SincellTE 2024

Aziza CAIDI

# Mapping, quality control and quantification



GUSTAVE / ROUSSY
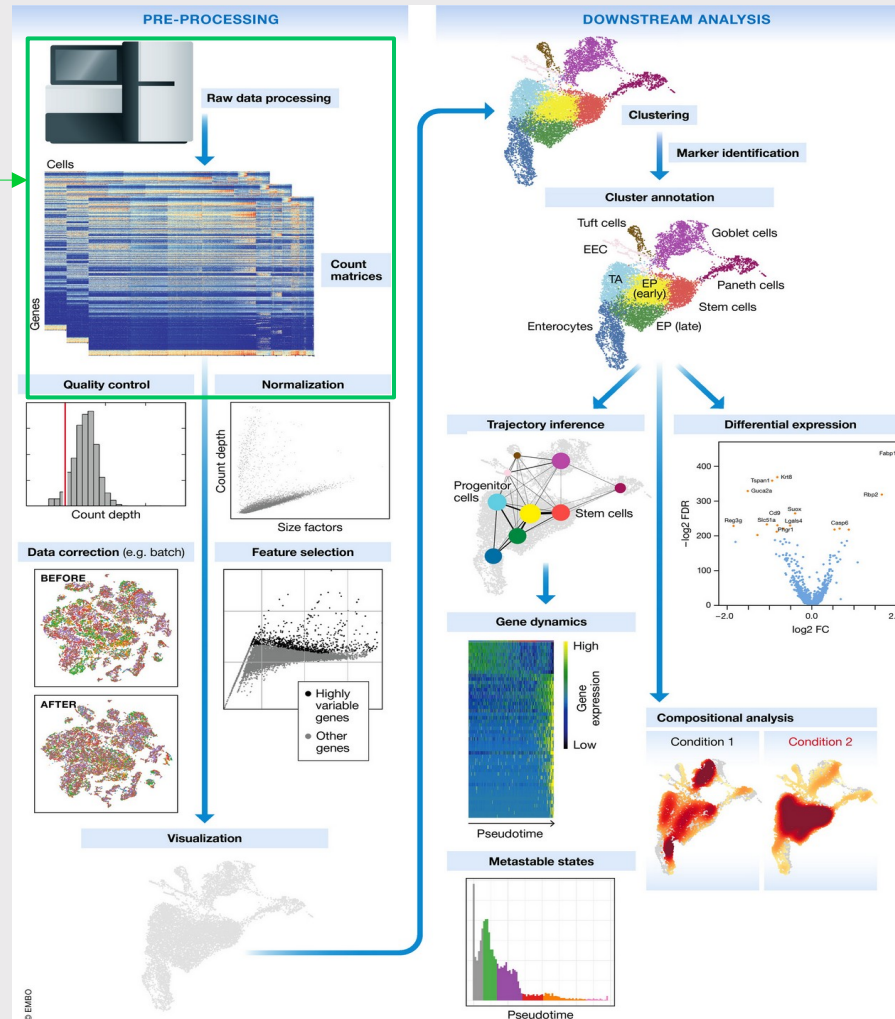CANCER CAMPUS
GRAND PARIS

# Main steps of single cell data processing



From Luecken and Theis, Mol Systems Biology 2019

# Main steps of single cell data processing
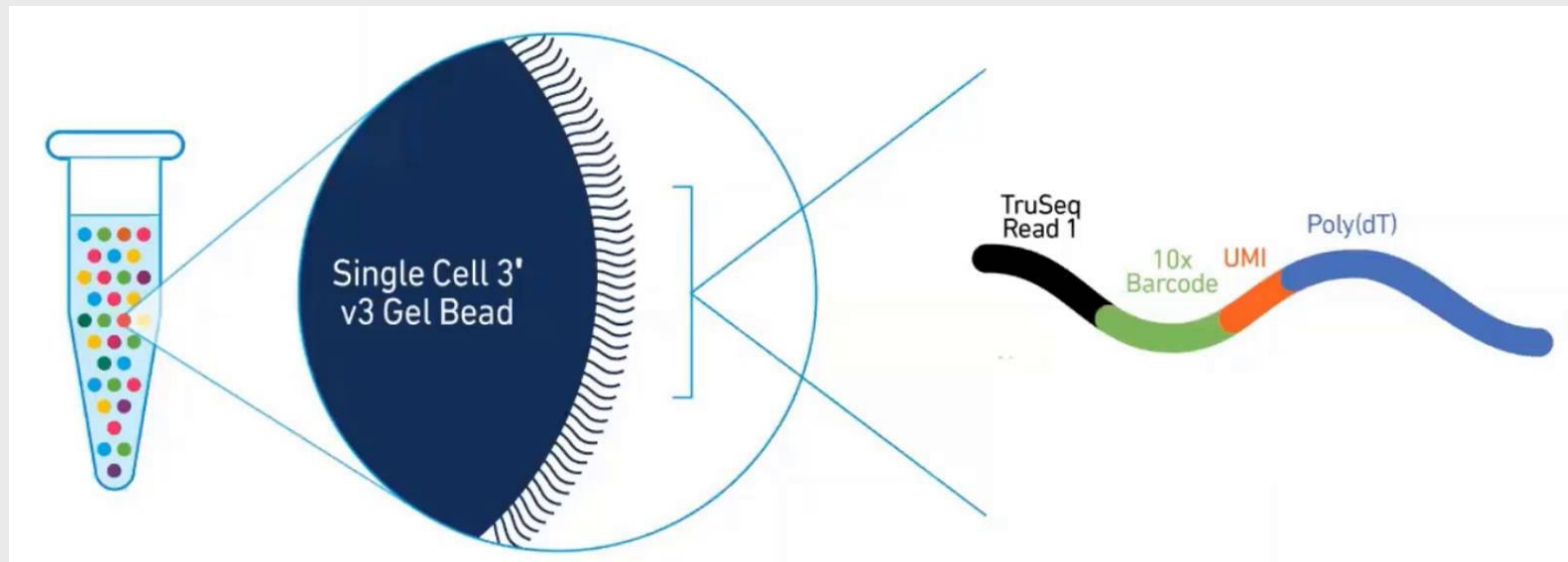


*this course*

From Luecken and Theis, Mol Systems Biology 2019
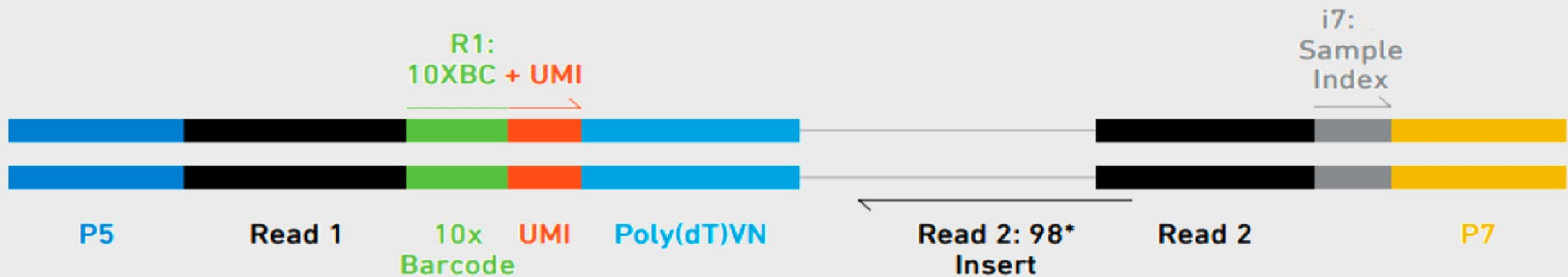
# The starting library

We will use a droplet-based library as an example.

# The starting library

We will use a droplet-based library as an example.



Read1: unique cell barcode (x nt) + UMI (y nt)

Read2: RNA 3' sequence

I7: sample index: determines which sample the read originated from

Cellular barcode: determines which cell the read originated from
Unique molecular identifier (UMI): determines which transcript molecule the read originated from

# Plan

- Demultiplexing: generating fastqs from bcl

- Quality Check

- Generating a gene x cell count matrix

# Demultiplexing

Convert BCL files (sequencer output) to fastq files
Most used tool : 10X's cellranger mkfastq a wrapper around
bcl2fastq

- Usual sample sheet

- You must know :
  - i7 (i5) index sequence
  - R1 and R2 lengths
  - (depends on technology, version...)

|  | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | [Header] | | | | | |
| 2 | IEMFileVersion | 5 | | | | |
| 3 | Investigator Name | MD | | | | |
| 4 | Experiment Name | sincellte | | | | |
| 5 | Date | 31/12/18 | | | | |
| 6 | Workflow | GenerateFASTQ | | | | |
| 7 | Application | NovaSeq FASTQ Only | | | | |
| 8 | Instrument Type | NovaSeq | | | | |
| 9 | Assay | Chromium SingleCell 10x | | | | |
| 10 | Index Adapters | Chromium SingleCell 10x Indexes (4x96 Indexes) | | | | |
| 11 | Description | PE26-98_SingleCell-10X | | | | |
| 12 | Chemistry | Default | | | | |
| 13 | [Reads] | | | | | |
| 14 | 26 | | | | | |
| 15 | 98 | | | | | |
| 16 | [Settings] | | | | | |
| 17 | [Data] | | | | | |
| 18 | Lane | Sample_ID | Sample_Name | index | Sample_Project | Description |
| 19 | 1 | SI-3A-A1_1 | sample1 | AAACGGCG | Chromium_20211119 | Homo_sapiens |
| 20 | 1 | SI-3A-A1_2 | sample1 | CCTACCAT | Chromium_20211119 | Homo_sapiens |
| 21 | 1 | SI-3A-A1_3 | sample1 | GGCGTTTC | Chromium_20211119 | Homo_sapiens |
| 22 | 1 | SI-3A-A1_4 | sample1 | TTGTAAGA | Chromium_20211119 | Homo_sapiens |
| 23 | | | | | | |

https://support.10xgenomics.com/single-cell-gene-expression/software/pipelines/latest/using/bcl2fastq-direct

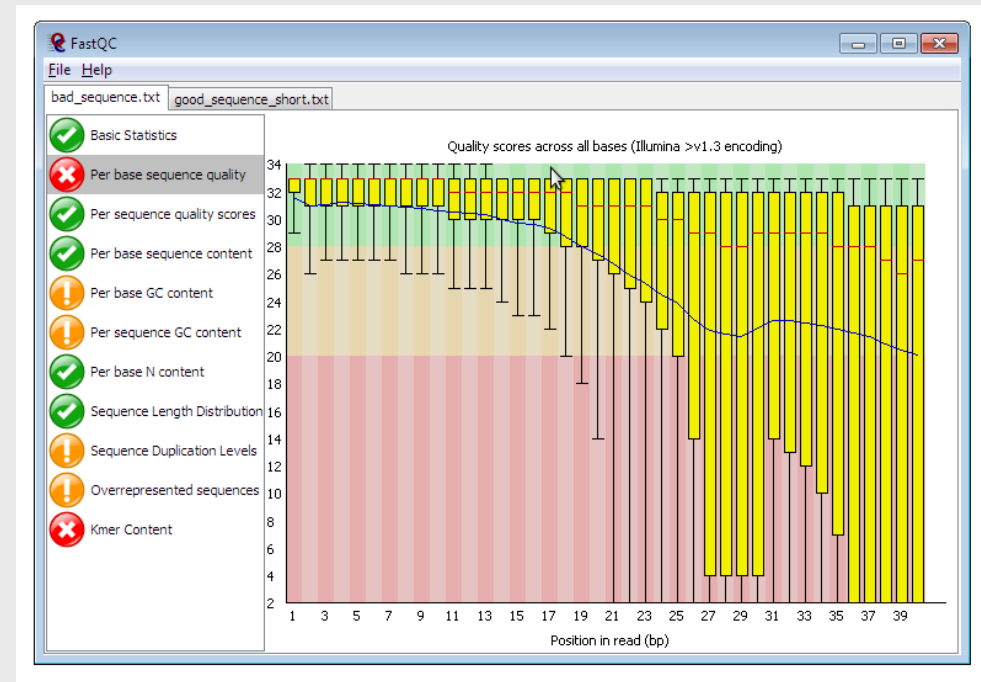GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

# Demultiplexing

Convert BCL files (sequencer output) to fastq files
Most used tool : 10X's cellranger mkfastq a wrapper around
bcl2fastq

- Usual sample sheet

- You must know :
  - i7 (i5) index sequence
  - R1 and R2 lengths
  - (depends on technology, version...)

- 10X: 1 index = 4 sequences ⇒ 4 lines

| Header | Sequence | Quality |

```
@HWI-ST227:389:C4WA2ACXX:7:1204:2272:59979
GGAGGAAGGTCCTCGCTCCTCTTTCATATAAGGGAAATGGCTGAAT
+
FFFFHHHHHHJIJJJJJJJJIJJJJIGIGIGGGIJJIJIJJJJJJIII
```

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

# Check reads quality : fastqc

- Performs various basic QC on reads

- For 10X scRNA datasets :
  - R1 (BC + UMI) : QC is mandatory. Watch out for Ns and highly repeated sequences
  - R2 : do as usual



https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

# Trimming

- If QC is not good:
  - Low base quality
  - Remaining adapter sequence
  - Homopolymer tailing
  - Low complexity

- Many tools to trim reads:
  - Trimmomatic (Bolger A.M. *et al.,* Bioinformatics (2014).
  - TrimGalore (Krueger F., https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/, unpublished 2012).
  - Cutadapt (Martin M., EMBnet.journal 2011)
  - Fastp (Chen *et al.,* Bioinformatics 2018).

- For single cell, like with xenome, apply to R2 file, then sync the R1 file.

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Check cross-species contaminations: FastQ Screen



https://www.bioinformatics.babraham.ac.uk/projects/fastq_screen/

- Quick mapping (bowtie2) of a subset of reads across multiple genomes and common contaminants: human, mouse, rat, E. coli, adapters, vectors...

- Identifies 1hit-1library, multi hits-1library, 1hit-multi libraries and multi hits-multi libraries
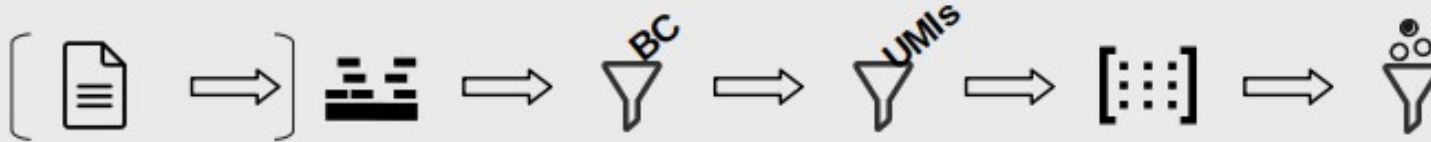
## Principle



- A classical scRNA-seq workflow contains four main steps:

    − Mapping the cDNA fragments to a reference

    − Assigning reads to genes

    − Assigning reads to cells (cell barcode demultiplexing)

    − Counting the number of unique RNA molecules (UMI deduplication).

## Principle



- Various tools have been developed:

  - **Cellranger**: 10X solution for 10X libraries only

  - **STARsolo**: an open source alternative to cellranger

  - **kallisto+bustools:** a pseudomapper and tool suite needing very little resources

  - (**Alevin**: a pseudomapper integrated with the salmon software)

## Cellranger



- A set of pipelines for single cell analysis

- Many languages + task scheduler Martian

- Aligner: STAR

- single cell gene expression: *cellranger count*



https://support.10xgenomics.com/single-cell-gene-expression/software/downloads/latest

**Reference Preparation**

- Human/mouse retained biotypes :
- Protein coding
- Long noncoding RNA
- Antisense
- All biotypes belonging to BCR/TCR (i.e. V/D/J) Genes
- All pseudogenes and small noncoding rnas are removed.

(note that older Cell Ranger reference versions do not include BCR/TCR Genes )

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

**Alignement**

- Cell Ranger further aligns confidently mapped exonic and intronic reads to annotated transcripts by examining their compatibility with the transcriptome

- Reads are classified based on whether they are exonic (light blue) or intronic (red) and whether they are sense or antisense (purple).

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

**Alignement**

- In Cell Ranger 7.0, by default, the cellranger-count and cellranger-multi pipelines will include intronic reads for whole transcriptome gene expression analysis –> recommended to maximize sensitivity
- Any reads that map in the sense orientation to a single are carried forward to UMI counting.
- Cell Ranger ignores antisense reads (purple).

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Reference Preparation

| Cell Ranger Reference | Species | Assembly/Annotation | Genes before filtering | Genes after filtering |
|---|---|---|---|---|
| 2020-A | human | GRCh38/GENCODE v32 | 60668 | 36601 |
| 2020-A | mouse | mm10/GENCODE vM23 | 55421 | 32285 |
| 3.0.0 | human | GRCh38/Ensembl 93 | 58395 | 33538 |
| 3.0.0 | human | hg19/Ensembl 87 | 57905 | 32738 |
| 3.0.0 | mouse | mm10/Ensembl 93 | 54232 | 31053 |
| 2.1.0 | mouse | mm10/Ensembl 84 | 47729 | 28692 |
| 1.2.0 | human | GRCh38/Ensembl 84 | 60675 | 33694 |
| 1.2.0 | human | hg19/Ensembl 82 | 57905 | 32738 |
| 1.2.0 | mouse | mm10/Ensembl 84 | 47729 | 27998 |

GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

**UMI Counting**

- Cell Ranger attempts to correct for sequencing errors in the UMI sequences by association ( Group confidently mapped reads -> Associate  UMIs differ by a single base (less confidently mapped) to their assigned group)

- Gene annotation with the most supporting reads is kept for UMI counting, and the other read groups are discarded

| Before Clustering | | | After Clustering | |
|---|---|---|---|---|
| **Umi** | **Count** | | **Umi** | **Count** |
| ATGGCGTT | 653 | ➔ | ATGGCGTT | 673 |
| ATGGCGTA | 12 | | | |
| ATGGCGTC | 8 | | | |
| CTGGCAAC | 403 | ➔ | CTGGCAAC | 406 |
| CTGGCGAC | 2 | | | |
| CTGGCTAC | 1 | | | |
| TACCGGAT | 42 | ➔ | TACCGGAT | 45 |
| TACAGGAT | 3 | | | |
| | | | | |
| sum reads | 1124 | | sum reads | 1124 |
| unique UMI | 8 | ➔ | unique UMI | 3 |

GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

# Building the count matrix

**UMI Counting**

- Cell Ranger attempts to correct for sequencing errors in the UMI sequences by association ( Group confidently mapped reads -> Associate  UMIs differ by a single base (less confidently mapped) to their assigned group)

- Gene annotation with the most supporting reads is kept for UMI counting, and the other read groups are discarded

- Aggregation: 1 BC+UMIs = 1 unique RNA molecule (filter excess)

- Finally, construct matrix with selected reads: *genes x barcodes*

GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

# Cellranger

## Outputs

## Cellranger

### 🙂

- Turnkey solution

```
cellranger count --id=count_hgmm_100_hg19_mm10 \
--transcriptome=/db/off_biomaj/10xgenomics/refdata-cellranger-hg19-and-mm10-3.0.0 \
--fastqs=../../Data/fastqs/original --sample=hgmm_100 --jobmode=local \
--localcores=4 --localmem=50 --expect-cells=100 --nosecondary
```

- Many QC-metrics, results summarized in 1 html.

- Some secondary analysis

- More complex experiences: VDJ analysis, feature-barcoding

- Versions for ATAC-Seq, TCR-seq and BCR-seq

### ☹

- Proprietary

- Analyze only 10X product (cannot customize BC and UMI)

- Has its own scheduler: hard to include in another pipeline

- Compatibility not guaranteed with all HPC managers

GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

## Cellranger

| Single Cell Gene Expression Solution | CR 7.1 | CR 7.0 | CR 6.1 | CR 6.0 | CR 5.0 | CR 4.0 | CR 3.1 | CR 3.0 | CR 2.2 |
|---|---|---|---|---|---|---|---|---|---|
| 3' Gene Expression v2 Libraries | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |
| 3' Gene Expression v3 Libraries | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |
| 3' Gene Expression v3 + Cell Surface Protein Libraries | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |
| 3' Gene Expression v3 + CRISPR Screening Libraries | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ |
| 3' Cell Surface Protein Libraries only | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ |
| Targeted Gene Expression | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |
| 3' Cell Multiplexing | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 3' LT (Low Throughput) | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ |
| 3' HT (High Throughput) | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |
| Fixed RNA Profiling | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ | ✗ |

Cell Ranger v8.0 introduces support for the analysis of GEM-X libraries.
Cell Ranger v7.2 is the last version to support the analysis of LT (low throughput) libraries.

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Kallisto/bustools

BC

UMIS

[:::]

_____
**kallisto**

_____
**bustools**

- Make use of the pseudo-aligner kallisto and the toolsuite bustools

- Very good time and memory performance.

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Kallisto/bustools



- Kallisto is a pseudo aligner: fast, low memory

- Working with a reference transcriptome, not genome

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Kallisto/bustools

- Kallisto is a pseudo aligner: fast, low memory

- Working with a reference transcriptome, not genome

- Principle:
  - reference chunked into k-mers ==> de Bruijn Graph

  - Reads chunked into k-mers and assigned to the
  - transcript(s) they overlap with

  - 1 read generally compatible with several transcripts:
  - proportion of transcripts computed by
  - Expectation Maximization from all reads

From Bray *et al.*, *Nat Biointechno* 2016

A very nice explanation of kallisto: https://bioinfo.iric.ca/fr/comprendre-comment-kallisto-fonctionne

## Kallisto/bustools

- Many technologies already accepted, the CB + UMI geometry is configurable

- Gives relative abundance, not absolute counts

- Output format in a specific, compressed format: bus instead of sam or bam files.

**Allows analysis of non 10X technologies**

From Melsted et al., Bioinformatics 2019

## Kallisto/bustools

- Next steps: bustools

Bus file + BC whitelist

*bustools correct*: correct and filter BC

*bustools sort*: sort results by BC, UMIs and gene

*bustools count*: correct and filter UMIs, construct matrix

raw gene x barcodes matrix

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

## Kallisto/bustools

- For modular pipeline construction

- Not proprietary

- Allows analysis of non 10X technologies

- The fastest and less resource consuming (can run on a laptop)

- Easy to include in a pipeline

- Compatible with HPC managers

- Not a turnkey solution

- No secondary analysis

- Gap with cellranger

- No Add sample_name and well range

GUSTAVE
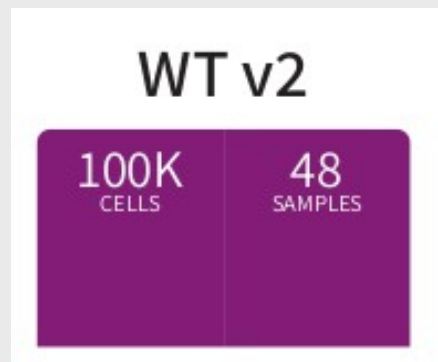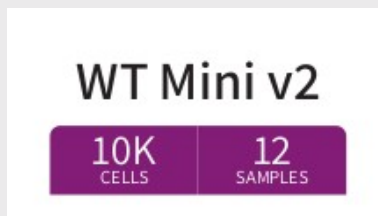ROUSSY
CANCER CAMPUS
GRAND PARIS

From Shulze Bruning *et al.*, *GigaScience* 2022

## Design Flexible Experiments that Scale

- Multiples samples are fixed and can be sequenced up to 6 months later

GUSTAVE ROUSSY
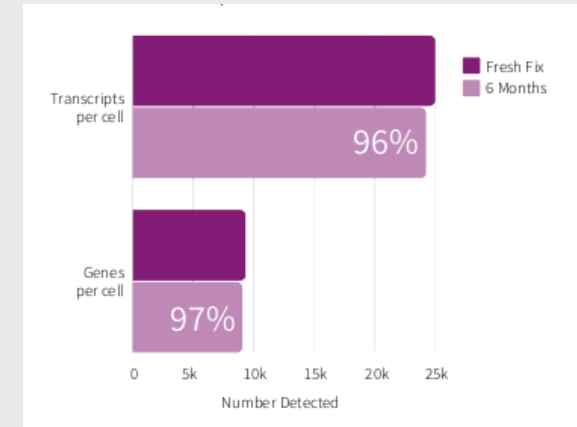CANCER CAMPUS
GRAND PARIS
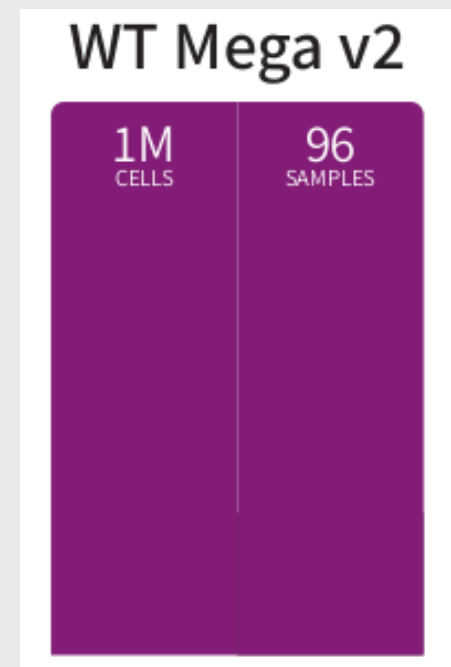
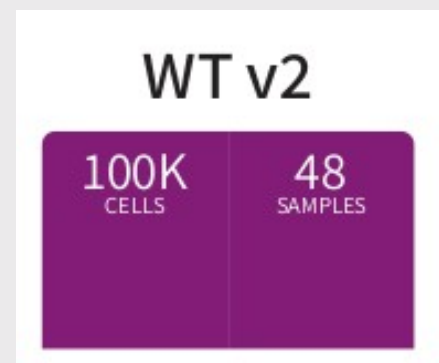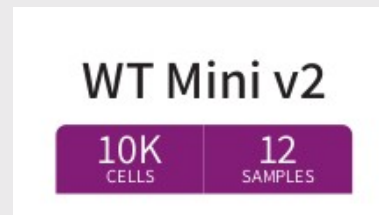## Design Flexible Experiments that Scale

- Multiples samples can be fixed sequenced up to 6 months later

- 3 kits are available
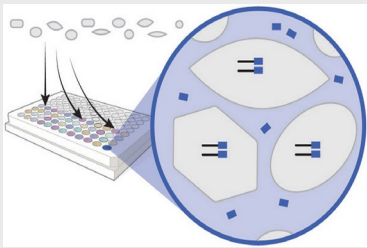
## Design Flexible Experiments that Scale

Multiples samples can be fixed sequenced up to 6 months later

- 3 kits are available
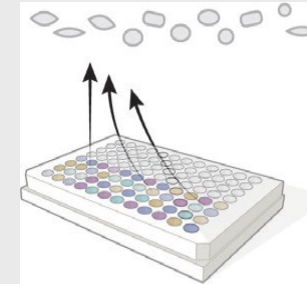
- four barcoding steps are required

**1**   Reverse Transcription

Split : Fixed cells/nuclei are distributed into wells, and the first sample-specific barcodes are added by in-cell reverse transcription.
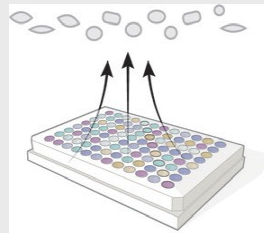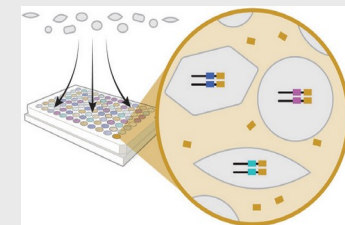


**2**

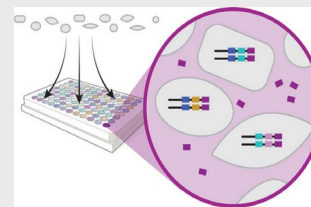Pool : All the cells are pooled together.



**3**   Ligation

Split : The pooled cells are distributed across a plate, and an in-cell ligation adds the second barcode.
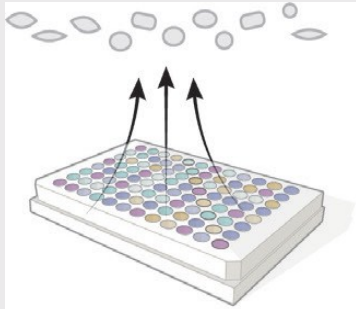


**4**   Pool : All the cells are pooled together.



**5**   Split | The pooled cells are again distributed across a plate, and a third barcode is added via in-cell ligation reaction.

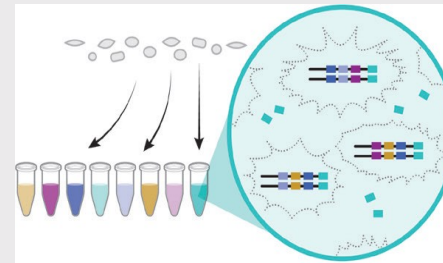GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS

**6**

Pool : All the cells are pooled together.



**7** Lysis and Library Prep

Split : The pooled cells are distributed across several sub-libraries then lysed. The fourth barcode is added via PCR.



**8** Sequencing with Illumina

Each transcript is assigned to a single cell based on a unique combination of barcodes.



**9** Data Analysis

ParseBiosciences-Pipeline.1.2.0.zip

```
# Create new environment with Python 3.10
conda create -n spipe conda-forge::python==3.10

# Activate your new environment
conda activate spipe
```

GUSTAVE
ROUSSY
CANCER CAMPUS
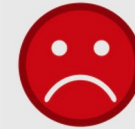GRAND PARIS

## ParseBiosciences-Pipeline

- Reference preparation : use split-pipe --mode mkref

- Aligner: STAR

- single cell gene expression : Parse count

- specify samples name for well ranges

- Running the pipeline for each library

- Combine libraries results

```
split-pipe \
    --mode  all \
    --chemistry v2 \
    --kit WT \
    --fq1 data/Parse_WT1_sublibrary_1_A1_S1_R1_001.fastq.gz \
    --fq2 data/Parse_WT1_sublibrary_1_A1_S1_R2_001.fastq.gz \
  --output_dir results/Parse_WT1_sublibrary_1_A1 \
    --genome_dir genomes/hg38 \
    --sample CABE048-Total_cells A1-A3 \
```

GUSTAVE / ROUSSY
CANCER CAMPUS
GRAND PARIS

- Turnkey solution

- Many QC-metrics, results summarized in 1 html.

- Some secondary analysis

- No empty droplet

- Versions for TCR-Seq and BCR-seq

- Batch effect reduced



- Proprietary

- Analyze only Parse product (cannot customize BC and UMI)

- Has its own scheduler: hard to include in another pipeline
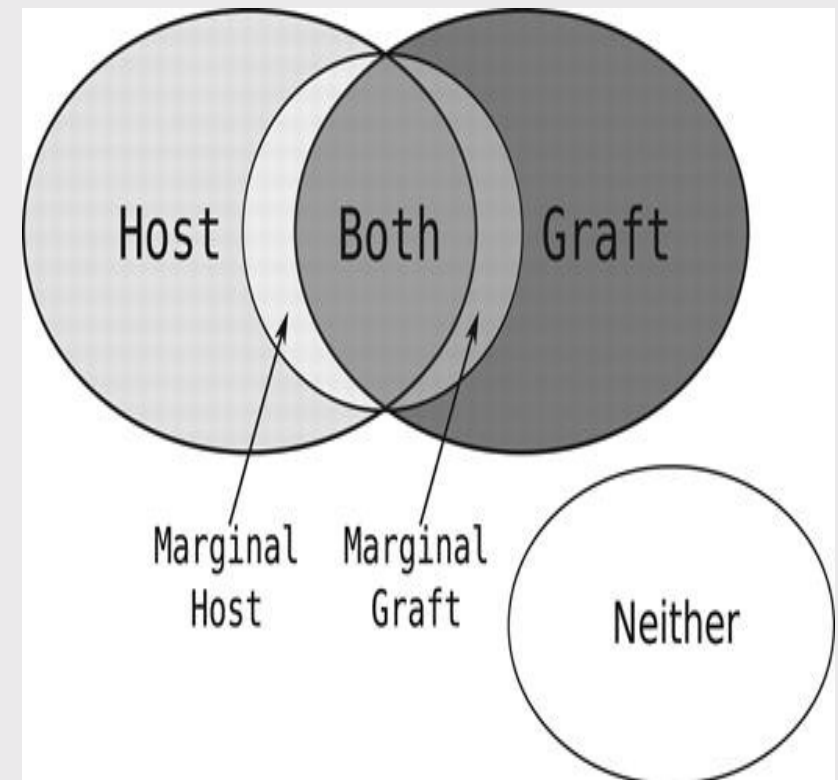
SincellTE 2024

Aziza CAIDI

# Thank you for your attention!

Thanks to Marine AGLAVE

GUSTAVE / ROUSSY
CANCER CAMPUS
GRAND PARIS

## Multiple species: Xenome

- For xenografts or contaminated samples

- 5 fastq files :
  - Graft
  - Host
  - Both
  - Neither
  - Ambiguous

- For single cell, apply to R2 only and sync R1: e.g. seqkit:
  - *seqkit seq* lists the selected read names.
  - *seqkit grep* filters R1 by keeping only reads in this list.
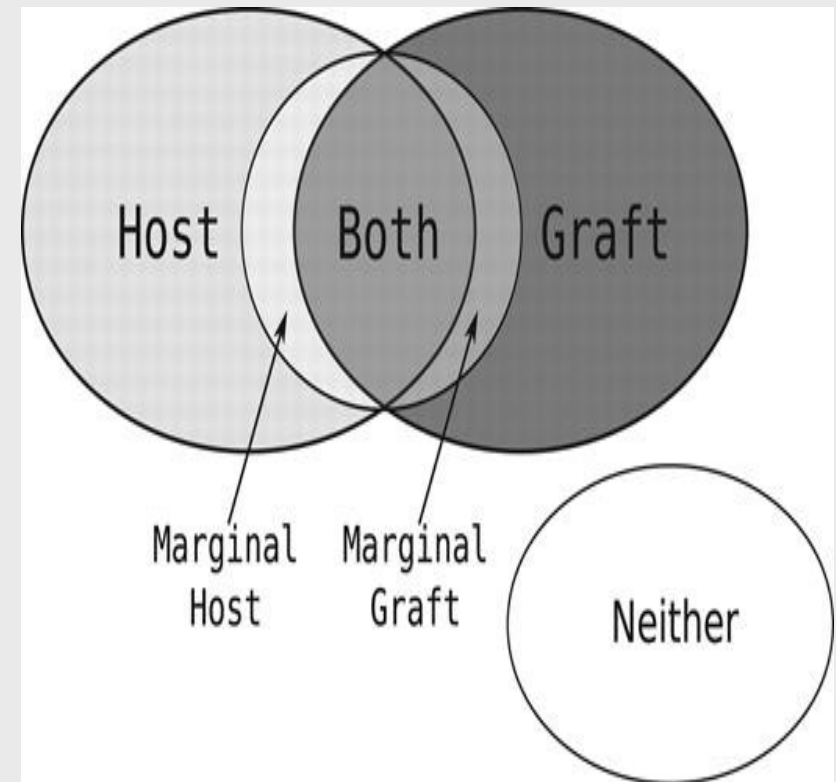  - *seqkit pair* pairs filtered R1 with R2.



https://github.com/data61/gossamer/blob/master/docs/xenome.md

GUSTAVE ROUSSY
CANCER CAMPUS
GRAND PARIS

## Multiple species: Xenome

- Xenome version is bugged: patch gossamer
- https://github.com/data61/gossamer


- Alternatives :
  - Xengsort (Zentgraf and Rahmann, S. Mol Biol 2021).
  - XenofilteR (Kluin *et al*, BMC Bioinfo 2018)
  - Bamcmp (Khandelwal *et al.*, MCR 2017).
  - XenoSplit: (https://github.com/goknurginer/XenoSplit Unpublished 2019).



https://github.com/data61/gossamer/blob/master/docs/xenome.md

GUSTAVE
ROUSSY
CANCER CAMPUS
GRAND PARIS