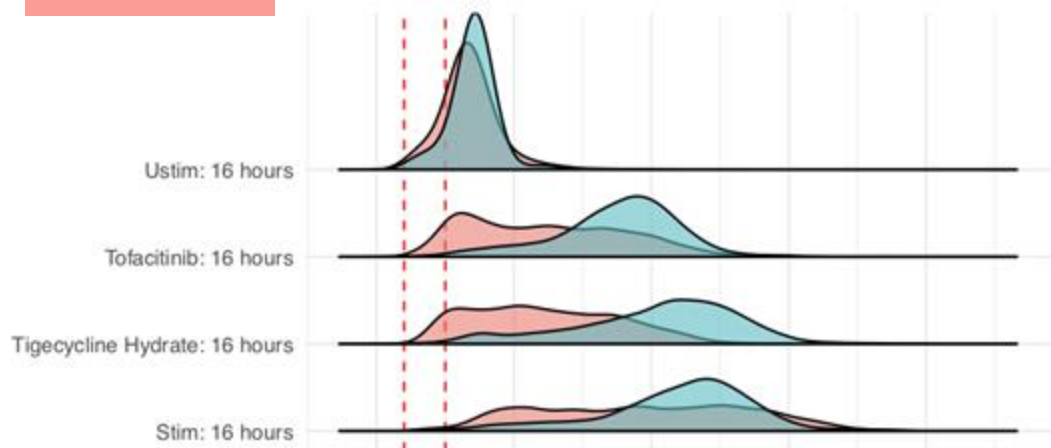
Better expression data

Expression intensity  
(UMI/cell)

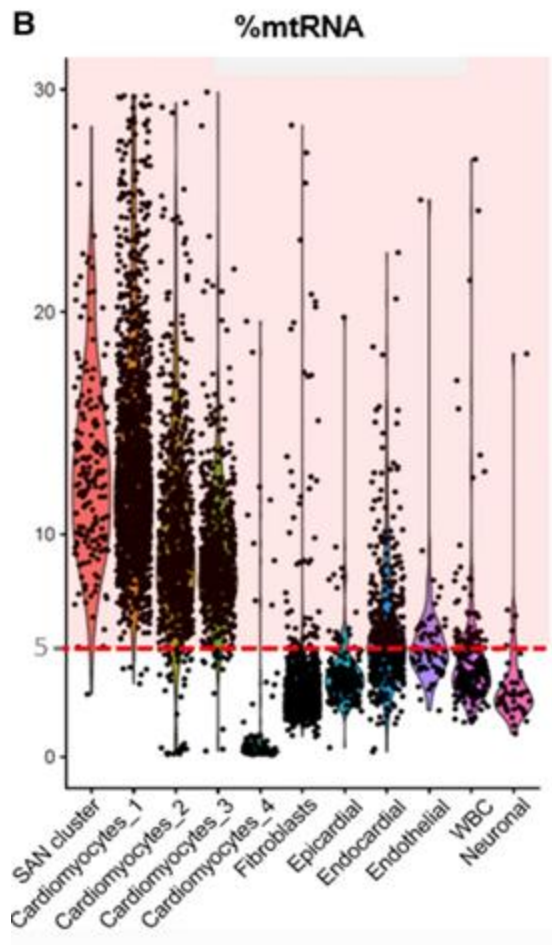


CD4 naive  
CD4 mem

Expressed genes  
(genes/cell)



plots: Ziyang Ke



AM Gallow, 2021

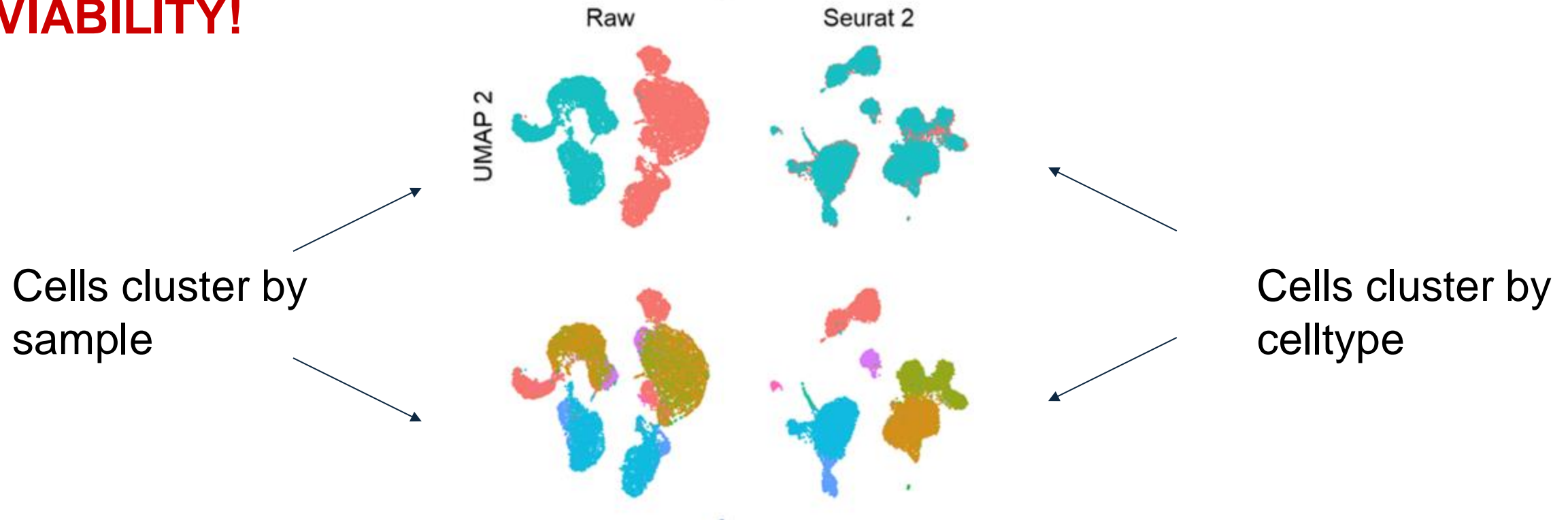


questions

end

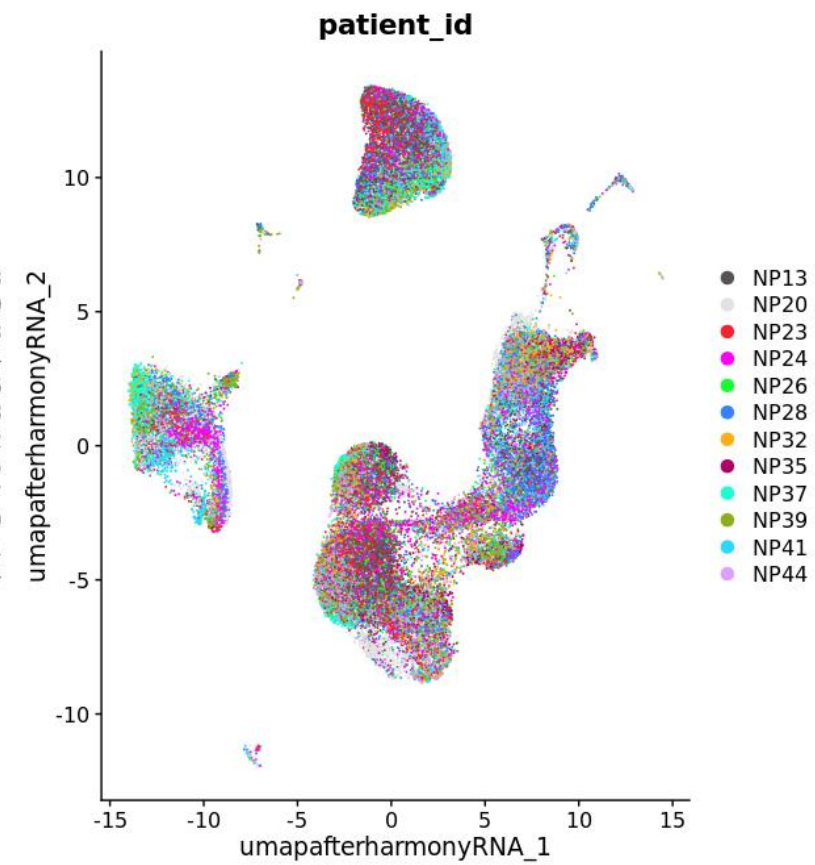
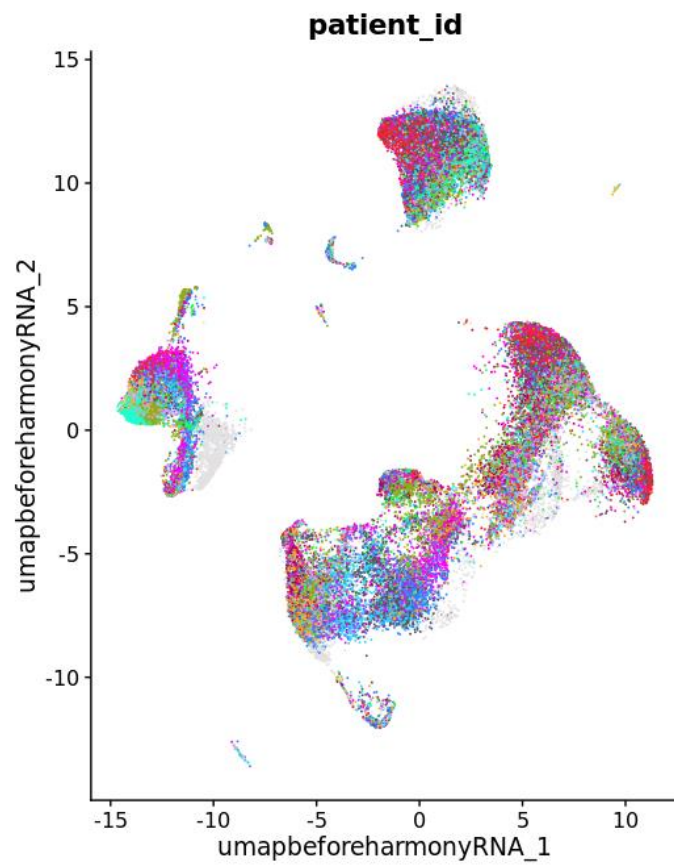
# Batch effects

**VIABILITY!**



Tran et al. Genome Biology (2020)

Seurat/Scanpy/etc, many methods



# QC and filtering → post-CellRanger

## Goals:

Exclusion of non-cells

Exclusion of dying cells

Exclusion of doublets

Debris - high background

**VIABILITY!**

Iterative process: QC filter → normalise → UMAP, annotate → redo  
QC filter → normalise → UMAP, annotate

Seurat/Scanpy/etc, software for doublet detection

# UMAP, cell annotation, DGE

## Goals:

Representation of relationships between cells in 2D

Identification of similar cells

Cell annotation

Downstream analysis

DimRed and clustering are dominated by the genes which vary most between cells → one can use just these genes

Genes to exclude from the analysis: TCRs, BCRs

Seurat/Scanpy/etc