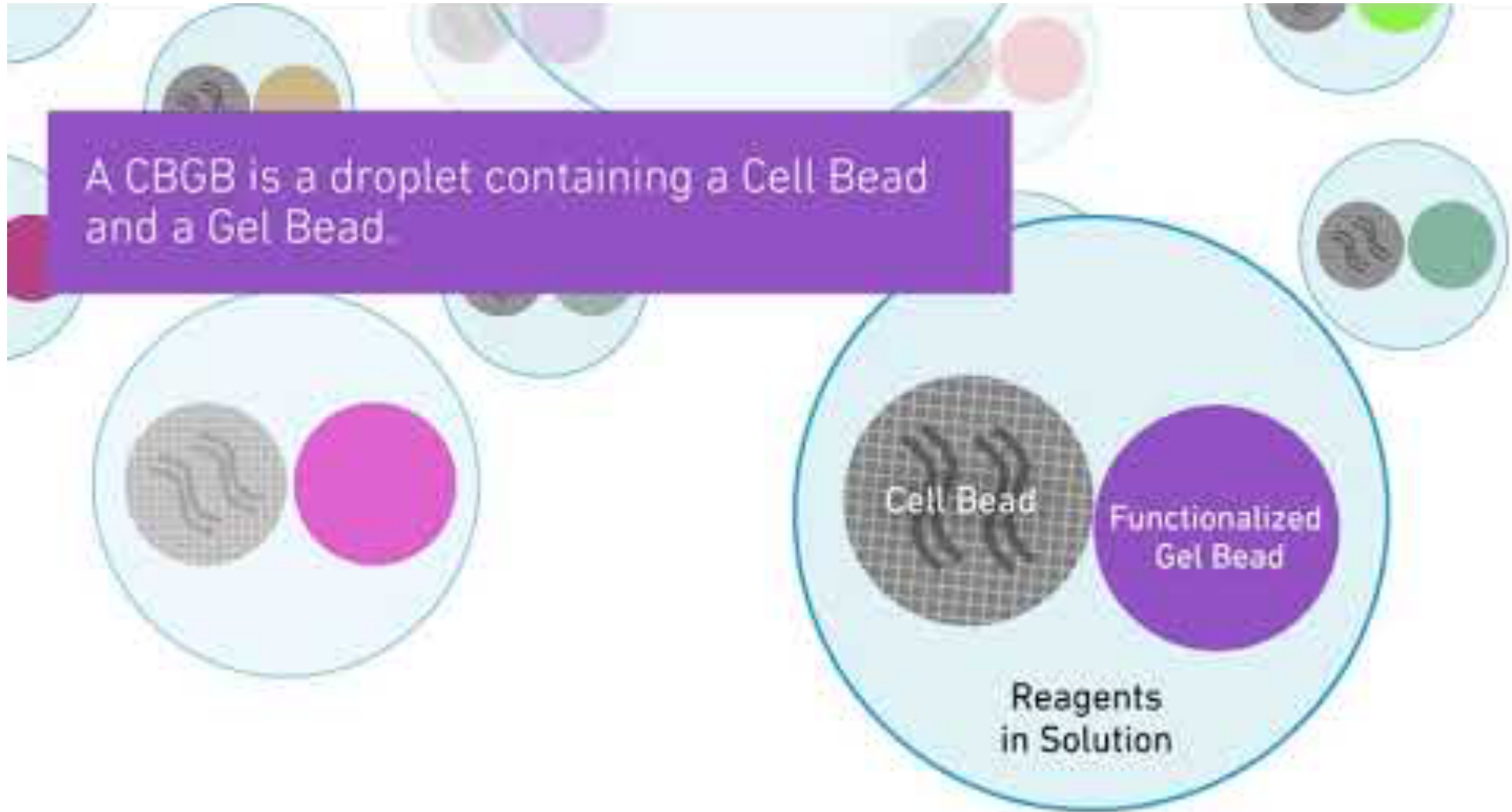*BTEP Presentation: scCNV + scATAC*

Keyur Talsania

CCR-SF Bioinformatics Group
Advanced Biomedical and Computational Sciences
Biomedical Informatics and Data Science (BIDS)  Directorate
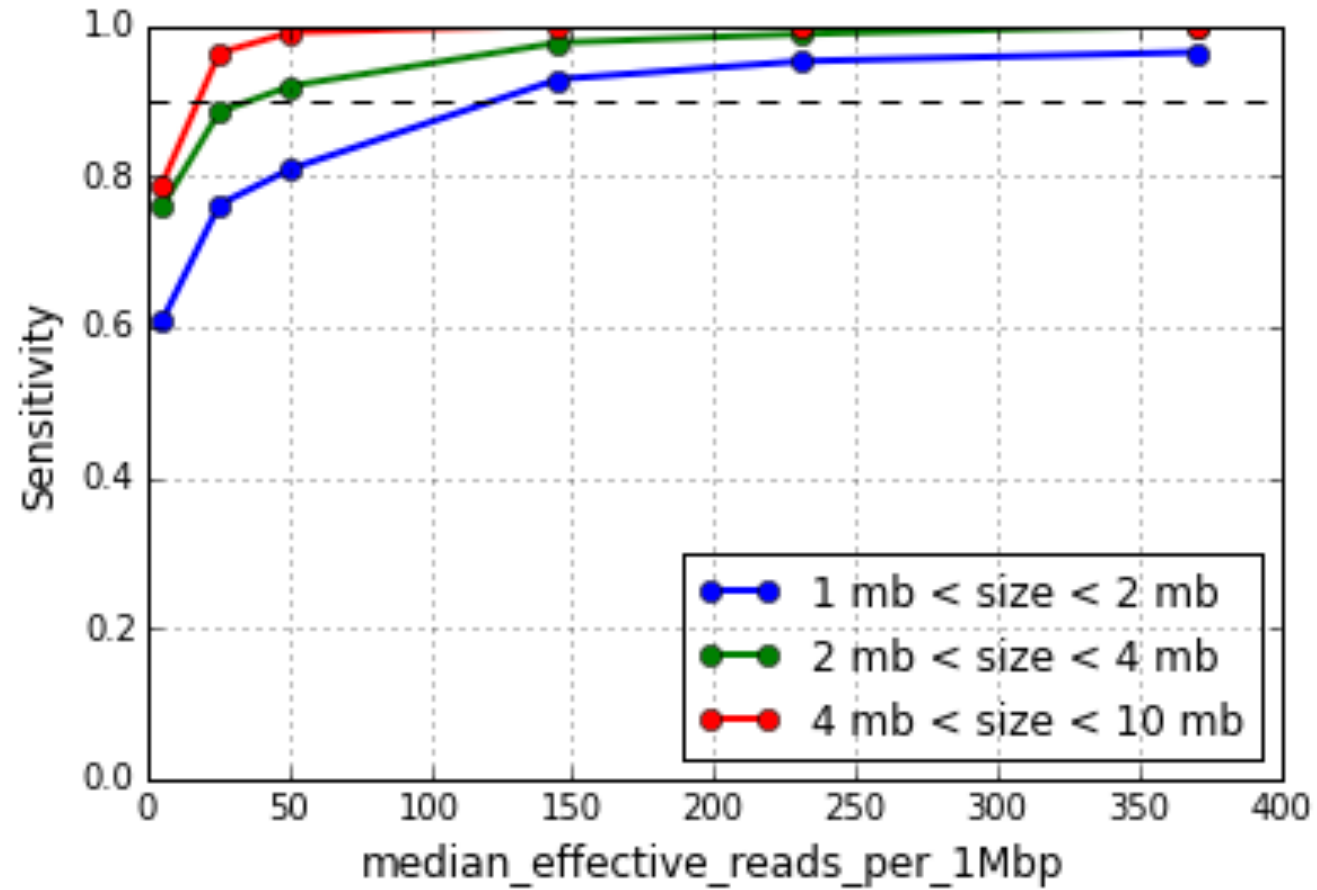Frederick National Laboratory for Cancer Research

# scCNV

- **Sequencing Requirement:**

" For an approximately diploid human sample we recommend a sequencing depth of 750,000 read-pairs per cell. At this depth, the metric median effective reads per 1Mbp is between 350-400, and we expect to be able to detect single cell copy number events in the size range 1-2 megabases (and upwards) with high sensitivity and positive predictive value. In groups of 10 or more cells we expect to be able to detect copy number events in the 100-200 kilobase (and upwards) with high sensitivity and positive predictive value."

# scCNV – 10x Genomics

| Read pairs per cell | Single cell CNV resolution (Mb) |
|---|---|
| 50K | 13 +/- 4 |
| 100K | 7 +/- 2 |
| 150K | 5 +/- 2 |
| 300K | 2.5 +/- 0.7 |
| 500K | 1.8 +/- 0.5 |
| 750K | 1.4 +/- 0.3 |

# scCNV – 10x Genomics
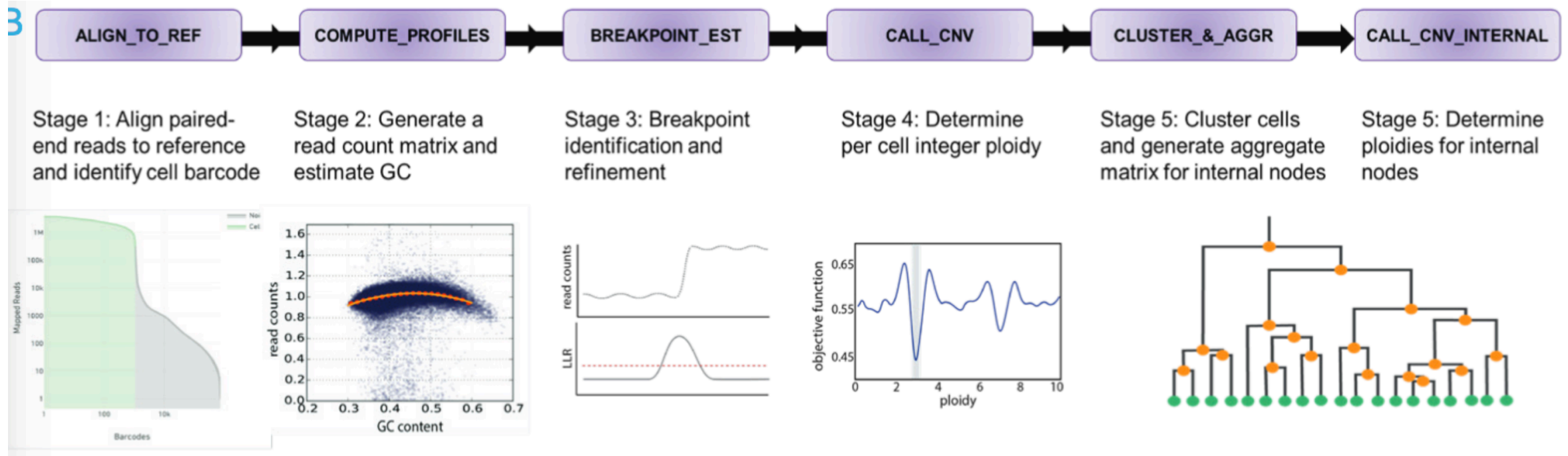
- Non-human organisms

    We expect high quality CNV detection when median effective reads per 1Mbp is in the range 350-400, and the results in the graph above will likely translate across organisms. This level of coverage can be achieved by scaling the recommended coverage of 1.5-2.0 million reads per cell by the ratio of the organism genome size to the human genome size.

- Very high ploidy samples

    For samples that contain cells with average ploidy significantly different from two, as is some times the case in cancer genomes, we recommend scaling the input coverage in proportion to the average ploidy / 2. In a tetraploid sample, for example, the extra coverage allows us to distinguish 4 -> 5 copy number changes and other n -> n+1 higher copy number transitions where the relative ploidy difference can be small.

- Non-human organisms

    We expect high quality CNV detection when median effective reads per 1Mbp is in the range 350-400, and the results in the graph above will likely translate across organisms. This level of coverage can be achieved by scaling the recommended coverage of 1.5-2.0 million reads per cell by the ratio of the organism genome size to the human genome size.

- Very high ploidy samples

    For samples that contain cells with average ploidy significantly different from two, as is some times the case in cancer genomes, we recommend scaling the input coverage in proportion to the average ploidy / 2. In a tetraploid sample, for example, the extra coverage allows us to distinguish 4 -> 5 copy number changes and other n -> n+1 higher copy number transitions where the relative ploidy difference can be small.

- Ploidy of 3 = 1.5m * (3/2) = 2.25m reads – 1.2m read pairs

- Ploidy of 3 = 1.5m * (4/2) = 3m reads – 1.5m read pairs

# 10x Genomis scCNV Pipeline



| ALIGN_TO_REF | COMPUTE_PROFILES | BREAKPOINT_EST | CALL_CNV | CLUSTER_&_AGGR | CALL_CNV_INTERNAL |

Stage 1: Align paired-end reads to reference and identify cell barcode

Stage 2: Generate a read count matrix and estimate GC

Stage 3: Breakpoint identification and refinement

Stage 4: Determine per cell integer ploidy

Stage 5: Cluster cells and generate aggregate matrix for internal nodes

Stage 5: Determine ploidies for internal nodes

# ScCNV: 10x scDNA Loupe

# ScCNV: 10x scDNA Loupe

# ScCNV: 10x scDNA Loupe

# ScCNV: 10x scDNA Loupe Demo

# *scATAC*

# scATAC – 10x Genomics

- How it works

- Results from 10x Genomics

- Other tools

- Data integration

# Results from 10x Genomics

- **cellranger-atac count** takes FASTQ files from `cellranger-atac mkfastq` and performs ATAC analysis, including:

  - Read filtering and alignment
  - Barcode counting
  - Identification of transposase cut sites
  - Detection of accessible chromatin peaks
  - Cell calling
  - Count matrix generation for peaks and transcription factors
  - Dimensionality reduction
  - Cell clustering
  - Cluster differential accessibility

# Summary of results

10x GENOMICS  **Cell Ranger ATAC**   Sample   Sequencing   Cells   Cell Clustering   Insert Sizes   Targeting   Library Complexity

## pbmc345 - jdoe's PBMC

| 581 | 13,822 | 91.0% |
|---|---|---|
| Estimated number of cells | Median fragments per cell | Fraction of fragments overlapping any targeted region |

**63.8%**

Fraction of transposition events in peaks

### Sample

| | |
|---|---|
| Sample ID | pbmc345 |
| Sample description | jdoe's PBMC |
| FASTQ path | /home/jdoe/HAWT7ADXX/outs/fastq_path |
| Pipeline version | cellranger-atac_1.0.0 |
| Reference path | …e/refdata-cellranger-atac-hg19_1.0.0 |

### Sequencing ⓘ

| | |
|---|---|
| Total number of read pairs | 47,276,182 |
| Fraction of read pairs with a valid barcode | 98.2% |
| Q30 bases in Read 1 | 94.9% |
| Q30 bases in Read 2 | 94.8% |
| Q30 bases in Barcode | 82.7% |
| Q30 bases in Sample Index | 89.1% |

# Summary of results



**Cells** ⓘ

| | |
|---|---|
| Estimated number of cells | 581 |
| Lower threshold on the number of fragments overlapping peaks per barcode to annotate barcode as cell | 635.00 |
| Median fragments per cell | 13,822 |
| Median fragments per non-cell barcode | 1 |

Cells - Sample pbmc345
jdoe's PBMC

Fragment Distribution - Sample pbmc345
jdoe's PBMC

**Cell Clustering** ⓘ

Cell Clustering (By Cluster) - Sample pbmc345
jdoe's PBMC

- Cluster 1 (160)
- Cluster 2 (125)
- Cluster 3 (118)
- Cluster 4 (102)
- Cluster 5 (76)

Cell Clustering (Colored by Depth) - Sample pbmc345
jdoe's PBMC

# Summary of results

## Insert Sizes ⑦

| | |
|---|---|
| Fragments in nucleosome-free regions | 51.1% |
| Fragments flanking a single nucleosome | 29.3% |



Insert Size Distribution – Sample pbmc345
jdoe's PBMC

# Summary of results

## Targeting ⓘ

| | |
|---|---|
| Fraction of fragments overlapping any targeted region | 91.0% |
| Fraction of fragments overlapping TSS | 56.6% |
| Fraction of fragments overlapping DNase HS regions | 86.6% |
| Fraction of fragments overlapping enhancer regions | 24.0% |
| Fraction of fragments overlapping promoter regions | 48.6% |
| Fraction of fragments overlapping blacklisted regions | 0.2% |
| Fraction of fragments overlapping called peaks | 70.2% |
| Enrichment score of transcription start sites | 13.71 |

| | |
|---|---|
| Fraction of total read pairs mapped confidently to genome (>30 mapq) | 89.2% |
| Fraction of total read pairs that are unmapped and in cell barcodes | 0.3% |
| Fraction of total read pairs in mitochondria and in cell barcodes | 0.4% |



Enrichment around TSS (normalized) - Sample pbmc345
jdoe's PBMC



Singlecell Targeting (Peaks) - Sample pbmc345
jdoe's PBMC

# Other tools

- scABC
- Destin
- ChromVAR
- Cicero
- CoupledNMF
- CisTopic
- Brockman
- SNAP-ATAC
- Signac

# Other tools - Destin

# ChromVAR

chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data

Alicia N Schep, Beijing Wu, Jason D Buenrostro ✉ & William J Greenleaf ✉

Technology

# Molecular Cell

## Cicero Predicts *cis*-Regulatory DNA Interactions from Single-Cell Chromatin Accessibility Data

### Graphical Abstract

### Authors

Hannah A. Pliner, Jonathan S. Packer, José L. McFaline-Figueroa, ..., Frank J. Steemers, Jay Shendure, Cole Trapnell

### Correspondence

shendure@uw.edu (J.S.), coletrap@uw.edu (C.T.)

### In Brief

Pliner et al. introduce Cicero, a software program to connect distal regulatory elements with target genes using single-cell ATAC-seq data. They find evidence that groups of co-accessible elements form chromatin hubs and undergo coordinated changes in histone marks that are predictive of changes in gene expression in skeletal muscle development.

# CoupledNMF

# Integrative analysis of single-cell genomics data by coupled nonnegative matrix factorizations

Zhana Duren[a,b,1], Xi Chen[a,b,1], Mahdi Zamanighomi[a,b,c,1], Wanwen Zeng[a,b,d], Ansuman T. Satpathy[c], Howard Y. Chang[c], Yong Wang[e,f], and Wing Hung Wong[a,b,c,2]

[a]Department of Statistics, Stanford University, Stanford, CA 94305; [b]Department of Biomedical Data Science, Stanford University, Stanford, CA 94305; [c]Center for Personal Dynamic Regulomes, Stanford University, Stanford, CA 94305; [d]Ministry of Education Key Laboratory of Bioinformatics, Bioinformatics Division and Center for Synthetic & Systems Biology, Department of Automation, Tsinghua University, 100084 Beijing, China; [e]Academy of Mathematics and Systems Science, Chinese Academy of Sciences, 100080 Beijing, China; and [f]Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, 650223 Kunming, China

# CisTopic

nature methods

## cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data

Carmen Bravo González-Blas [1,2,3], Liesbeth Minnoye [1,2,3], Dafni Papasokrati [1,2], Sara Aibar [1,2], Gert Hulselmans [1,2], Valerie Christiaens [1,2], Kristofer Davie [1,2], Jasper Wouters [1,2] and Stein Aerts [1,2]*

# CisTopic

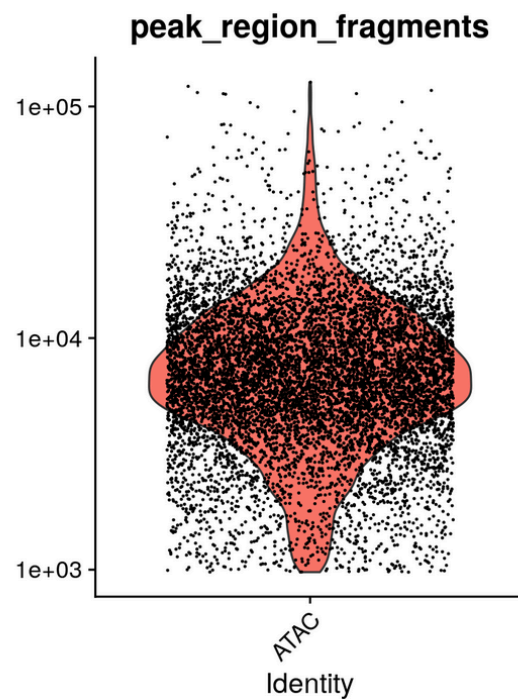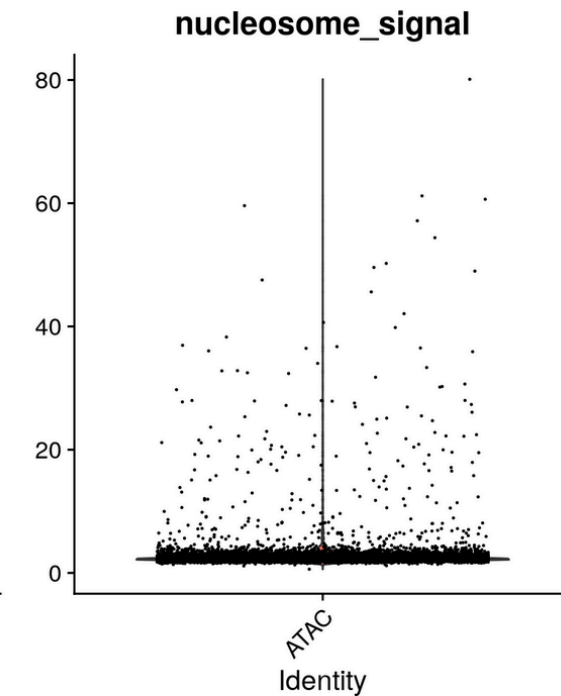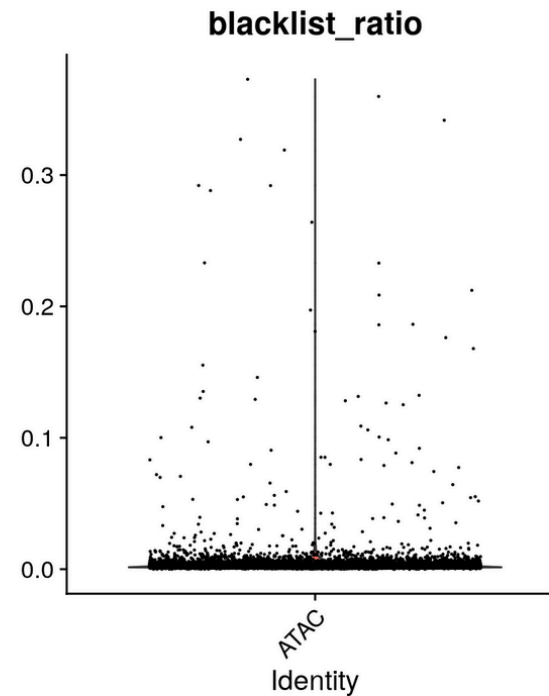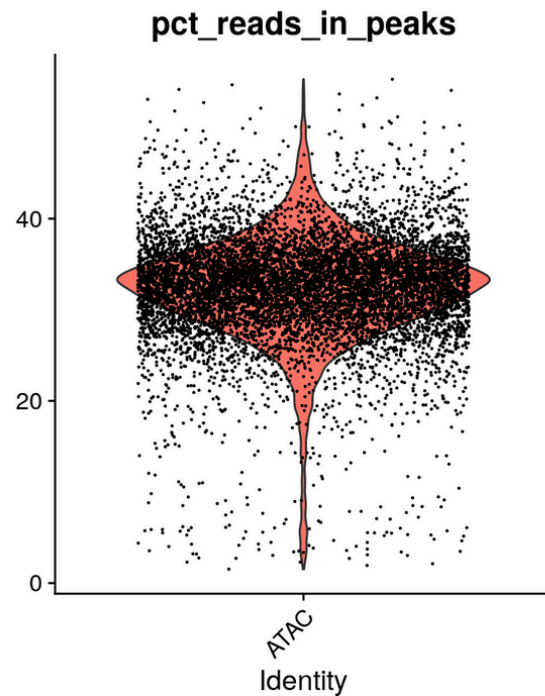## cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data

Carmen Bravo González-Blas [1,2,3], Liesbeth Minnoye [1,2,3], Dafni Papasokrati [1,2], Sara Aibar [1,2], Gert Hulselmans [1,2], Valerie Christiaens [1,2], Kristofer Davie [1,2], Jasper Wouters [1,2] and Stein Aerts [1,2]*

# CisTopic

## cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data

Carmen Bravo González-Blas [1,2,3], Liesbeth Minnoye [1,2,3], Dafni Papasokrati [1,2], Sara Aibar [1,2], Gert Hulselmans [1,2], Valerie Christiaens[1,2], Kristofer Davie [1,2], Jasper Wouters [1,2] and Stein Aerts [1,2*]

# CisTopic

# Signac

- Signac is an extension of [Seurat](#) for the analysis, interpretation, and exploration of single-cell chromatin datasets.

- Calculating single-cell QC metrics

- Dimensional reduction, visualization, and clustering

- Identifying cell type-specific peaks

- Visualizing 'pseudo-bulk' coverage tracks

- Integration of multiple single-cell ATAC-seq datasets

- Integration with single-cell RNA-seq datasets

- Motif enrichment analysis


- Integration with other single cell tools (Harmony, Cicero, Chromvar)

**Signac - QC**

https://satijalab.org/signac/index.html

# Signac - QC

# Signac- Normalization and linear dimensional reduction

- Normalization: Signac performs term frequency-inverse document frequency (TF-IDF) normalization.
  - two-step normalization procedure,
    - normalizes across cells to correct for differences in cellular sequencing depth & across peaks to give higher values to more rare peaks.

- Feature selection: Binary nature of scATAC-seq data makes it challenging to perform 'variable' feature selection, as we do for scRNA-seq.
  - Instead, use only the top $n\%$ of features (peaks) for dimensional reduction,
  - or remove features present in less that $n$ cells with the FindTopFeatures function.

- Dimensional reduction: Singular value decomposition (SVD) on the TD-IDF normalized matrix, using the features (peaks) selected above. This returns a low-dimensional representation of the object (for users who are more familiar with scRNA-seq, you can think of this as analogous to the output of PCA)

# Signac- Clustering

# Signac- Cicero: *gene activity matrix*

# scATAC + scRNA integration

# Signac- scATAC + scRNA using Seurat

# Signac- scATAC + scRNA using Seurat

# Signac- scATAC + scRNA using Harmony

# Signac- scATAC + scRNA using Conos
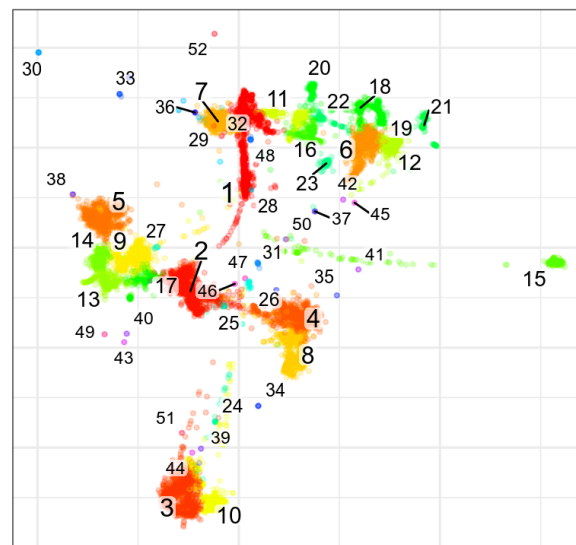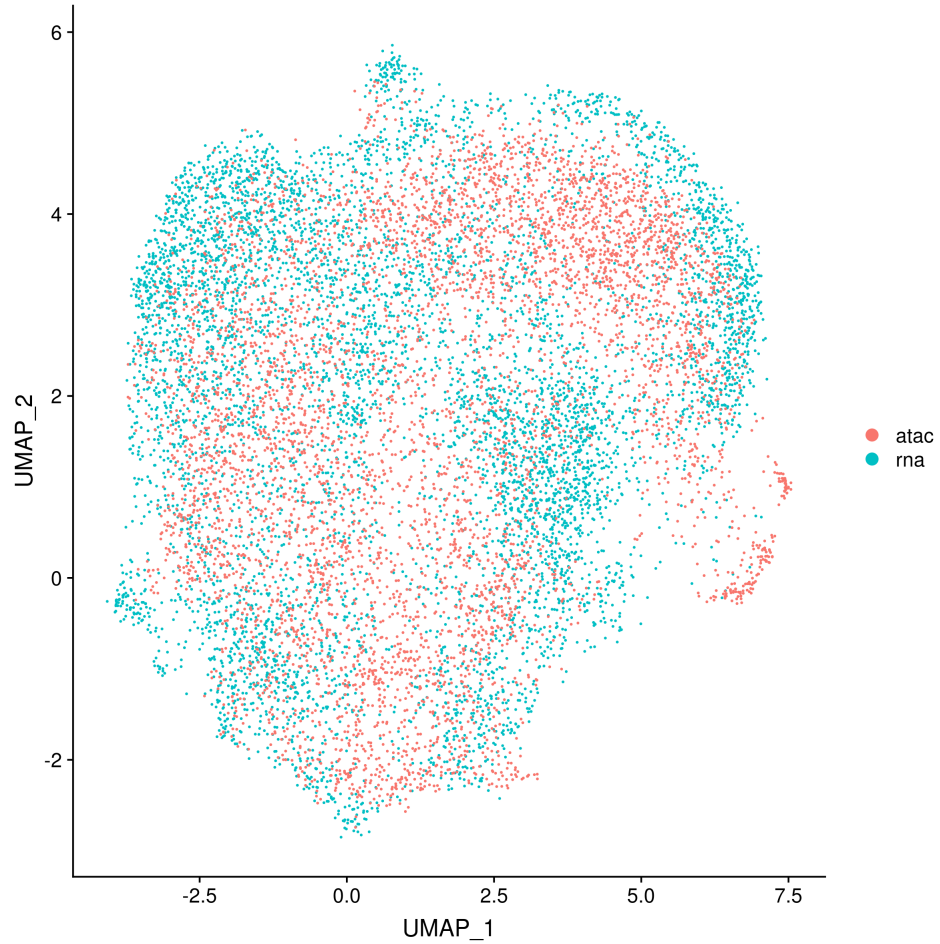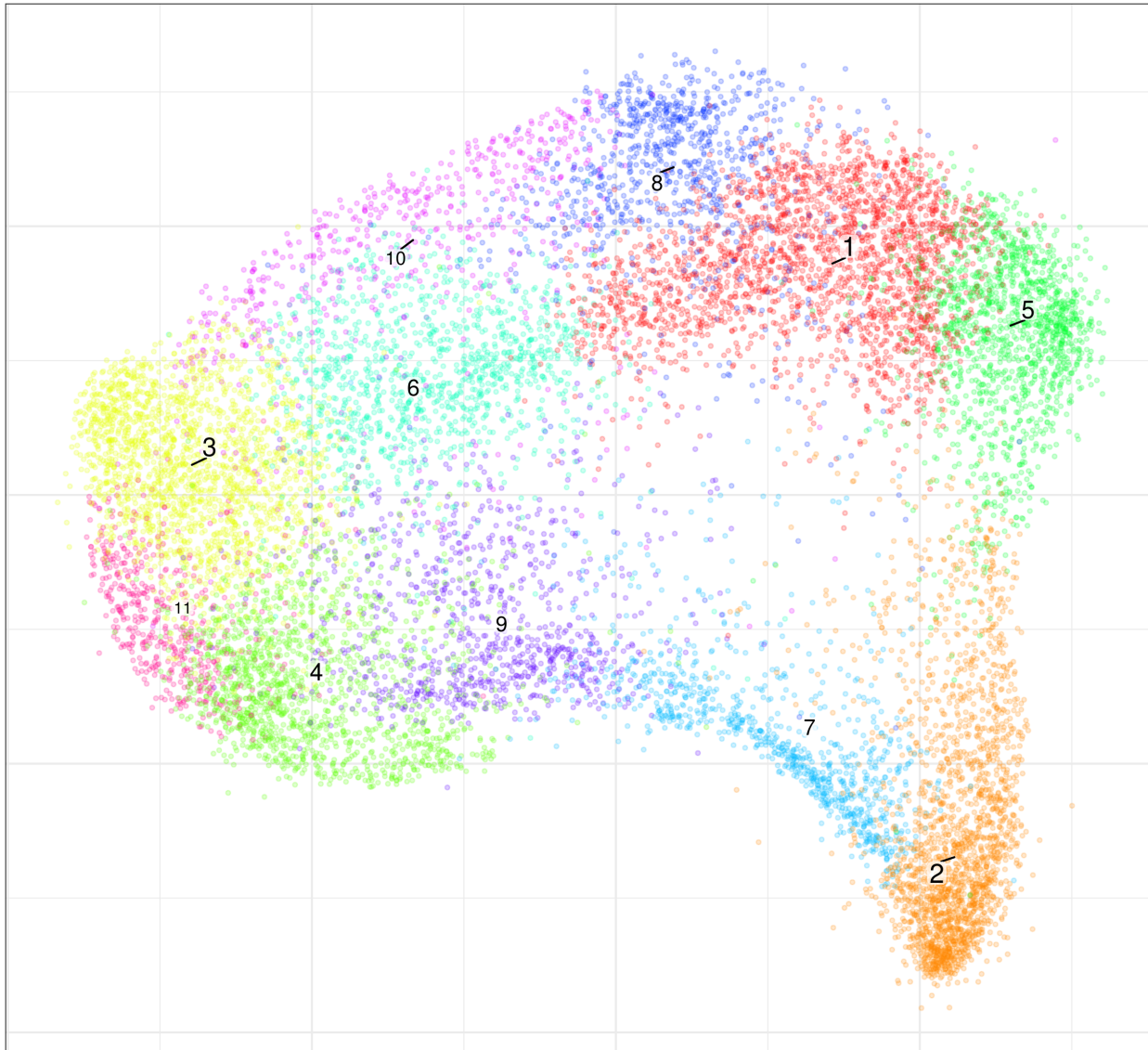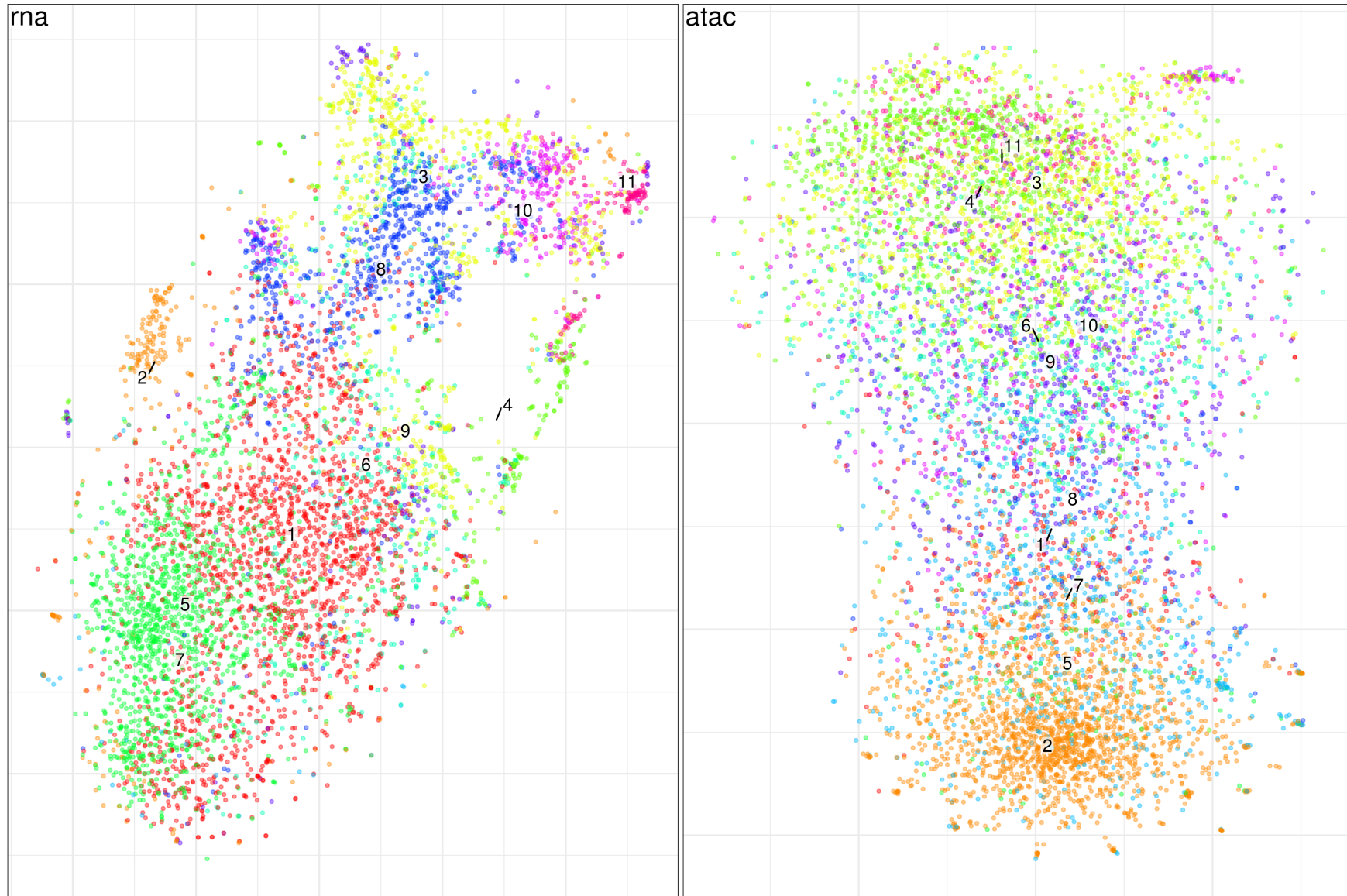
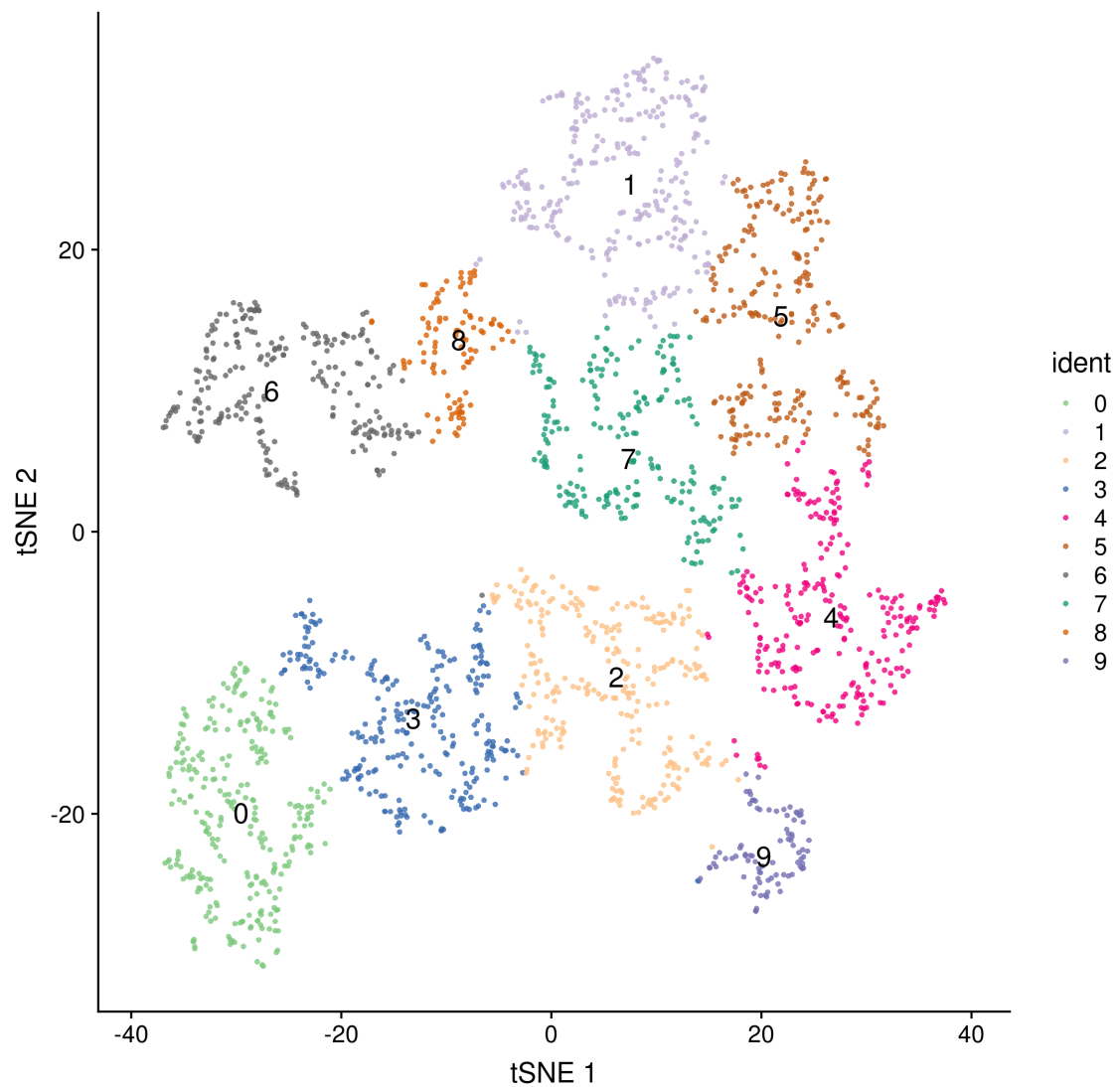# Only Seurat vs Seurat + Conos
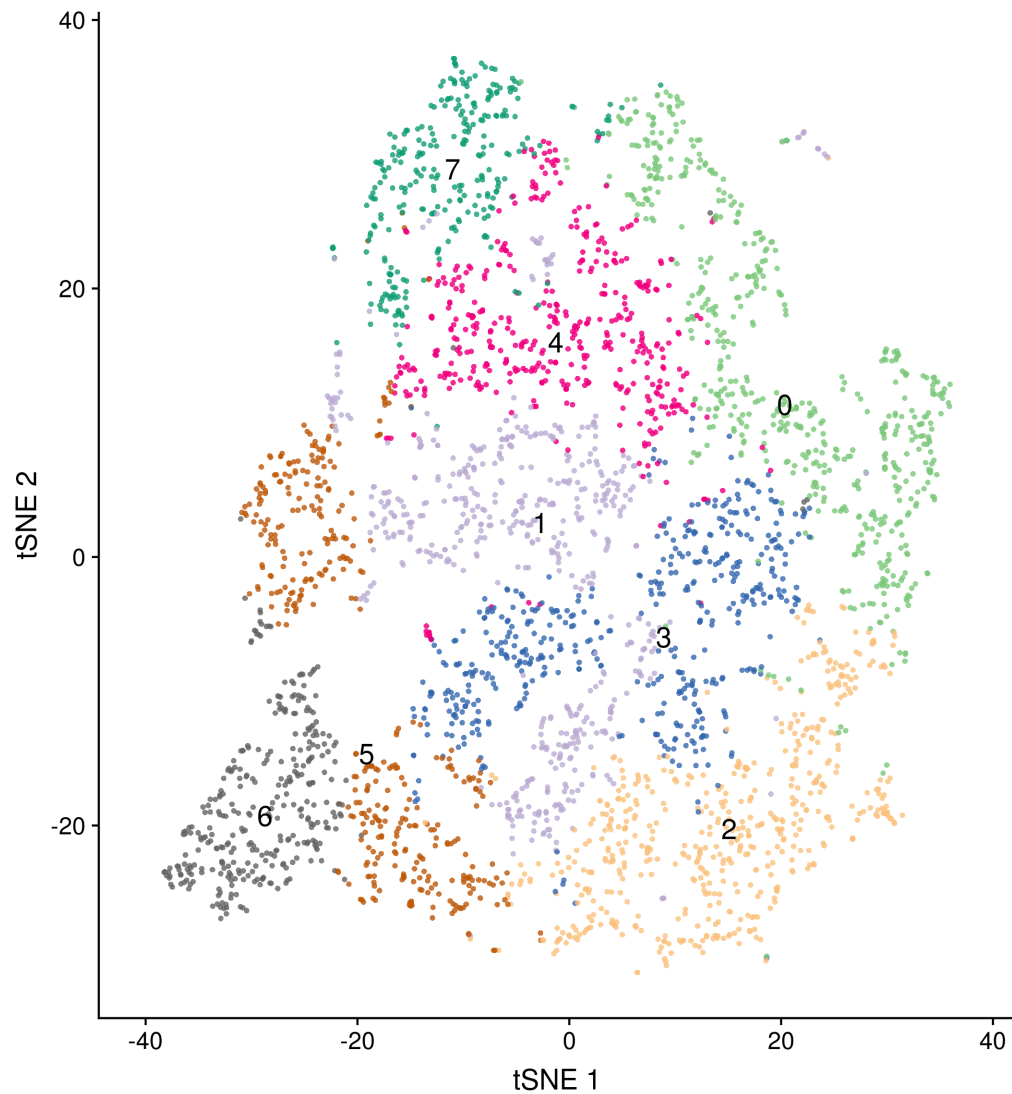
# Seurat + Conos
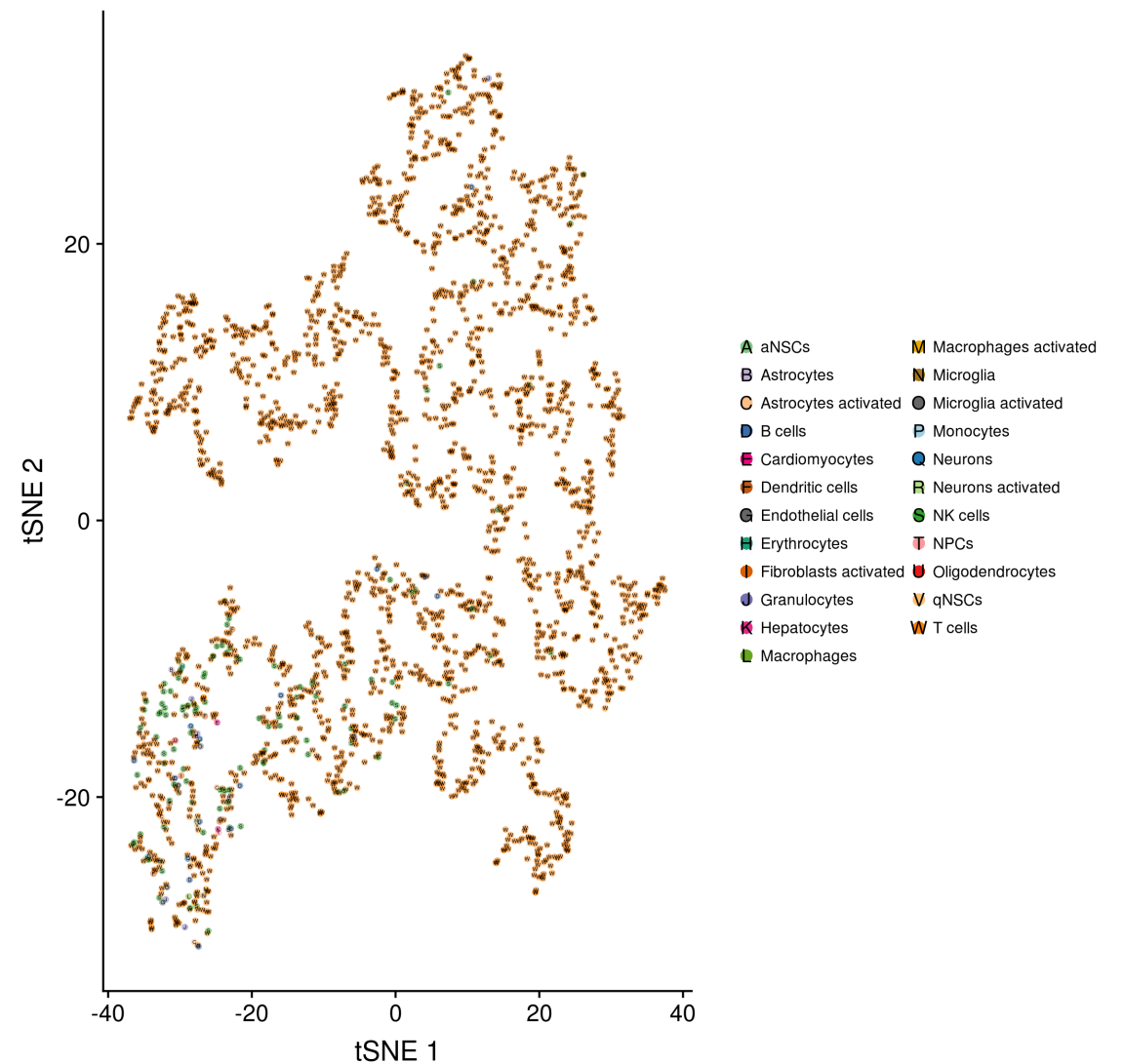


conos clusters

platform

# Seurat + Conos

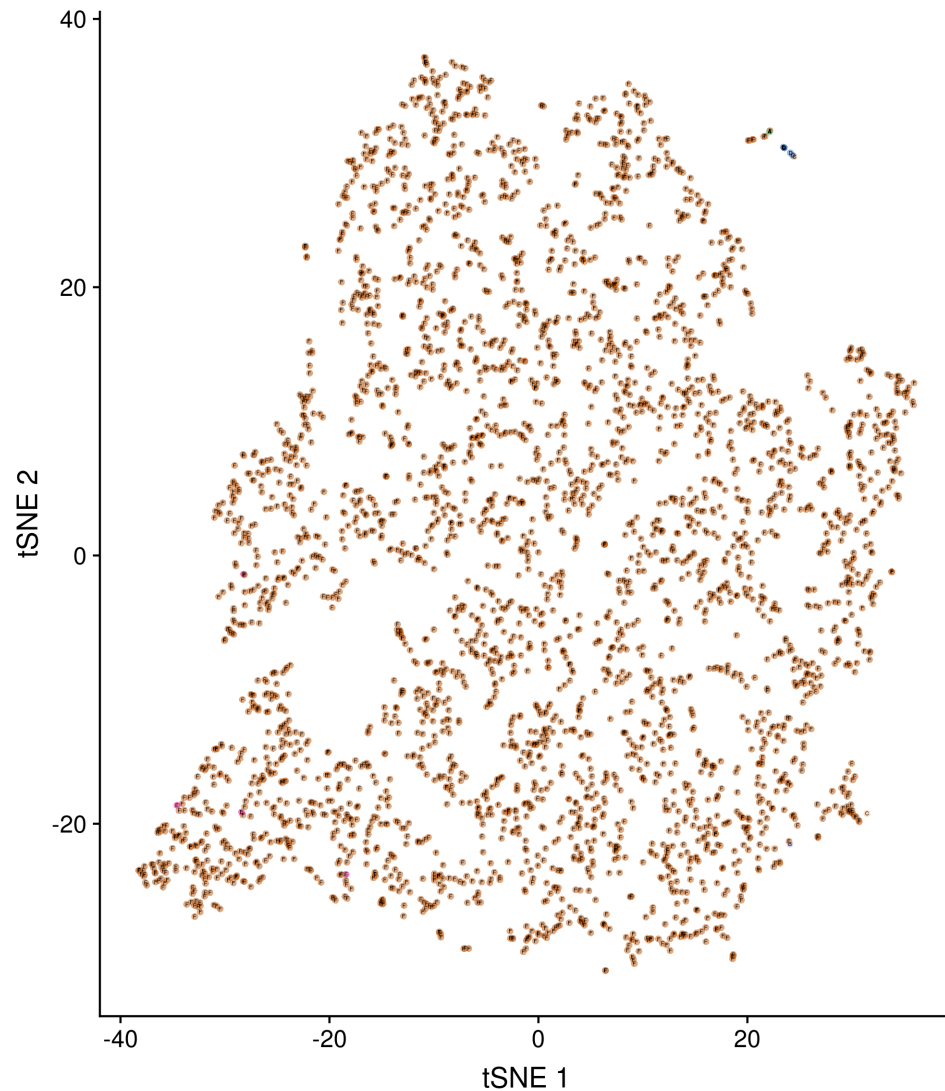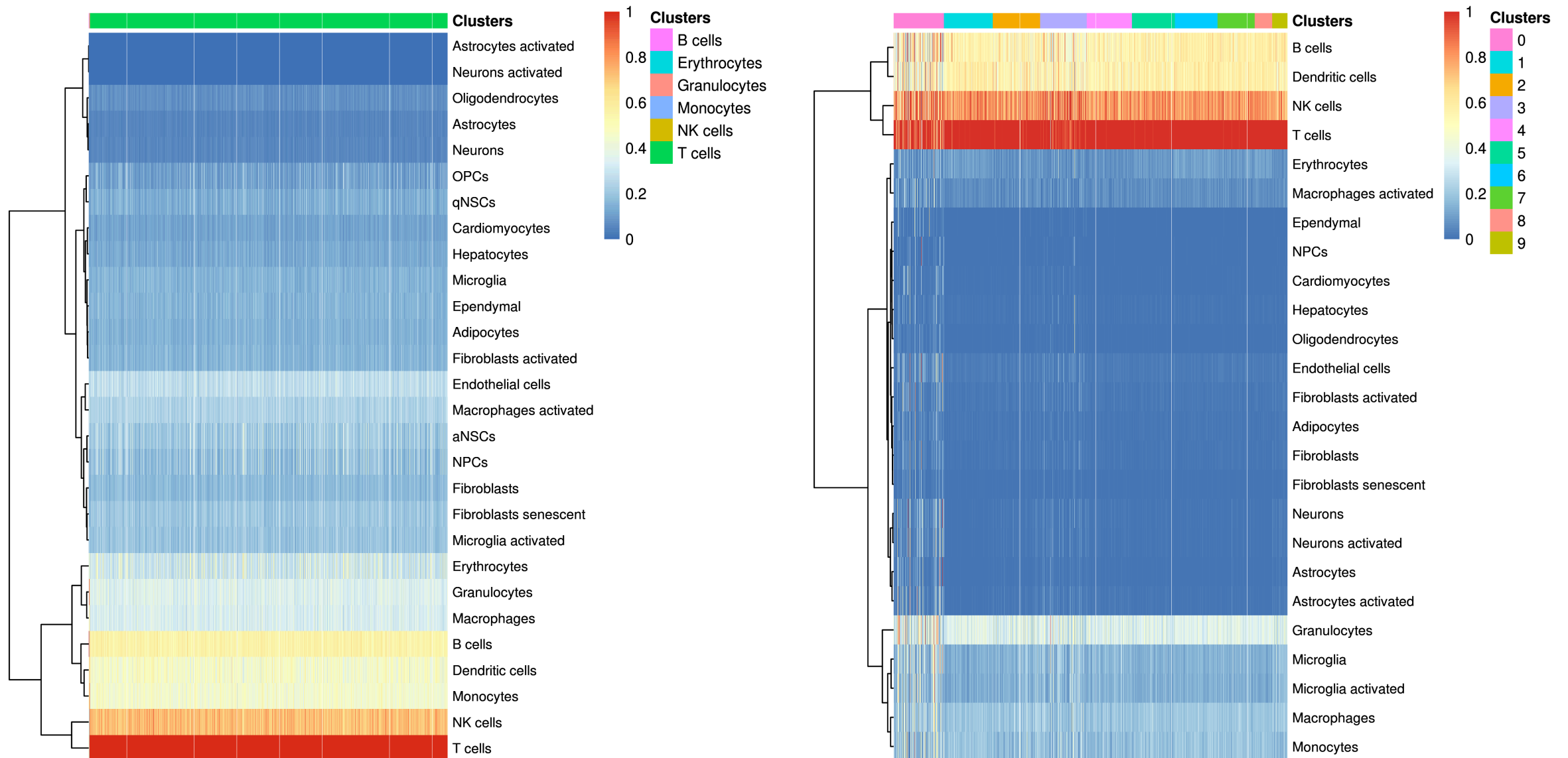# scRNA + scATAC - Compare

# scRNA + scATAC – Compare - SingleR

# scRNA + scATAC – Compare - SingleR

*Question?*