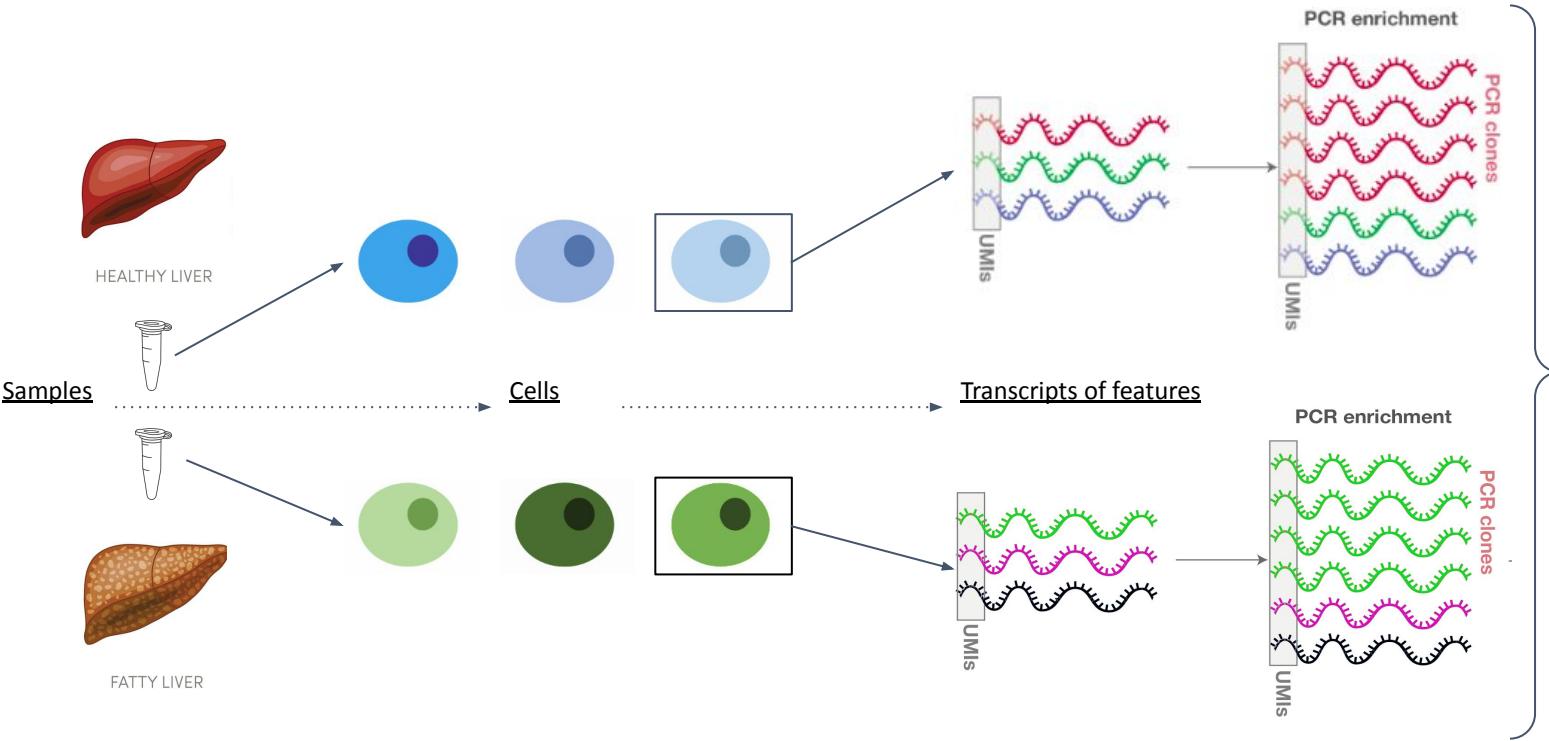

From Raw Read Data to Cell Clustering: Cell ranger and Seurat

David Carbonez

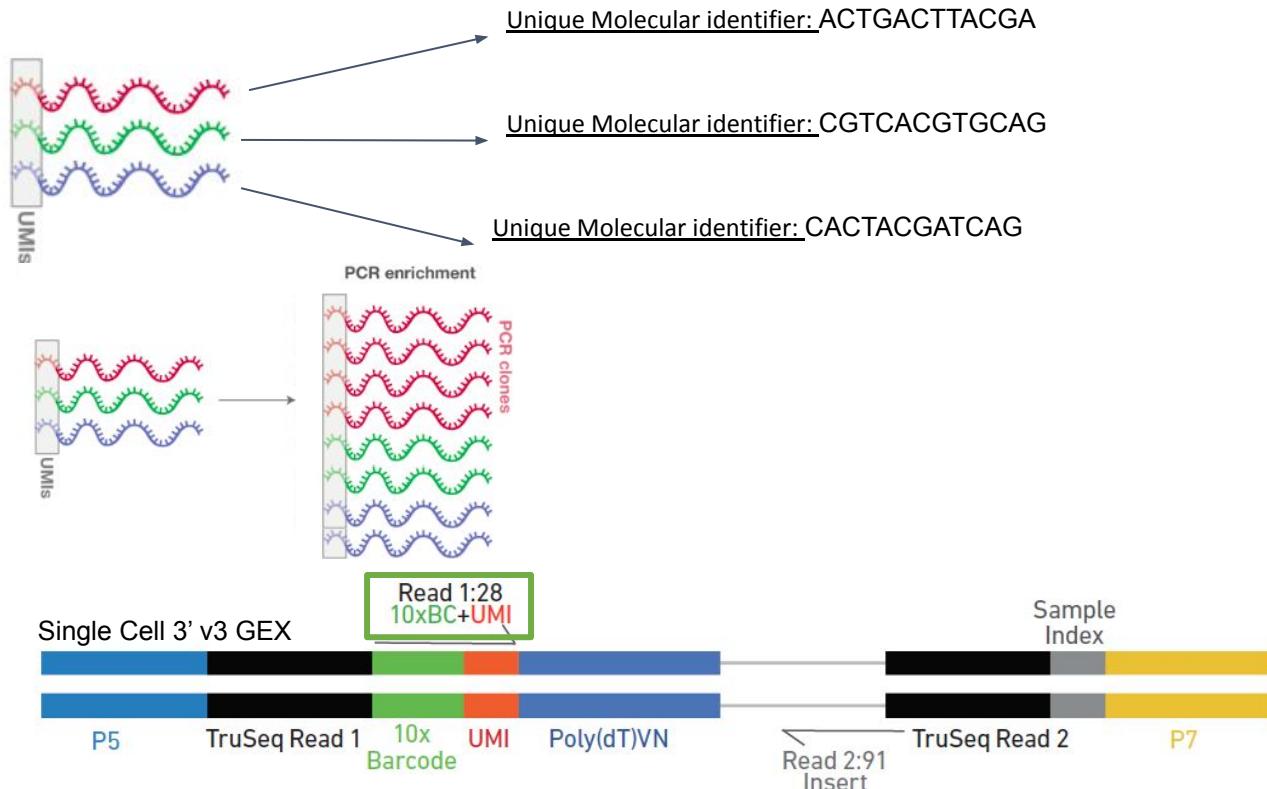
From Raw read to Cell Clustering

1. Part 1: Preparation of the data
 - a. Introduction & Terminology
 - b. The anatomy of a 10x Construct
 - c. Sequencing
 - d. Demultiplexing & Counting
2. Part 2: Analysis of the data

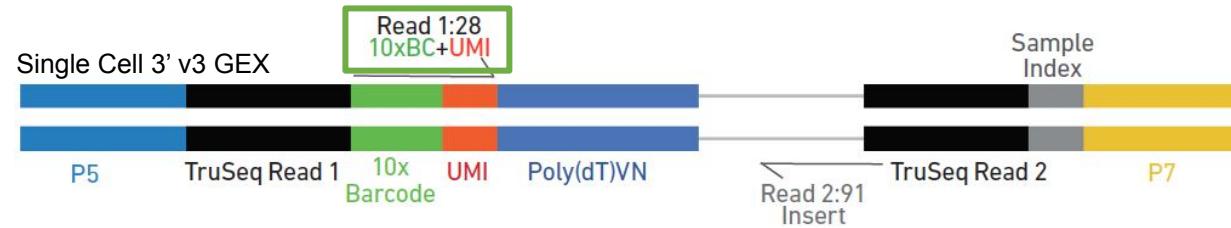
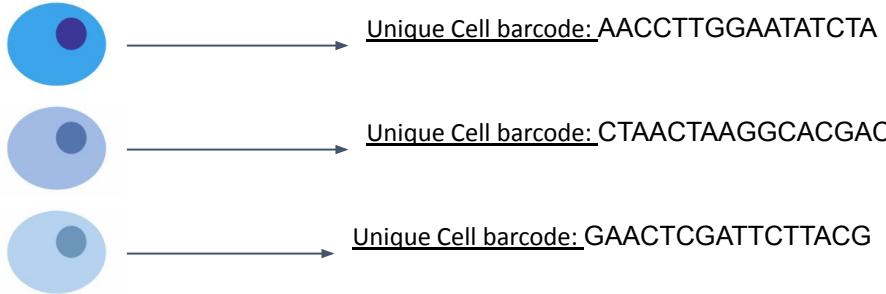
Introduction



UMI

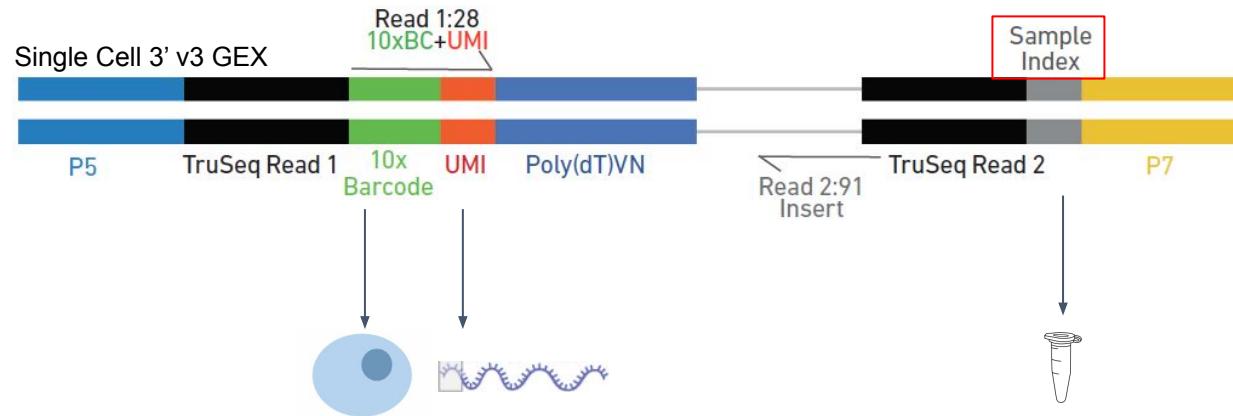
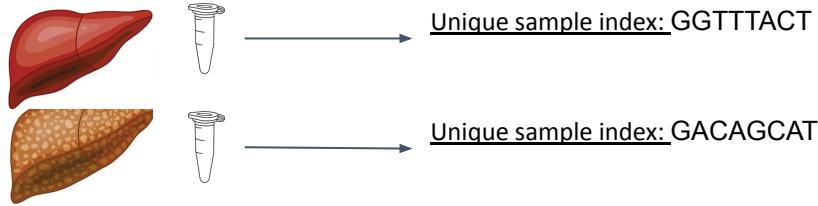


Cell Barcode



Sample index

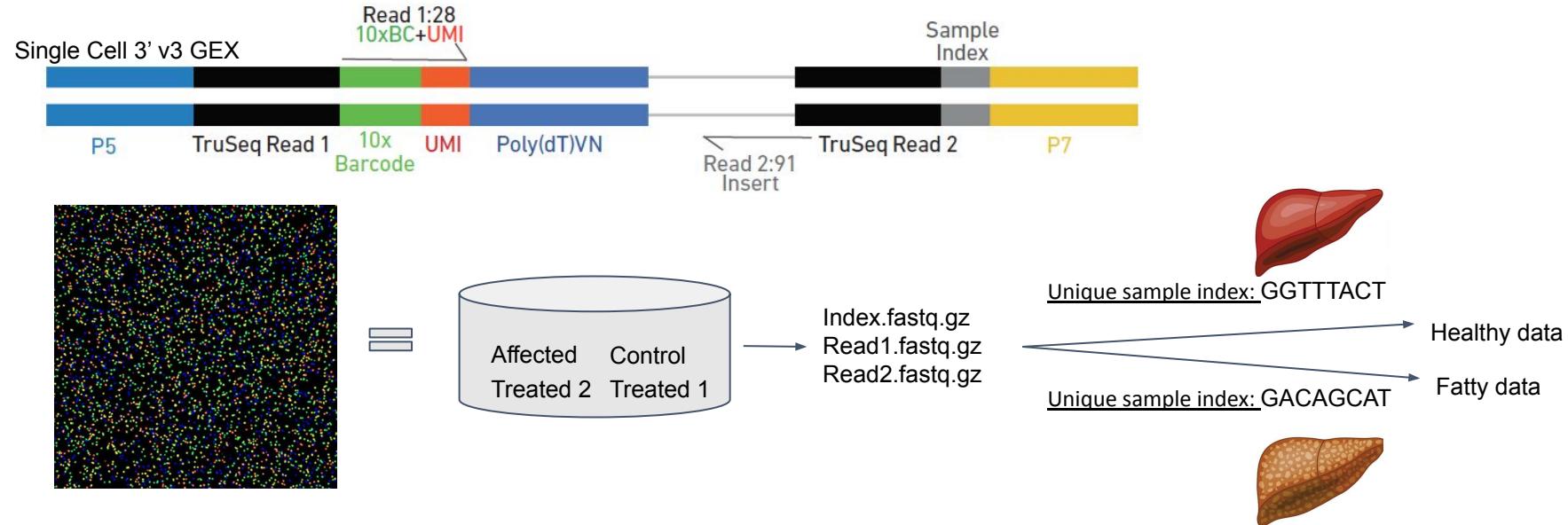
Samples



Sequencing



BioIT analysis: Demultiplexing



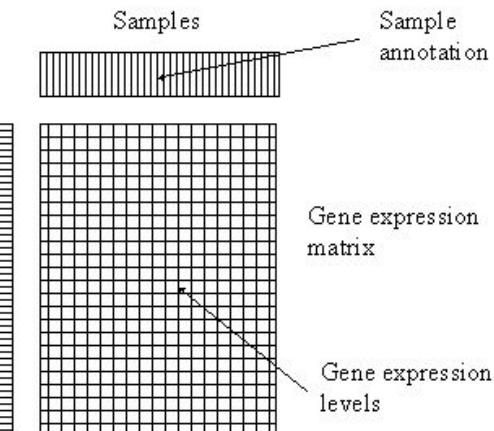
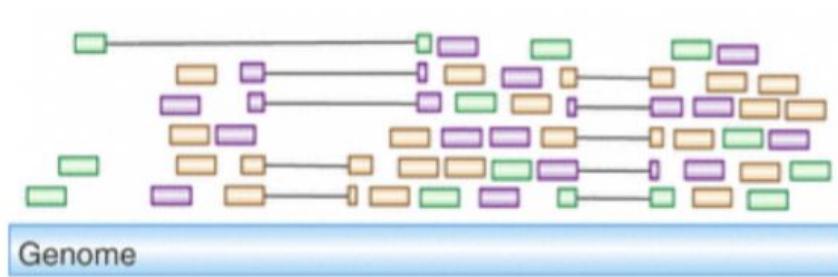
BioIT analysis: Counting

Healthy_I1.fastq.gz
Healthy_R1.fastq.gz
Healthy_R2.fastq.gz

Sample indices (8bp)

→ 10x Barcode (16bp) + UMI (12 bp)
mRNA (91 bp)

Identify sample
Identify cell + Transcript
Identify gene



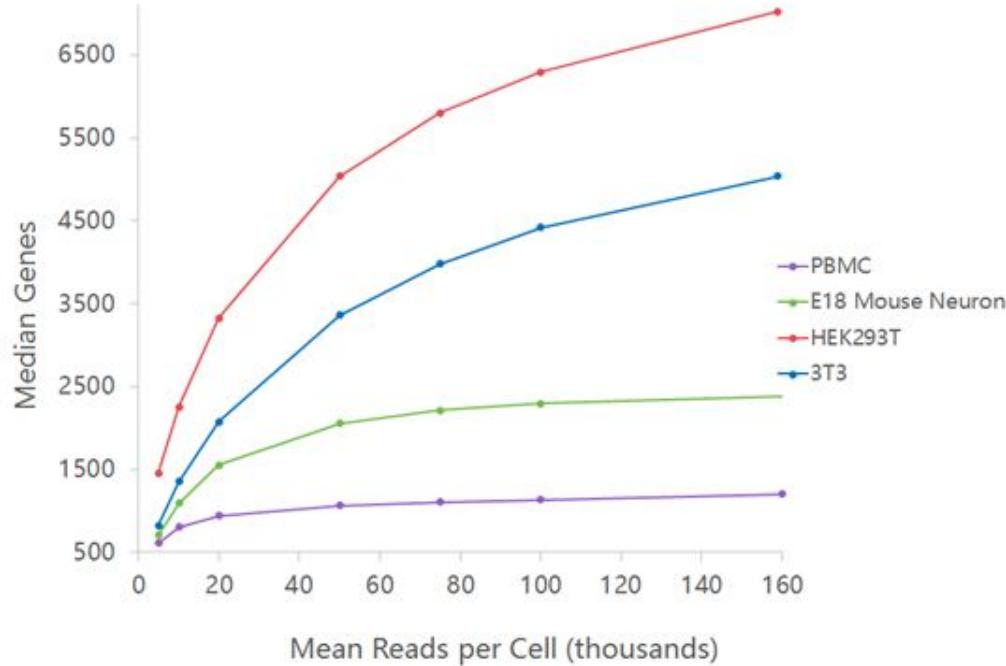
Sequencing requirements

How many reads/cell?



Saturation =

'a measure of the fraction of library complexity that was sequenced in a given experiment'



Sequencing requirements

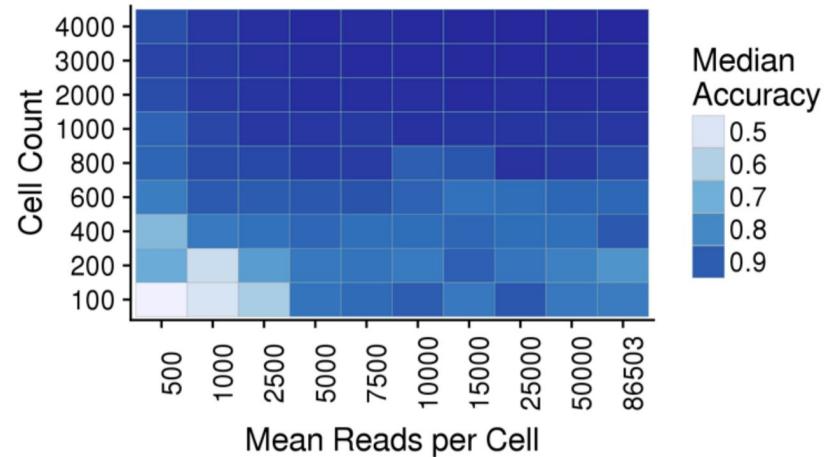
1. Overview of the cell types

mRNA content ↑ → required reads/cell ↑

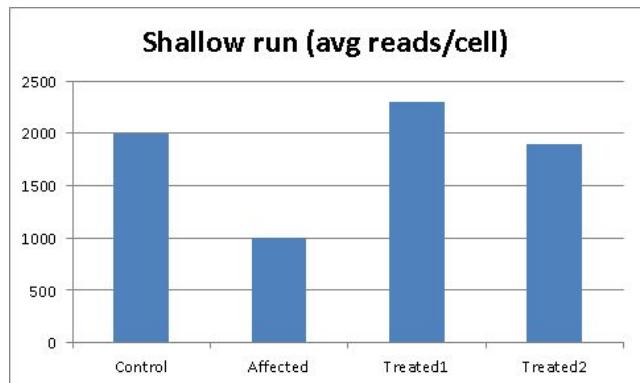
RARE cell types: Cell count ↓ → required reads/cell ↑

→ Usually ~25k reads/cell

2. Detect as many genes as possible



BioIT analysis: Shallow sequencing



Estimated Number of Cells

3,868

Mean Reads per Cell

4,939

Median Genes per Cell

722

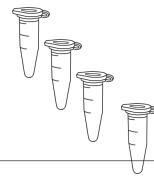
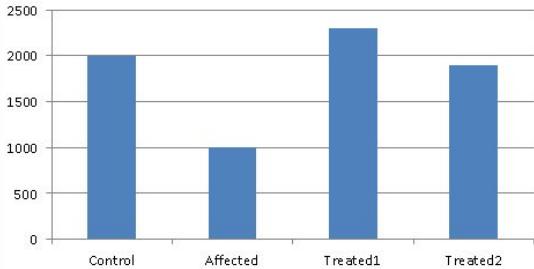
Sequencing

Number of Reads	19,107,635
Valid Barcodes	97.7%
Sequencing Saturation	4.7%
Q30 Bases in Barcode	97.2%
Q30 Bases in RNA Read	92.5%
Q30 Bases in Sample Index	96.3%
Q30 Bases in UMI	96.9%

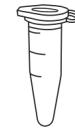


BioIT analysis: Deep sequencing

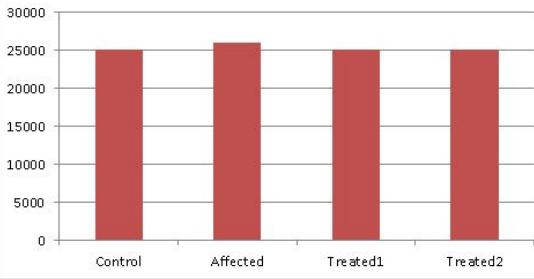
Shallow run (avg reads/cell)



Adjust pooling



Deep run (avg reads/cell)



Estimated Number of Cells

4,824

Mean Reads per Cell

57,119

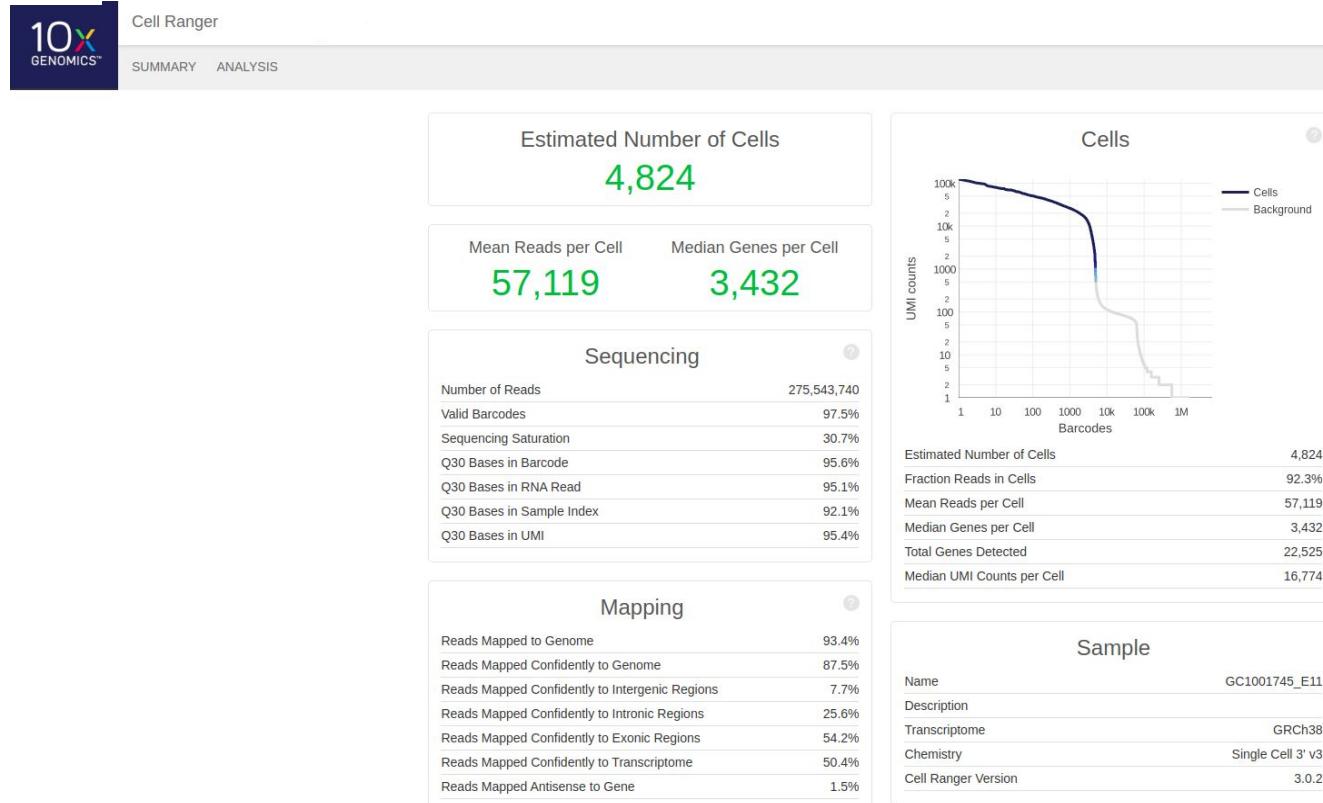
Median Genes per Cell

3,432

Sequencing

Number of Reads	275,543,740
Valid Barcodes	97.5%
Sequencing Saturation	30.7%
Q30 Bases in Barcode	95.6%
Q30 Bases in RNA Read	95.1%
Q30 Bases in Sample Index	92.1%
Q30 Bases in UMI	95.4%

Web summary

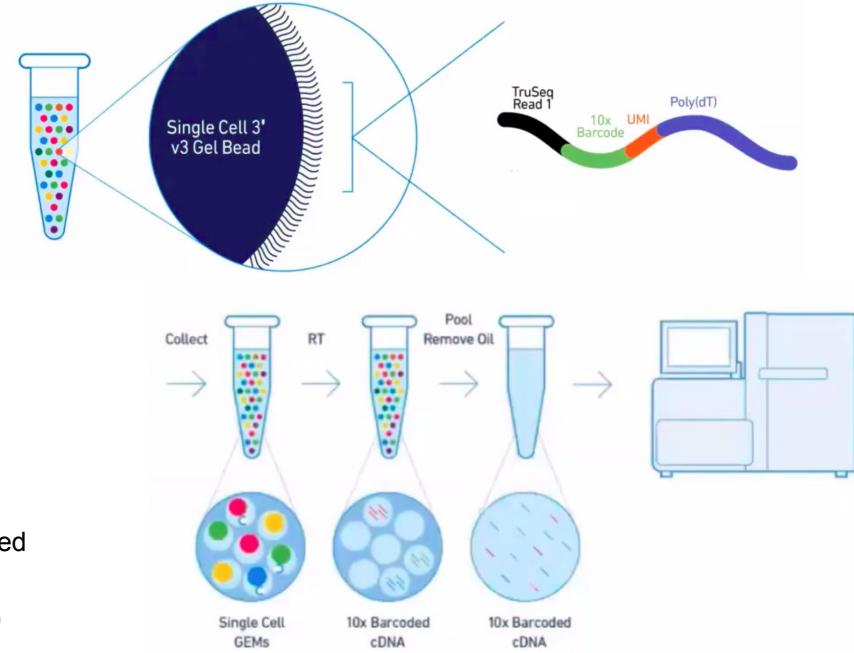
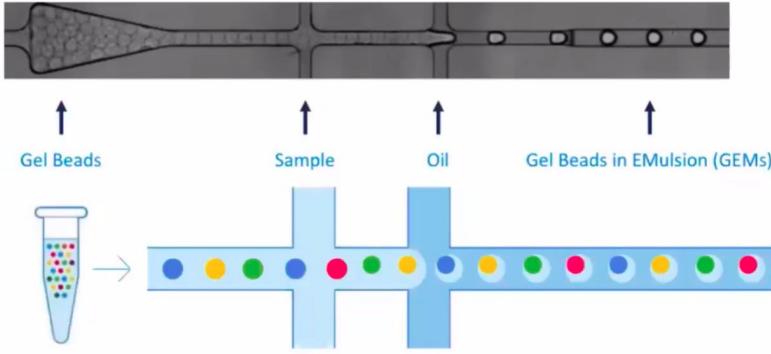


Analysing the data itself through Seurat

Analysing the data itself through Seurat

1. Part 1: Preparation of the data
2. Part 2: Analysis of the data
 - a. Filtering
 - b. Feature selection & Dimension reduction
 - c. Clustering
 - d. Identification
 - e. Subsetting
 - f. Integration
 - g. Further analysis

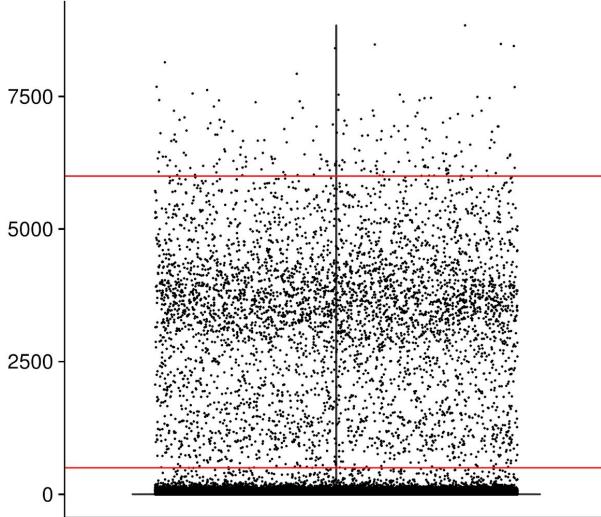
Recap: Library Prep 10x



BioIT analysis: Filtering

nFeature_RNA

Lower limit: 500 – Upper limit: 6000



Multiplet



Correct



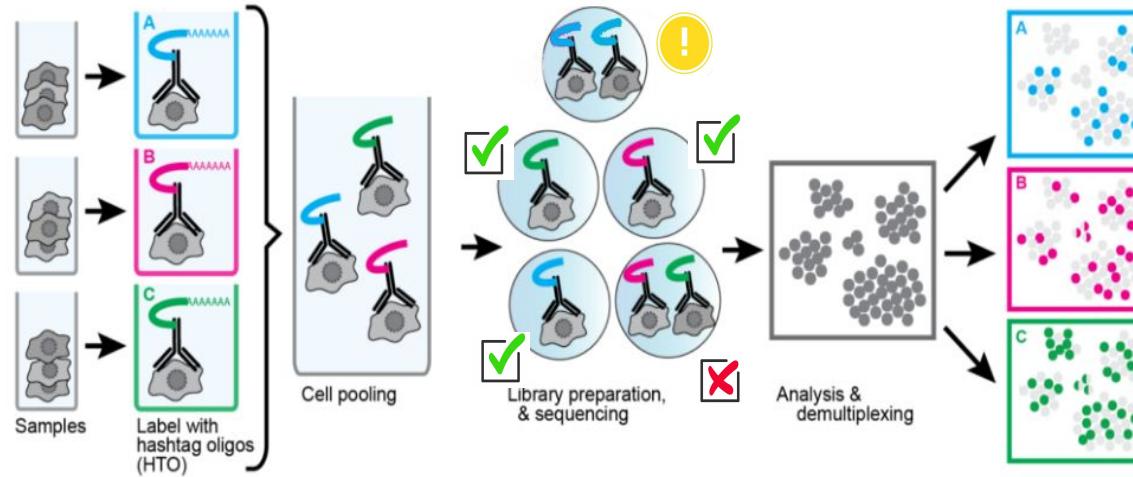
Empty GEM



Multiplet Rate (%)	# of Cells Loaded	# of Cells Recovered
~0.4%	~870	~500
~0.8%	~1700	~1000
~1.6%	~3500	~2000
~2.3%	~5300	~3000
~3.1%	~7000	~4000
~3.9%	~8700	~5000
~4.6%	~10500	~6000
~5.4%	~12200	~7000
~6.1%	~14000	~8000
~6.9%	~15700	~9000
~7.6%	~17400	~10000

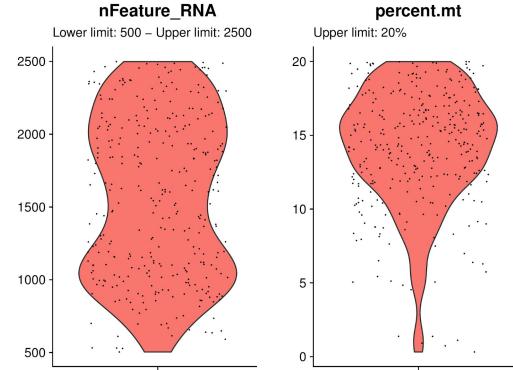
BioIT analysis: Multiplet detection

Cell hashing

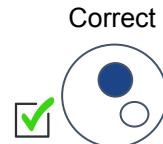
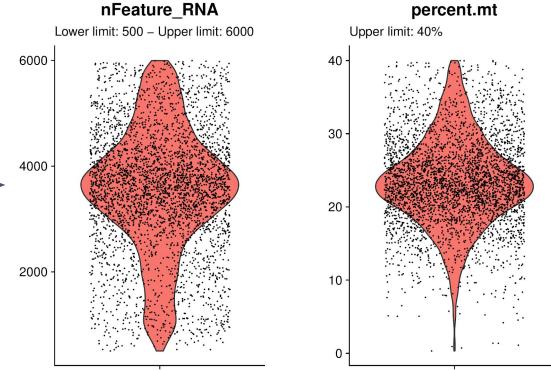


BioIT analysis: Filtering

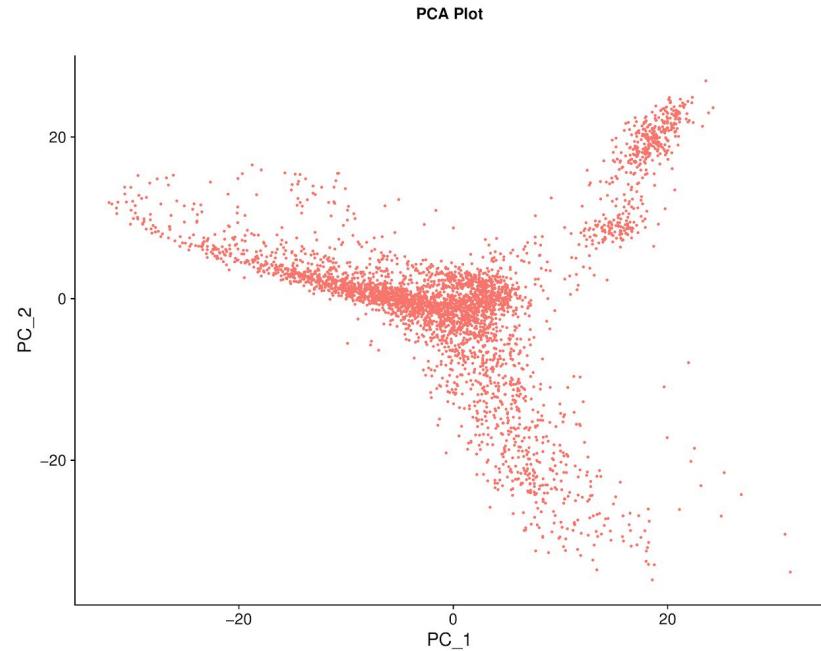
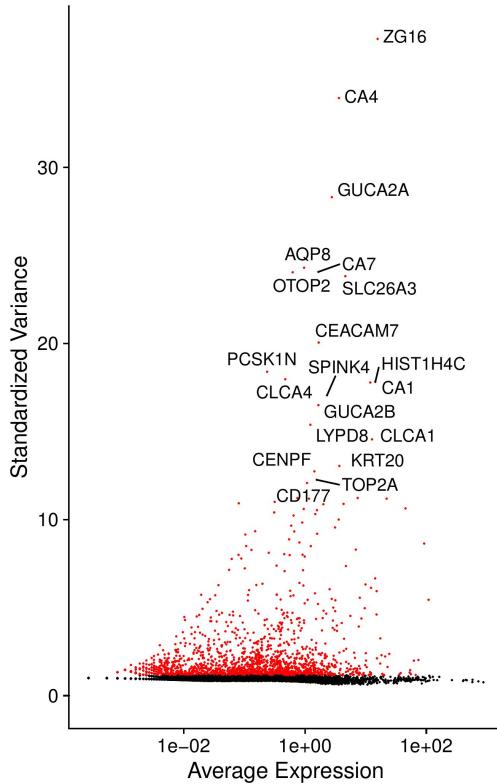
First pass



Second pass



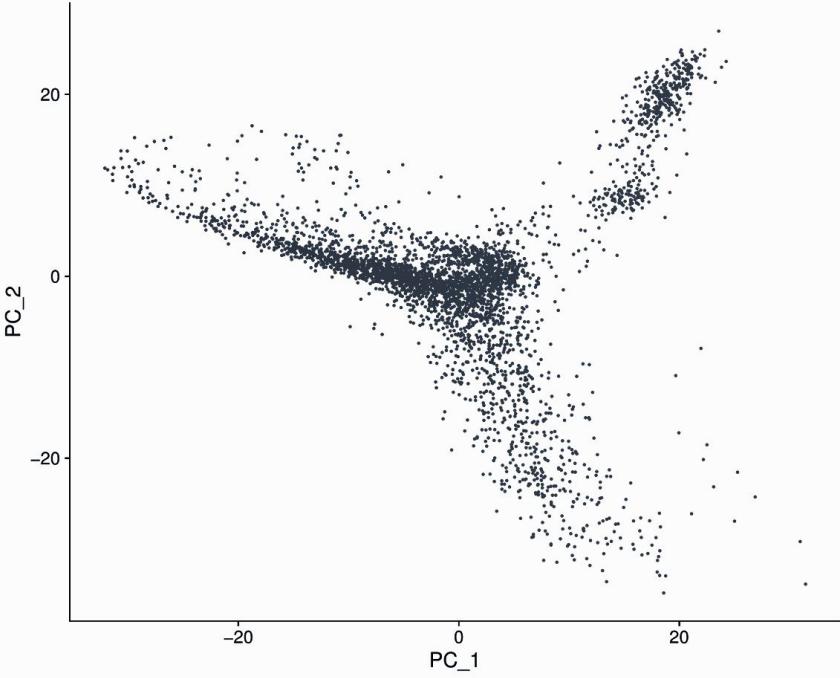
BioIT analysis: Feature selection & PCA



BioIT analysis: PCA

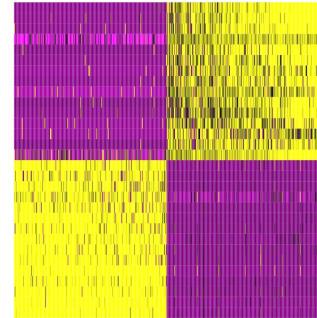
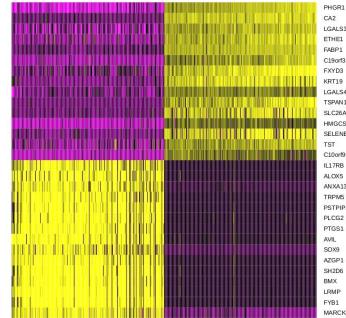
Linear reduction

PCA Plot
Coloured by cluster



PC_1

PC_2



C10orf99

TST

SLEENP1

HMGCS2

SLC25A2

TSPAN1

LGALS4

KRT19

FXYD3

C10orf99

FABP1

ETHE1

LGALS3

CA2

PHGR1

IL17RB

ALOX5

ANXA13

TRPM5

PSTPIP2

PLCG2

PTGS1

AVL

SOX9

AZGP1

SH2D6

BMX

LRMP

FYB1

MARCKSL1

HMGCB

PBK

TVMS

MAD2L1

CENPF

UBE2C

SMC2

BIRC5

CENPW

TPX2

CDK1

H2AFZ

NUSAP1

TOP2A

MK167

AZGP1

MATK

SH2D7

RASSF6

PIK3CG

HPGDS

HSPB

PLX2

PTGS1

TRPM5

LRMP

PSTPIP2

BMX

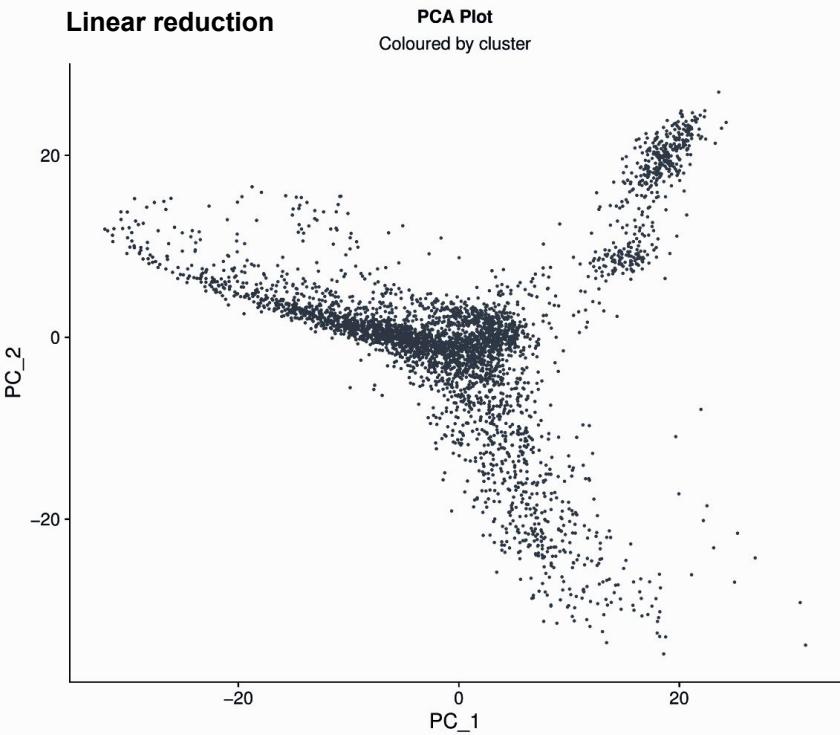
SH2D6

FYB1

PC_1

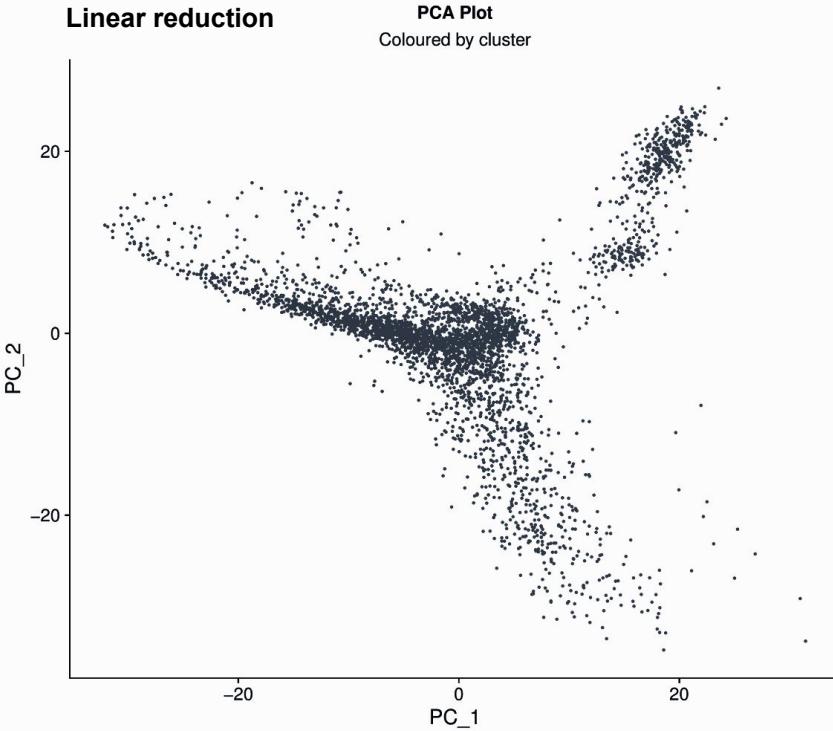
PC_2

BioIT analysis: Clustering

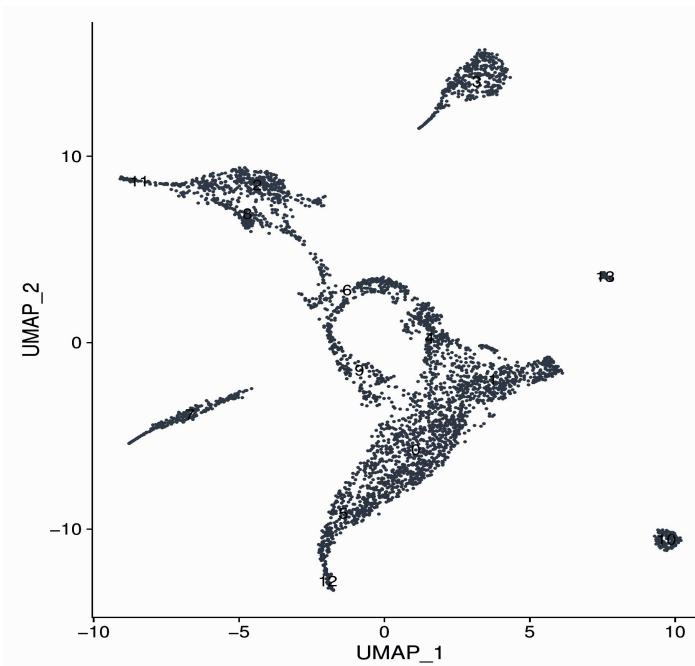


BioIT analysis: Clustering

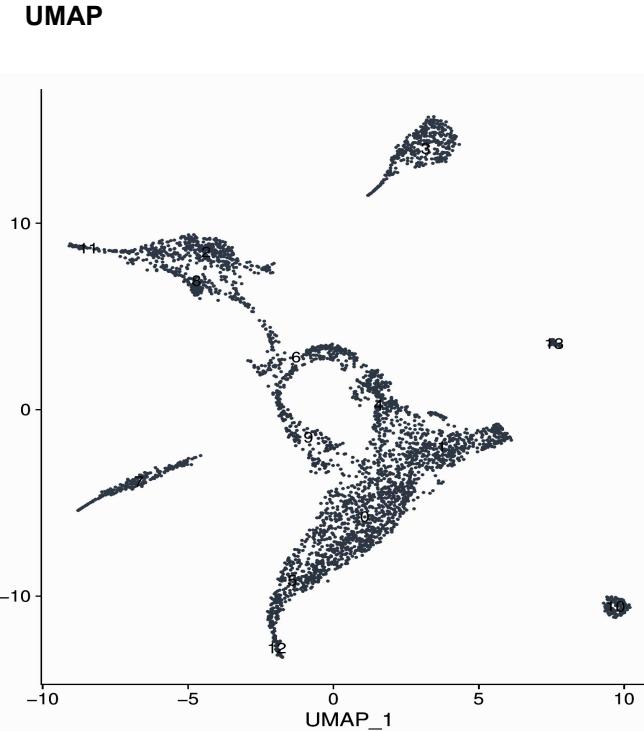
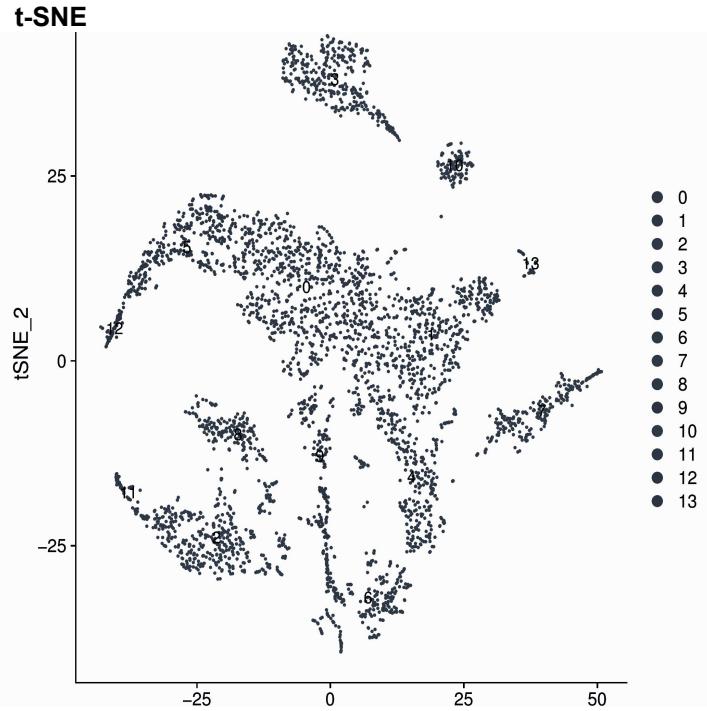
Linear reduction



Nonlinear reduction

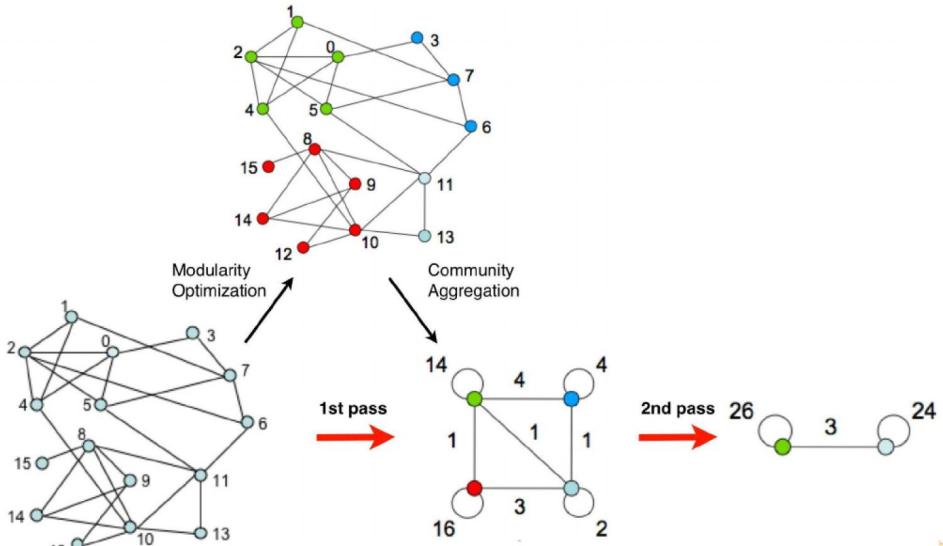


BioIT analysis: Clustering

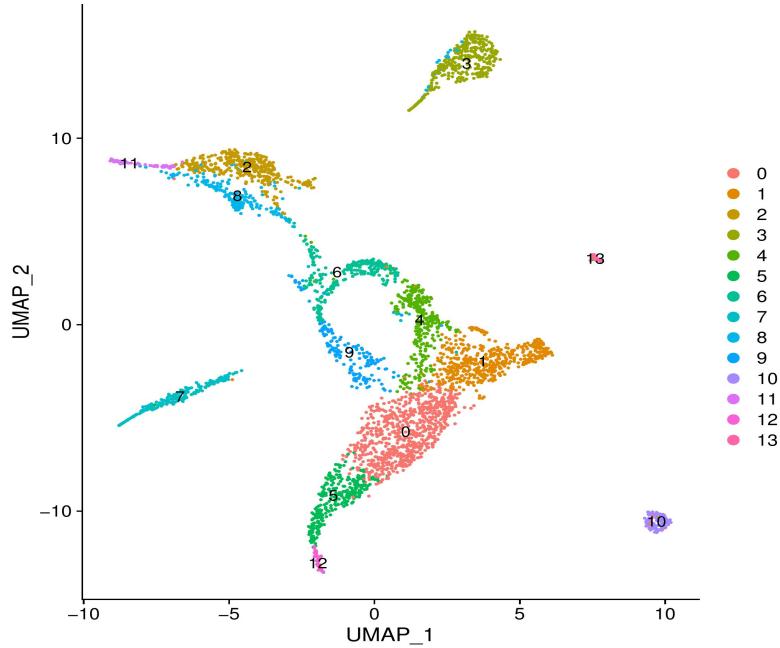


BioIT analysis: Clustering

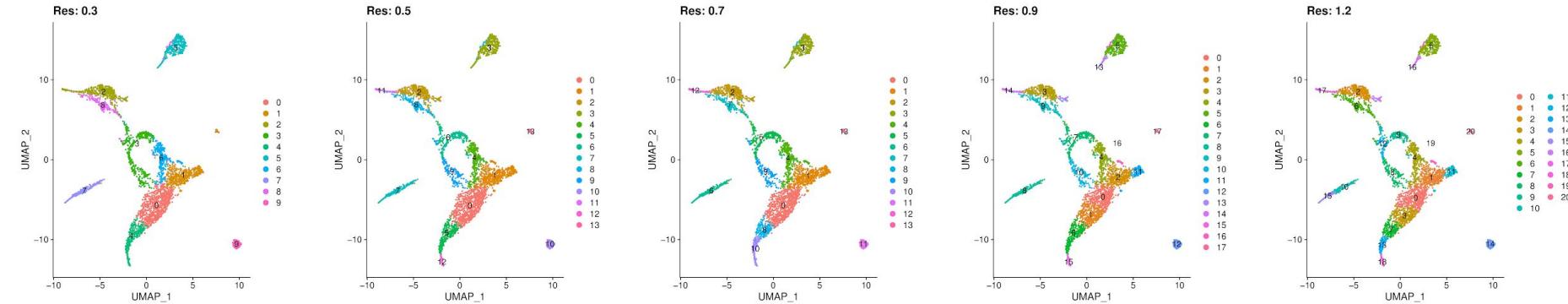
Graph-based



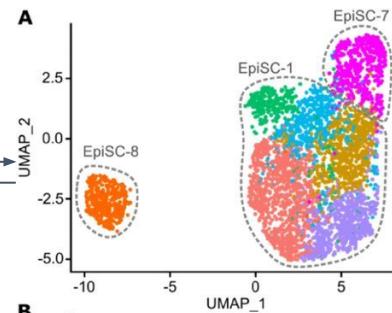
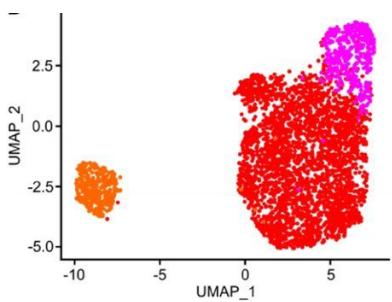
Nonlinear reduction



BioIT analysis: Clustering

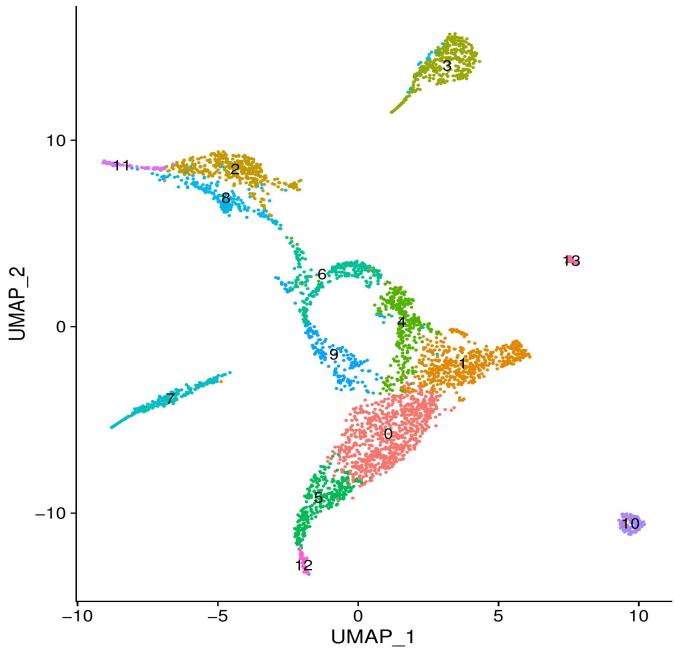


Low resolution

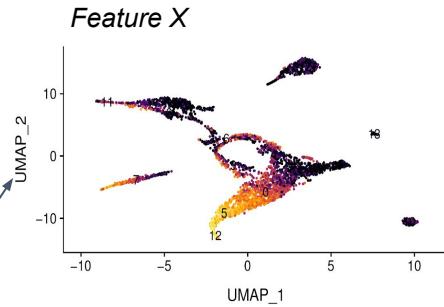


High resolution

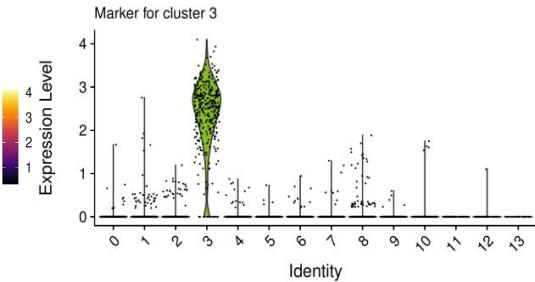
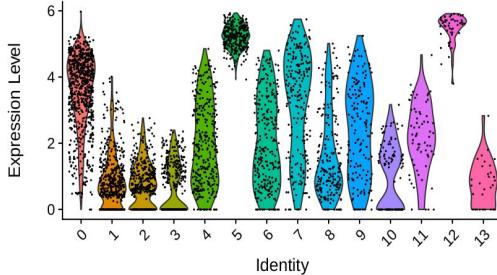
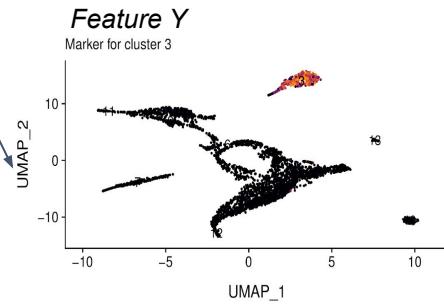
BioIT analysis: Identification



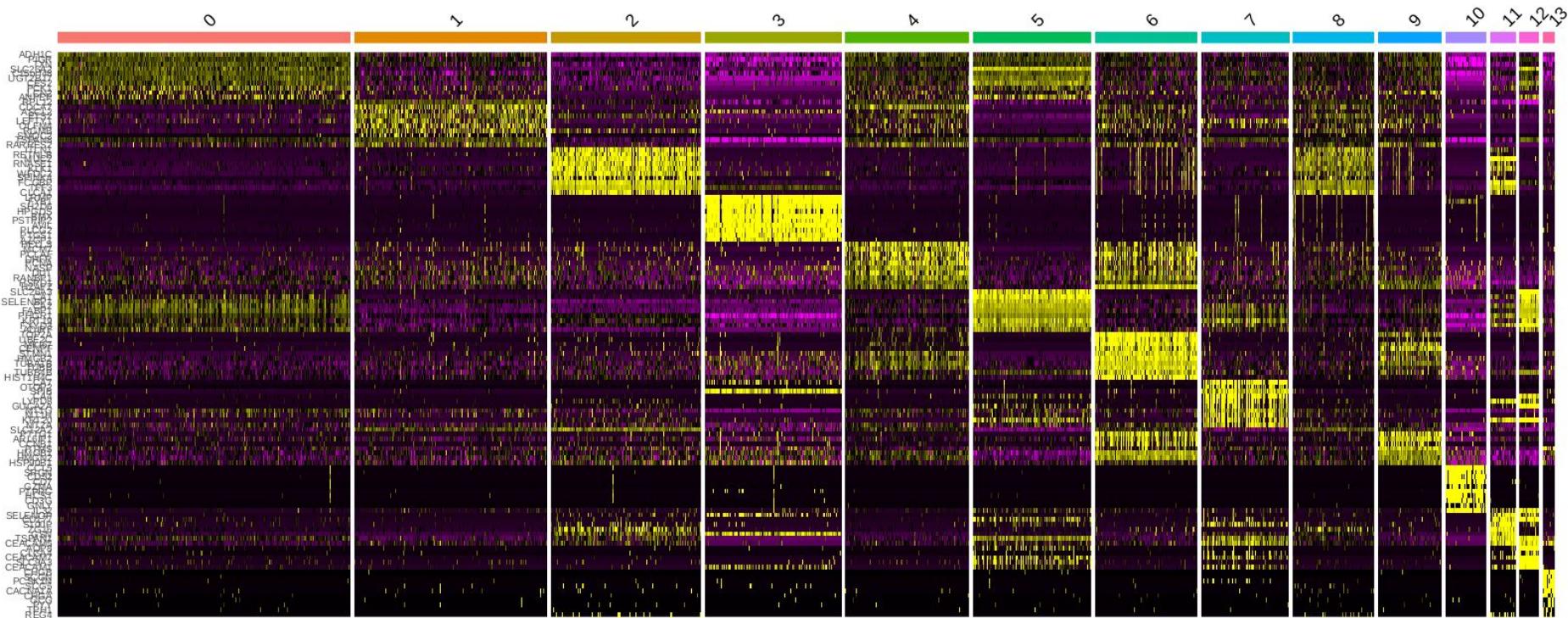
Visualize genes typically expressed in cell type



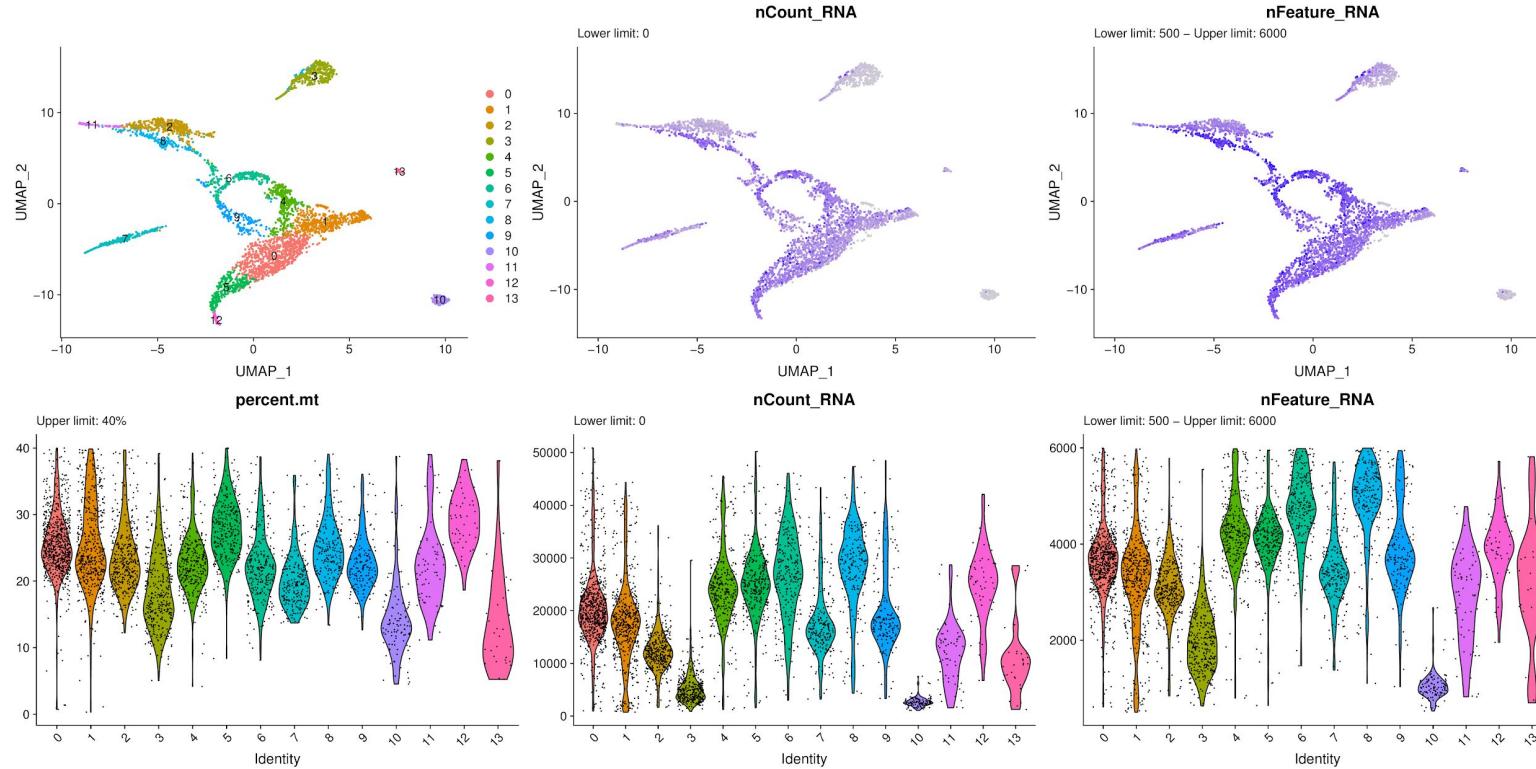
Find markers expressed for each cluster



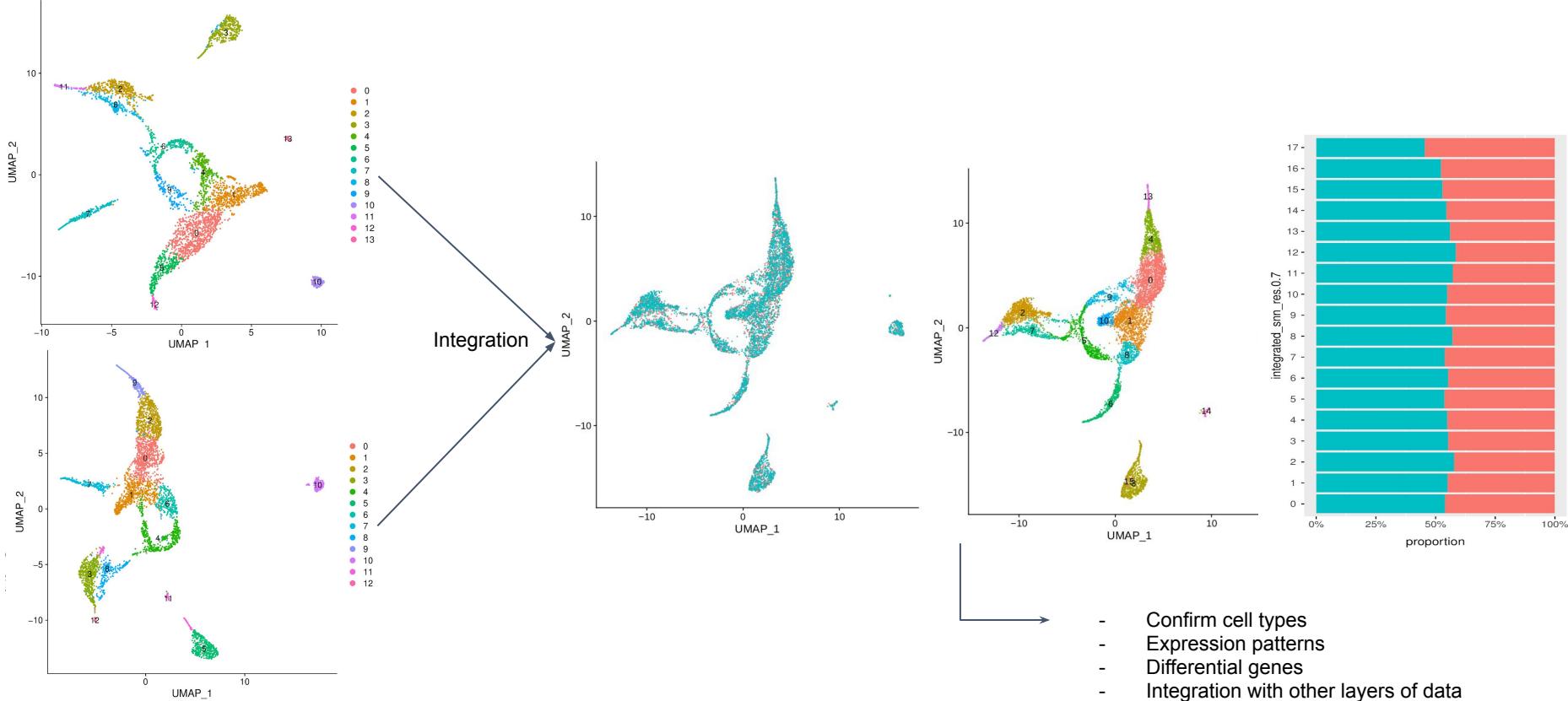
BioIT analysis: Identification



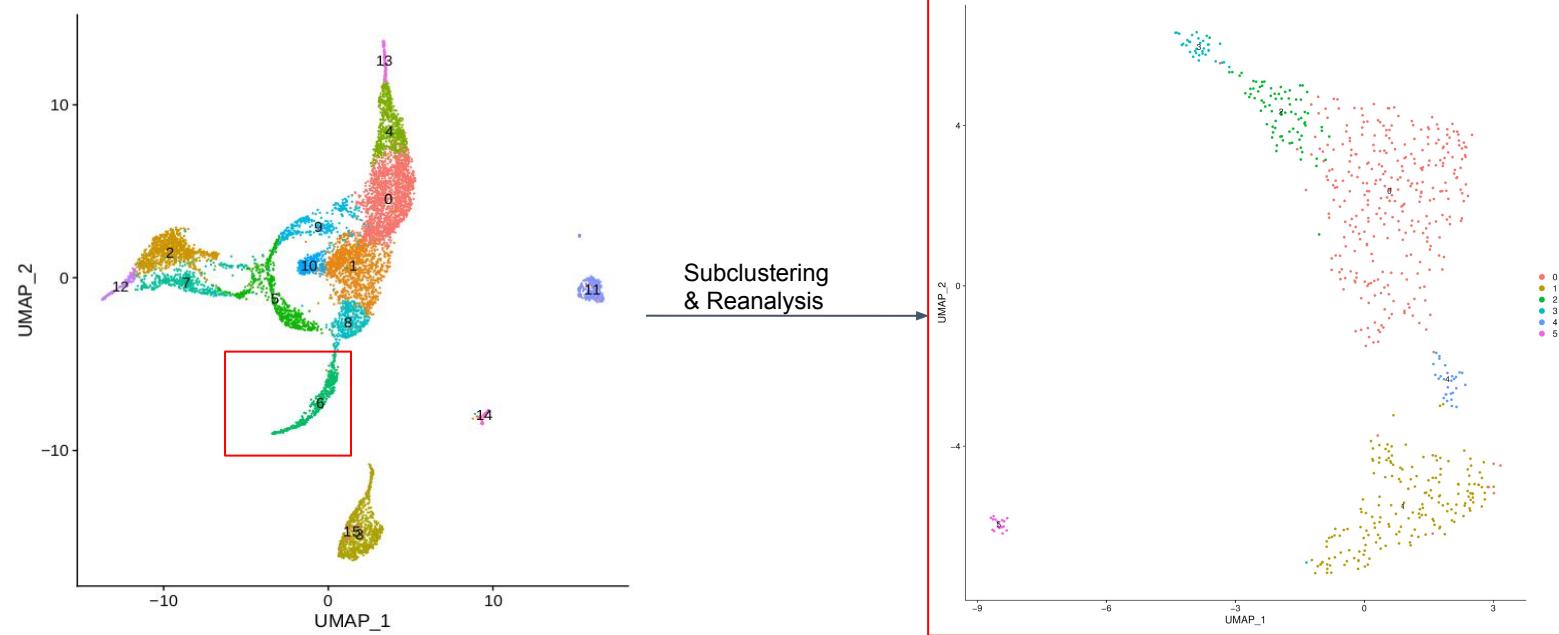
BioIT analysis: Subsetting



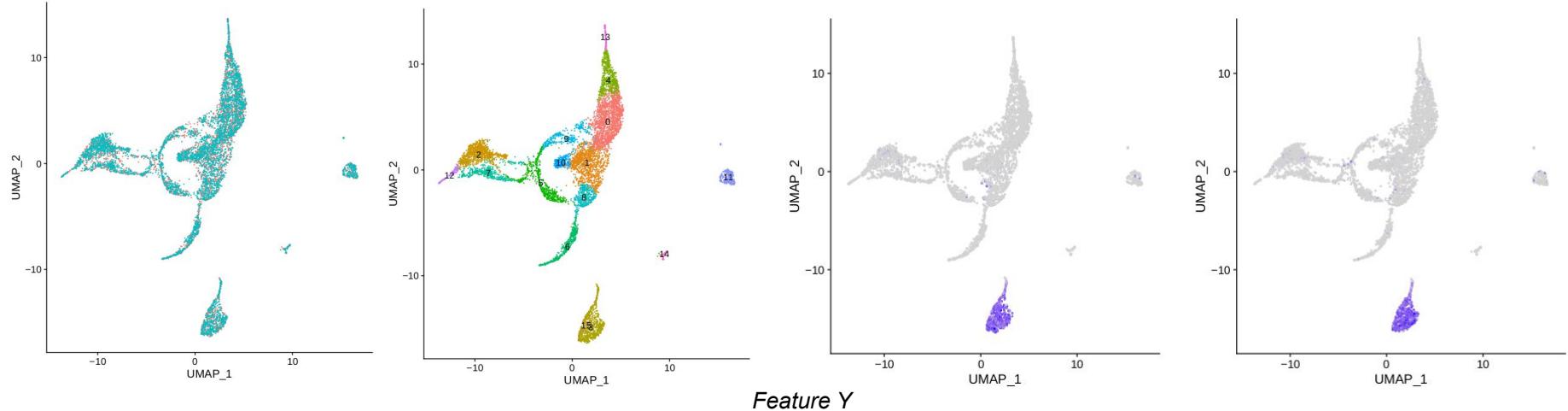
BioIT analysis: Integration



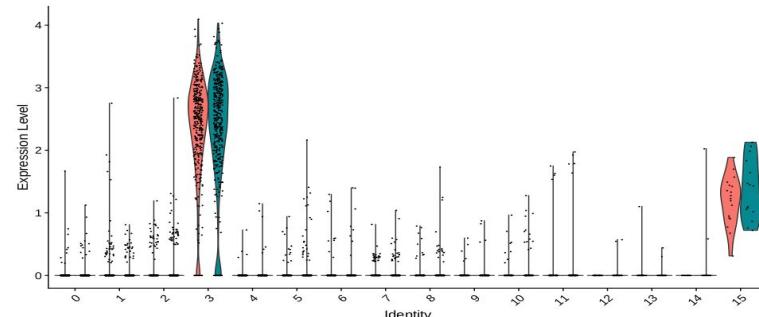
BioIT analysis: Subsetting



BioIT analysis: Integration

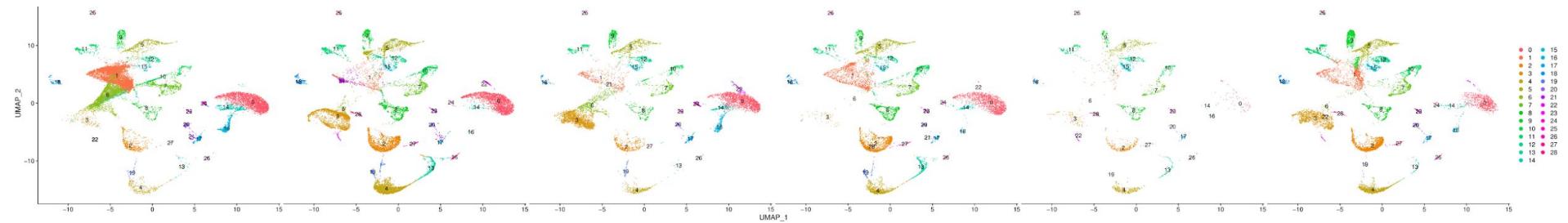
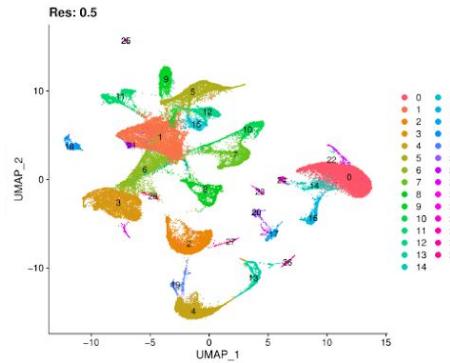
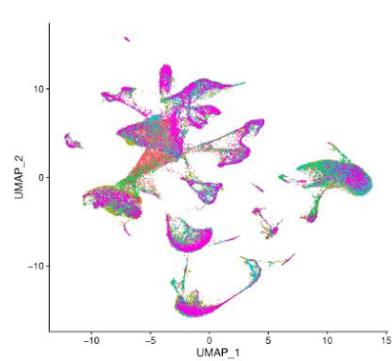


Feature Y



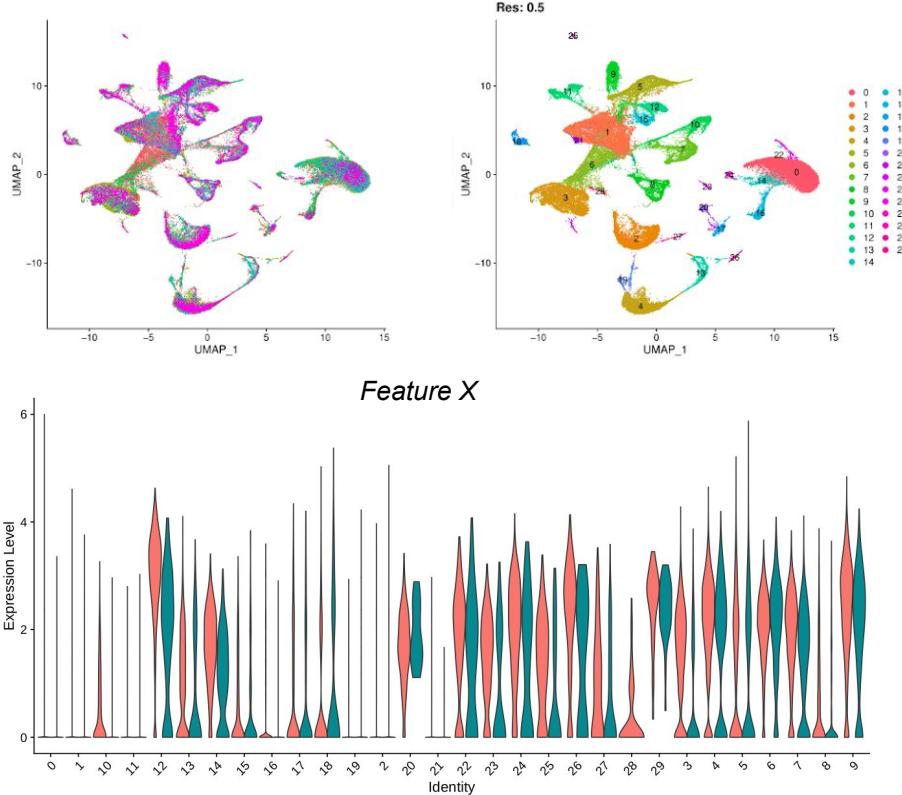
BioIT analysis: Scalability

~78k cells

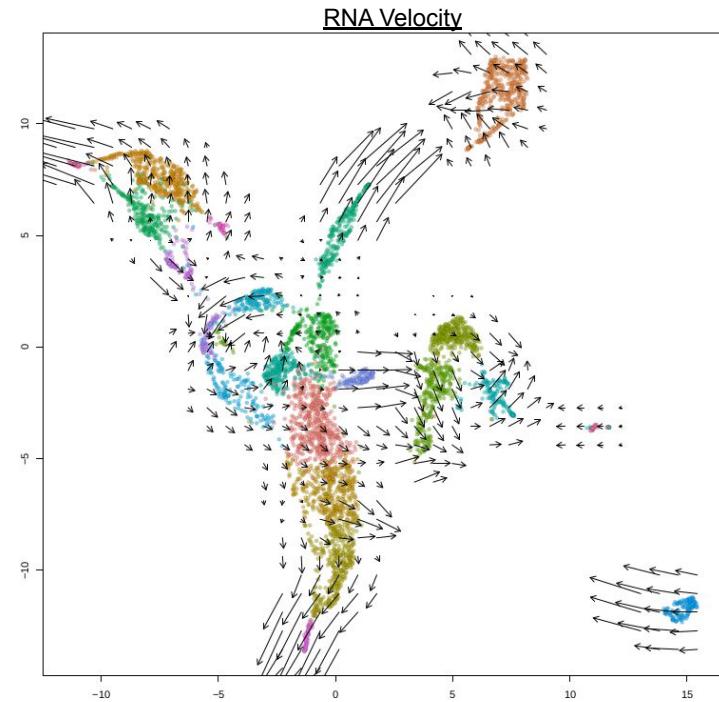
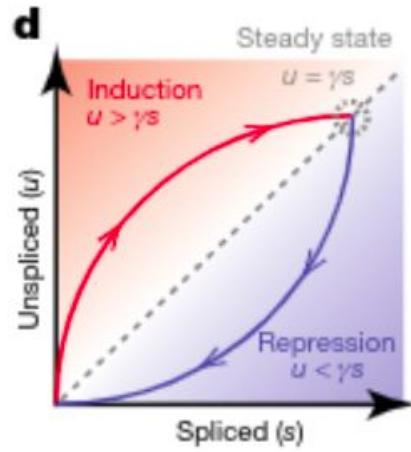
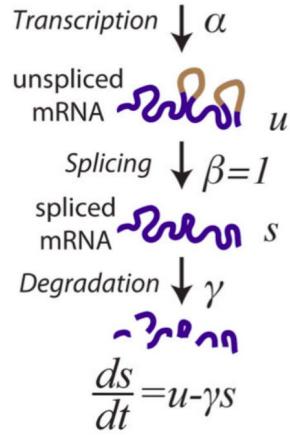


BioIT analysis: Scalability

~78k cells

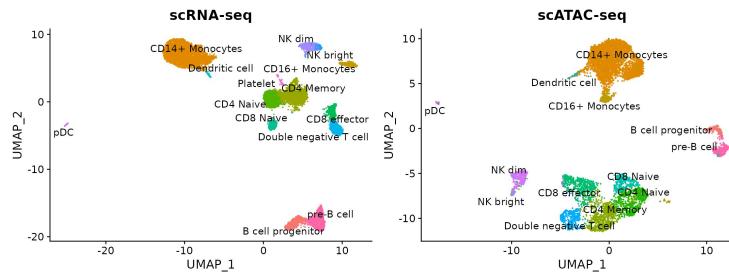
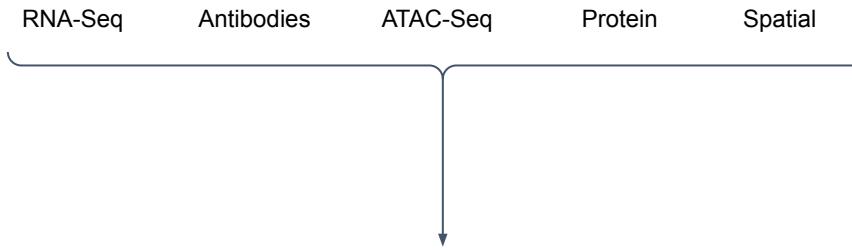


BioIT analysis: Further analysis



BioIT analysis: Further analysis

Multimodal analyses

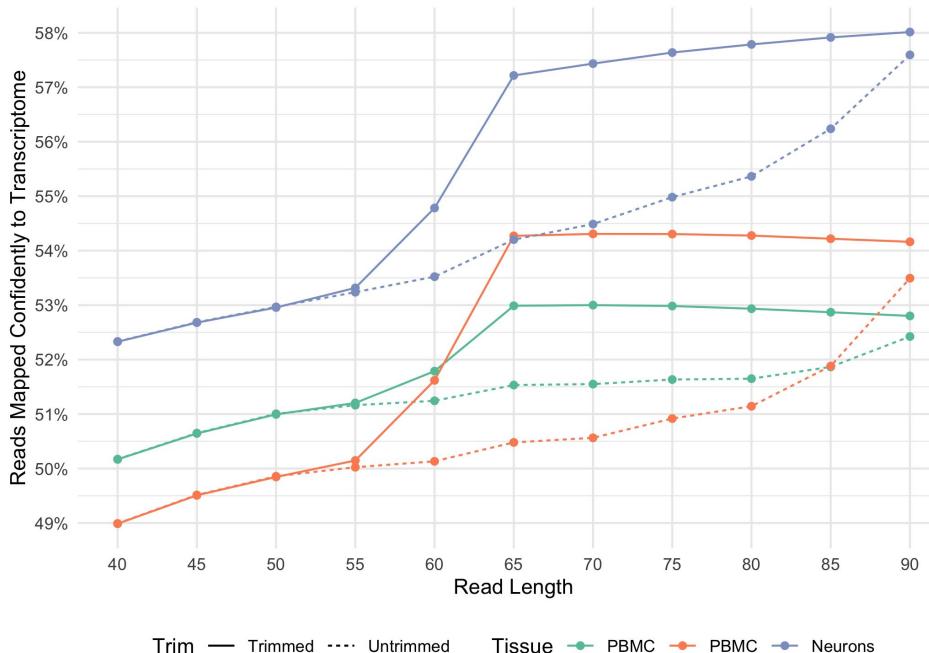


Thank you for your attention!

Questions?

info@genomicscore.be

Additional figures:



Cell calling algorithm

- 1) First pass: (High RNA content cells)
 - a) Cutoff based on UMI count
 - i) Take 99th percentile of UMI counts as m
 - ii) All GEMs with nUMI > (m/10) → mark as cells
- 2) Second pass (Lower RNA content cells)
 - a) Take set of GEMs which appear empty (low nUMI)
 - b) Create a model of their RNA profile
 - c) RNA profile of all GEM's uncalled in the first pass is compared to this RNA profile
 - d) If strong disagreement with the profile → mark as cells

