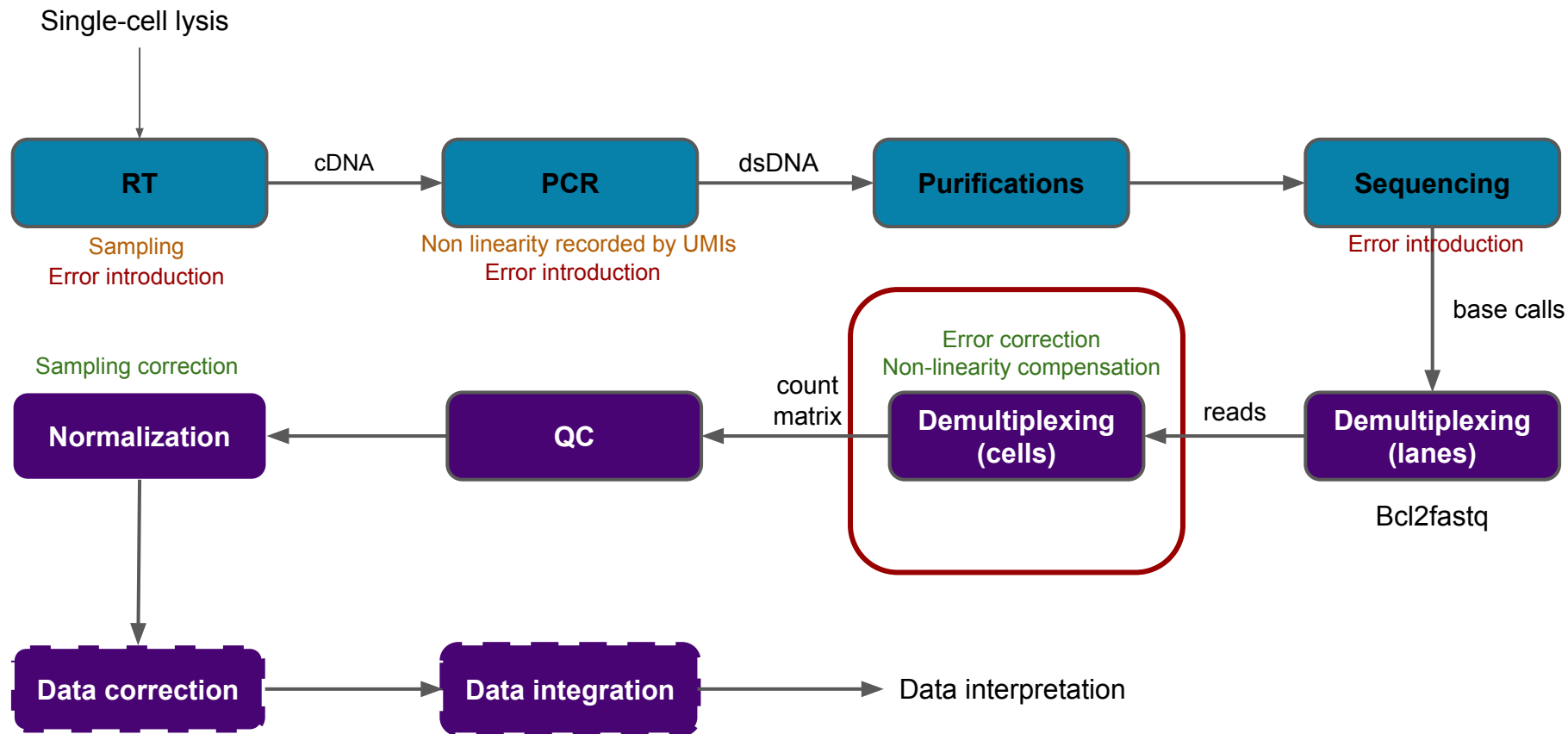# Demultiplexing
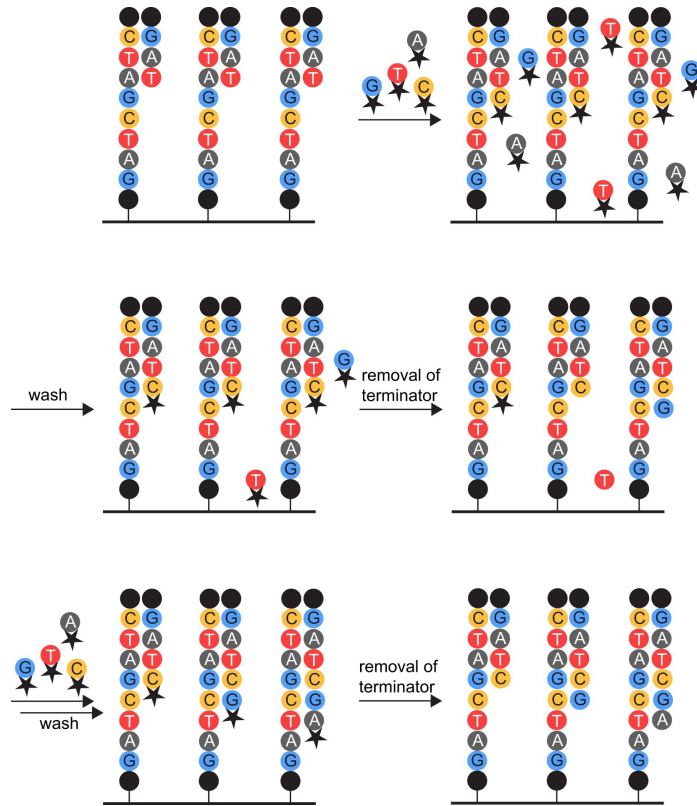# from FASTQ files to count matrix

05/11/20 - Gaël Blivet

# Different demultiplexing tools with different strategies

- Which one should I use to process unusual -different from 10X- read structure?

- Do these different strategies have an impact on the next stages of the analysis?

- Can it be wise to go back to multiplexing to integrate different analyzes?

# What error rate are we talking about?

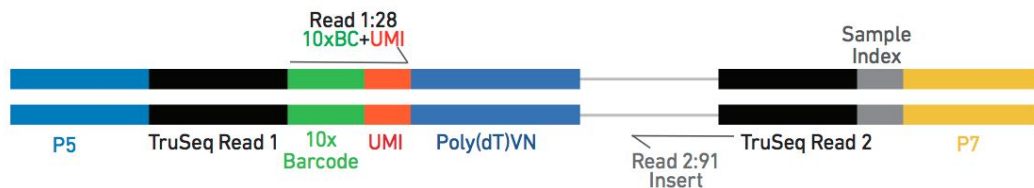| Step | RT | PCR | Sequencing |
|---|---|---|---|
| **Enzyme** | Superscript III (MMLV-RT) | Phusion, Q5 | ? |
| **Error rate / base** | ~$10{-}5$ | ~$10{-}6$ to ~$10{-}7$ | ~$10{-}3$ + phasing |
| **References** | *Invitrogen; Orton RJ, BMC Genomics. 2015; Potter J, Focus. 2003* | NEB | *Marinier 2015, BMC Bioinformatics; Orton RJ. BMC Genomics. 2015, Tilo Buschmann 2016* |

# Origin of phasing effects

*Pfeiffer, F., Gröber, C., Blank, M. et al. Systematic evaluation of error rates and causes in short samples in next-generation sequencing. Sci Rep 8, 10950 (2018).*
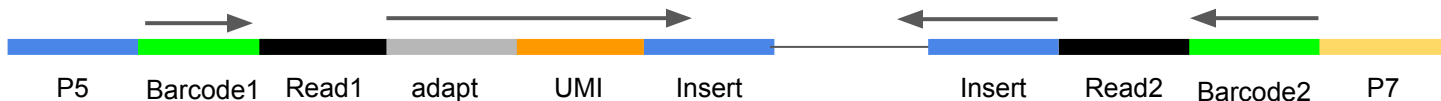
# Protocol specific read structure with similar info

- Cell barcode (BC)
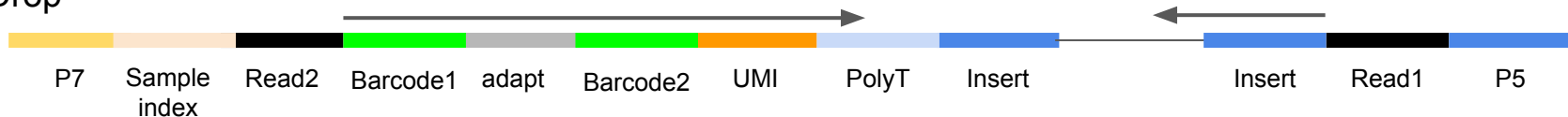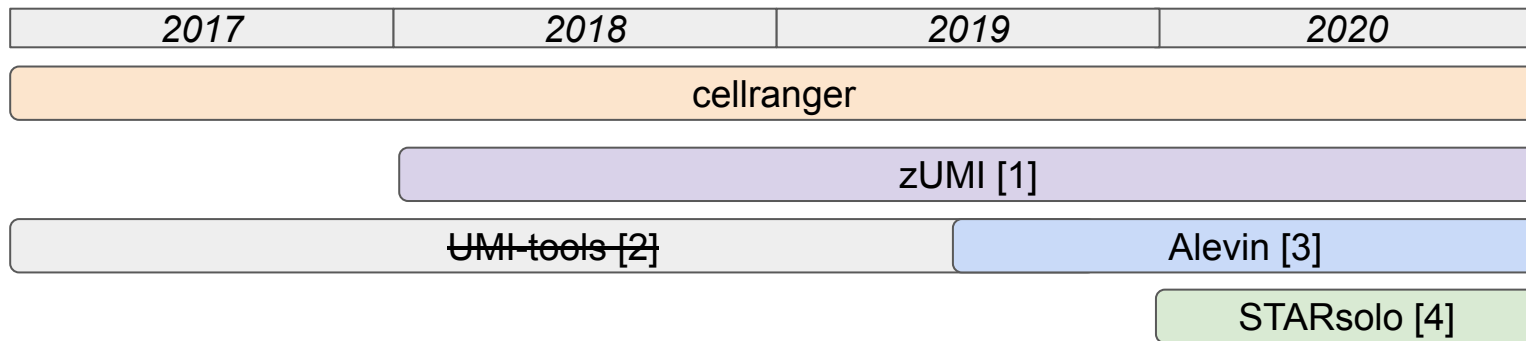- UMI (Unique Molecular identifier)
- Insert (gene)

# Different demultiplexing tools...

- Cellranger is:
    - Not flexible: can't analyse something else than fastq generated by 10x
    - Slow
    - Resource consuming (memory, space storage?)
    - Black box
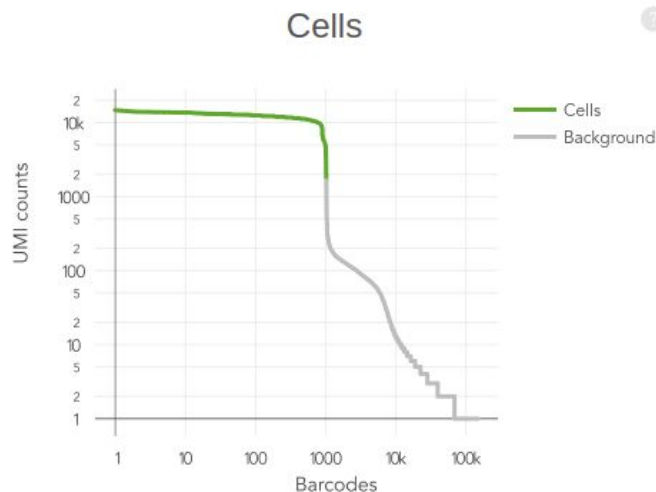- (non-exhaustive) alternatives

| 2017 | 2018 | 2019 | 2020 |
|---|---|---|---|

cellranger

zUMI [1]

~~UMI-tools [2]~~          Alevin [3]

STARsolo [4]

# ...with different strategies

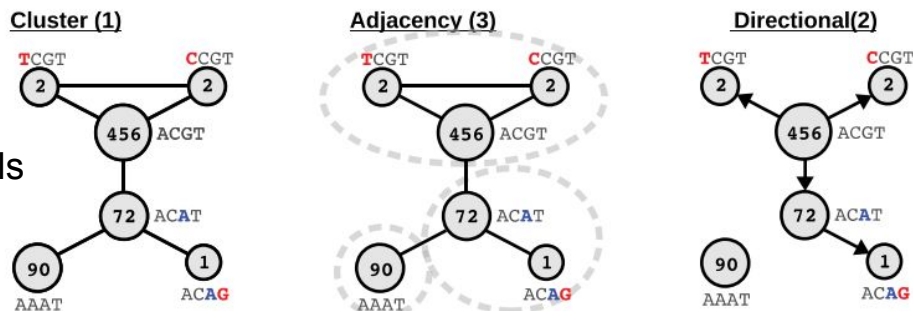|  | Flexible str. | BC correction | UMI correction | Mapping | Speed est. |
|---|---|---|---|---|---|
| **Cellranger** | 10X compatible | Hamming dist. based | No? | STAR | 16h (4 CPUs) |
| **alevin** | A bit, more to come | 1 error (sub or indel) | Graph-based | Pseudo-alignment | 20' (4 CPUs) |
| **STARsolo** | Yes | Hamming dist. based | Various modes | STAR | 35' (3 CPUs) |
| **zUMI** | Yes | No? | Hamming dist based | STAR | Fast |
| **LBCpipe** | Yes | Seqlev or Hamming dist. | Shift correction | Rsubread | 22h (4 CPUs) |

# Barcode extraction

- Pattern flexibility
- Error correction strategies
  - Error type: sub vs indels, shifts
  - Metric: Hamming, Levenshtein, Seqlev, Alevin specific...
  - [LBCpipe] Shift correction
- Barcode database filtering
  - "Knee plot" approach
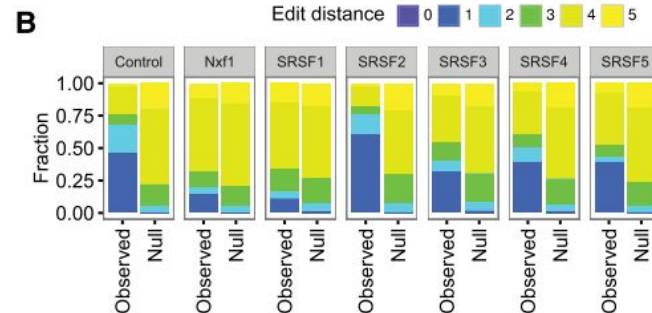  - [alevin] second pass with trained classifier

# UMI deduplication

- Pattern flexibility
- UMI deduplication strategies
  - No deduplication
  - Hamming distance
  - Graph-based approach
- [alevin] Multimapping strategy (transcript-level)
- [LBCpipe] shift correction
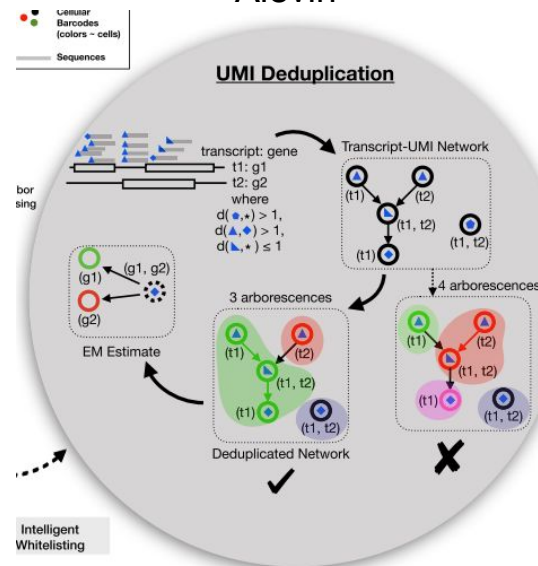
UMI-tools

Alevin

10

# Insert mapping

- Bulk mapping considerations + 3' bias
- [alevin] Pseudo-alignment + decoys
- [cellranger] Biotype discard
- Gene/transcript level => impact on UMI correction
- Resources: computing time, memory footprint

# Litterature comparison

- Comparison of alevin, cellranger and EmptyDrops:
  https://lazappi.github.io/phd-thesis/5-analysis.html#pre-processing
- Alevin compares to cellranger in its paper [3]:
  - Better for genes with a lower "sequence uniqueness"
  - Cellranger discards multimappers? (not true anymore since cellranger3?)

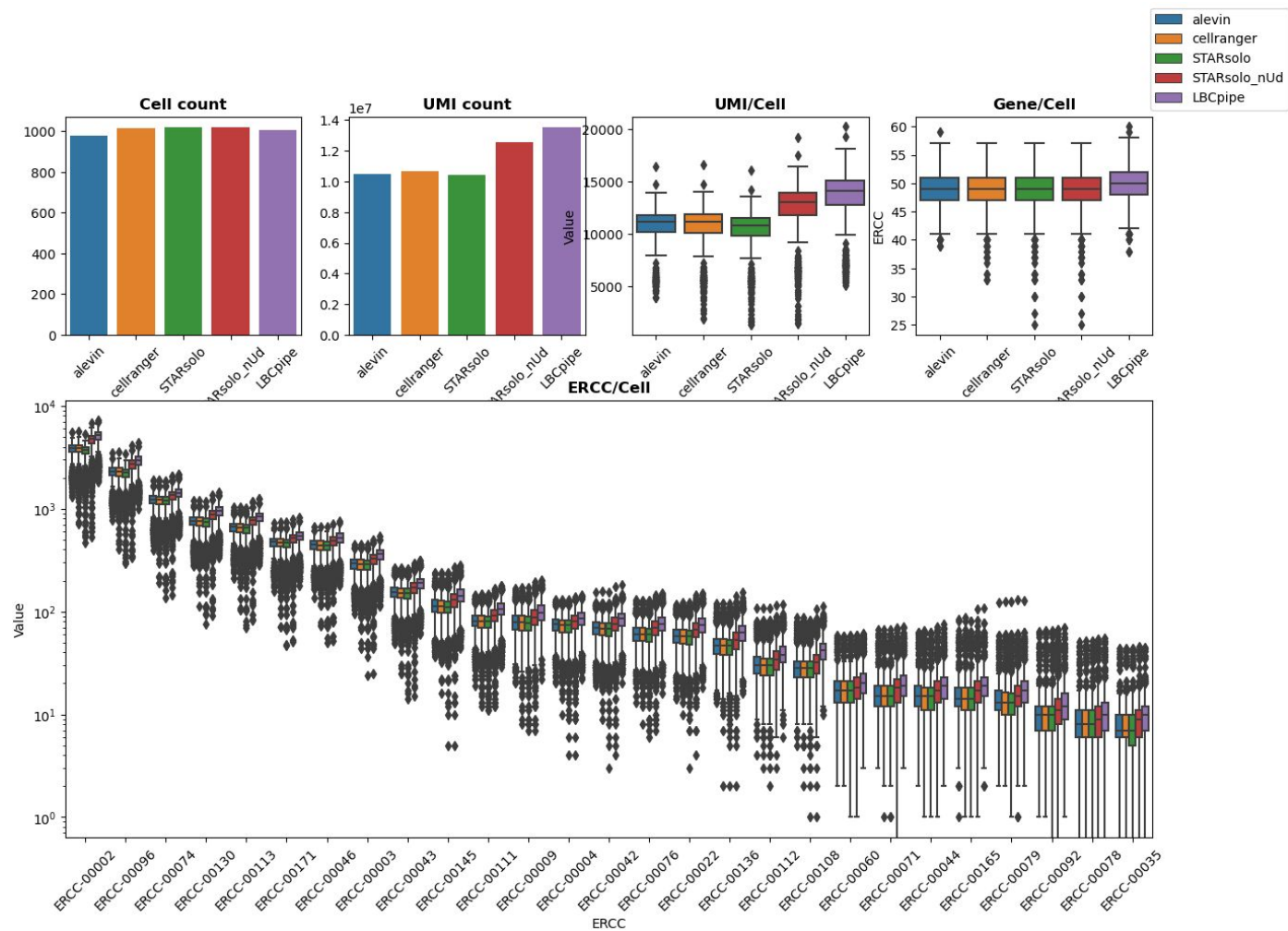# Different demultiplexing tools with different strategies

- Which one should I use to process unusual -different from 10X- read structure?

- Do these different strategies have an impact on the next stages of the analysis?

- Can it be wise to go back to multiplexing to integrate different analyzes?

# Considered strategy

1. Benchmarking on a spike-ins dataset
   - No reference genome issues
   - No biological variations (cell degradation, cell types/states, etc)
   - Only technical variations
   - Dropout exploration
2. Benchmarking on a classical single-cell dataset
   - Reference genome management
3. Benchmarking on a multi-organism single-cell dataset
   - Quality check through secondary analysis
4. Datasets integration
   - From matrices vs from FASTQ
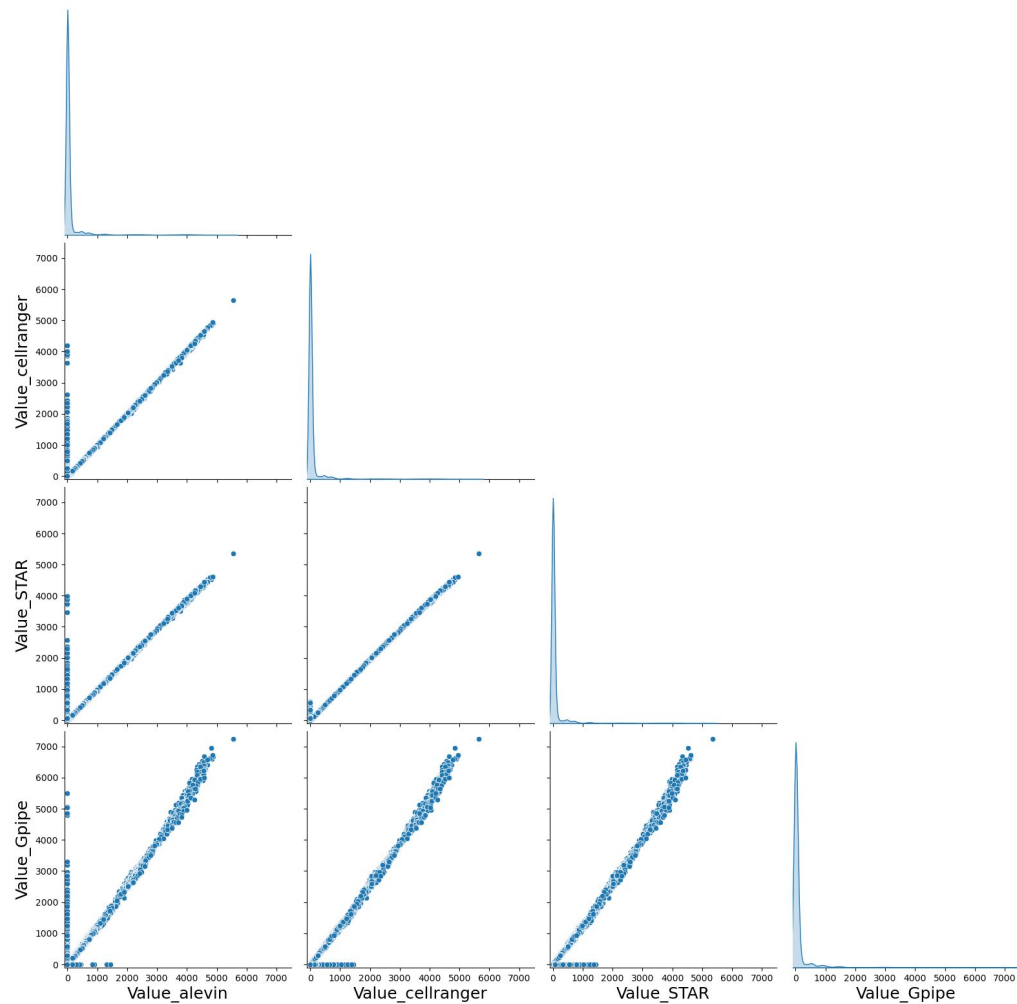   - Quality check through secondary analysis

# Dataset 1

- 10X ERCC 1k cells
    - [Massively parallel digital transcriptional profiling of single cells](#), Zheng *et al.* 2017, Nature communications
    - [Listed in 10X datasets](#)
    - 10X v2 chemistry
    - 14bp BC + 10bp UMI
    - 1k droplets
    - 92 ERCCs (spike-ins) in various expected quantities
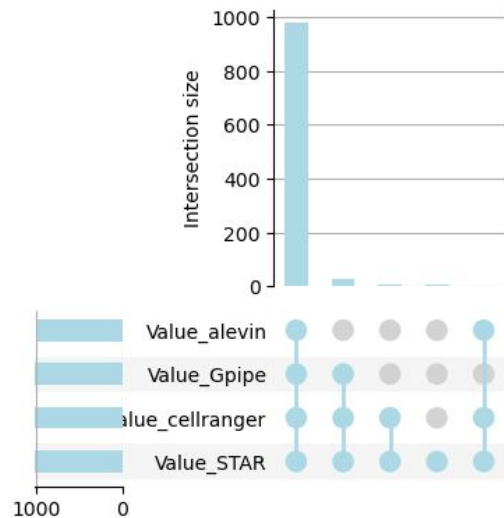- No reference genome compatibility issues
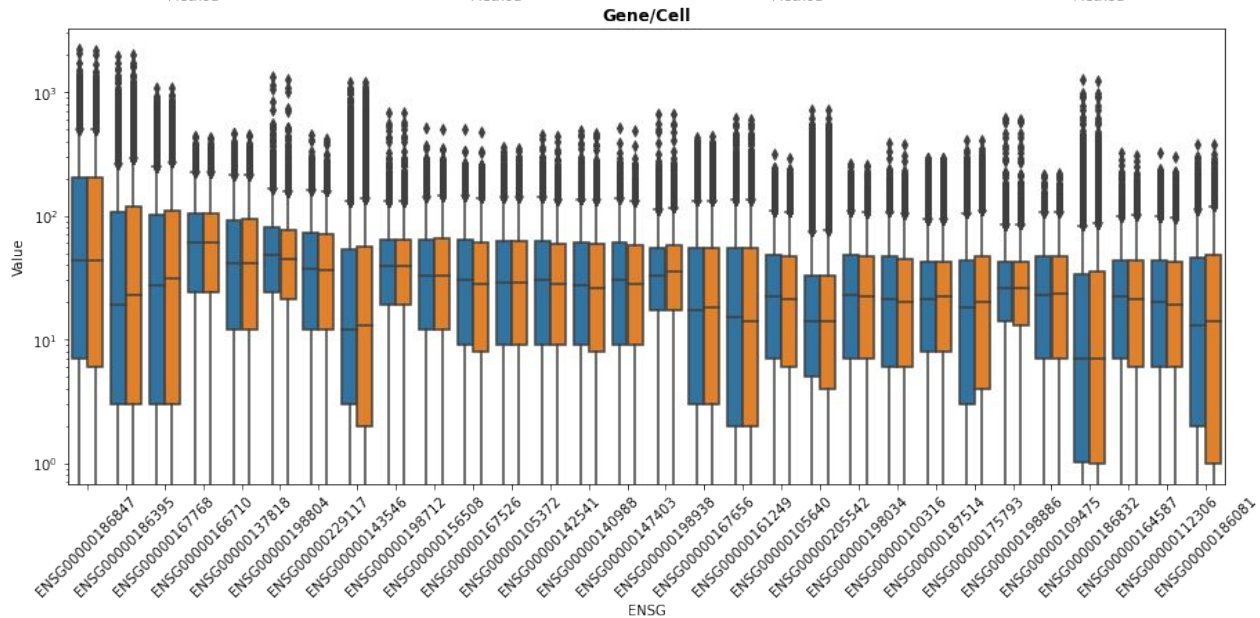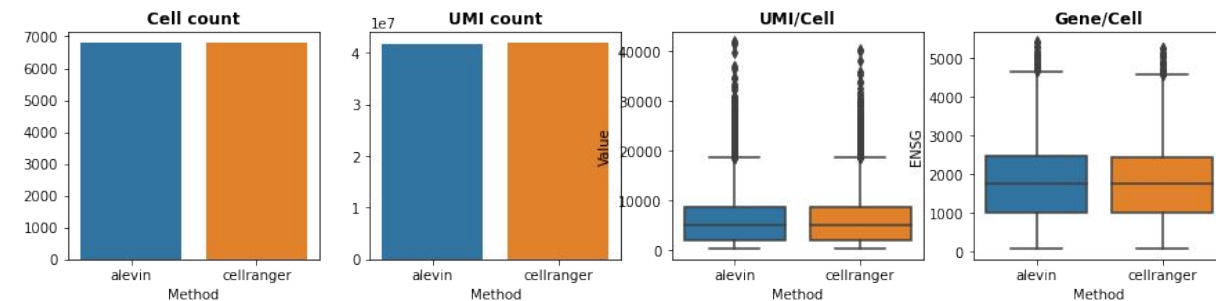
# Counts correlation

# Intersection of detected cells

# Dataset 2

- [Multimodal Analysis of Composition and Spatial Architecture in Human Squamous Cell Carcinoma](#), Ji *et al.*, Cell 2020
  - One "Normal" skin sample, ~320M reads
  - Chromium v2 chemistry
  - 16 bp BC + 10bp UMI
- Reference genome issues
  - Built from Ensembl 101
  - Considering gene IDs (not names)

# First impressions

- Cellranger can be replaced without results quality loss
  - to go faster
  - to handle unusual read structure
  - STARsolo seems to be the best candidates for these 2 features
- Differences between common demultiplexing tools seem to be marginal though different strategies are applied
- Need to investigate relevancy of shift correction and UMI correction

# Perspectives

- Study the impact of these different strategies on secondary analysis
  - On dataset2
  - On a reference dataset from "[Benchmarking single-cell RNA-sequencing protocols for cell atlas projects](...)", Mereu *et al.*, Nature biotechnology 2020
    - ~3k cells
    - A lots of protocols (n=13: 10X, Quartz-seq3, plate-based, etc)
    - Control mix of human, murine and canine cells
- Datasets integration
  - From matrices vs from FASTQ
  - Quality check through secondary analysis

# Tools references

- [1] [zUMIs - A fast and flexible pipeline to process RNA sequencing data with UMIs](), Parekh *et al.*, GigaScience 2018
- [2] [UMI-tools: Modeling sequencing errors in Unique Molecular Identifiers to improve quantification accuracy](), Smith *et al.*, Genome Research 2017
- [3] [Alevin efficiently estimates accurate gene abundances from dscRNA-seq data](), Srivastava *et al.*, Genome Biology 2019
- [4] [https://github.com/alexdobin/STAR/blob/master/docs/STARsolo.md](https://github.com/alexdobin/STAR/blob/master/docs/STARsolo.md), Dobin, 2020