

Single-cell RNA-Seq with kallisto

Vasilis Ntranos

Outline

- 10x Chromium dataset



10X GENOMICS®

10X GENOMICS®

PRODUCTS TECHNOLOGY COMPANY CAREERS

Support > Single Cell > Datasets

SE

3k PBMCs from a Healthy Donor

Chromium Demonstration Dataset by Cell Ranger 1.1.0

Peripheral blood mononuclear cells (PBMCs) from a healthy donor (same donor as pbmc6k).

PBMCs are primary cells with relatively small amounts of RNA (~1pg RNA/cell).

- 2,700 cells detected
- Sequenced on Illumina NextSeq 500 with ~69,000 reads per cell
- 98bp read1 (transcript), 8bp I5 sample barcode, 14bp I7 GemCode barcode and 10bp read2 (UMI)
- Analysis run with --cells=3000

Published on May 26, 2016

[View Summary](#)

Show batch download instructions

Input Files	Size	md5sum
FASTQs	17.38 GB	7999bb457af4e11c57f96bf97e4ee645

t-SNE projection of Cells Colored by k-means Clustering

Cluster	Cells
1	267
2	470
3	348
4	1,387
5	1
6	10
7	1
8	216

Outline

- 10x Chromium dataset



10X GENOMICS®

- Single-cell workflow using kallisto
github.com/pachterlab/scRNA-Seq-TCC-prep/

jupyter notebook/python scripts for clustering based on transcript compatibility counts



10X GENOMICS®

Support > Single Cell > Datasets

3k PBMCs from a Healthy Donor

Chromium Demonstration Dataset by Cell Ranger 1.1.0

Peripheral blood mononuclear cells (PBMCs) from a healthy donor (same donor as pbmc6k).

PBMCs are primary cells with relatively small amounts of RNA (~1pg RNA/cell).

- 2,700 cells detected
- Sequenced on Illumina NextSeq 500 with ~69,000 reads per cell
- 98bp read1 (transcript), 8bp I5 sample barcode, 14bp I7 GemCode barcode and 10bp read2 (UMI)
- Analysis run with --cells=3000

Published on May 26, 2016

[View Summary](#)

Show batch download instructions

Input Files	Size	md5sum
FASTQs	17.38 GB	7999bb457af4e11c57f96bf97e4ee645

t-SNE projection of Cells Colored by k-means Clustering

t-SNE2

t-SNE1

- 1 - 267 cells
- 2 - 470 cells
- 3 - 348 cells
- 4 - 1,387 cells
- 5 - 1 cells
- 6 - 10 cells
- 7 - 1 cells
- 8 - 216 cells

Outline

- 10x Chromium dataset

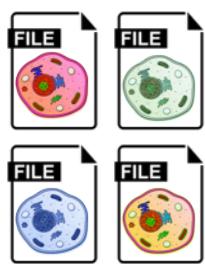


10X GENOMICS®

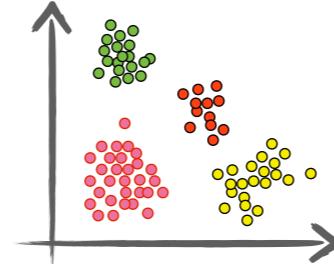
- Single-cell workflow using kallisto
github.com/pachterlab/scRNA-Seq-TCC-prep/

jupyter notebook/python scripts for clustering based on transcript compatibility counts

Single Cell Reads



Single Cell Clusters



Single Cell Datasets

▼ Chromium Demonstration

- 1:1 Mixture of Fresh Frozen Human (HEK293T) and Mouse (NIH3T3) Cells
- 33k PBMCs from a Healthy Donor
- 3k PBMCs from a Healthy Donor
- 6k PBMCs from a Healthy Donor

▼ Single Cell 3' Paper: Zheng et al. 2016 (29 datasets)

Datasets for the manuscript [Zheng et al., "Massively parallel digital transcriptional profiling of single cells"](#). We encourage you to download the data here, as the BAM files deposited in the SRA database have had the cell barcode tags removed. We are working with NCBI to resolve this issue.

- 293T Cells
- 293T and 3T3 Cell Mixture
- 50%:50% Donor B: Donor C PBMC Mixture
- 50%:50% Jurkat:293T Cell Mixture
- 90%:10% Donor B: Donor C PBMC Mixture
- 99%:1% Donor B: Donor C PBMC Mixture
- 99%:1% Jurkat:293T Cell Mixture
- AML027 Post-transplant BMMCs
- AML027 Pre-transplant BMMCs
- AML035 Post-transplant BMMCs
- AML035 Pre-transplant BMMCs
- CD14+ Monocytes
- CD19+ B Cells
- CD34+ Cells
- CD4+ Helper T Cells
- CD4+/CD25+ Regulatory T Cells
- CD4+/CD45RA+/CD25- Naive T cells
- CD4+/CD45RO+ Memory T Cells
- CD56+ Natural Killer Cells
- CD8+ Cytotoxic T cells
- CD8+/CD45RA+ Naive Cytotoxic T Cells
- ERCC (1k GEMS, 1:10 Dilution)
- Fresh 68k PBMCs (Donor A)
- Frozen BMMCs (Healthy Control 1)
- Frozen BMMCs (Healthy Control 2)
- Frozen PBMCs (Donor A)
- Frozen PBMCs (Donor B)
- Frozen PBMCs (Donor C)
- Jurkat Cells

Outline

- 10x Chromium dataset

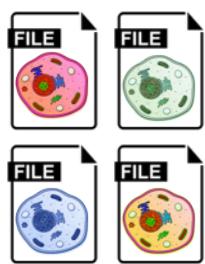


10X GENOMICS®

- Single-cell workflow using kallisto
- github.com/pachterlab/scRNA-Seq-TCC-prep/

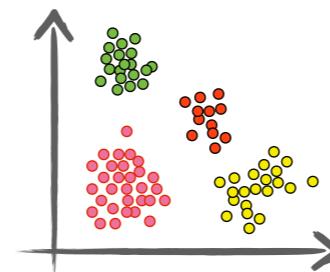
jupyter notebook/python scripts for clustering based on transcript compatibility counts

Single Cell Reads



```
>_ kallisto  
pseudoalign
```

Single Cell Clusters



Single Cell Datasets

Fluidigm C1

▼ Chromium Demonstration

- 1:1 Mixture of Fresh Frozen Human Lung Adenocarcinoma Cells (H3T3) Cells
- 33k PBMCs from a Healthy Donor A
- 3k PBMCs from a Healthy Donor B
- 6k PBMCs from a Healthy Donor C



▼ Single Cell 3' Paper: Zheng et al.

Datasets for the manuscript [Zheng et al. \(2012\). Single-cell transcriptome analysis reveals cell-to-cell heterogeneity in gene expression and lineage relationship in the mouse thymus](#). We encourage you to download the datasets. All BAM files deposited in the SRA database have had the cell barcode tag removed. Please contact us if you are working with NCBI to resolve this issue.



- 293T Cells
- 293T and 3T3 Cell Mixture
- 50%:50% Donor B: Donor C PBMC Mixture
- 50%:50% Jurkat:293T Cell Mixture
- 90%:10% Donor B: Donor C PBMC Mixture
- 99%:1% Donor B: Donor C PBMC Mixture
- 99%:1% Jurkat:293T Cell Mixture
- AML027 Post-transplant BMMCs
- AML027 Pre-transplant BMMCs
- AML035 Post-transplant BMMCs
- AML035 Pre-transplant BMMCs
- CD14+ Monocytes
- CD19+ B Cells
- CD34+ Cells
- CD4+ Helper T Cells
- CD4+/CD25+ Regulatory T Cells
- CD4+/CD45RA+/CD25- Naive T cells
- CD4+/CD45RO+ Memory T Cells
- CD56+ Natural Killer Cells
- CD8+ Cytotoxic T cells
- CD8+/CD45RA+ Naive Cytotoxic T Cells
- ERCC (1k GEMS, 1:10 Dilution)
- Fresh 68k PBMCs (Donor A)
- Frozen BMMCs (Healthy Control 1)
- Frozen BMMCs (Healthy Control 2)
- Frozen PBMCs (Donor A)
- Frozen PBMCs (Donor B)
- Frozen PBMCs (Donor C)
- Jurkat Cells

Outline

- 10x Chromium dataset

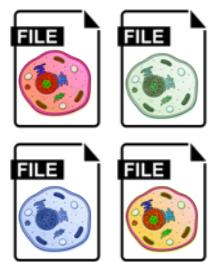


10X GENOMICS®

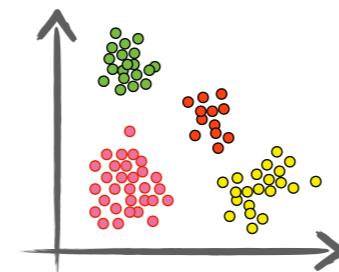
- Single-cell workflow using kallisto
github.com/pachterlab/scRNA-Seq-TCC-prep/

jupyter notebook/python scripts for clustering based on transcript compatibility counts

Single Cell Reads



Single Cell Clusters



Single Cell Datasets

Fluidigm C1

▼ Chromium Demonstration

- 1:1 Mixture of Fresh Frozen Human Lung Cells and 293T Cells
- 33k PBMCs from a Healthy Donor
- 3k PBMCs from a Healthy Donor
- 6k PBMCs from a Healthy Donor

▼ Single Cell 3' Paper: Zheng et al. 2015

Datasets for the manuscript Zheng et al. 2015 "Single-cell transcriptome analysis reveals global transcriptional profiling of single cells". We encourage you to download the datasets. All BAM files deposited in the SRA database have had the cell barcode tag removed. Please contact us working with NCBI to resolve this issue.

FLUIDIGM™

- 293T Cells
- 293T and 3T3 Cell Mixture
- 50%:50% Donor B: Donor C PBMC Mixture
- 50%:50% Jurkat:293T Cell Mixture
- 90%:10% Donor B: Donor C PBMC Mixture
- 99%:1% Donor B: Donor C PBMC Mixture
- 99%:1% Jurkat:293T Cell Mixture
- AML027 Post-transplant BMMC
- AML027 Pre-transplant BMMC
- AML035 Post-transplant BMMC
- AML035 Pre-transplant BMMC

Drop-seq

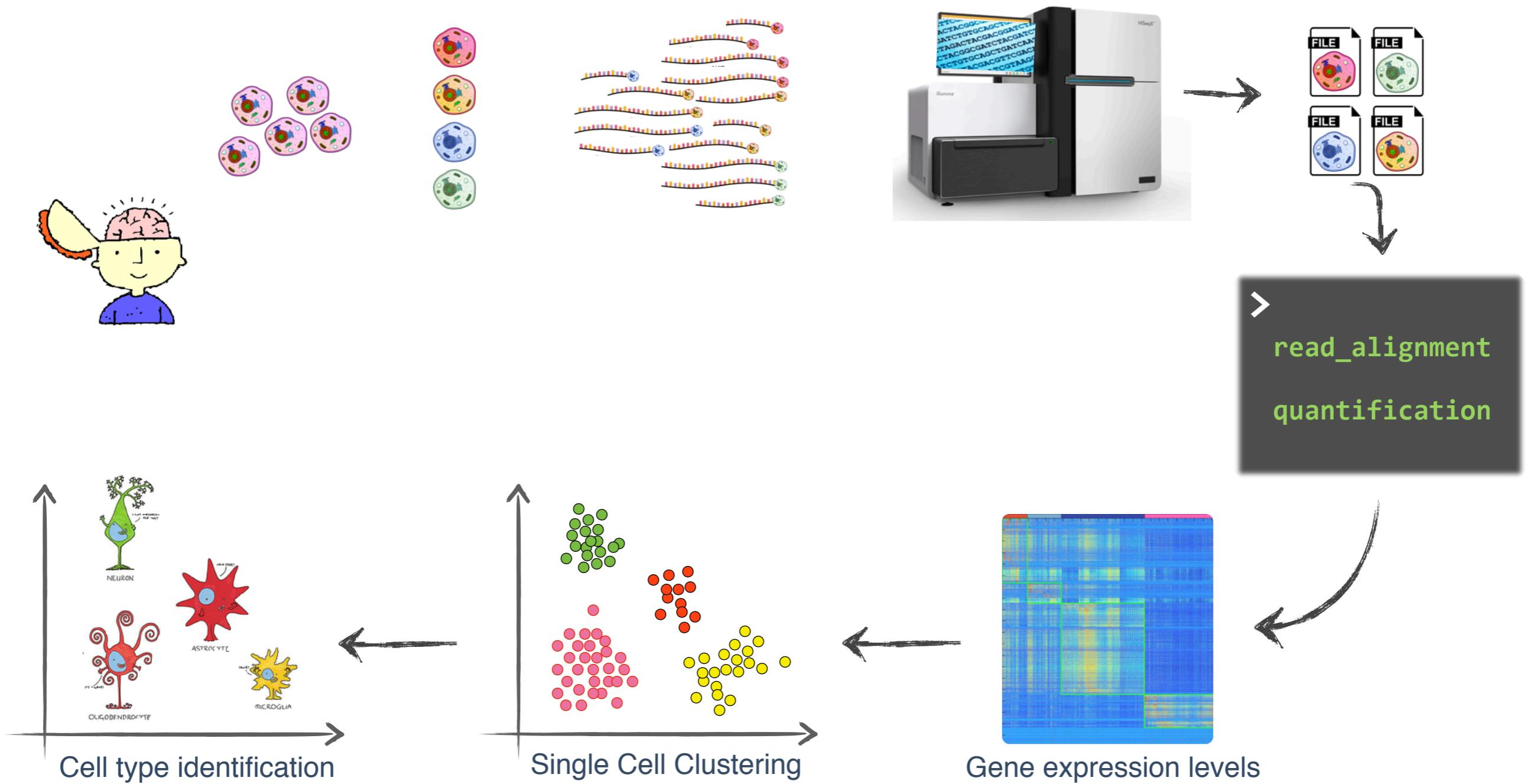
Drop-seq single cell analysis

A diagram illustrating the Drop-seq workflow. On the left, a photograph shows a laboratory setup with a microscope and a Drop-seq instrument. To the right, a schematic shows a flow: "Cells" are mixed with "Distinctly barcoded beads" and then loaded into a "Drop-seq" device. The process involves "Emulsion droplet formation" and "PCR amplification". Below the schematic, a sequence of three icons shows a cell being captured by a bead, the bead becoming a droplet, and finally a cluster of droplets representing a single cell library.

Macosko et al. 2015

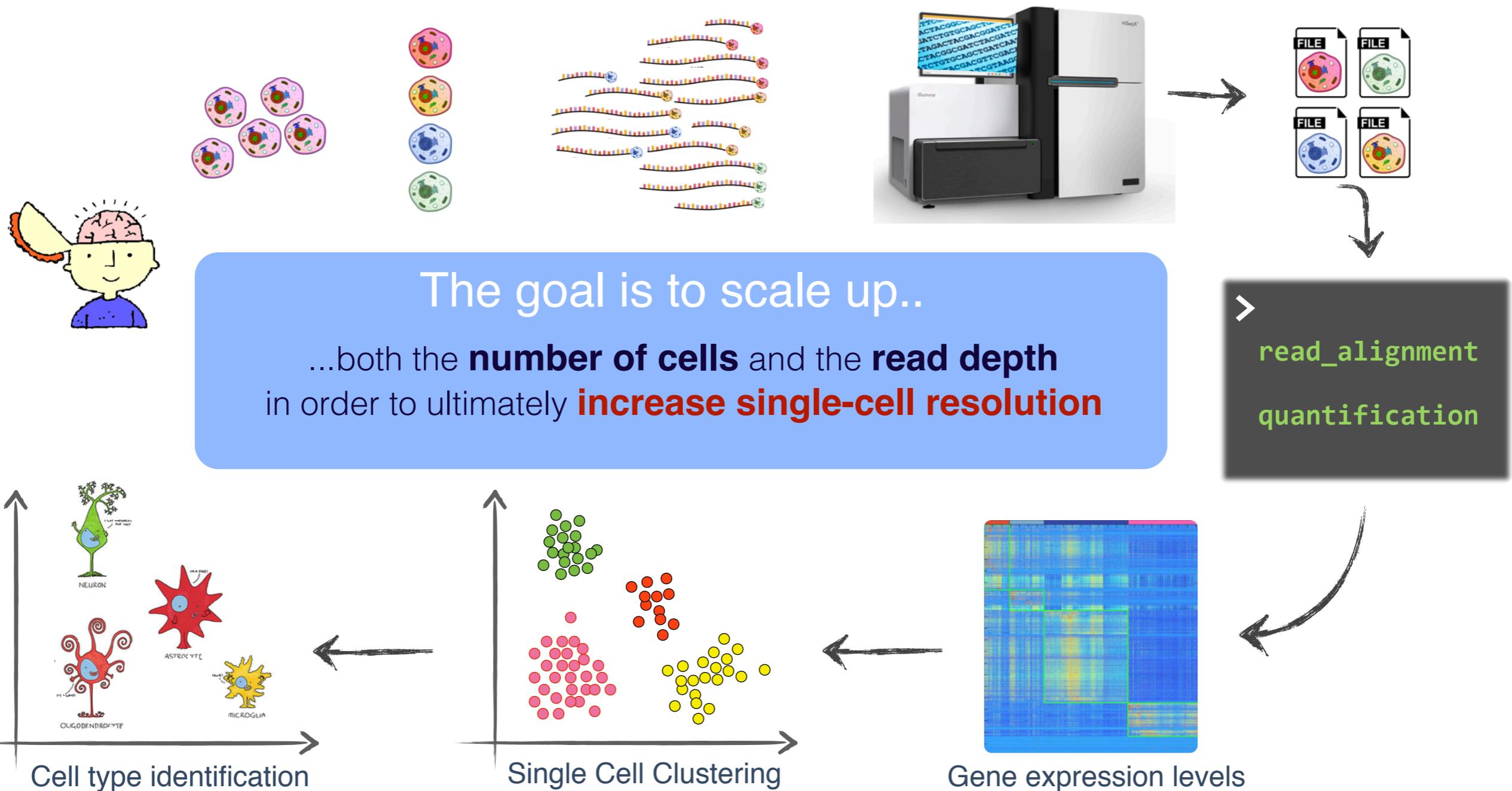
Introduction — Single-Cell RNA-Seq

Standard workflow:



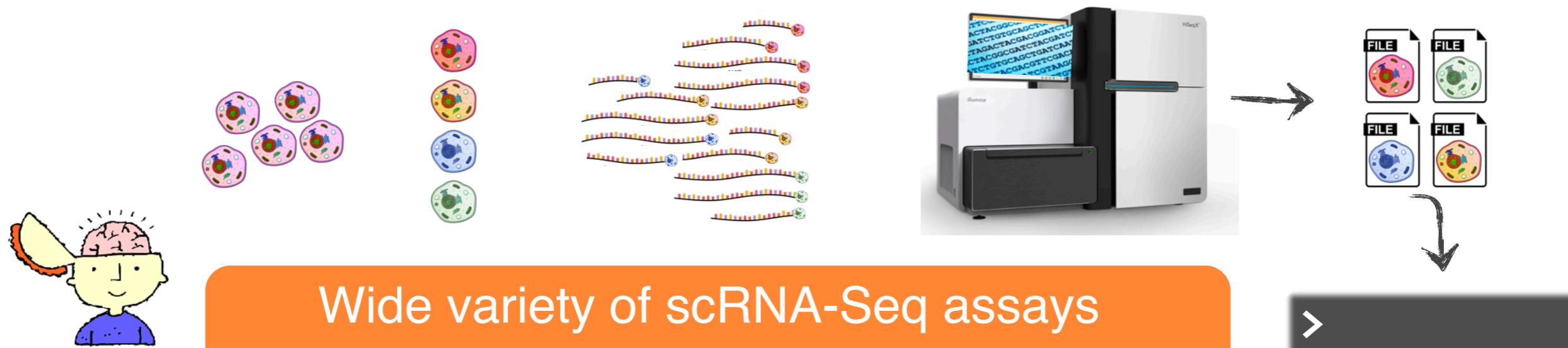
Introduction — Single-Cell RNA-Seq

Standard workflow:



Introduction — Single-Cell RNA-Seq

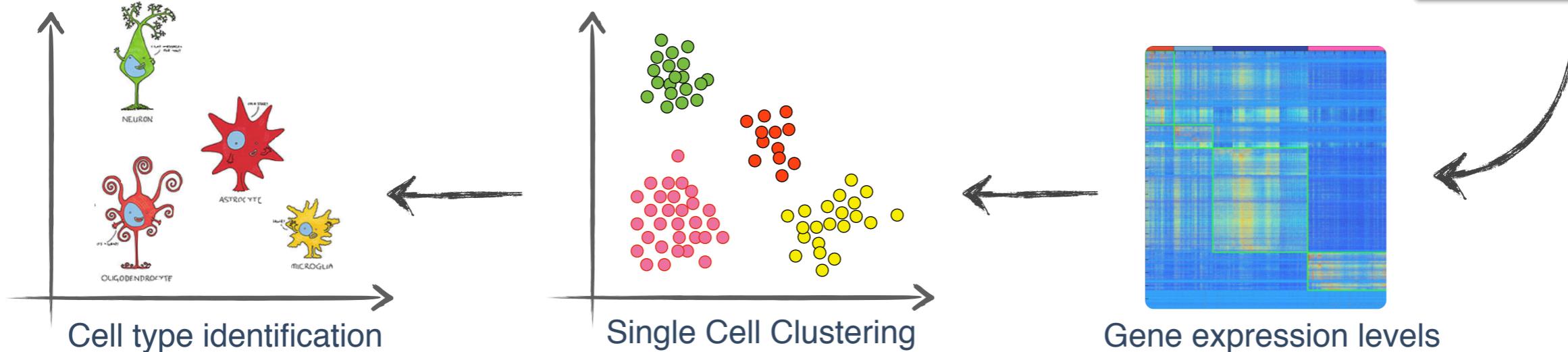
Standard workflow:



Wide variety of scRNA-Seq assays

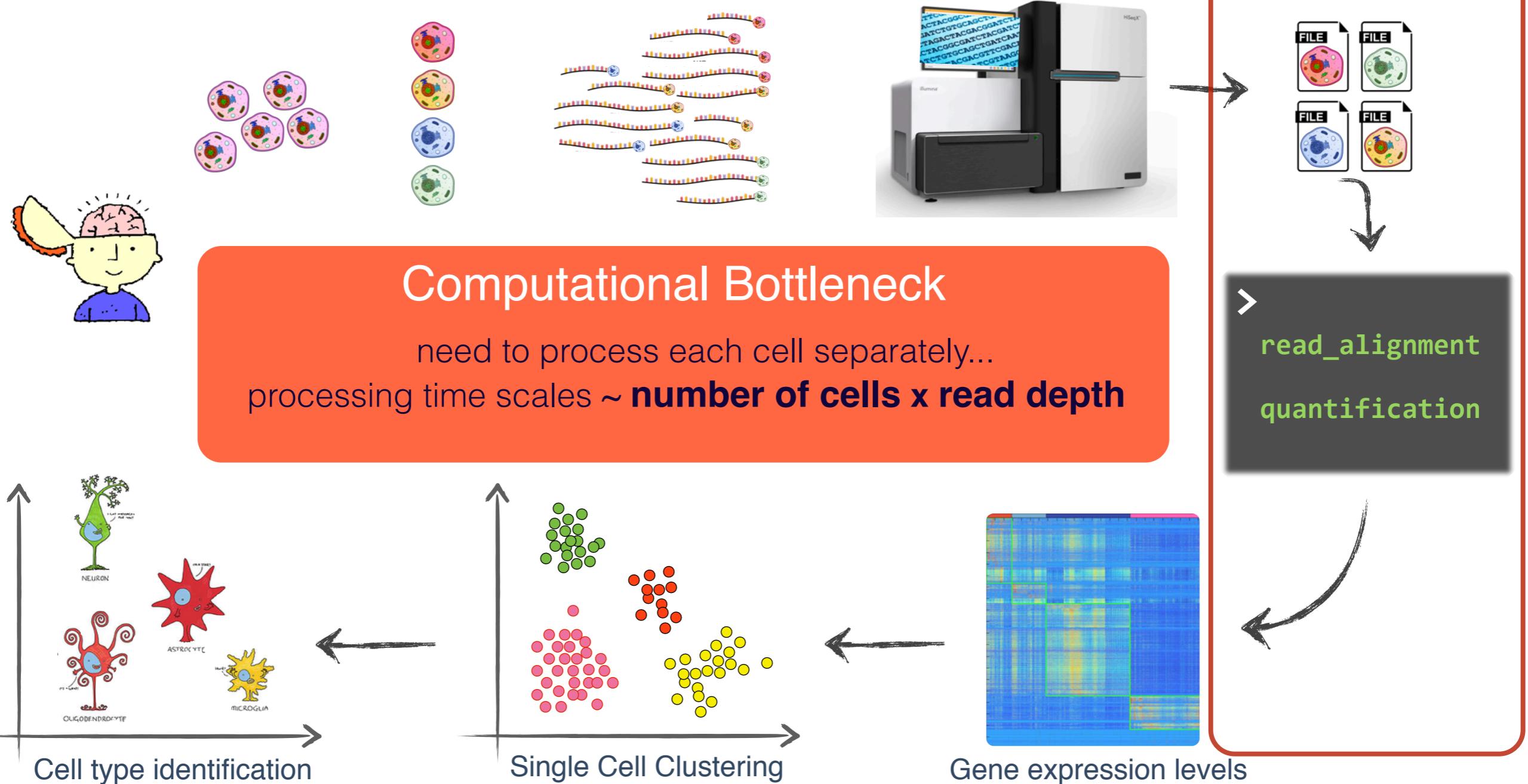
... technologies with different characteristics,
need different **models**, difficult to **compare** results...
lack of “**universality**”

>
read_alignment
quantification



Introduction — Single-Cell RNA-Seq

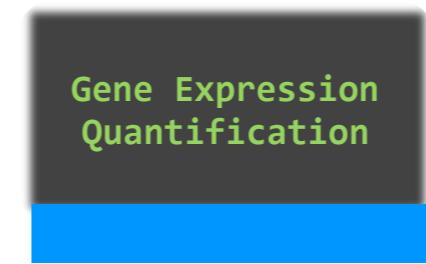
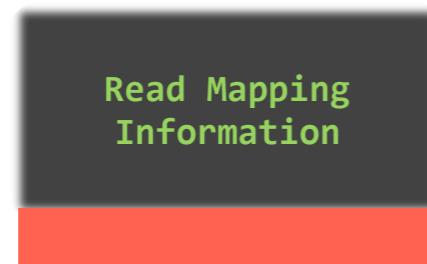
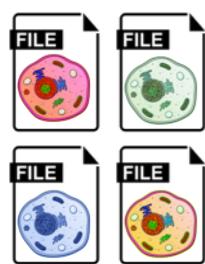
Standard workflow:



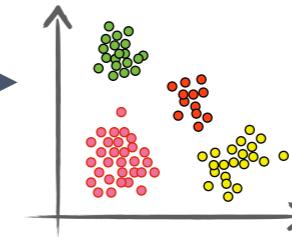
Big Picture

Standard analysis workflow:

Single Cell Reads



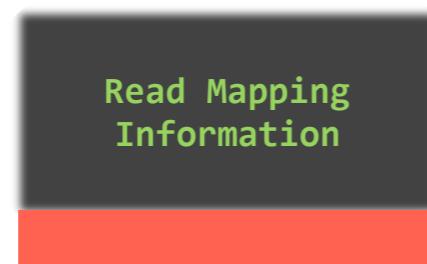
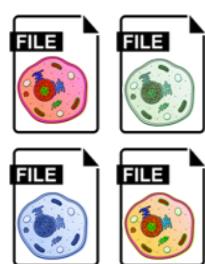
Clustering



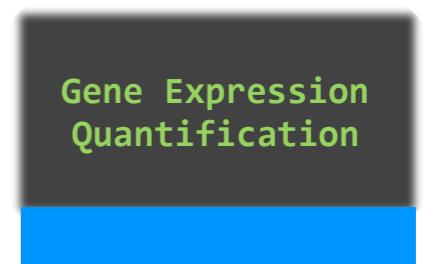
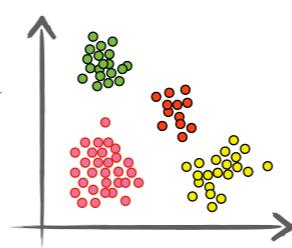
extract further information

Our proposed workflow:

Single Cell Reads



Clustering

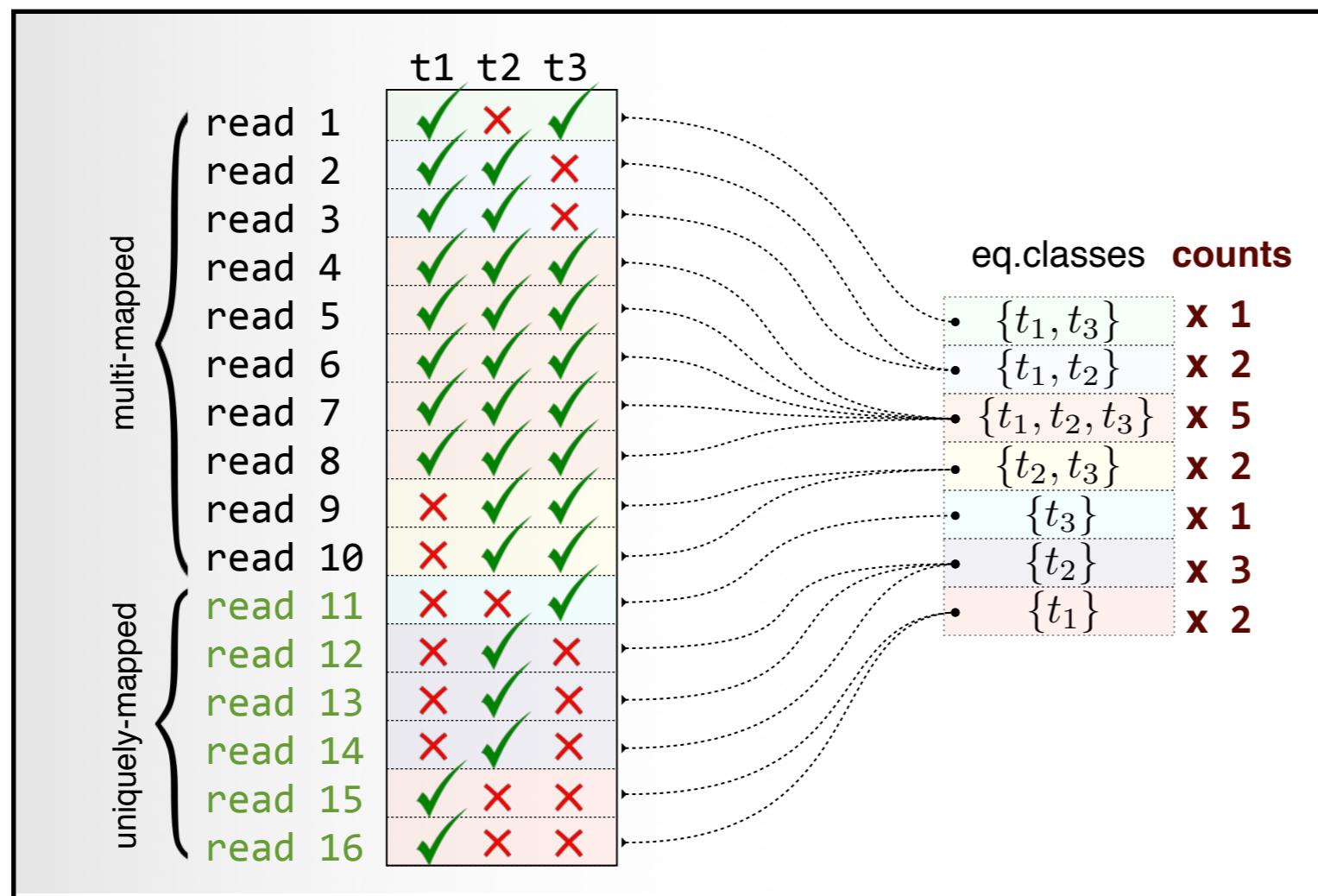


- Ntranos, V., Kamath, G.M., Zhang, J., Pachter, L., Tse, D.,N., “**Fast and accurate single-cell RNA-seq analysis by clustering of transcript-compatibility counts**,” **Genome Biology**, May 2016

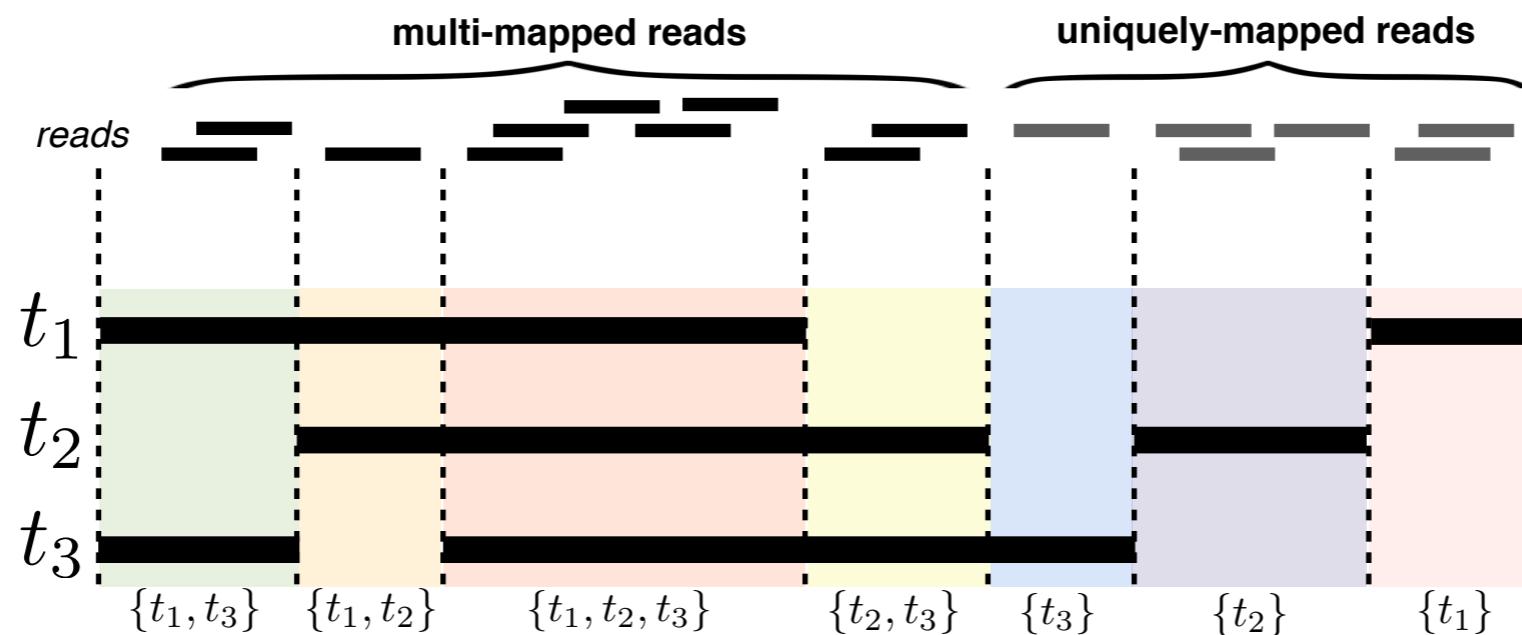
Equivalence Classes and Counts

Definitions:

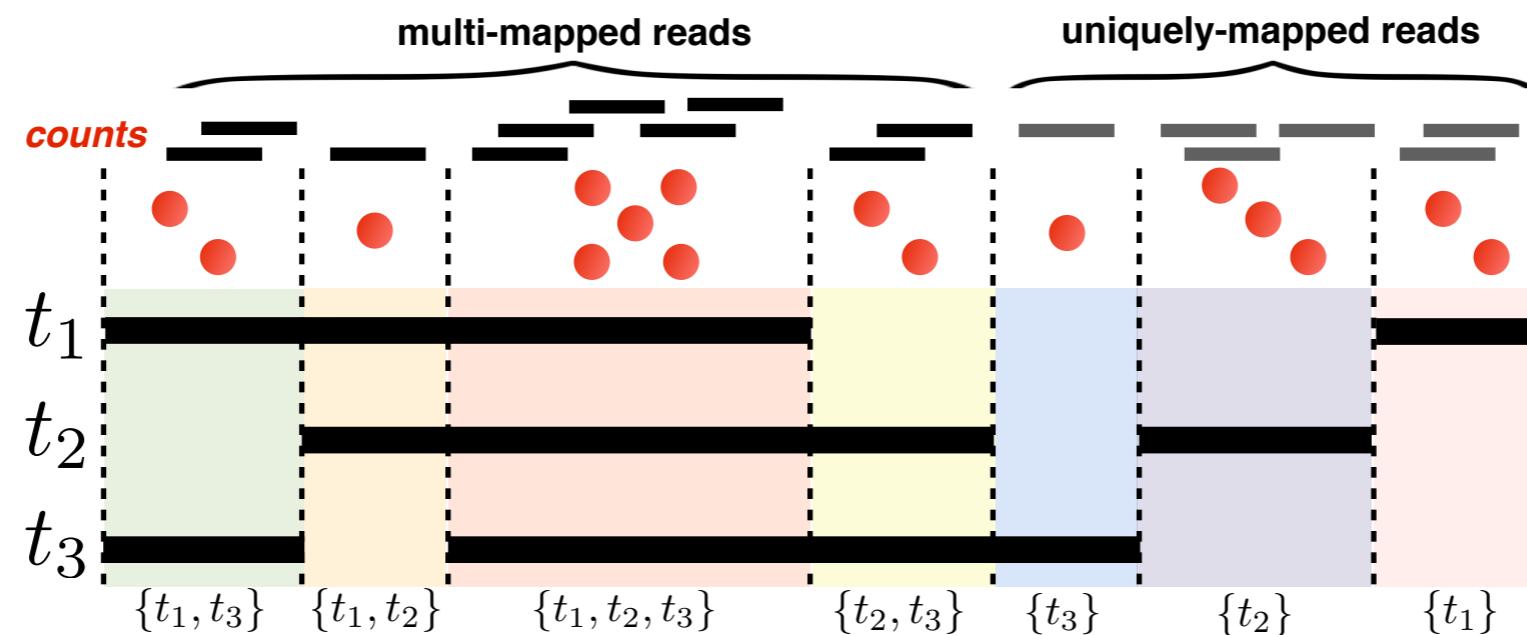
- Each read is *compatible* with one or more transcripts (i.e., the transcripts it could possibly have come from).
- An **equivalence class** is a group of reads that are *compatible* with the same set of transcripts.
- The **transcript-compatibility counts** are the *total number of reads* in each equivalence class.



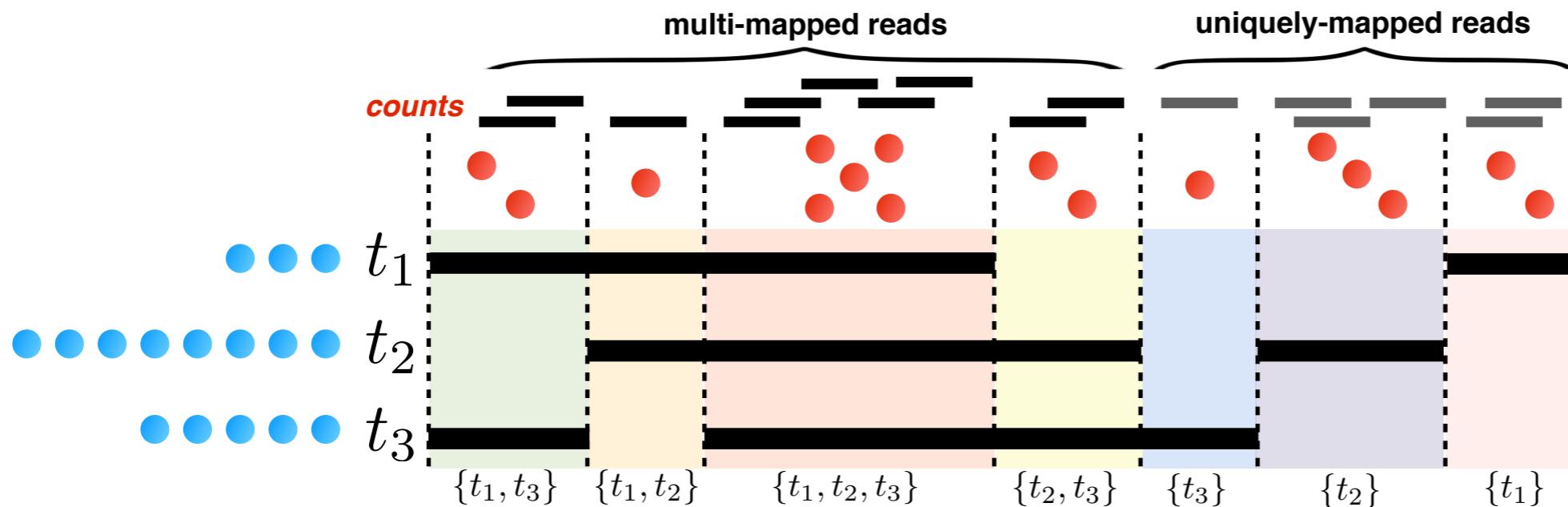
Pseudoalignment



Pseudoalignment



TCC Model



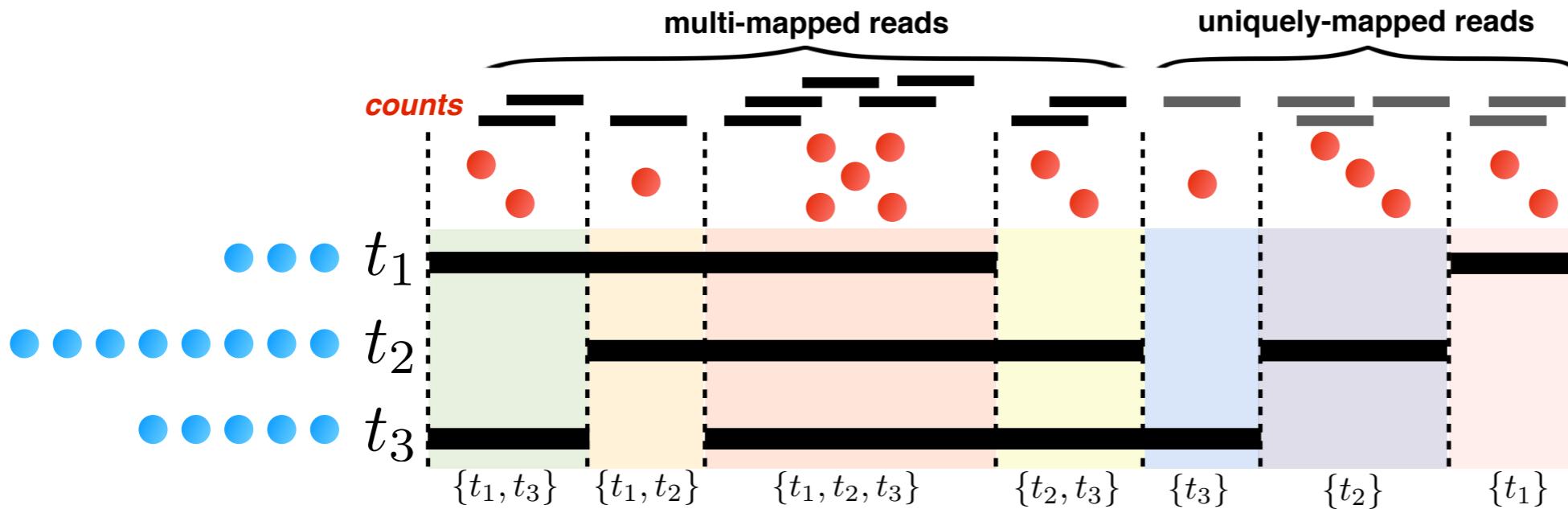
Underlying Model (ground truth):

$$p_t^* = \text{transcripts}$$

$$P_{S|T}^* = \begin{matrix} \text{transcripts} \\ \text{equiv. classes} \end{matrix}$$

$$p_S^* = \begin{matrix} \text{equiv. classes} \\ \text{classes} \end{matrix}$$

TCC Model



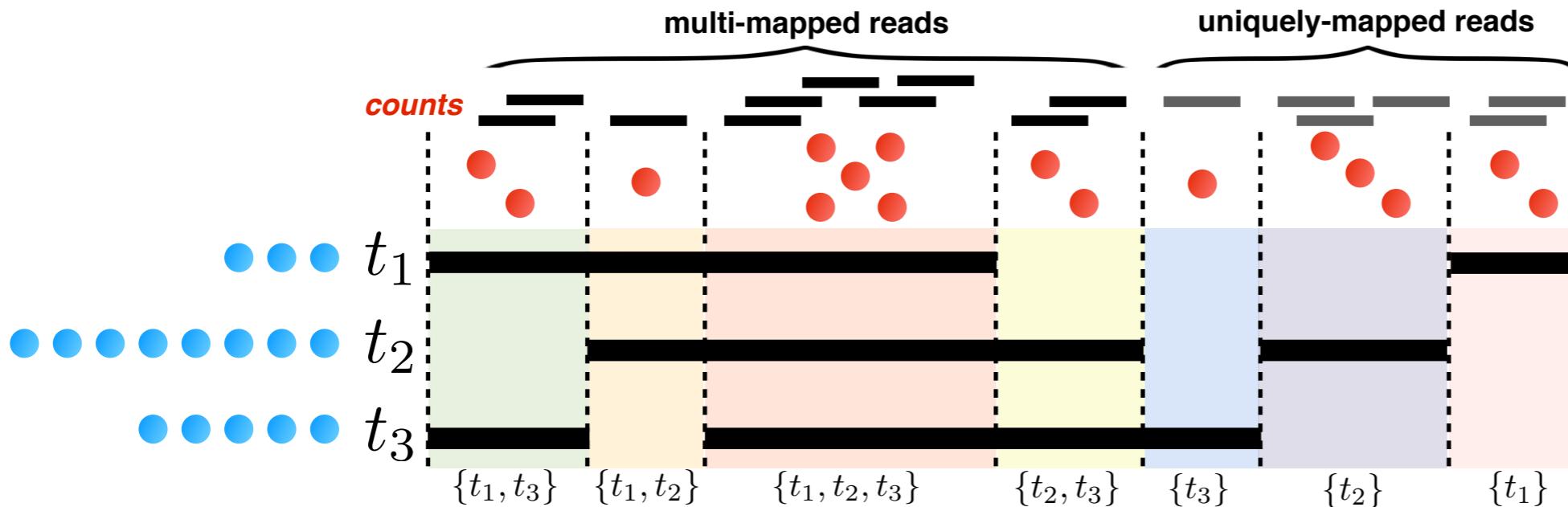
Underlying Model (ground truth):

$$p_S^* = \text{equiv. classes} \quad P_{S|T}^* \quad p_t^*$$

transcripts

equiv. classes

TCC Model

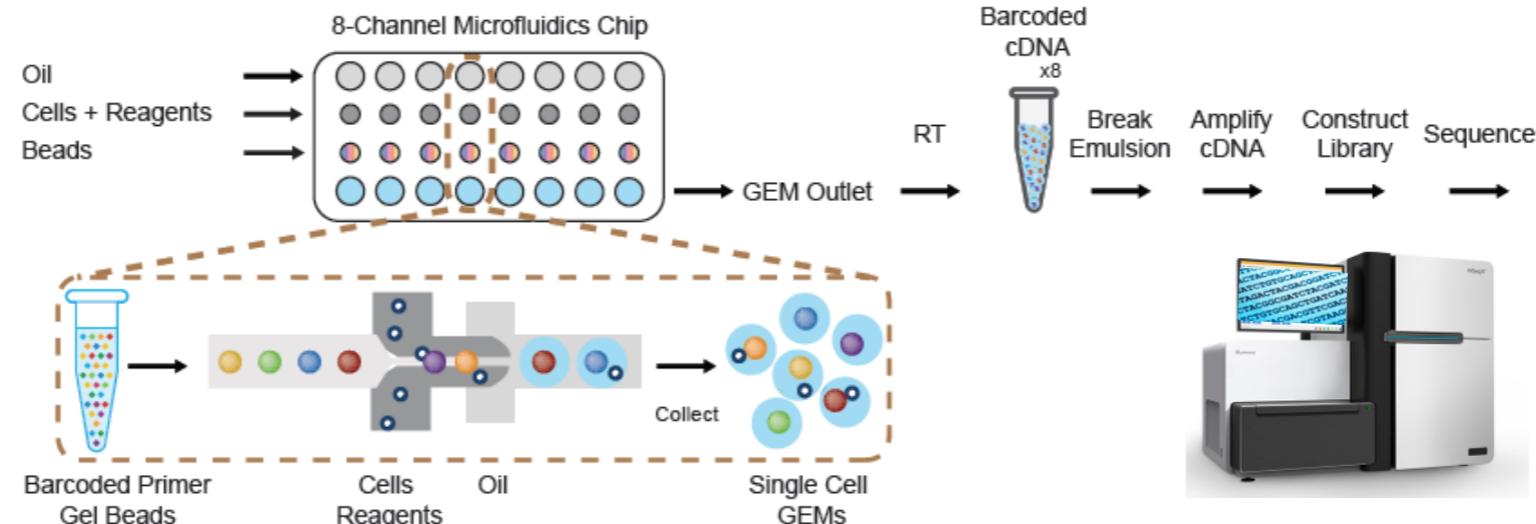


Underlying Model (ground truth):

$$\begin{array}{ccc} \text{cells} & \text{transcripts} & \text{cells} \\ P_S^* & = & P_{S|T}^* \\ \text{equiv. classes} & & \text{equiv. classes} \\ & & P_t^* \end{array}$$

TCC Model

for example:



10X GENOMICS®

Underlying Model (ground truth): 10x Chromium

$$\text{equiv. classes } P_S^* = \text{equiv. classes } P_{S|T}^* \text{ cells } P_t^*$$

TCC Model



Underlying Model (ground truth):

Fluidigm C1

$$P_S^* = P_{S|T}^*$$

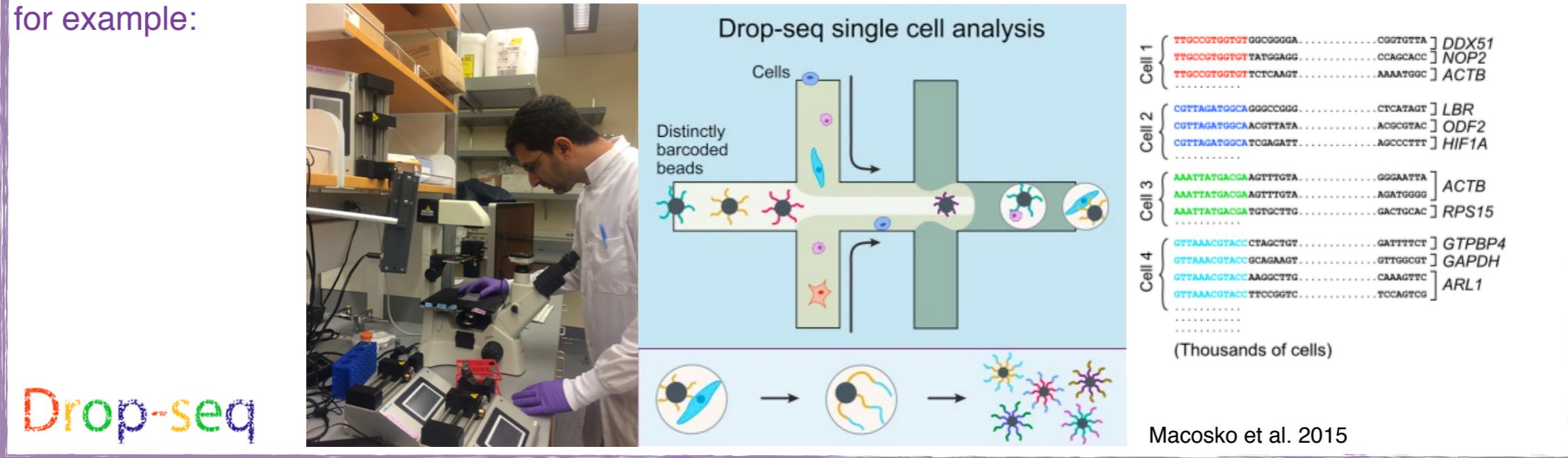
cells
equiv. classes

transcripts
equiv. classes

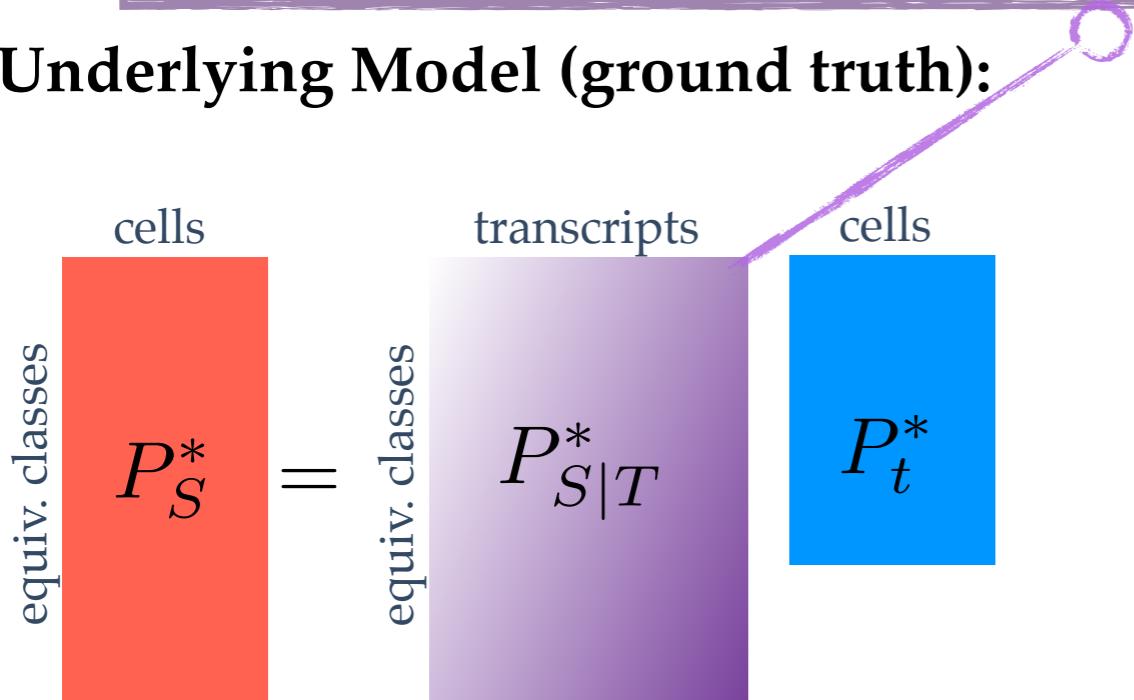
cells
 P_t^*

TCC Model *

for example:

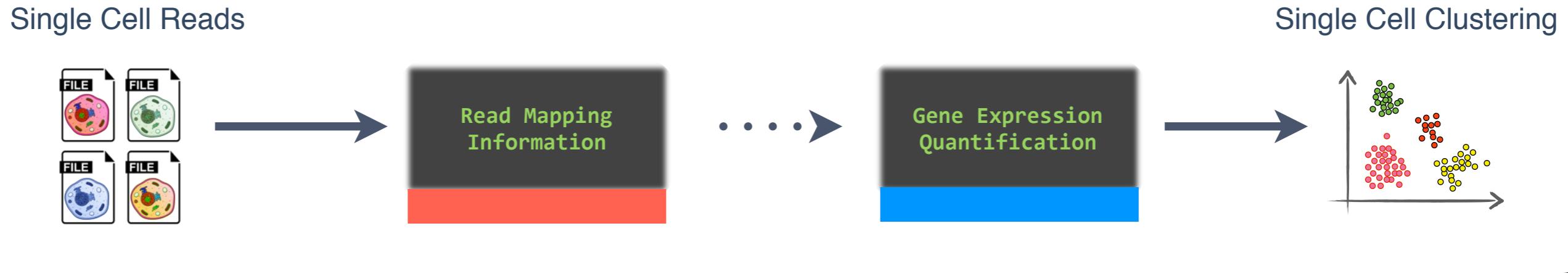


Underlying Model (ground truth): Drop-Seq



*Disclaimer: All characters and events in this slide, even those based on real people, are entirely fictional.

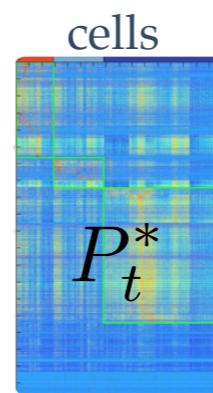
TCC — Single-cell processing workflow



main idea 

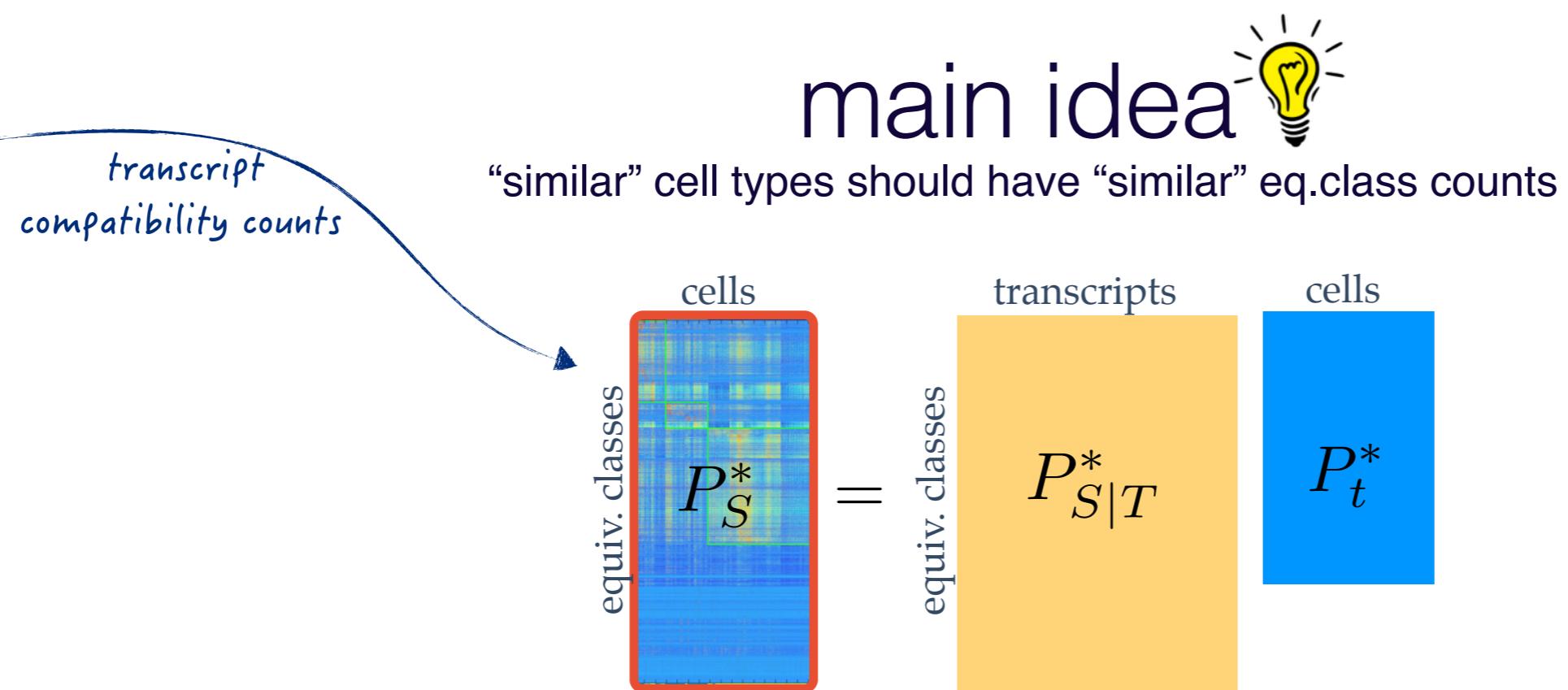
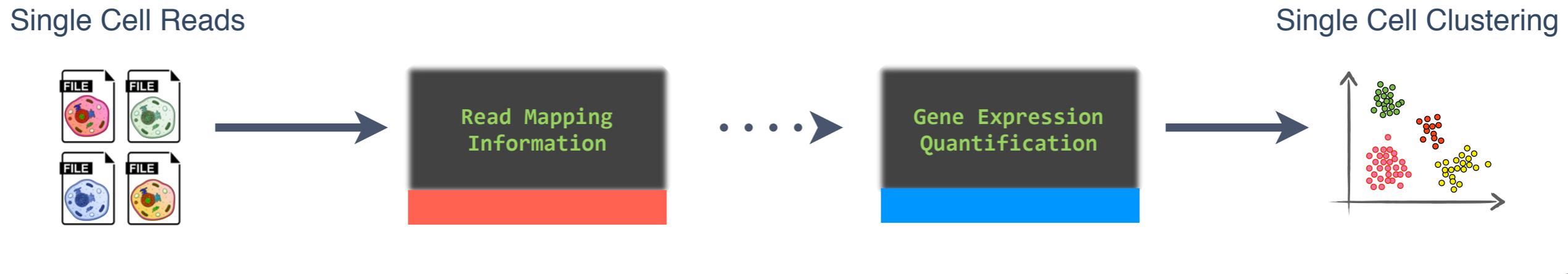
“similar” cell types should have “similar” eq.class counts

$$\begin{matrix} \text{cells} \\ \text{equiv. classes} \end{matrix} P_S^* = \begin{matrix} \text{transcripts} \\ \text{equiv. classes} \end{matrix} P_{S|T}^*$$

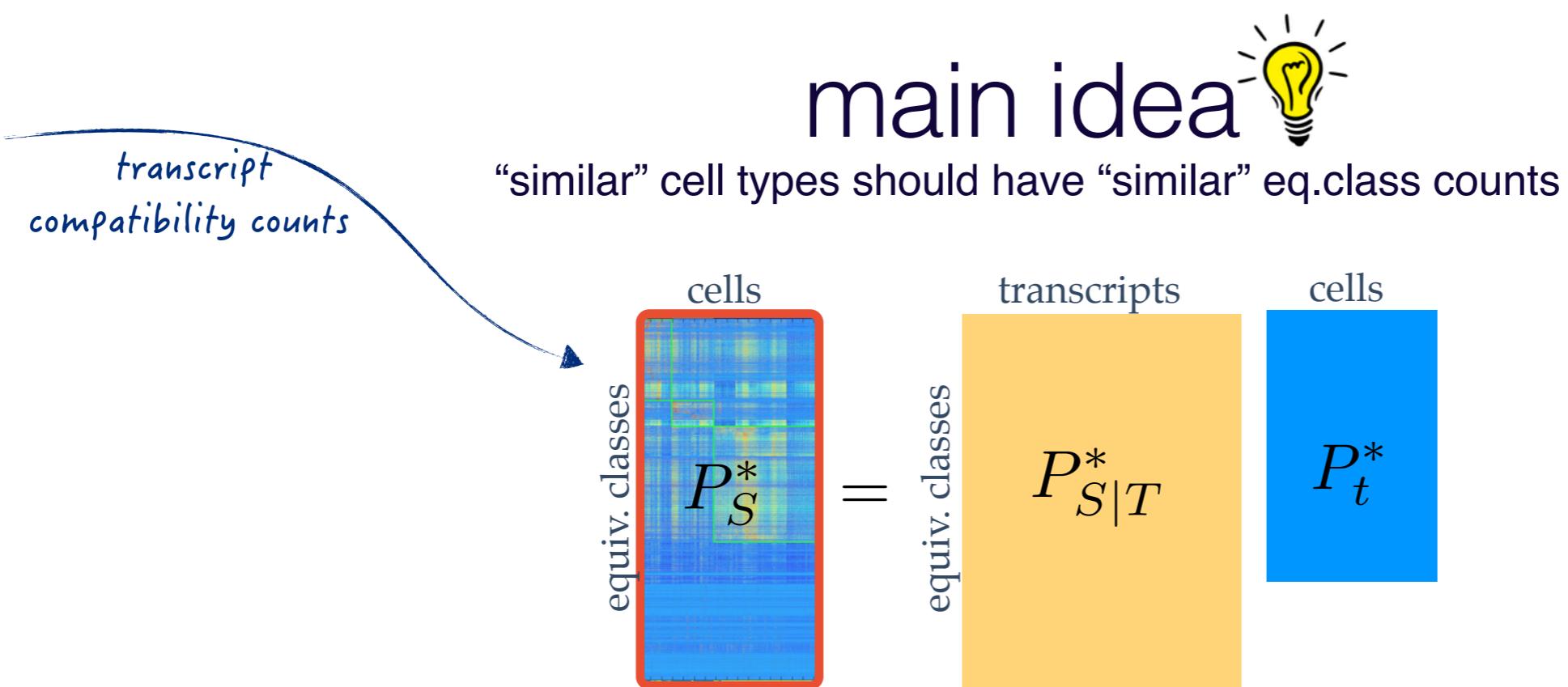
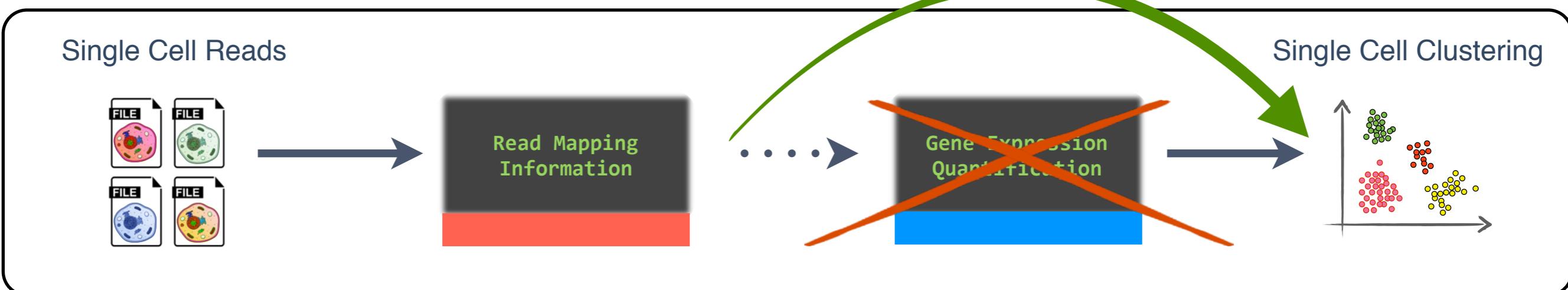


isoform / gene
expression profiles

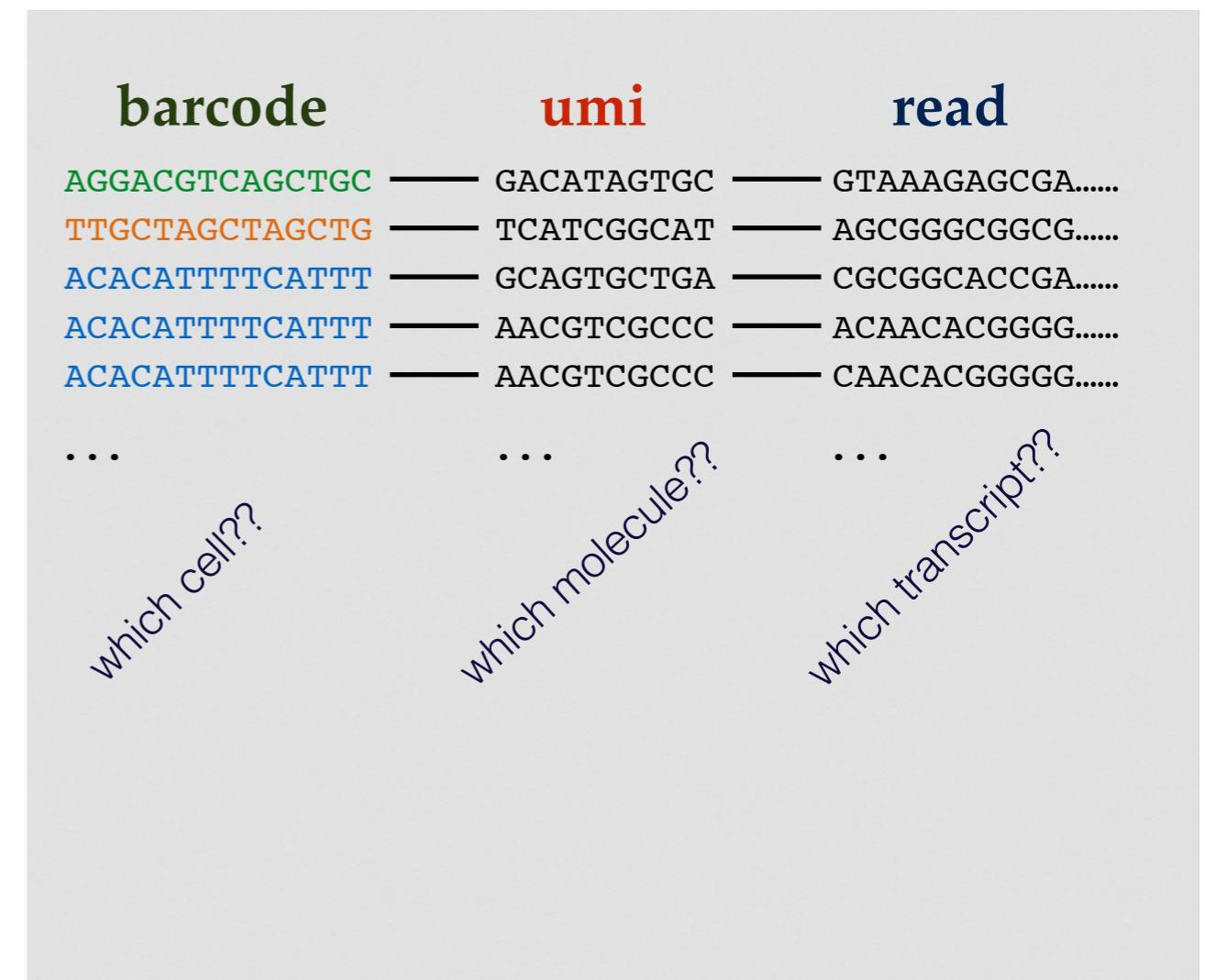
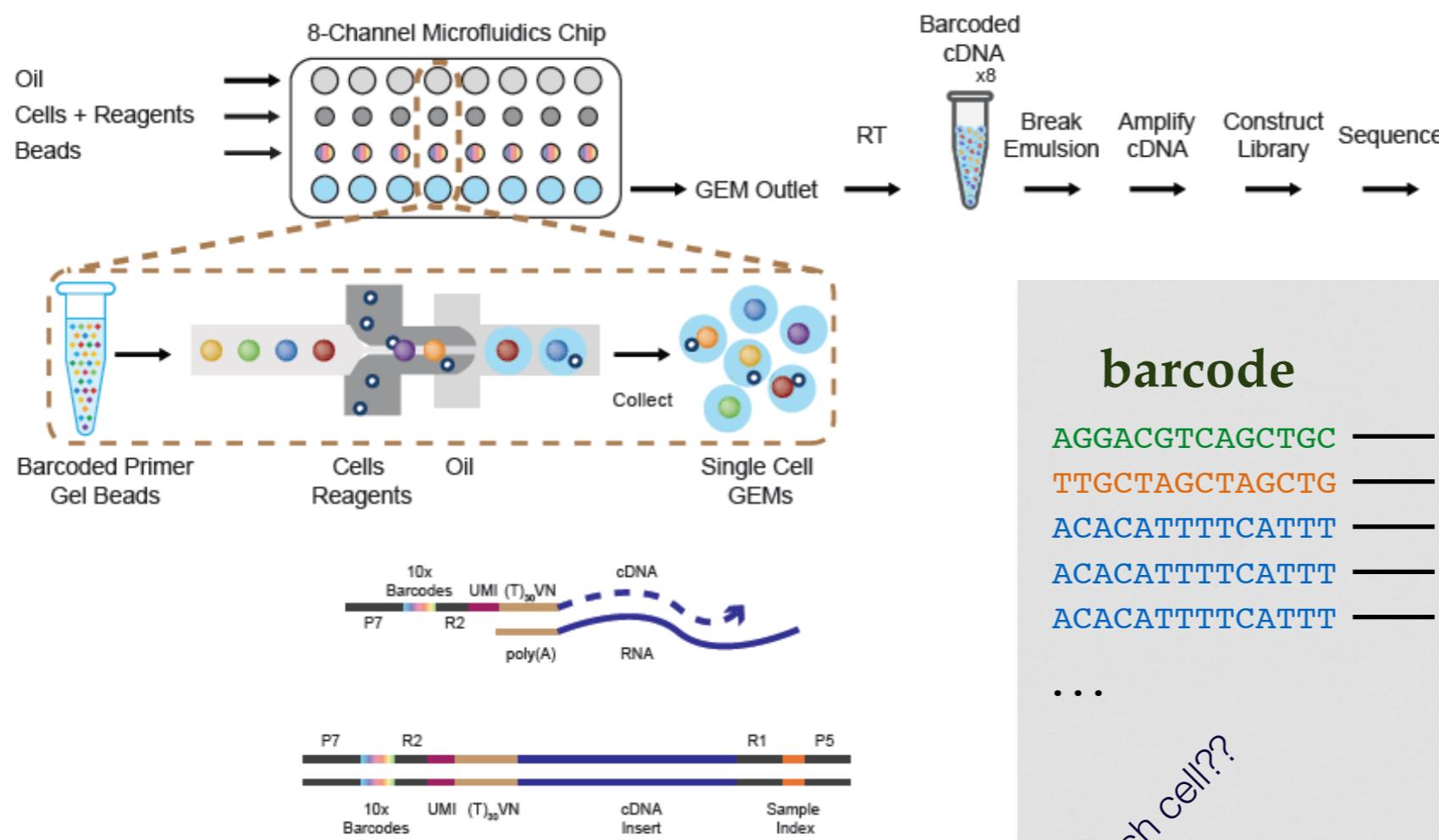
TCC — Single-cell processing workflow



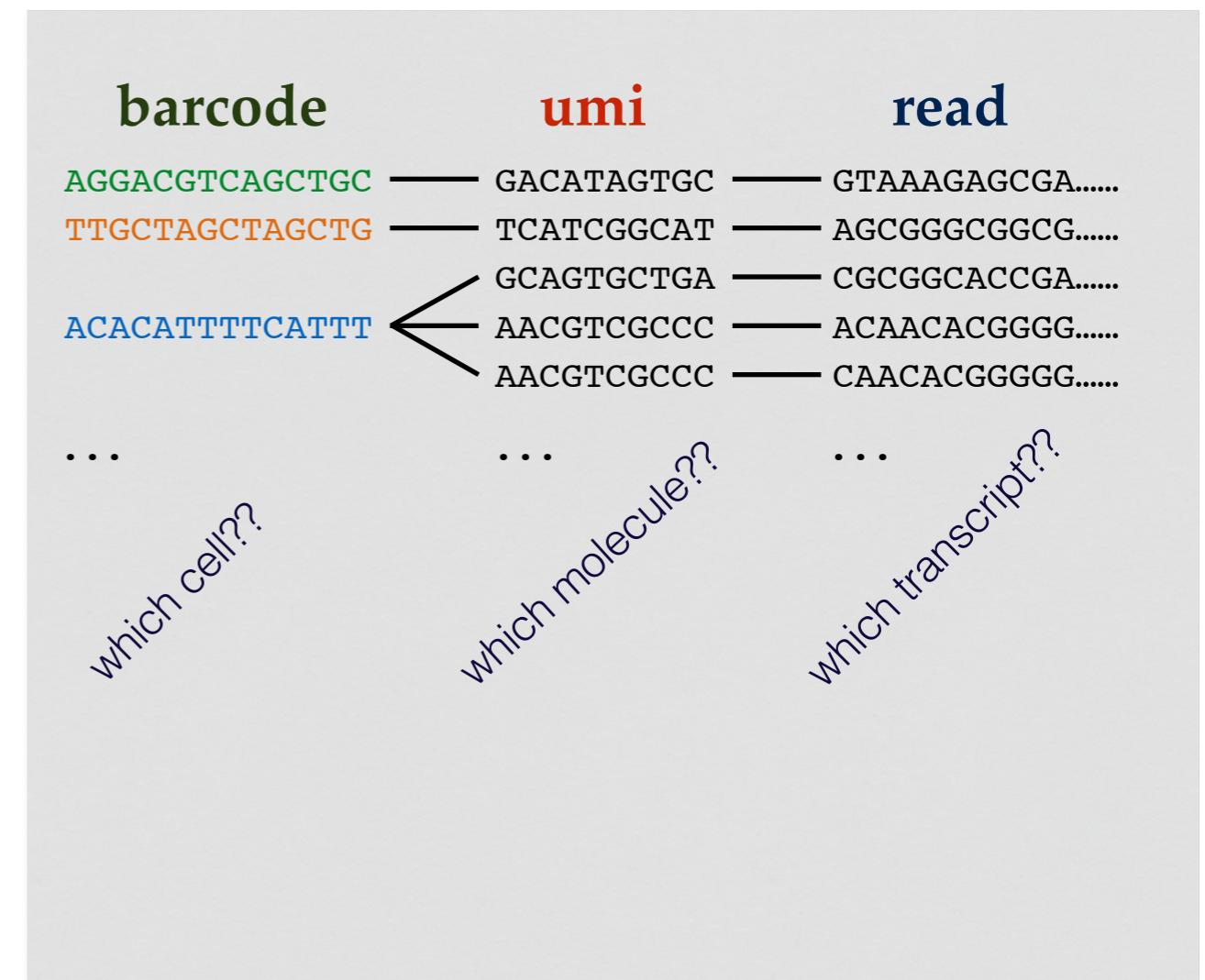
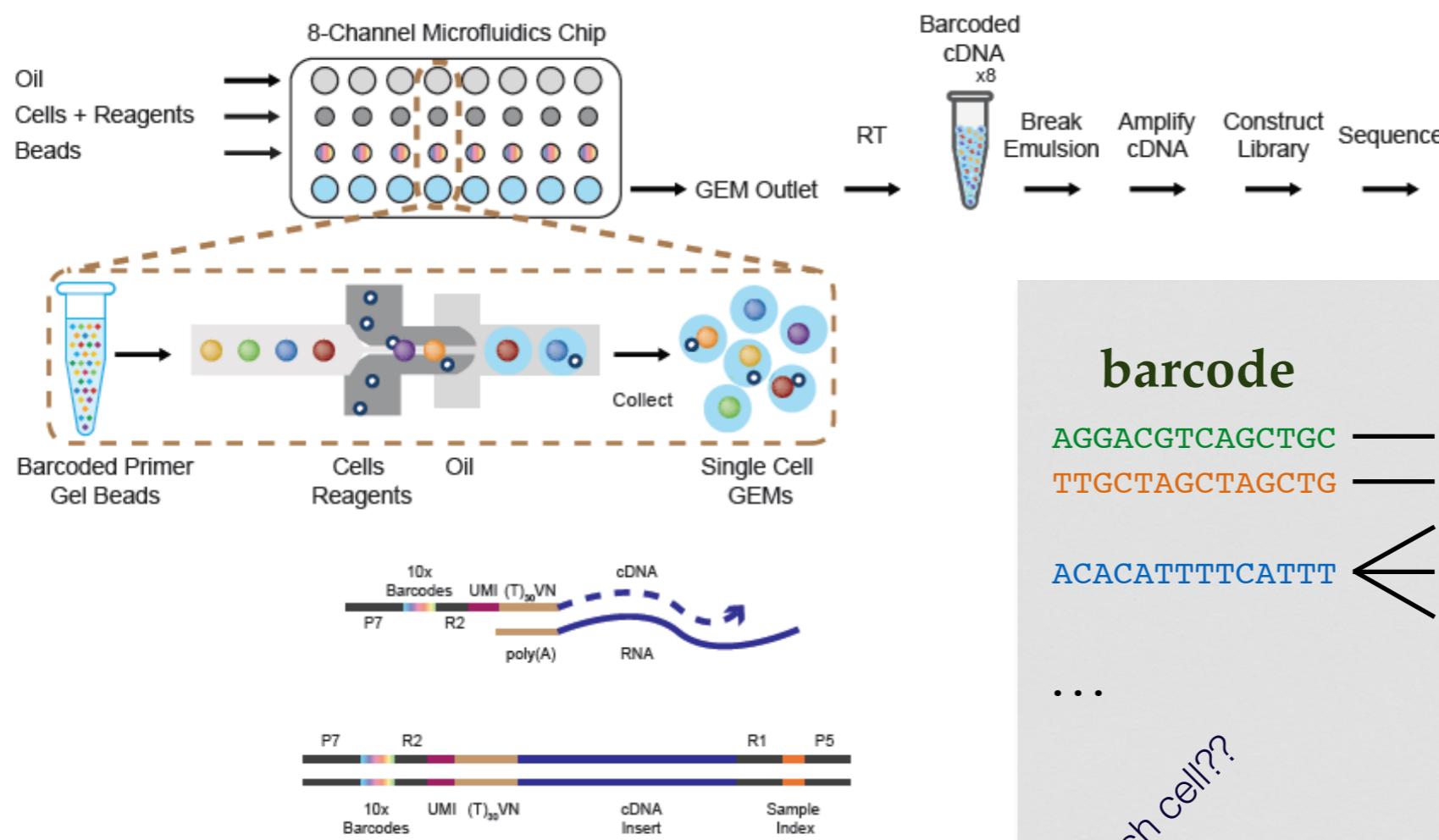
TCC — Single-cell processing workflow



Overview — 10x Chromium data



Overview — 10x Chromium data



TCC pipeline for 10x data



SI-3A-A10 (sample index) --> [ACAGCAAC, CGCAATT, GAGTTGCG, TTTCGCGA]

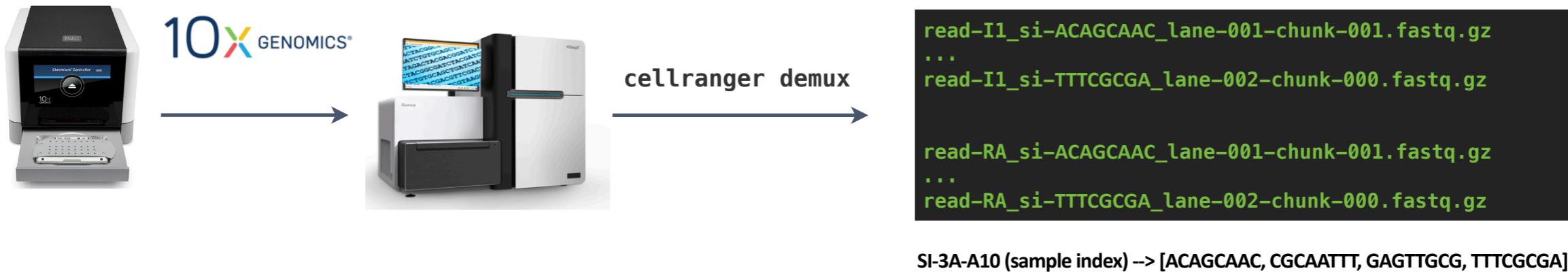
github.com/pachterlab/scRNA-Seq-TCC-prep/

v0.2

Python scripts and jupyter notebooks to:

- Detect and Correct Cell Barcodes
- Generate single-cell fastqs
- Get distinct UMI-TCC matrix
- Prep for Clustering and Analysis

TCC pipeline for 10x data



github.com/pachterlab/scRNA-Seq-TCC-prep/

v0.2

Python scripts and jupyter notebooks to:

- Detect and Correct Cell Barcodes
- Generate single-cell Fastqs
- Get distinct UMI-TCC matrix
- Prep for Clustering and Analysis

notebooks/

10xGet_cell_barcodes.ipynb

10xResults.ipynb

source/

10xPrep_Data.py

error_correct_and_split.py

compute_TCCs.py

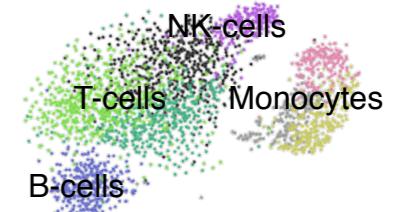
prep_TCC_matrix.py

Step 0: Find a 10xDataset, e.g.,

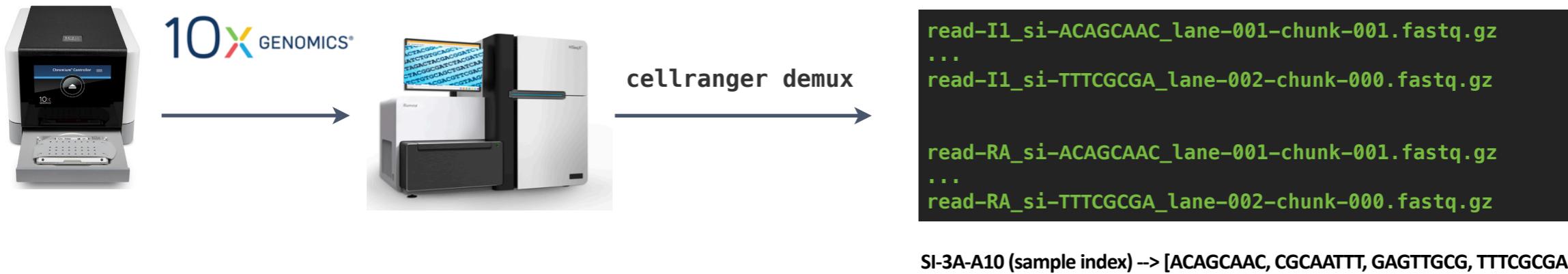
[Peripheral Blood Mononuclear Cells \(PBMCs\) from a Healthy Donor](#)

Step 1: In your project dir, copy the **source/** files and the two notebooks and edit a **config.json** file:

```
{  
    "NUM_THREADS": 8,  
    "WINDOW": [500, 5000],  
    "SOURCE_DIR": "/path/to/source/",  
    "BASE_DIR": "/path/to/fastq_files/",  
    "sample_idx": ["ATCGCTCC", "CCGTACAG", "GATAGGTA", "TGACTAGT"],  
    "SAVE_DIR": "/output/dir/for/save_data/",  
    "dmin": 4,  
    "BARCODE_LENGTH": 14,  
    "OUTPUT_DIR": "/output/dir/for/singlecell_fastqs/",  
    "kallisto":{  
        "binary": "/path/to/kallisto",  
        "index": "/path/to/kallisto_index.idx",  
        "TCC_output" : "/output/dir/for/TCC_output/"  
    }  
}
```



TCC pipeline for 10x data



github.com/pachterlab/scRNA-Seq-TCC-prep/

v0.2

Python scripts and jupyter notebooks to:

- Detect and Correct Cell Barcodes
- Generate single-cell Fastqs
- Get distinct UMI-TCC matrix
- Prep for Clustering and Analysis

notebooks/

10xGet_cell_barcodes.ipynb

10xResults.ipynb

source/

10xPrep_Data.py

error_correct_and_split.py

compute_TCCs.py

prep_TCC_matrix.py

Step 0: Find a 10xDataset, e.g.,

[Peripheral Blood Mononuclear Cells \(PBMCs\) from a Healthy Donor](#)

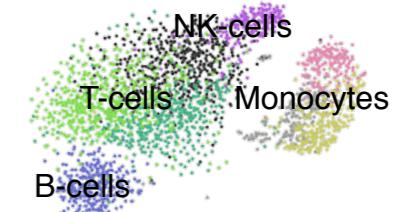
Step 1: In your project dir, copy the **source/** files and the two notebooks and edit a **config.json** file:

TCC pipeline:

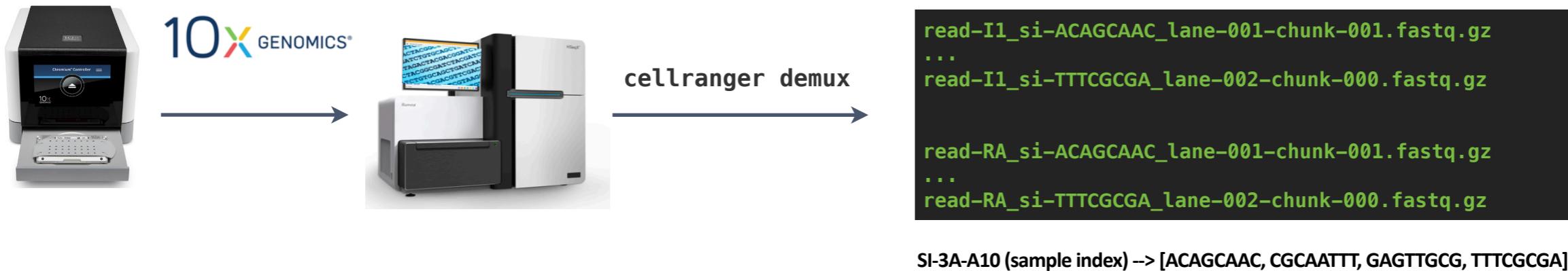
10xGet_cell_barcodes.ipynb

python source/**10xPrep_Data.py** config.json

10xResults.ipynb



TCC pipeline for 10x data



github.com/pachterlab/scRNA-Seq-TCC-prep/

v0.2

Python scripts and jupyter notebooks to:

- Detect and Correct Cell Barcodes
- Generate single-cell Fastqs
- Get distinct UMI-TCC matrix
- Prep for Clustering and Analysis

notebooks/

10xGet_cell_barcodes.ipynb

10xResults.ipynb

source/

10xPrep_Data.py

error_correct_and_split.py

compute_TCCs.py

prep_TCC_matrix.py

Step 0: Find a 10xDataset, e.g.,

[Peripheral Blood Mononuclear Cells \(PBMCs\) from a Healthy Donor](#)

Step 1: In your project dir, copy the **source/** files and the two notebooks and edit a **config.json** file:

TCC pipeline:

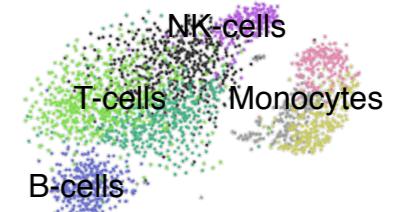
10xGet_cell_barcodes.ipynb

python source/error_correct_and_split.py config.json

python source/compute_TCCs.py config.json

python source/prep_TCC_matrix.py config.json

10xResults.ipynb



TCC pipeline for 10x data



github.com/pachterlab/scRNA-Seq-TCC-prep/

v0.2

Python scripts and jupyter notebooks to:

- Detect and Correct Cell Barcodes
- Generate single-cell Fastqs
- Get distinct UMI-TCC matrix
- Prep for Clustering and Analysis

notebooks/

10xGet_cell_barcodes.ipynb
10xResults.ipynb

source/

10xPrep_Data.py
error_correct_and_split.py
compute_TCCs.py
prep_TCC_matrix.py

Step 0: Find a 10xDataset, e.g.,

[Peripheral Blood Mononuclear Cells \(PBMCs\) from a Healthy Donor](#)

Step 1: In your project dir, copy the **source/** files and the two notebooks and edit a **config.json** file:

TCC pipeline:

```
10xGet_cell_barcodes.ipynb
python source/error_correct_and_split.py config.json
generate single cell files
python source/compute_TCCs.py config.json
kallisto pseudoalignment of all cells
python source/prep_TCC_matrix.py config.json
10xResults.ipynb
pwise dist, clustering and analysis
```

