

---

# Analysing single-cell CUT&Tag in tumors

And tools related to single-cell epigenomics analysis

25 / 10 / 2022

Pacôme Prompsy

# Plan

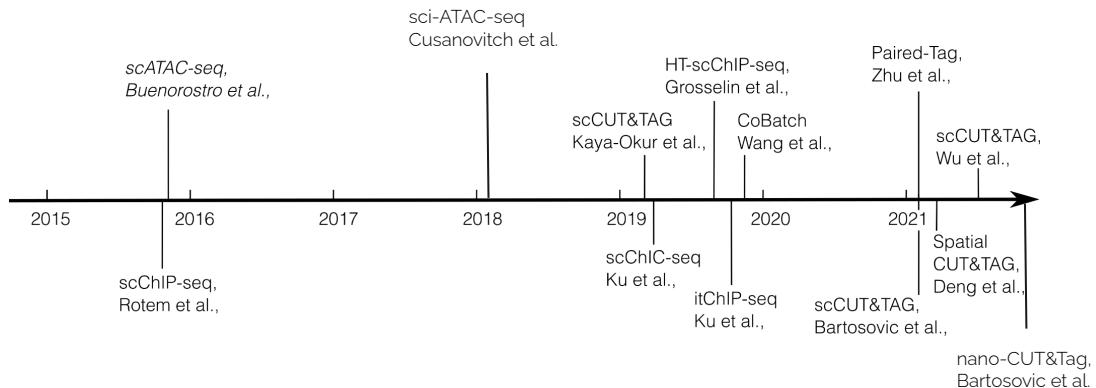
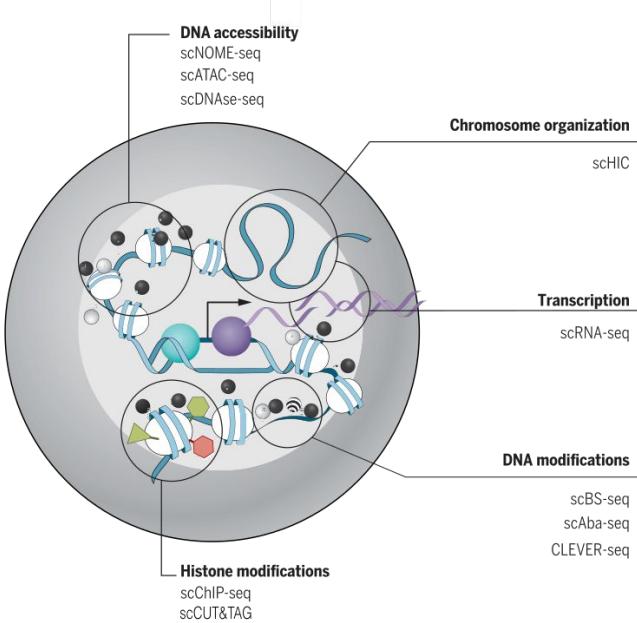
---

- 1) Introducing single-cell epigenomics & challenges in the analysis
  - 2) Quick comparison of scChIP-seq and scCUT&Tag
  - 3) Best practices for single-cell histone modification analysis
  - 4) Analysis of single-cell CUT&Tag data of Juxtatumoral tissue of triple negative breast cancer mouse using ArchR & ChromSCape
- 
- 5) **(If enough time)** Iterative Differential Clustering

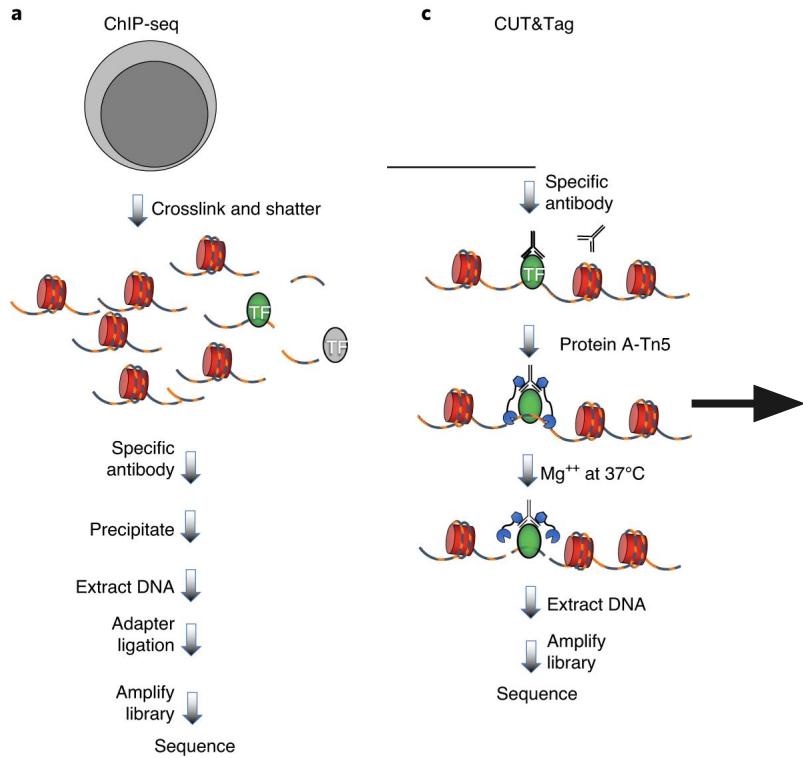
---

# **1 - Introducing single-cell epigenomics & challenges in the analysis**

# Existing Technologies



# ChIP-seq vs CUT&Tag

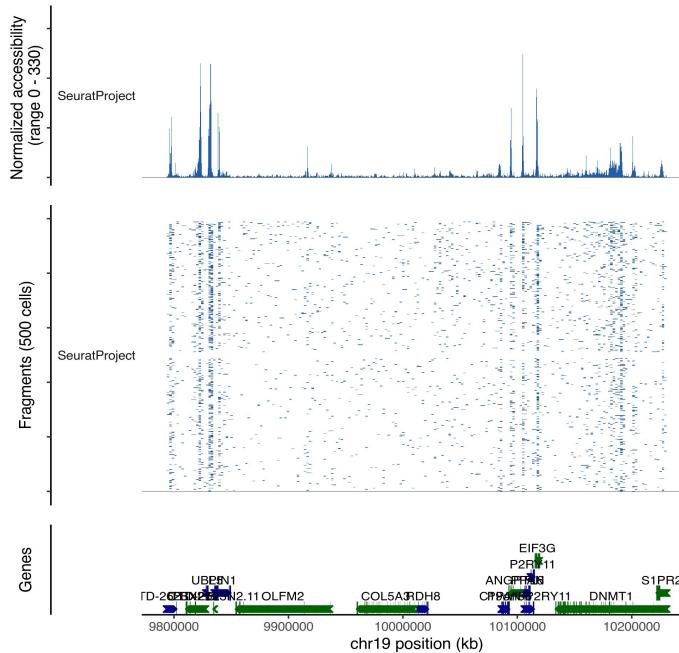


- Simpler protocol
- Uses **transposase** enzyme to both cut and add primers
- Better ratio signal vs noise
- Might be biases toward open chromatin (as coming from ATAC-seq)
- Protocol still needs to be optimized !

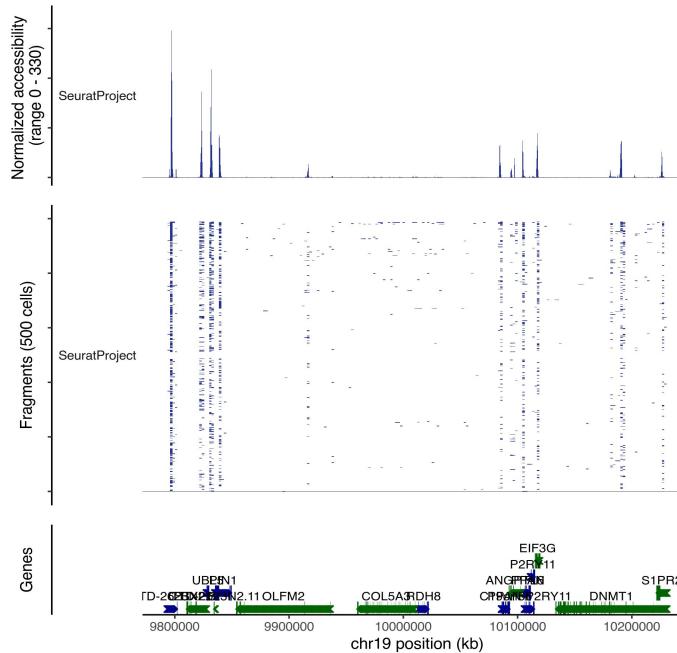
# ChIP-seq vs CUT&Tag



scChIP-seq H3K4me3



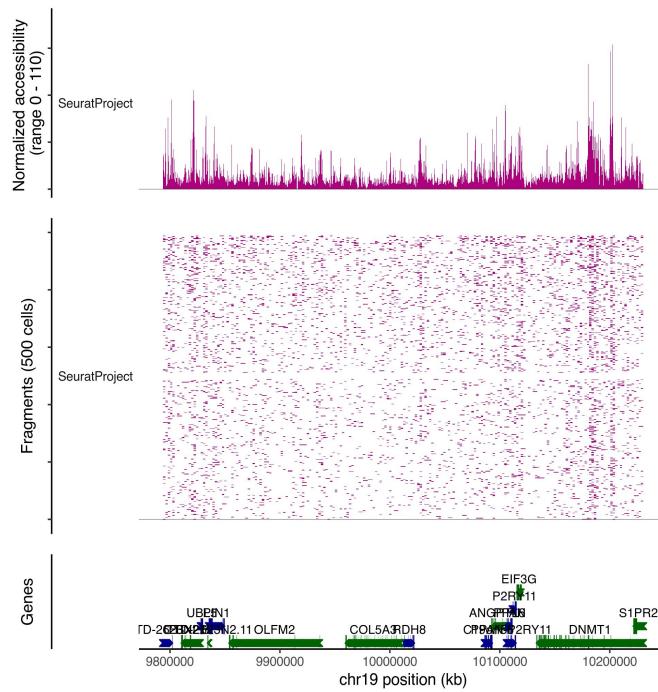
scCutTag H3K4me3



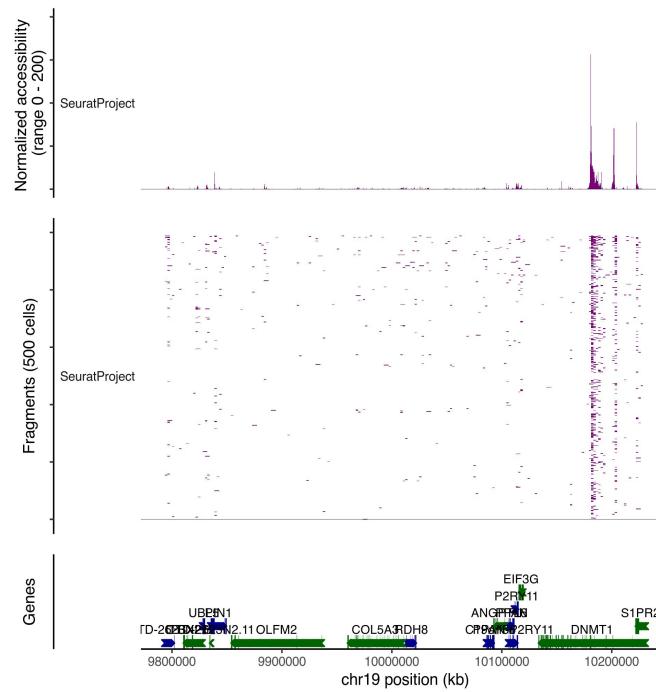
# ChIP-seq vs CUT&Tag



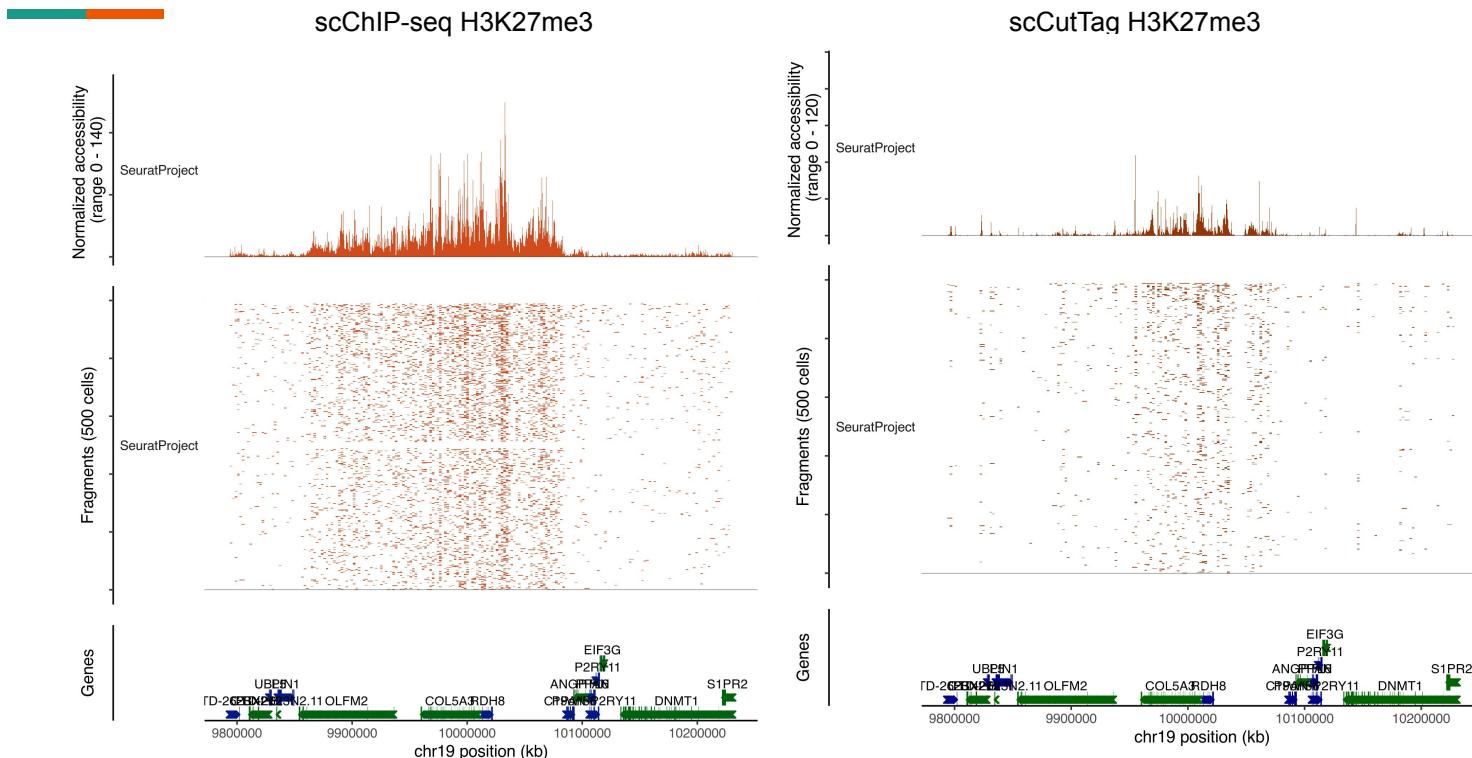
scChIP-seq H3K4me1



scCutTag H3K4me1

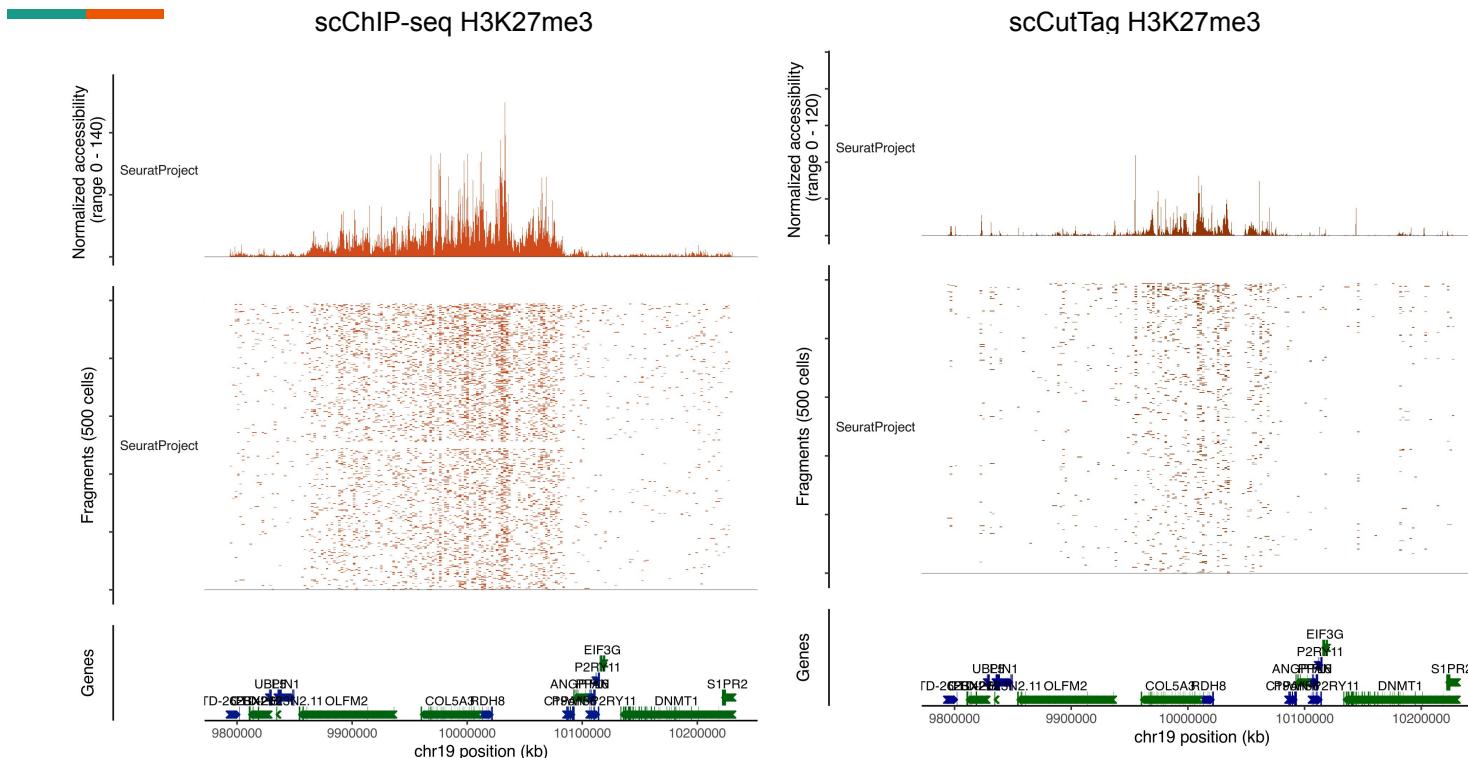


# ChIP-seq vs CUT&Tag

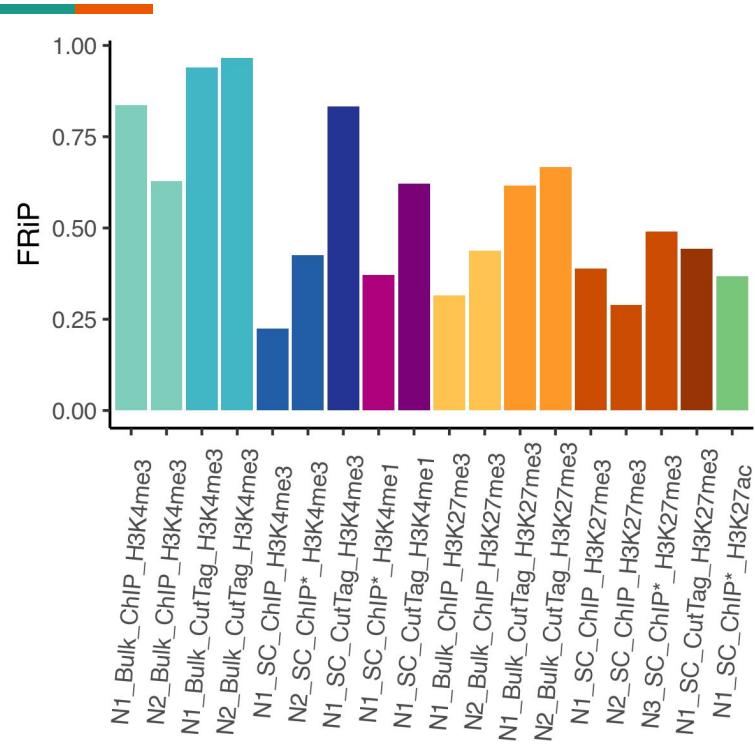


Beware of bias of the Transposase towards open chromatin

# ChIP-seq vs CUT&Tag



Beware of bias of the Transposase towards open chromatin



- The FRIP is always higher in CUT&TAG compared to ChIP-seq: ~33% increase for bulk H3K4me3 ; ~2-fold increase for scH3K4me3 & H3K4me1 & bulk H3K27me3
- The FRIP is constantly higher in bulk experiments than in single-cell experiments (except scChIP H3K27me3\*).
- The FRIP is higher in H3K4me3 compared to H3K4me1, H3K27ac and H3K27me3.

---

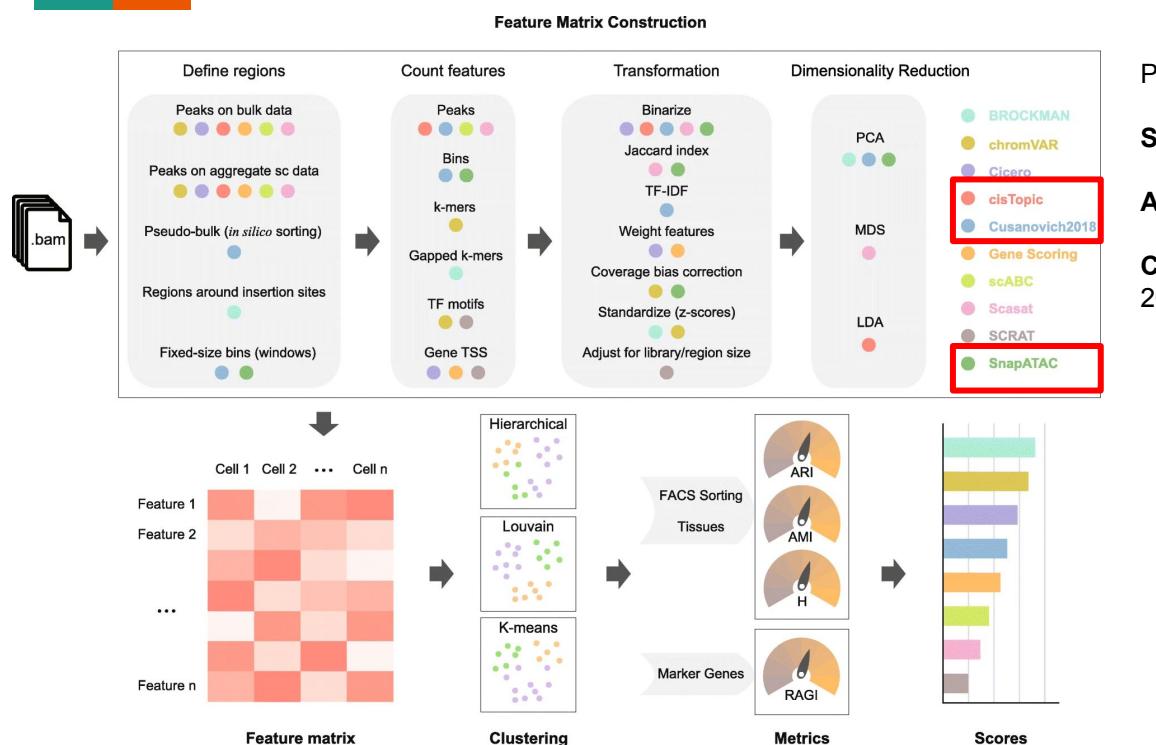
## **2 - Best practices for single-cell histone modification analysis**

# Challenges of scHistone (vs scRNA)

---

- 2 or 4 copies of DNA per cells vs 100s to 1000s of RNA per cell  
→ need to **strongly amplify signal** → lots of duplicates  
  
→ ‘Read recovery’ of 0.5 to 5%, to be compared to **70-95%** in scRNA using 10X.
- Median number of reads per cell is **50x to 100x** lower than in scRNA (**500-1000** reads/cell vs **50,000-100,000** reads/cell).
- There is no ‘natural feature’ (e.g. genes) to create the matrix. Features are usually one of **genomic bins, peaks or gene body + TSS**.

# Existing tools (dedicated to single-cell ATAC-seq)

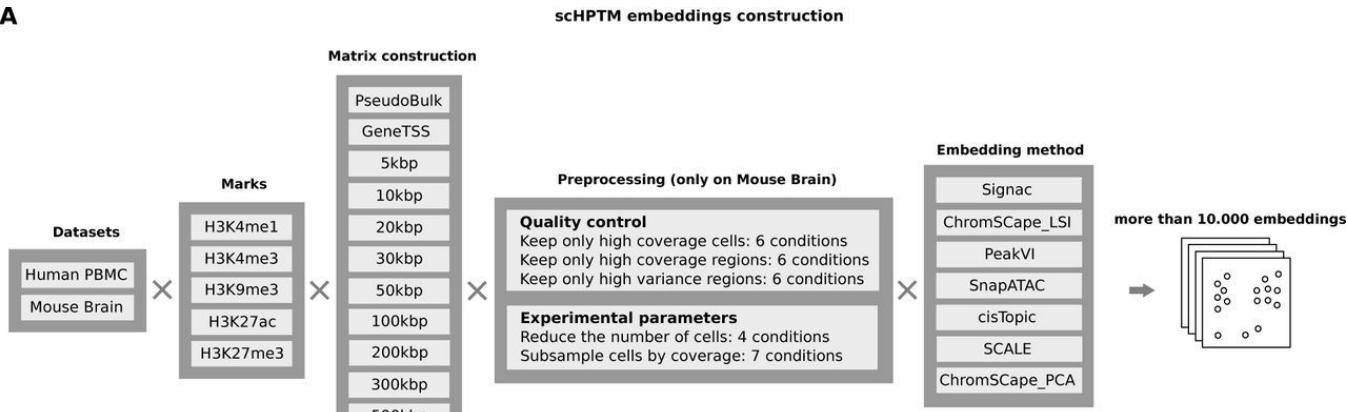


# Best practices for single-cell histone modification analysis

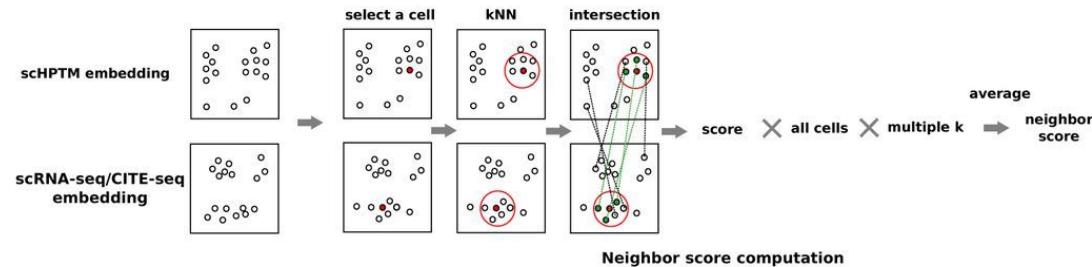


Félix Raimundo

A



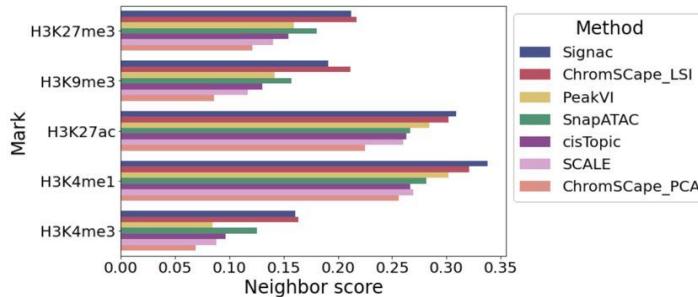
B



**bioRxiv**  
THE PREPRINT SERVER FOR BIOLOGY

# Best practices for single-cell histone modification analysis

A

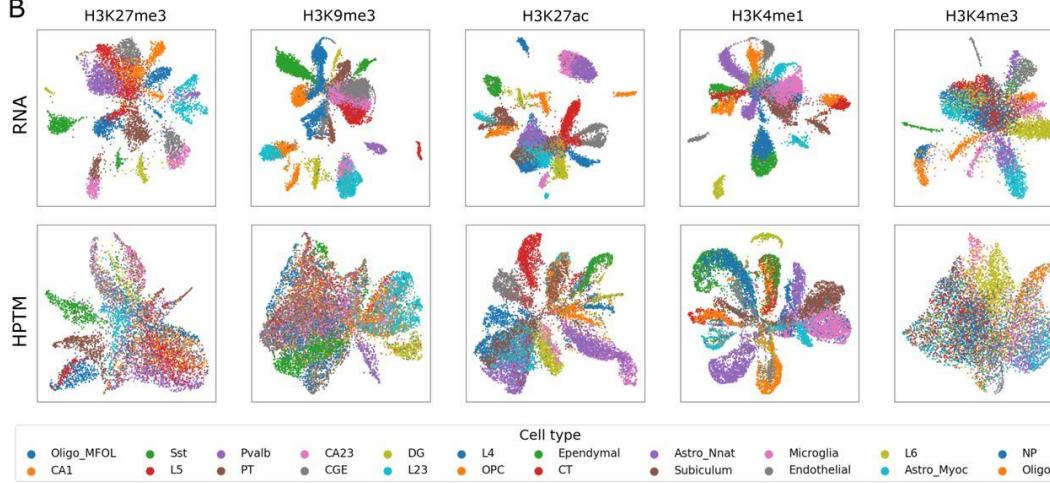


Globally, Latent Semantic Indexing ([LSI](#)) generated the best embeddings.

*This is simply a ‘TF-IDF’ log-normalisation followed by a regular PCA. 1st component is removed.*

*First used in scATAC by Cusanovich et al., 2018*

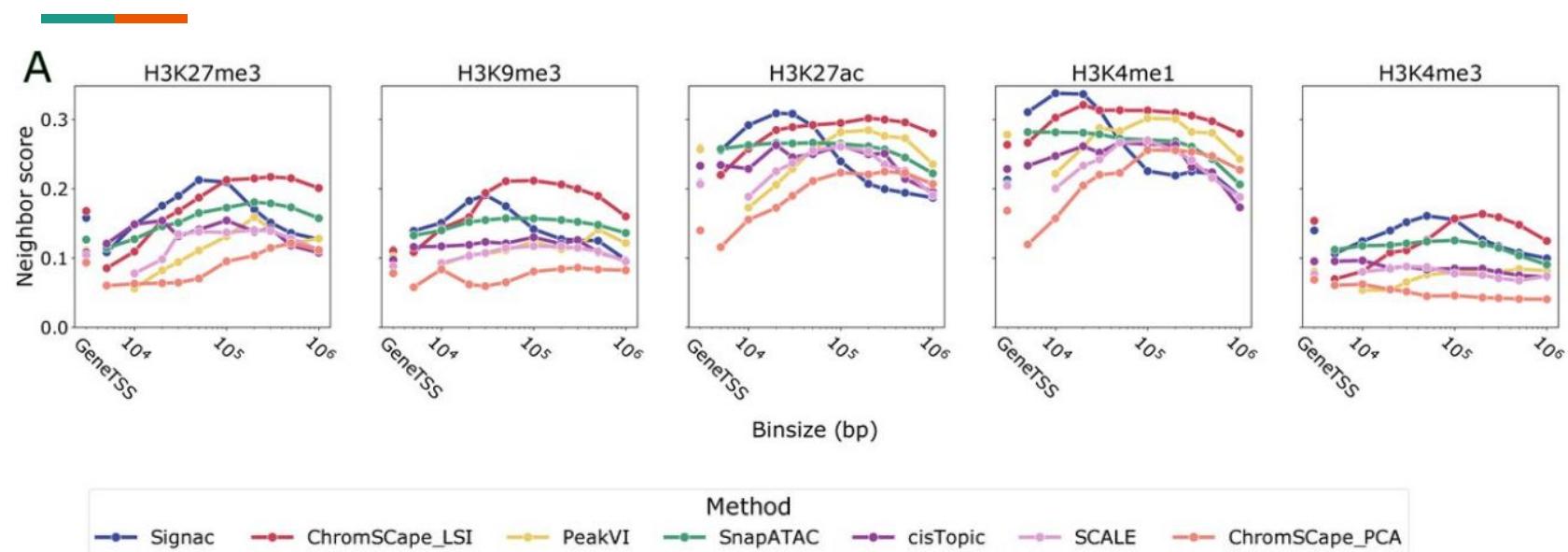
B



## Paired-Tag

Joint profiling of histone modifications and transcriptome in single cells from mouse brain, Zhu et al., 2021

# Choice of feature is highly important



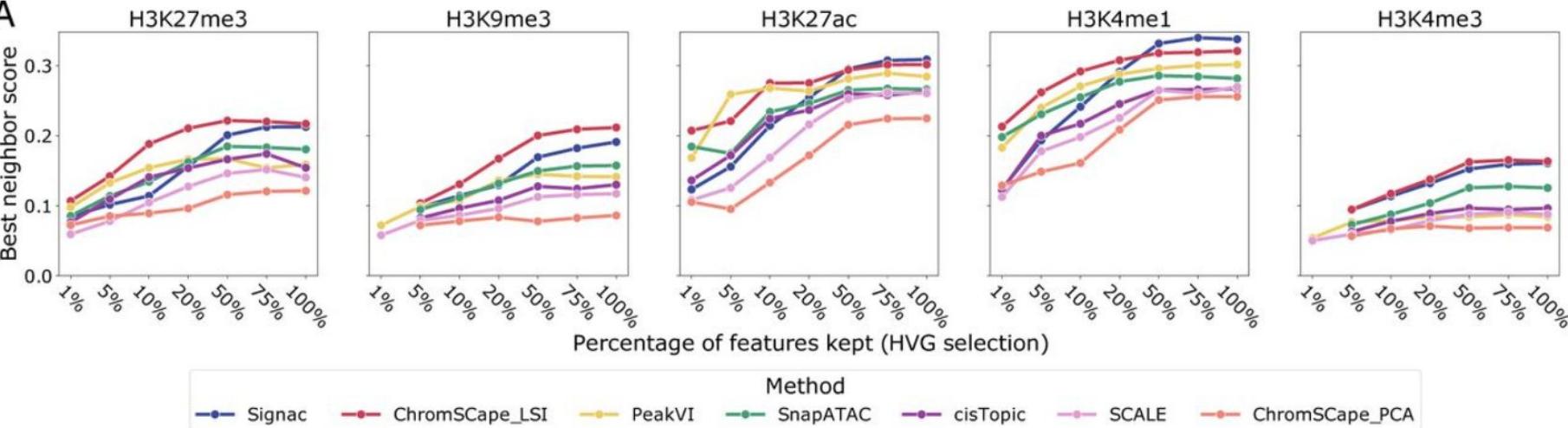
The matrix construction step strongly impacts performances.

Large bins (50kbp - 200kbp) perform overall best, even for active mark known to accumulate in sharp peaks.

Counting on genes and peaks constantly performed worse than genomic bins.

# Feature selection always reduce performances

A

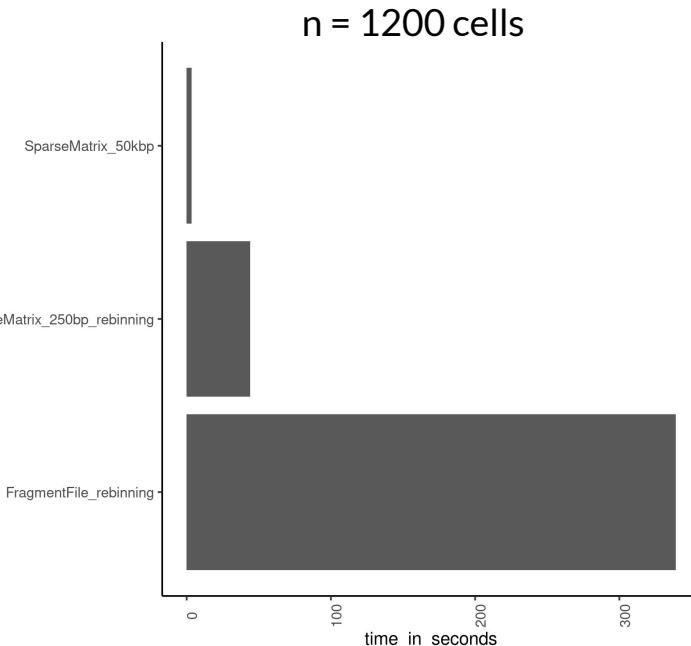


Removing features by either coverage or variance, such as in scRNA-seq always impact negatively performances.

# Digression of input formats and feature choice

---

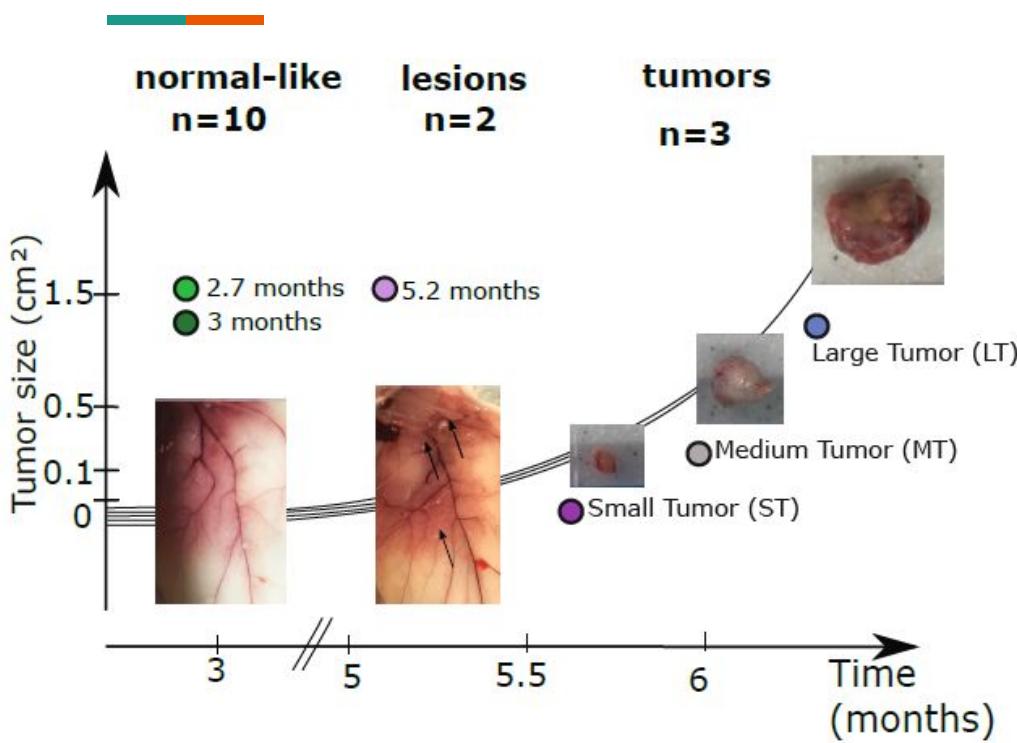
- **Sparse Matrix Format - Processed** - 10X scRNA like (genomic bins, genes or peaks)
  - **pros:** light, quick to read
  - **cons:** cannot re-count to recreate other features, create coverage files and call peaks
- **Fragment File Format - Raw** - BED like
  - **pros:** can re-count on any feature type, create coverage files and call peaks
  - **cons:** long time to read and re-count
- **Sparse Matrix 250-bp format - Semi Raw** - 10X scRNA like
  - + **pros:** light and fast to read, can re-bin quickly to create coverage files and call peaks
  - + **cons:** is an approximation of the real data, not precise (e.g. for TF motif analysis)



---

# **3 - Analysis of Breast Cancer Juxta-tumoral tissue in mouse**

# Analysis of Breast Cancer Juxtatumoral tissue in mouse



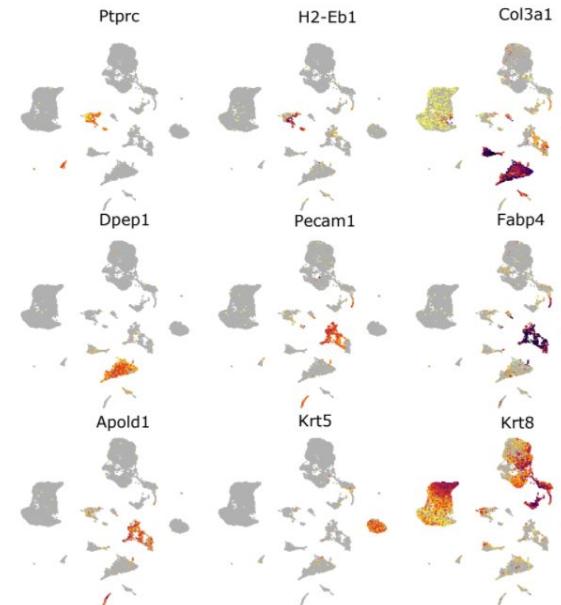
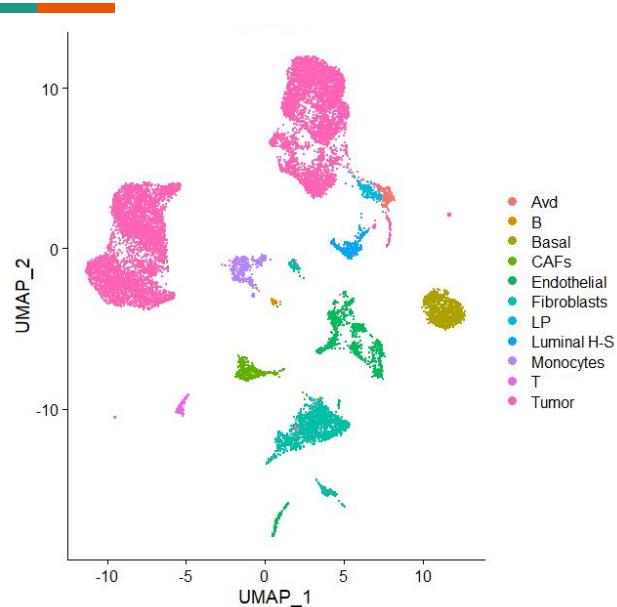
## Abbreviations:

- LP: Luminal Progenitor
- Avd: Alveolar differentiated cells
- ST: Small tumor; MT: Medium tumor;
- LT: Large tumor
- Juxta: Juxta-tumoral : part of tissue close to the tumor site

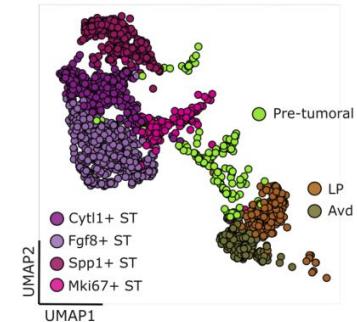


Camille Landragin

# Identification of pre-tumoral LP using scRNA-seq



Identified ‘pre-tumoral’  
Luminal Progenitors in  
juxta-tumoral tissue

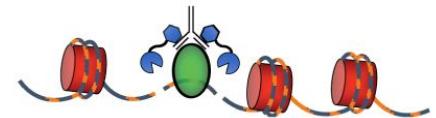


→ Can we find the same cell types as in scRNA-seq ?

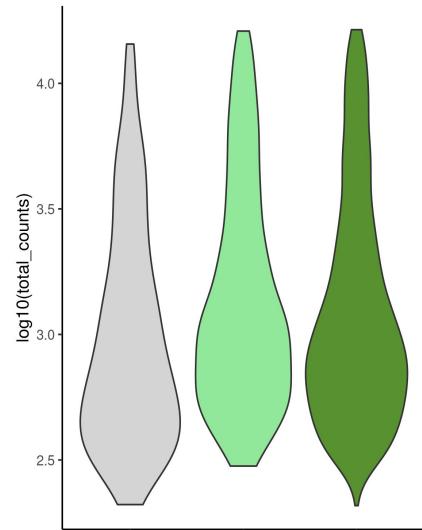
→ Can we identify a ‘LP pre-tumoral’ population in our juxta tissue, based on H3K4me1 landscape only ?

# scH3K4me1 of Juxta and Control tissues

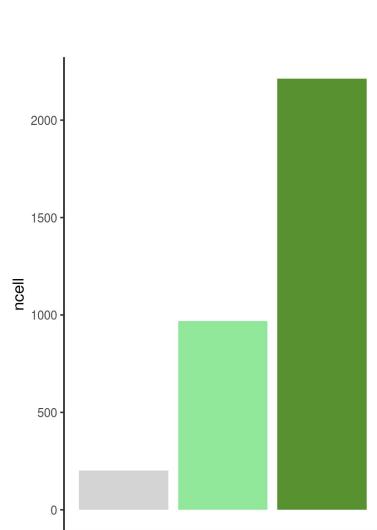
scCUT&Tag + 10X



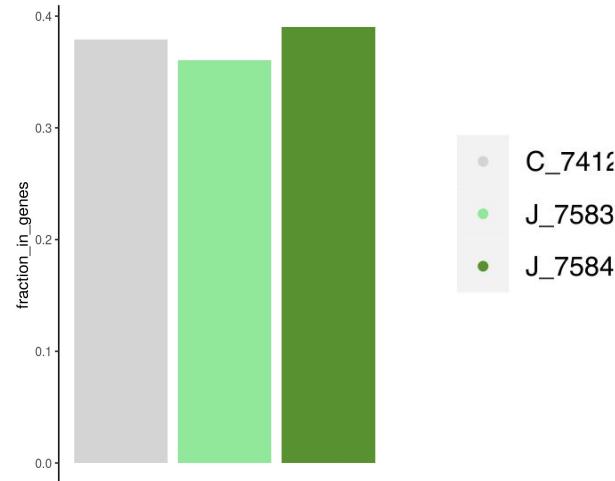
log10 total reads



Number of cells



Fraction of reads around genes  
(genebody + 1000bp)



C\_7412

J\_7583

J\_7584

Genomic bins of 20kbp - No feature removal - Cells > 200 unique fragments & remove top 5%

# ChromSCape - UMAPs

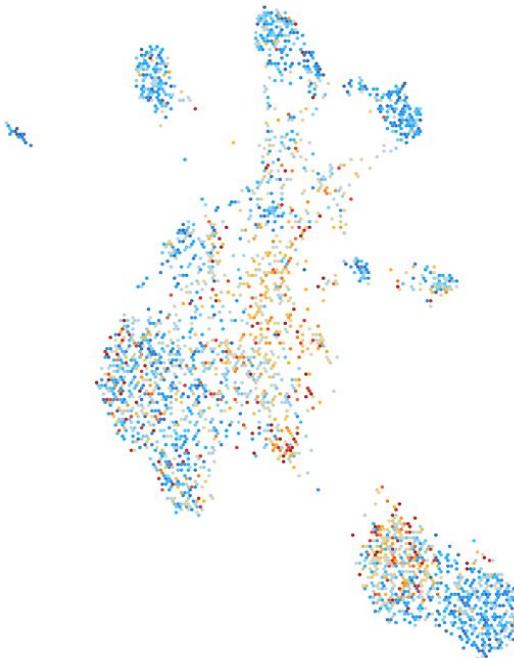
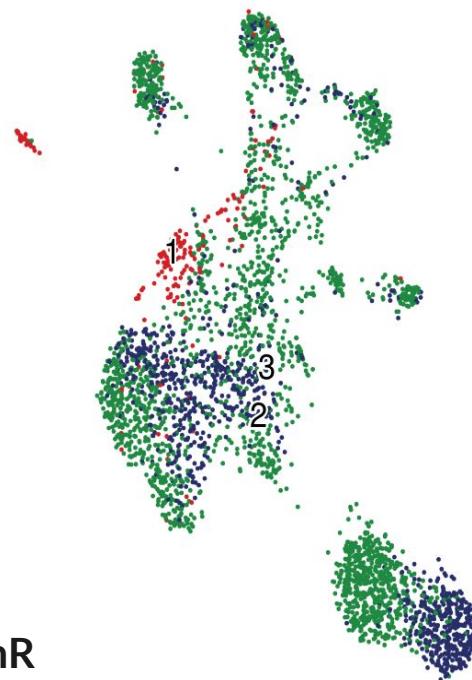
C\_7412\_  
J\_7583\_  
J\_7584\_



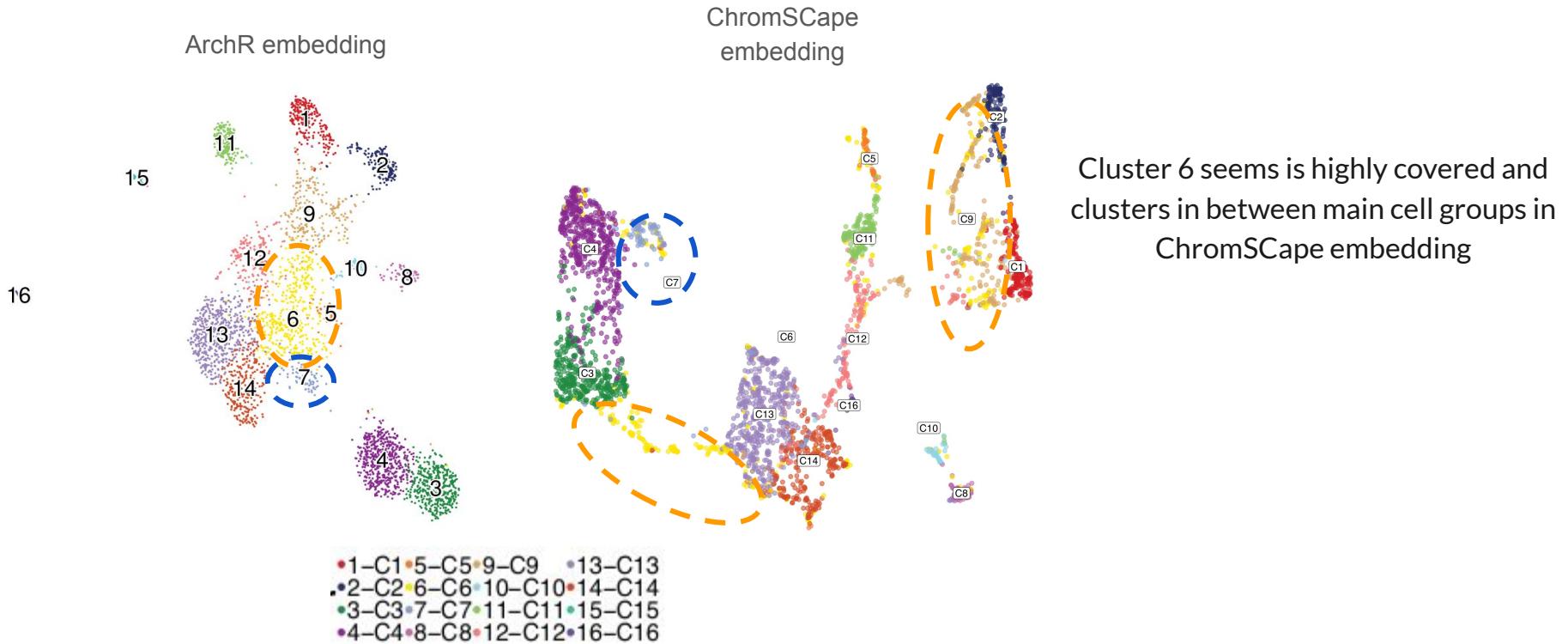
ChromSCape

# ArchR - UMAPs

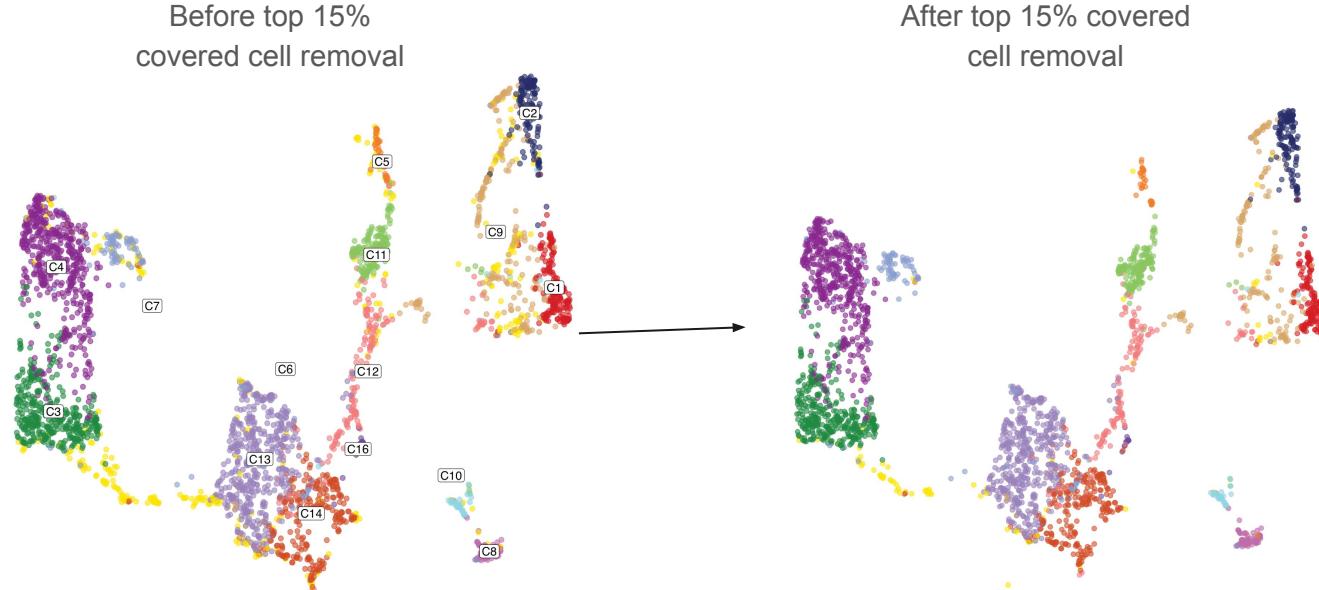
- C\_7412\_mm\_m02y22\_H3K4me1\_250
- J\_7583\_mm\_m06y22\_H3K4me1\_1\_250
- J\_7584\_mm\_m06y22\_H3K4me1\_250



# Identifying a 'strange' cluster



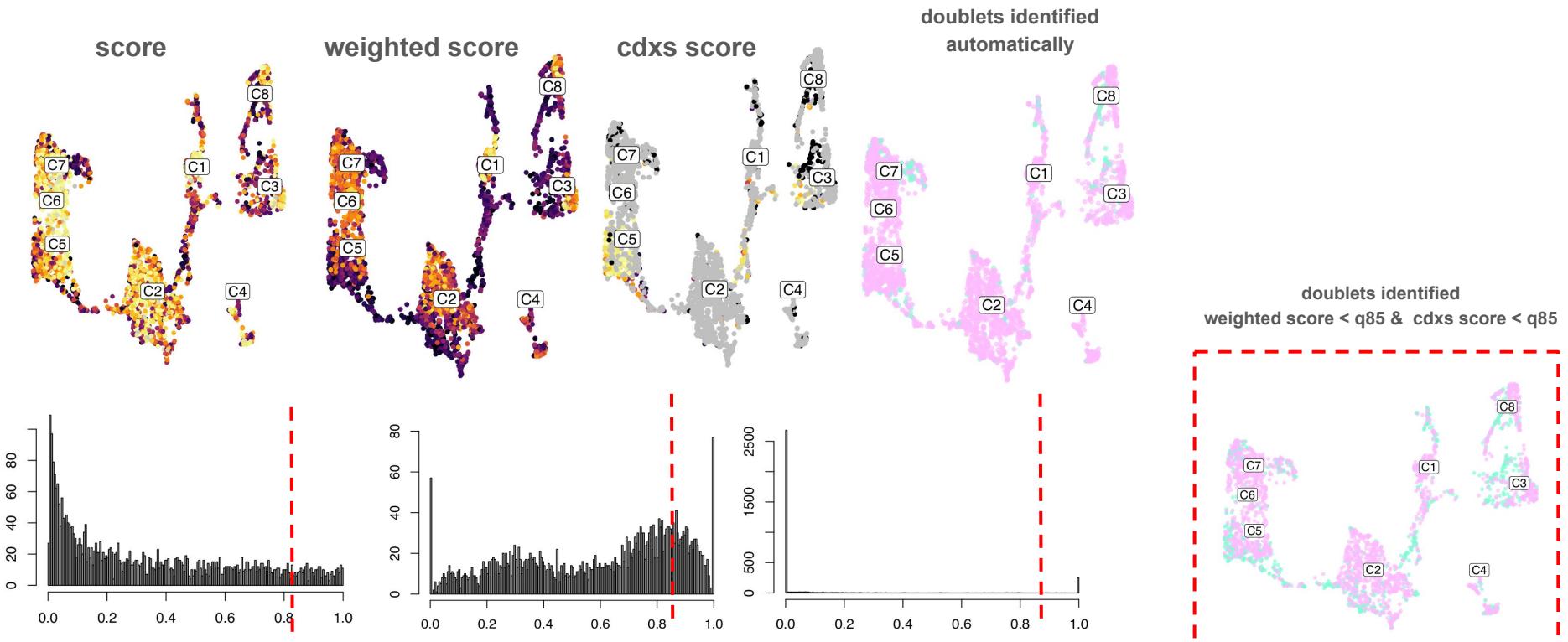
# Removing highly covered cells remove specifically cluster 6



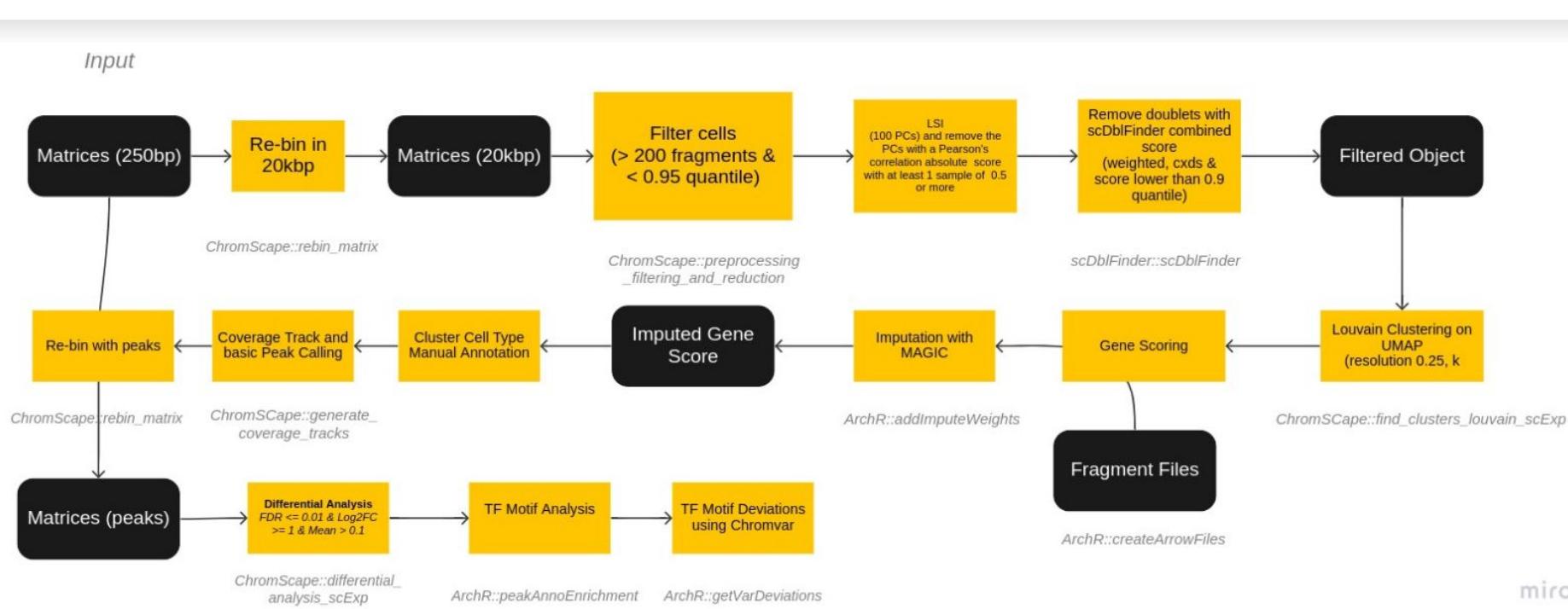
Removing additionally the top 15% of cells removes in priority cells from cluster 6 (80% vs ~8.5% in average)

- 1-C1 • 5-C5 • 9-C9 • 13-C13
- 2-C2 • 6-C6 • 10-C10 • 14-C14
- 3-C3 • 7-C7 • 11-C11 • 15-C15
- 4-C4 • 8-C8 • 12-C12 • 16-C16

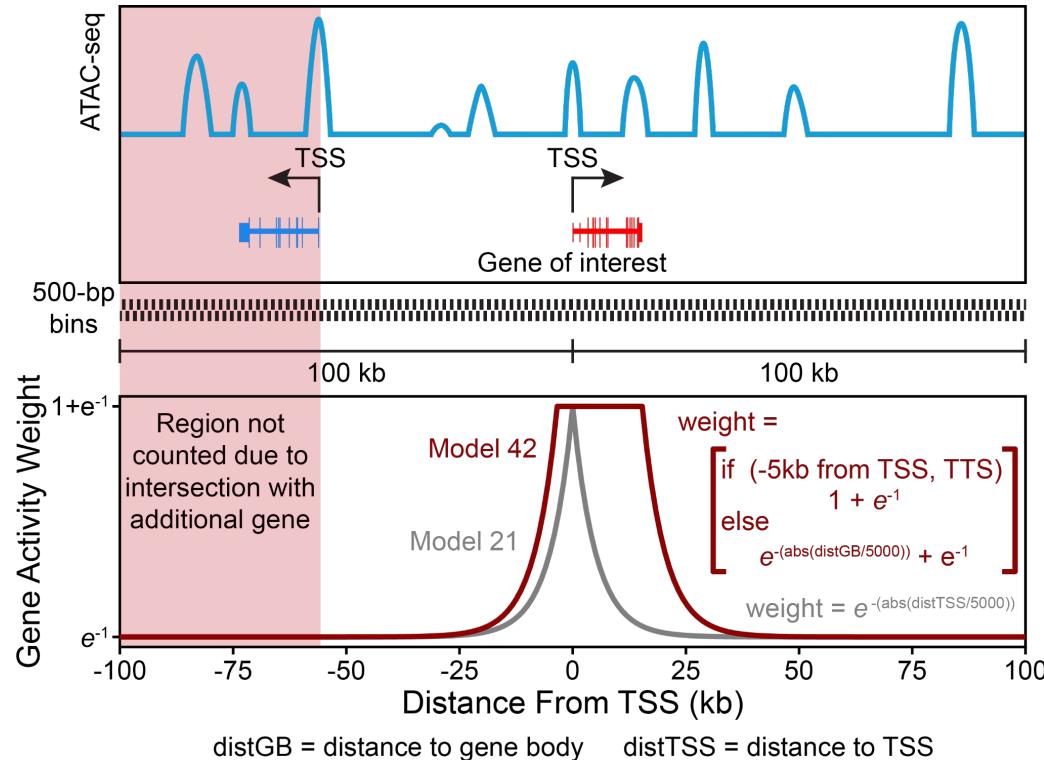
# scDblFinder



# Final 'Roadmap' of the scH3K4me1 analysis

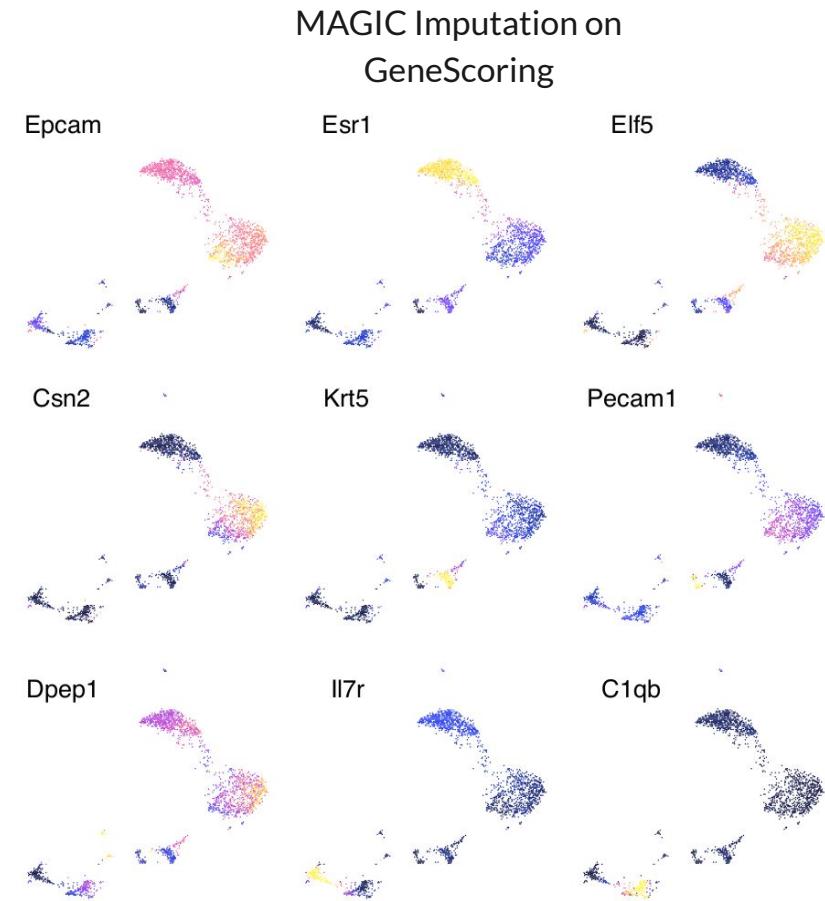
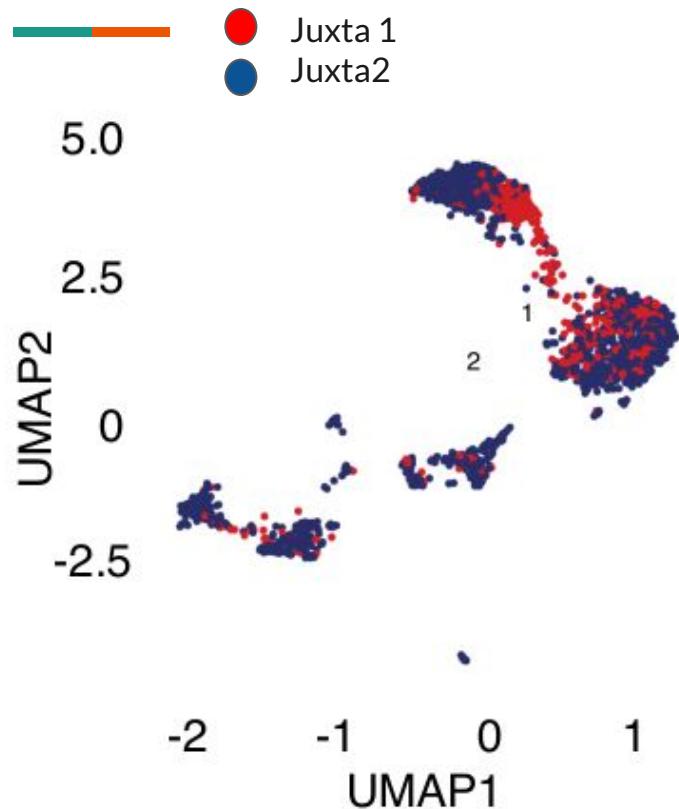


# Calculating Gene Scores with ArchR

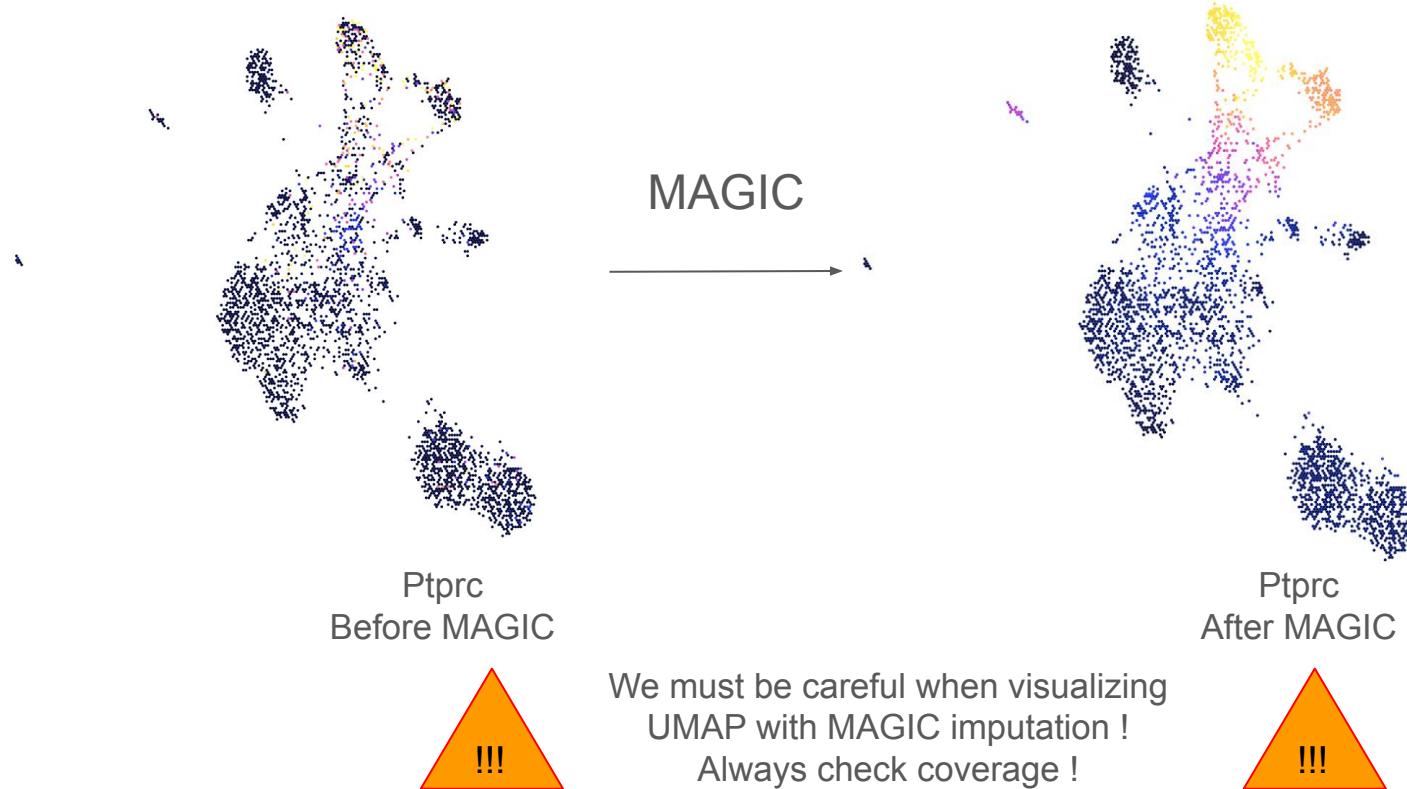


→ ArchR has a specific model to count reads around genes

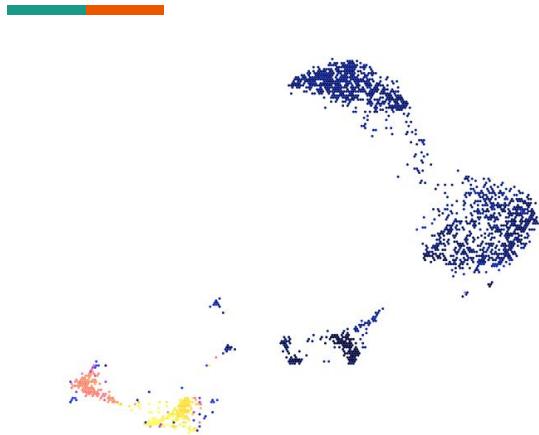
# H3K4me1 of Juxta tissues



# Comparison MAGIC vs Coverage



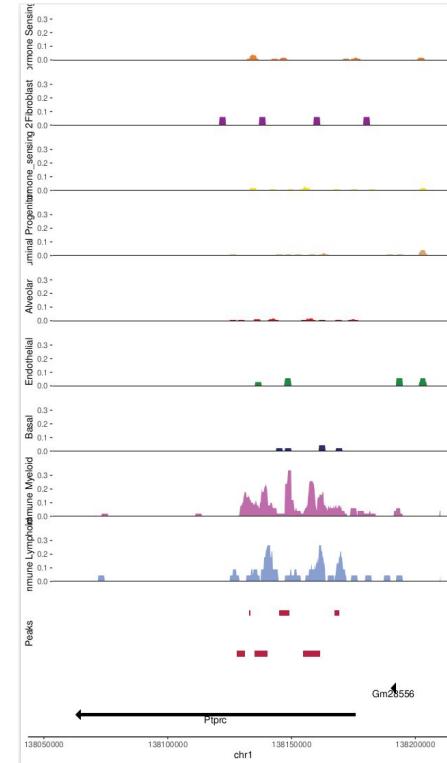
# Comparison MAGIC vs Coverage



Ptprc  
After MAGIC

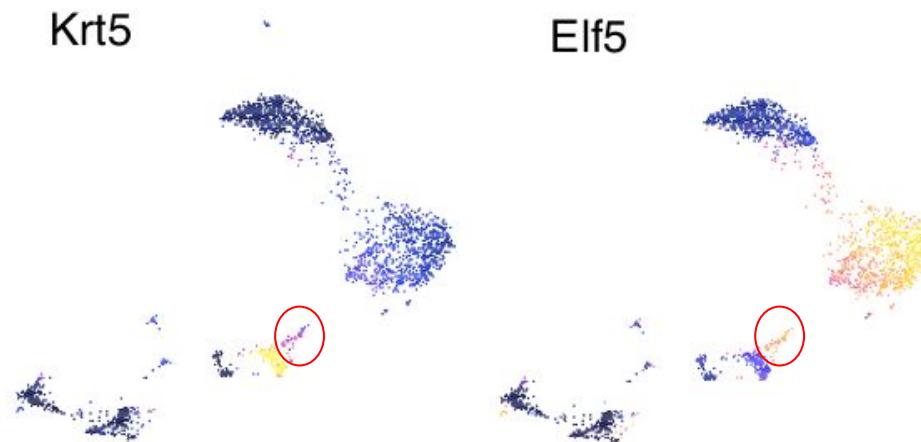


We must be careful when visualizing  
UMAP with MAGIC imputation !  
Always check coverage !

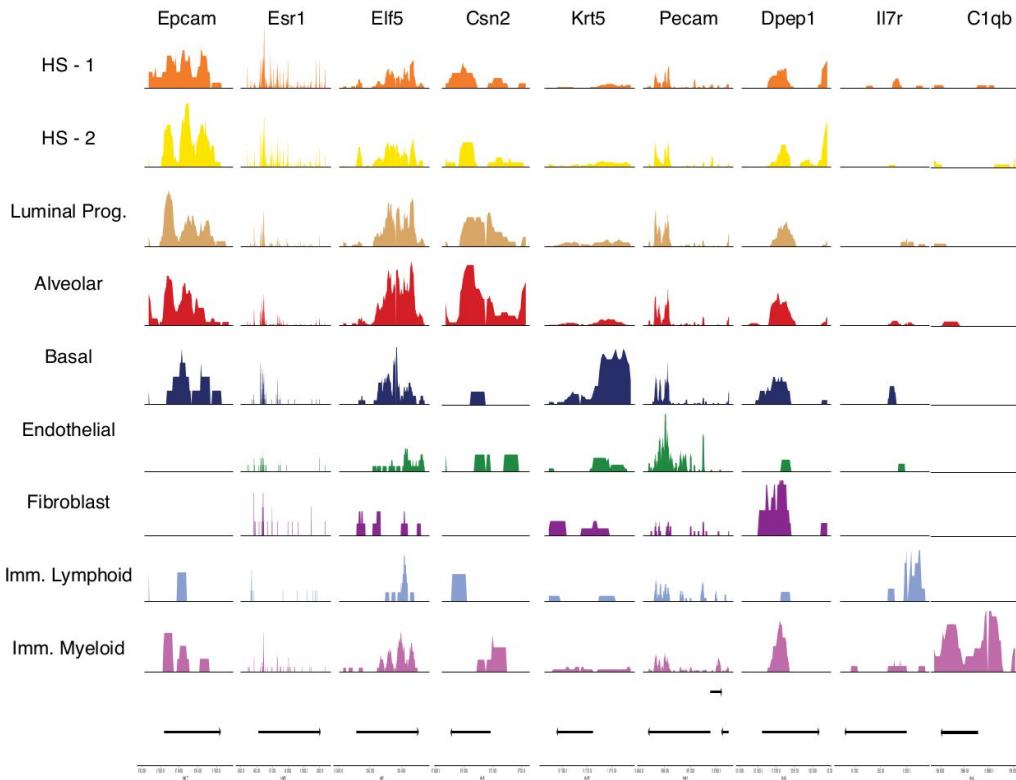
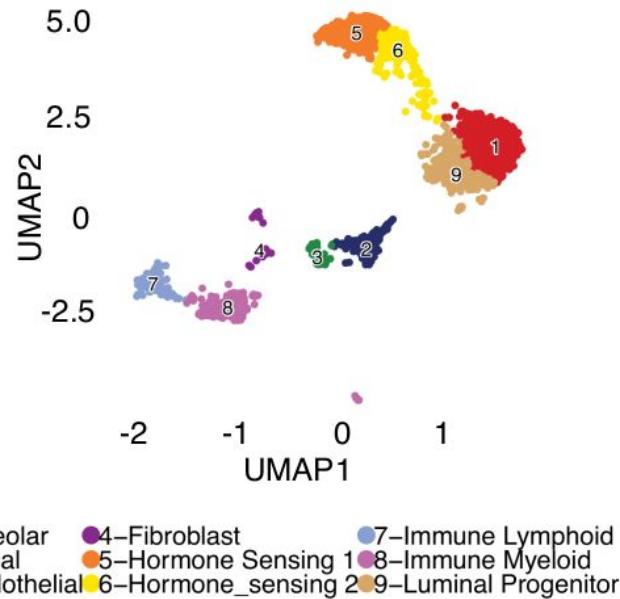


# A potentially interesting population between LP and basal

---

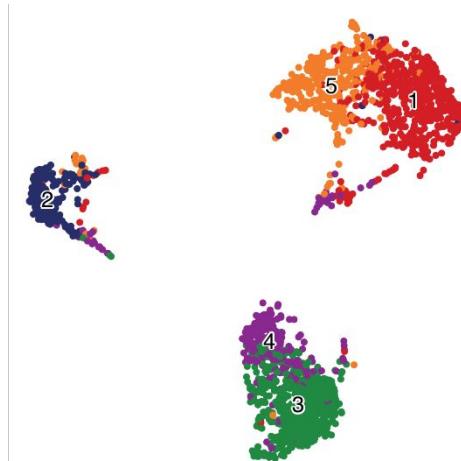


# H3K4me1 of Juxta tissues

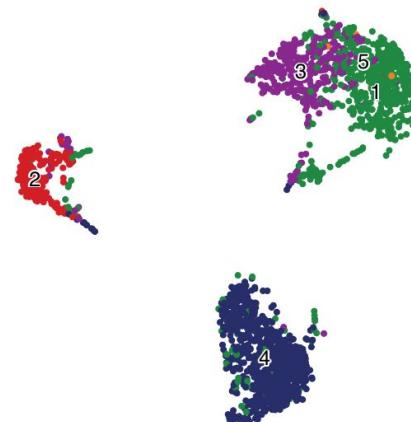


# Label transfer of scRNA annotation on the epithelial compartment using ArchR

Annotation of Epithelial Cells



Label Transfer from scRNA

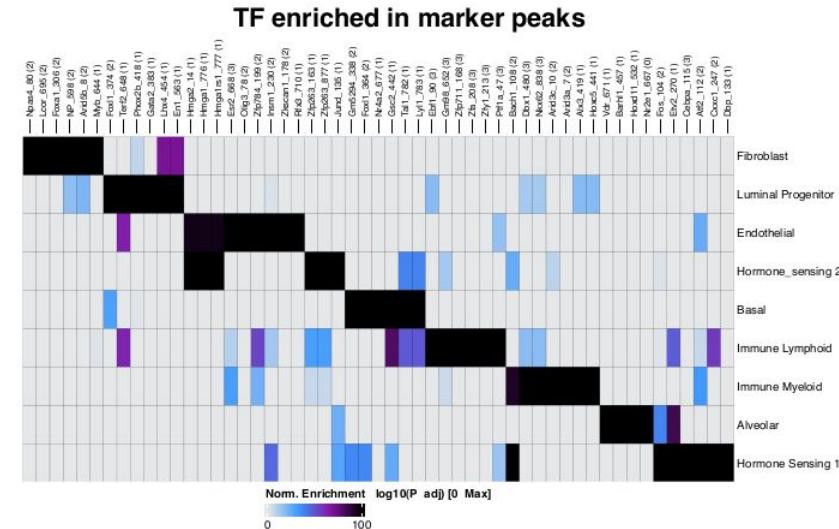
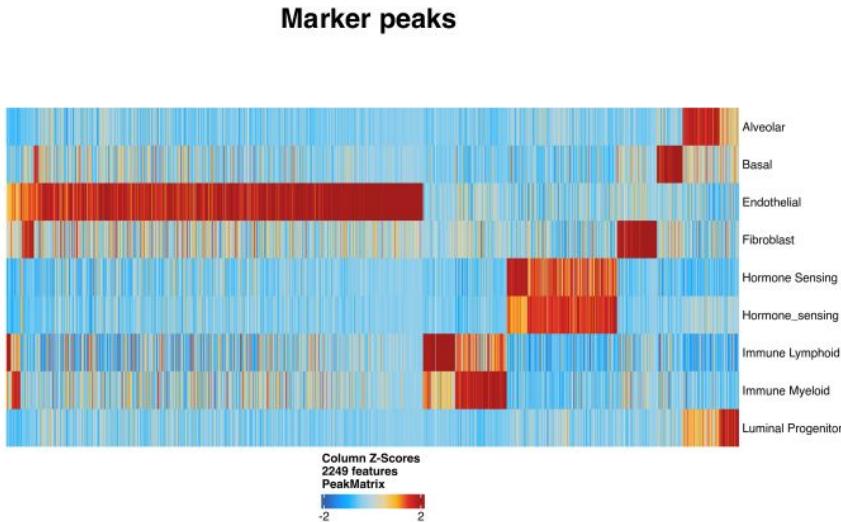


● 1-Alveolar  
● 2-Basal  
● 3-Hormone Sensing 2  
● 4-Luminal H-S  
● 5-Luminal Progenitor

● 2-Basal ● 4-Luminal H-S ● 1-Avd ● 3-LP ● 5-P16+ Pre-lesional

Best ways of integrating

# H3K4me1 of Juxta tissues



Problem → size of peaks might be too large to enrich relevant TF

# Conclusions

---

- The H3K4me1 landscape defines with good granularity the cell types in juxta-tumoral tissue
- Using different tools can be a generally good idea for exploration of a new dataset for example to identify doublets
- To overcome sparsity one can use larger bins, imputation and cluster coverage

---

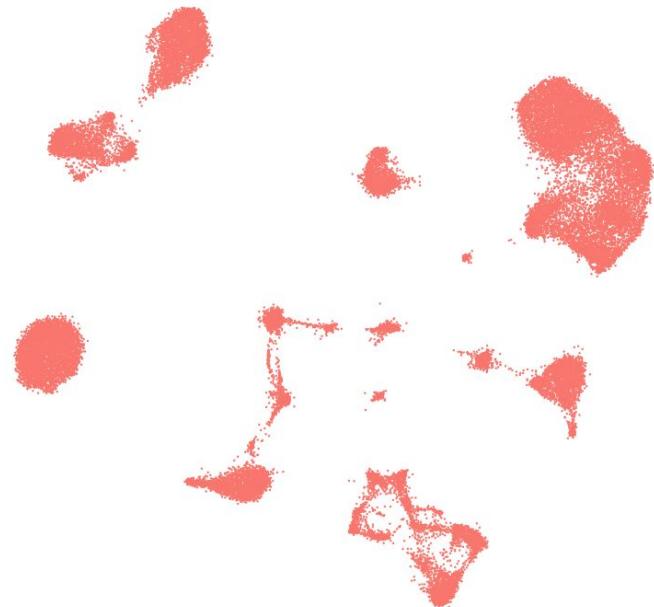
## 4 - Iterative Differential Clustering

*Biologically relevant unsupervised clustering*

Tool in development, looking for feedback !

# What is a 'cell cluster' ?

---



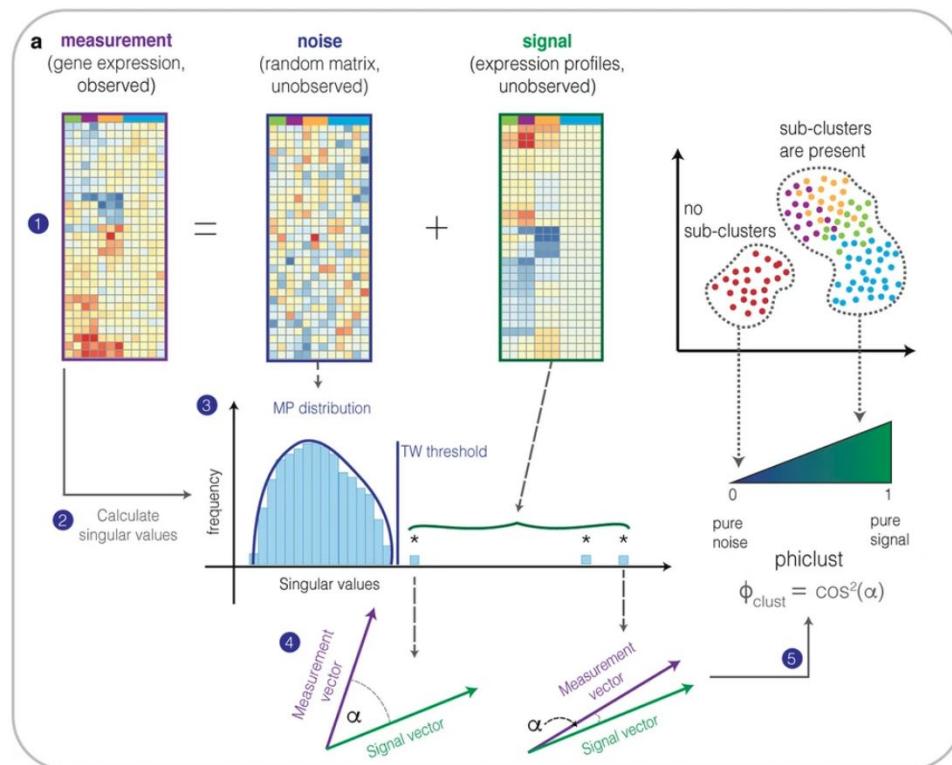
# What is a 'cluster' ? Phiclust

**Phiclust:** a clusterability measure for single-cell transcriptomics reveals phenotypic subpopulations (Mircea et al., 2022)

→ Gives a score cell clusters with non-random substructure using random matrix theory and PCA

→ The score enable users to decide if the cluster can be safely divided into additional subclusters or not

**cons:** ressource intensive (time x memory)



# What is a 'cluster' ? TooManyCells / TooManyPeaks

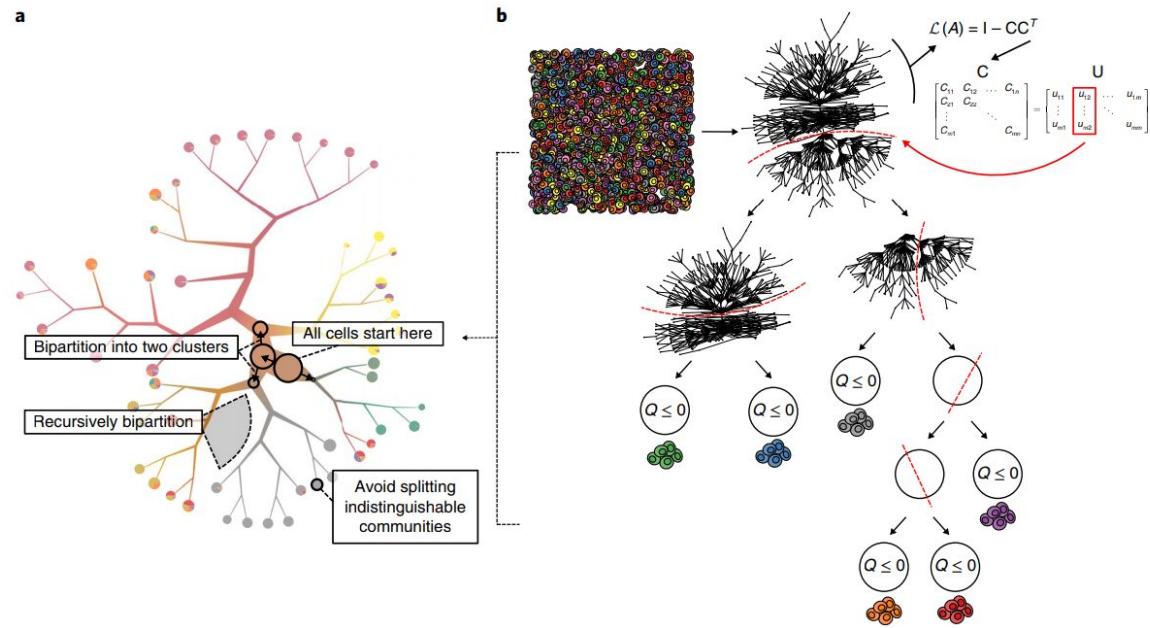
**TooManyCells** identifies and visualizes relationships of single-cell clades  
(Schwartz et al., 2020)

→ Hierarchical bipartition spectral clustering, runs LSI at each step.

→ Uses 'Newman–Girvan' modularity ( $Q$ ) as a stopping criteria instead of an optimization parameter

→ Allows detection of rare populations not found by other algorithms

**cons:** tends to produce high number of very small clusters



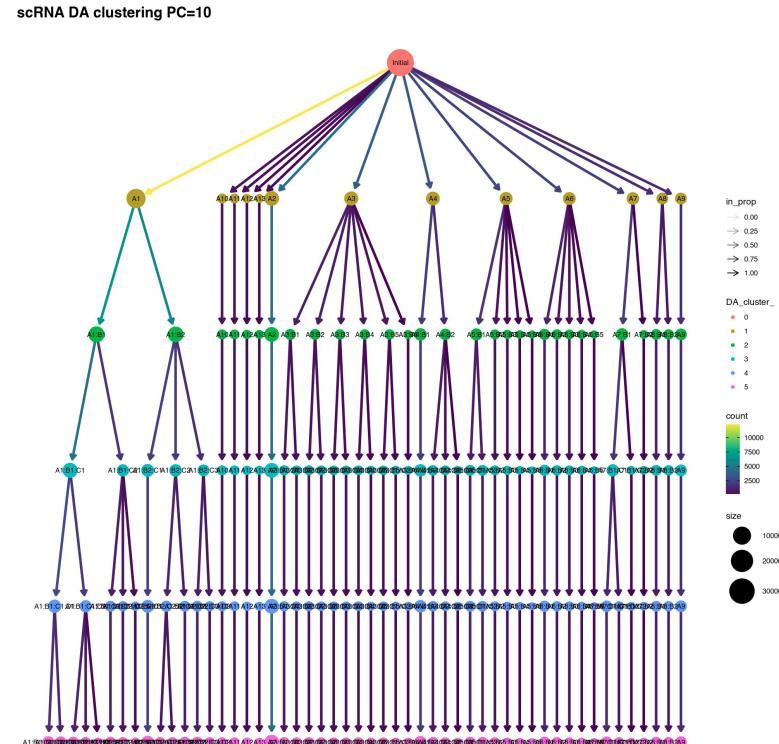
# What is a 'cluster' ? Clustree

---

*Clustree enables visualization of clusters at different resolutions.*

→ If the clusters are stable when changing the resolution, this should mean that the clusters are stable

**cons:** Do not indicate the number of clusters to choose.



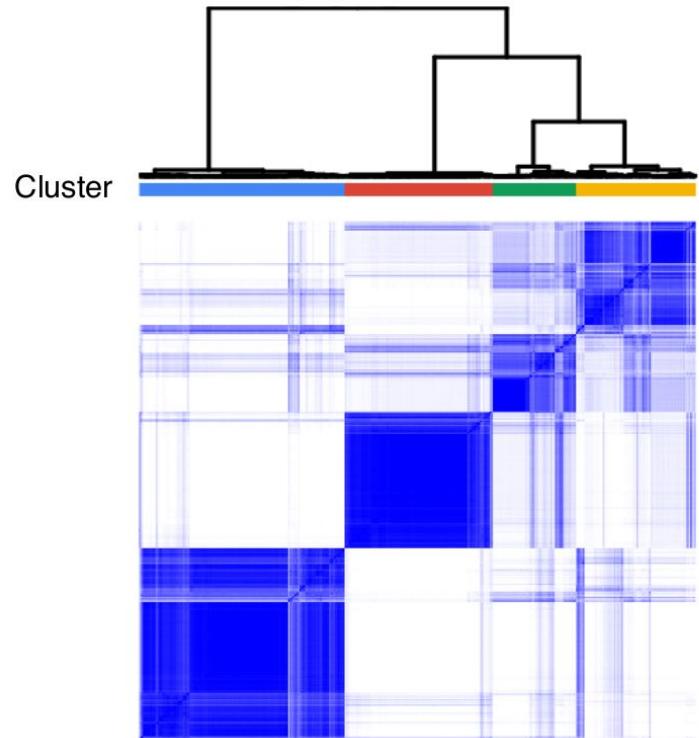
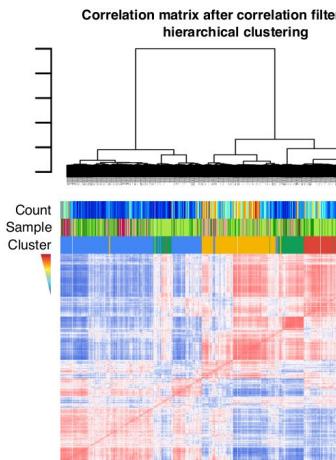
# What is a 'cluster' ?

## Consensus Clustering

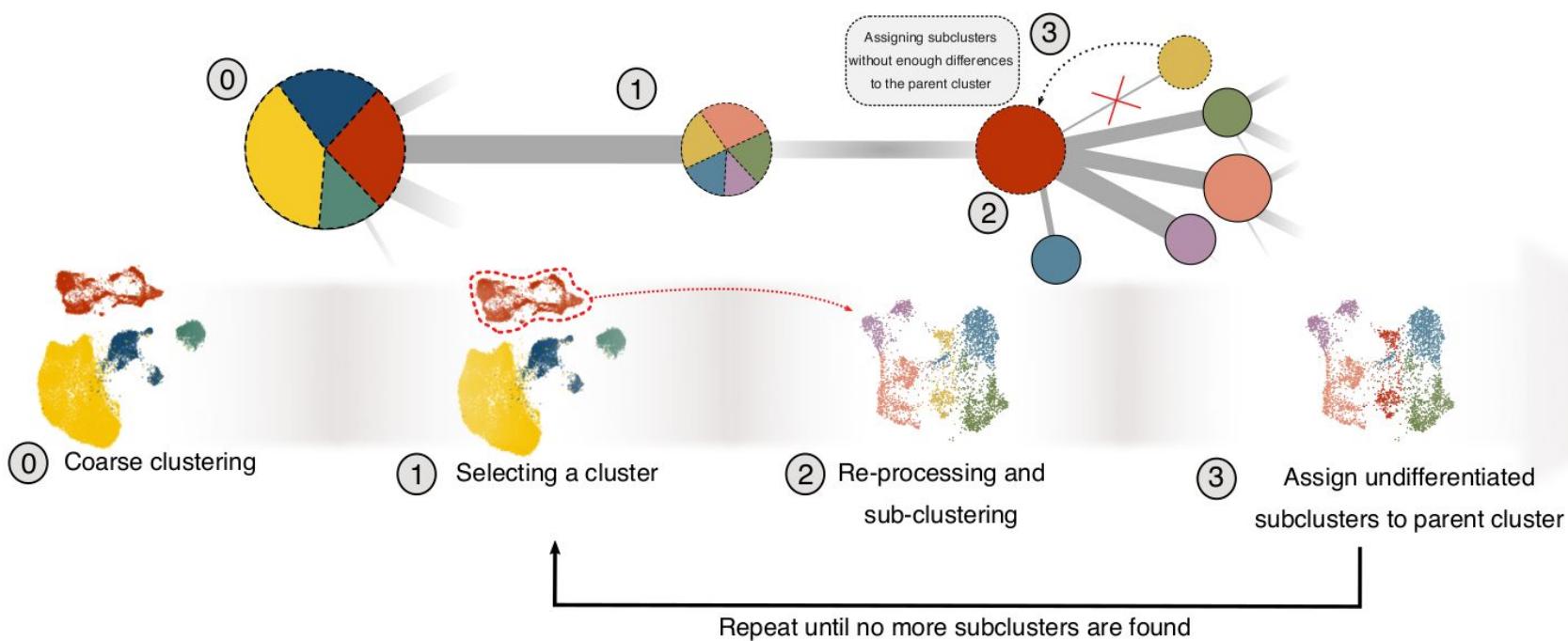
**ConsensusClusterPlus** or **SC3** enables to bootstrap clustering to assess

→ If the clusters are stable when subsampling the cells, this should mean that the clusters are stable

**cons:** ressource intensive (time x memory)

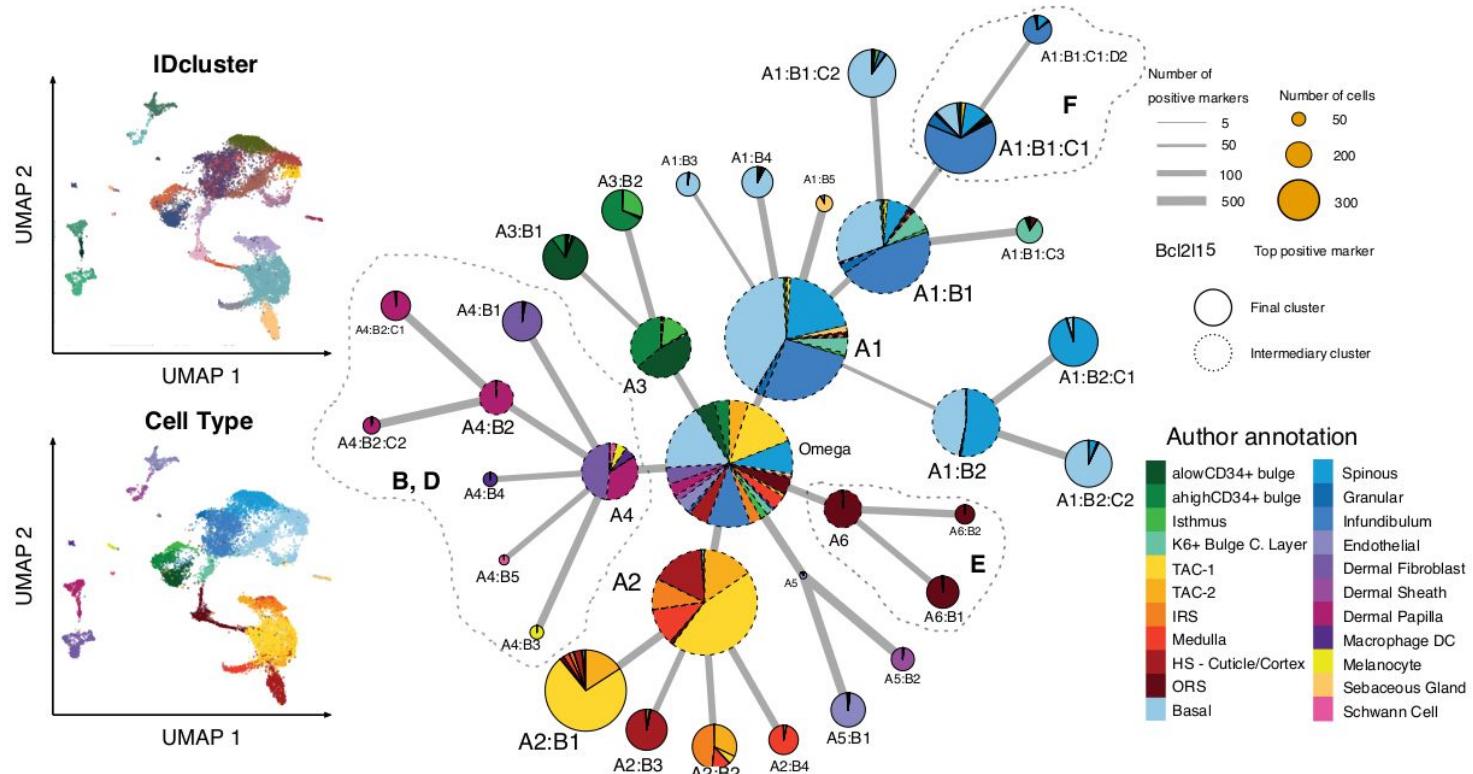
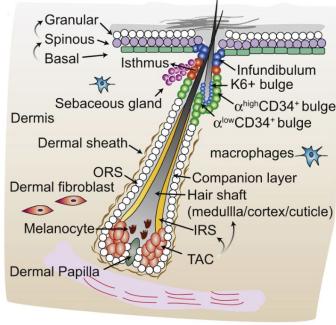


# What is a 'cluster' ? Iterative Differential Clustering (IDclust)

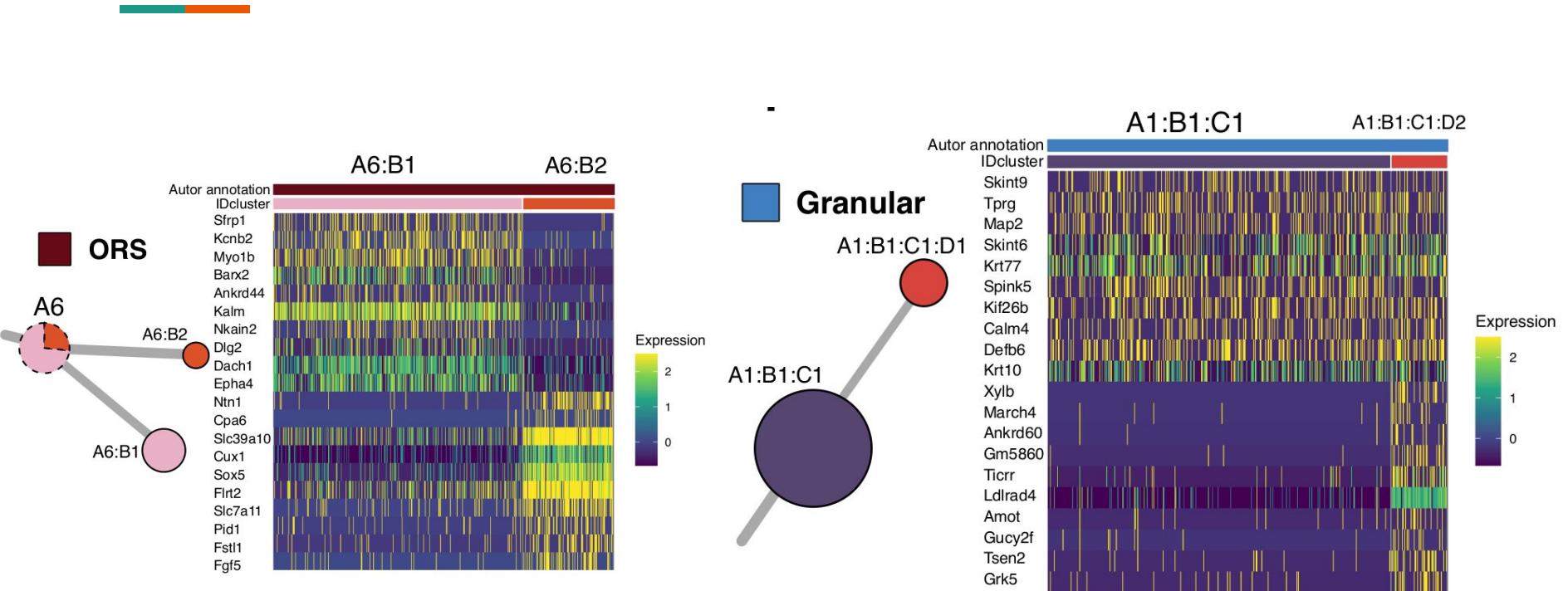


From abstract hyperparameters to biologically relevant parameters: logFC, q.value and % activation.

# SHARE-seq of mouse skin - clustering scRNA

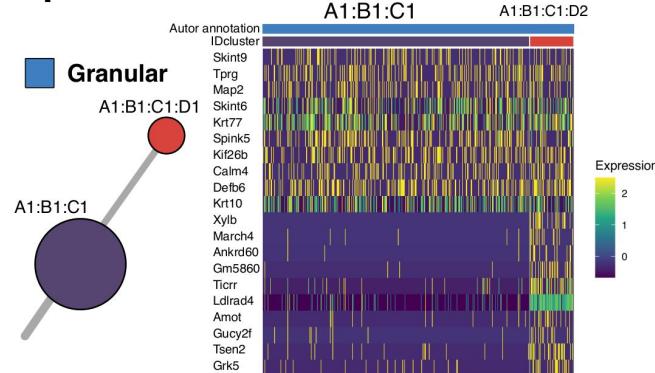
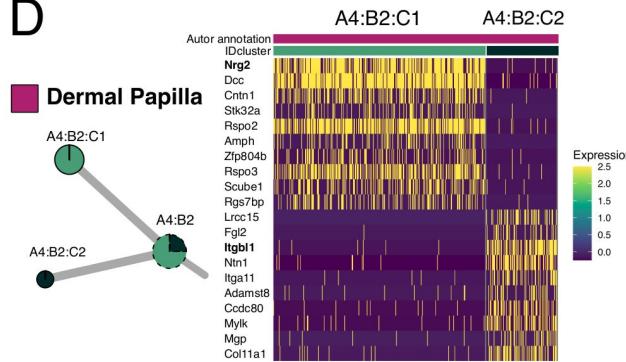


# Different hierarchies of clusters

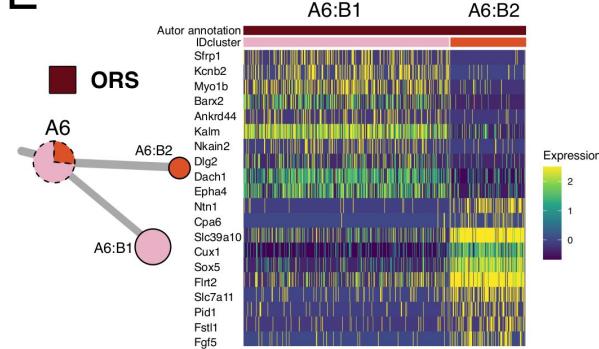


# Different hierarchies of clusters

D

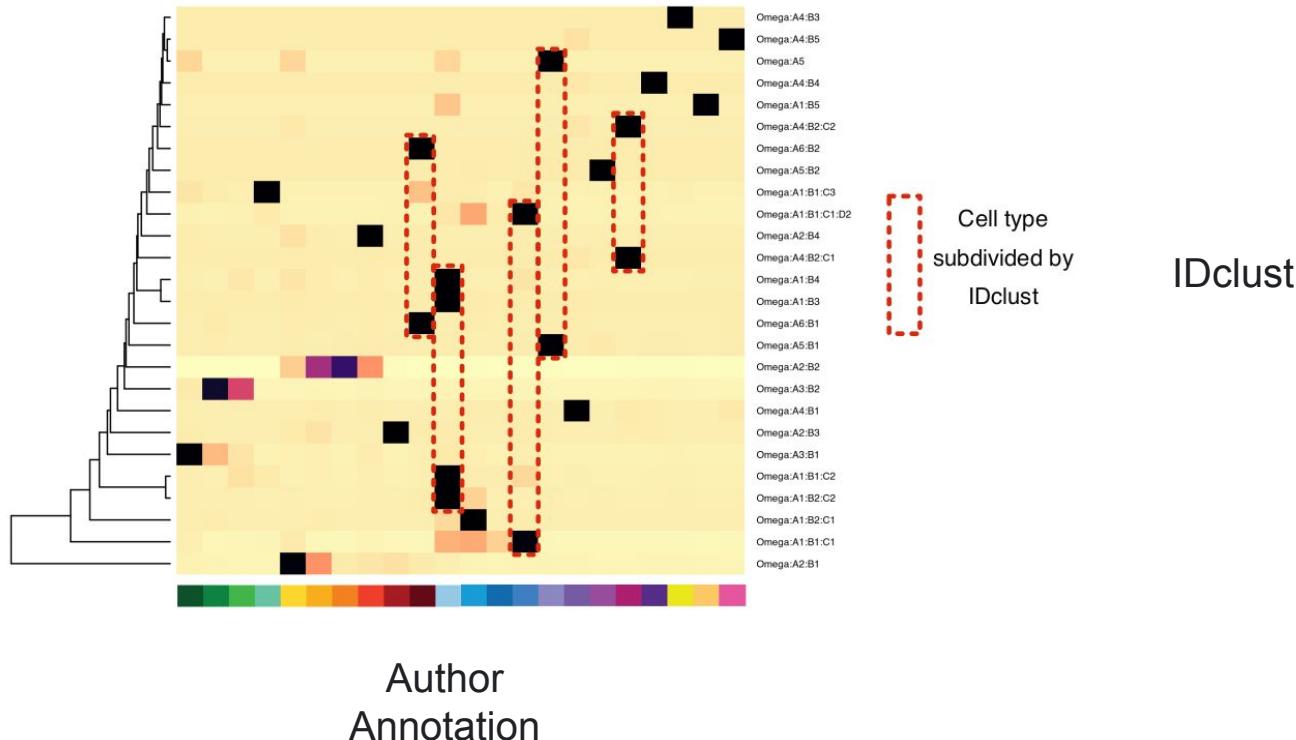


E

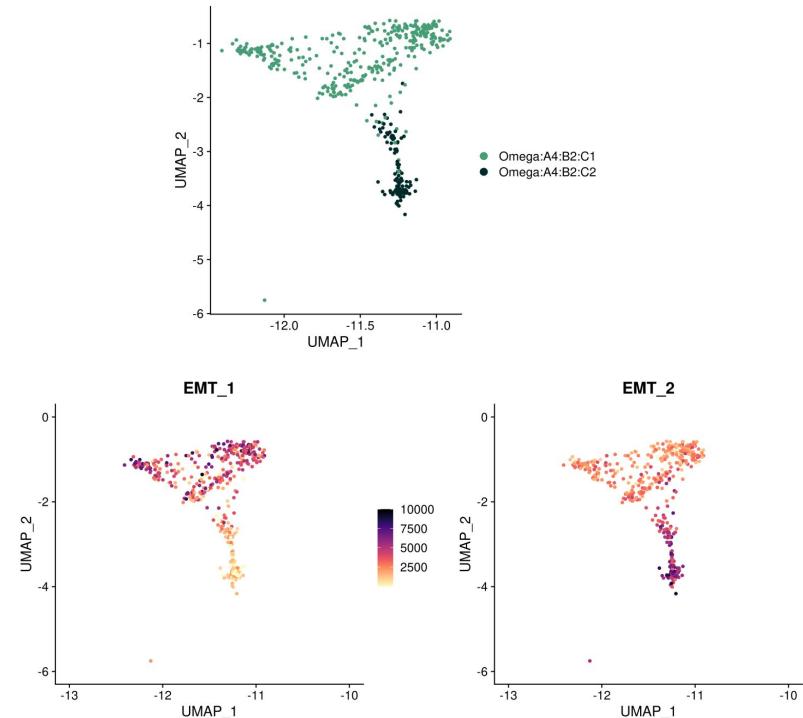
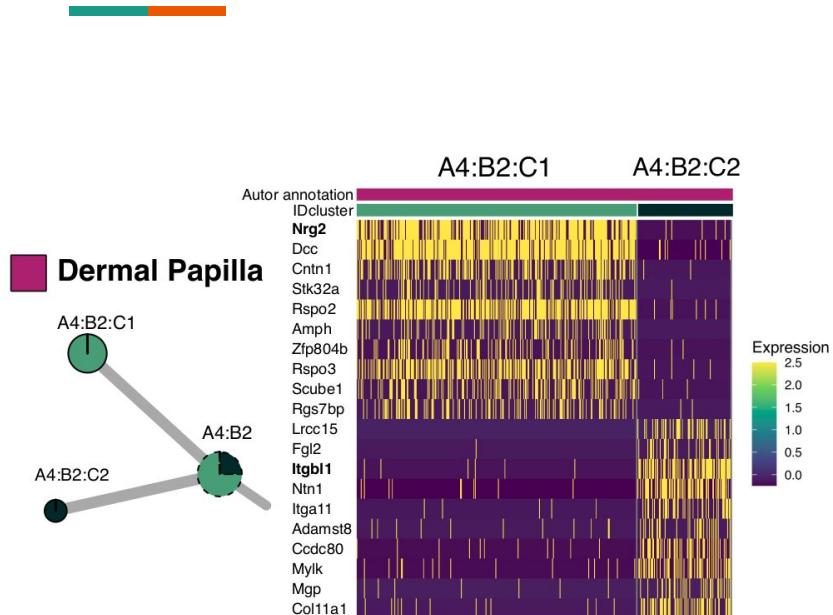


# Dividing some annotated clusters into multiple

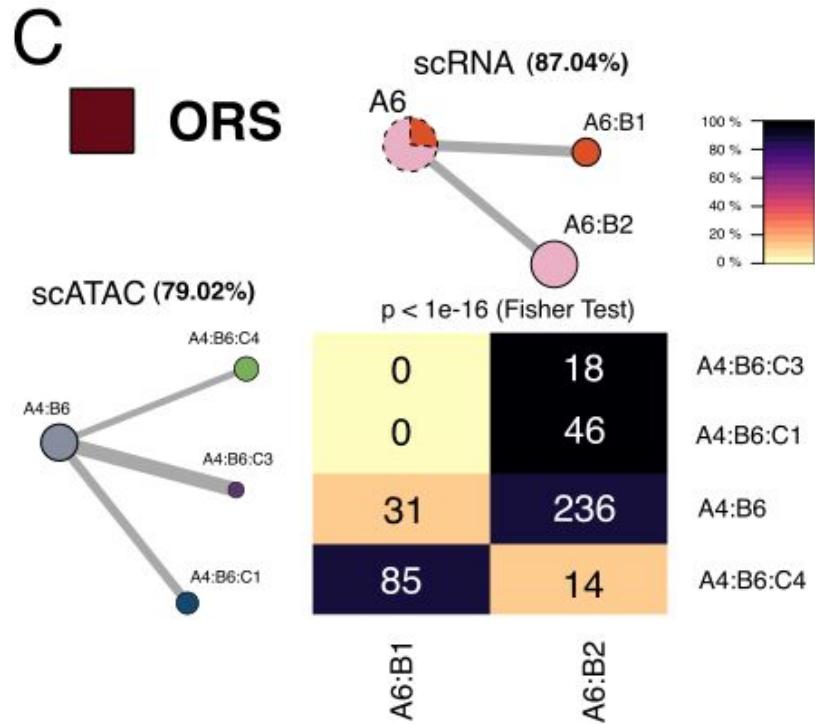
███████



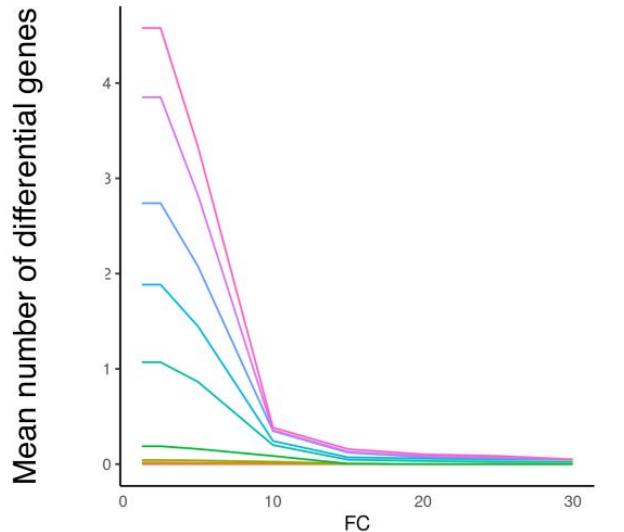
# Two different EMT submodules activated in Dermal Papilla



## IDclust can identify new consistent clusters in both modality (RNA & ATAC)

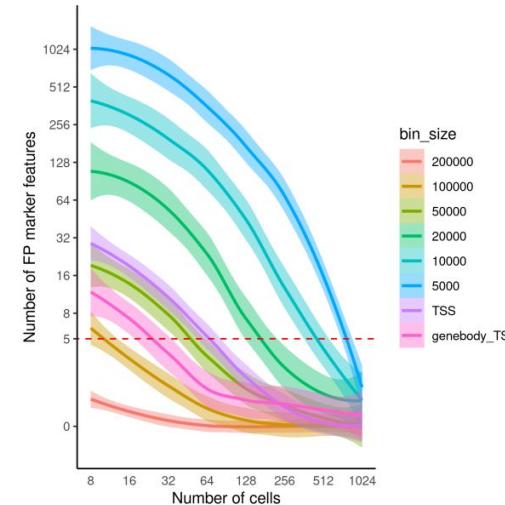


# Benchmarking the differential analysis parameters (epigenomics)



Adjusted p-value

- 1e-05
- 1e-04
- 0.001
- 0.01
- 0.05
- 0.1
- 0.15
- 0.2
- 0.25



- Running Differential analysis below 50 cells in the cluster produces high number of false positives

# Conclusions

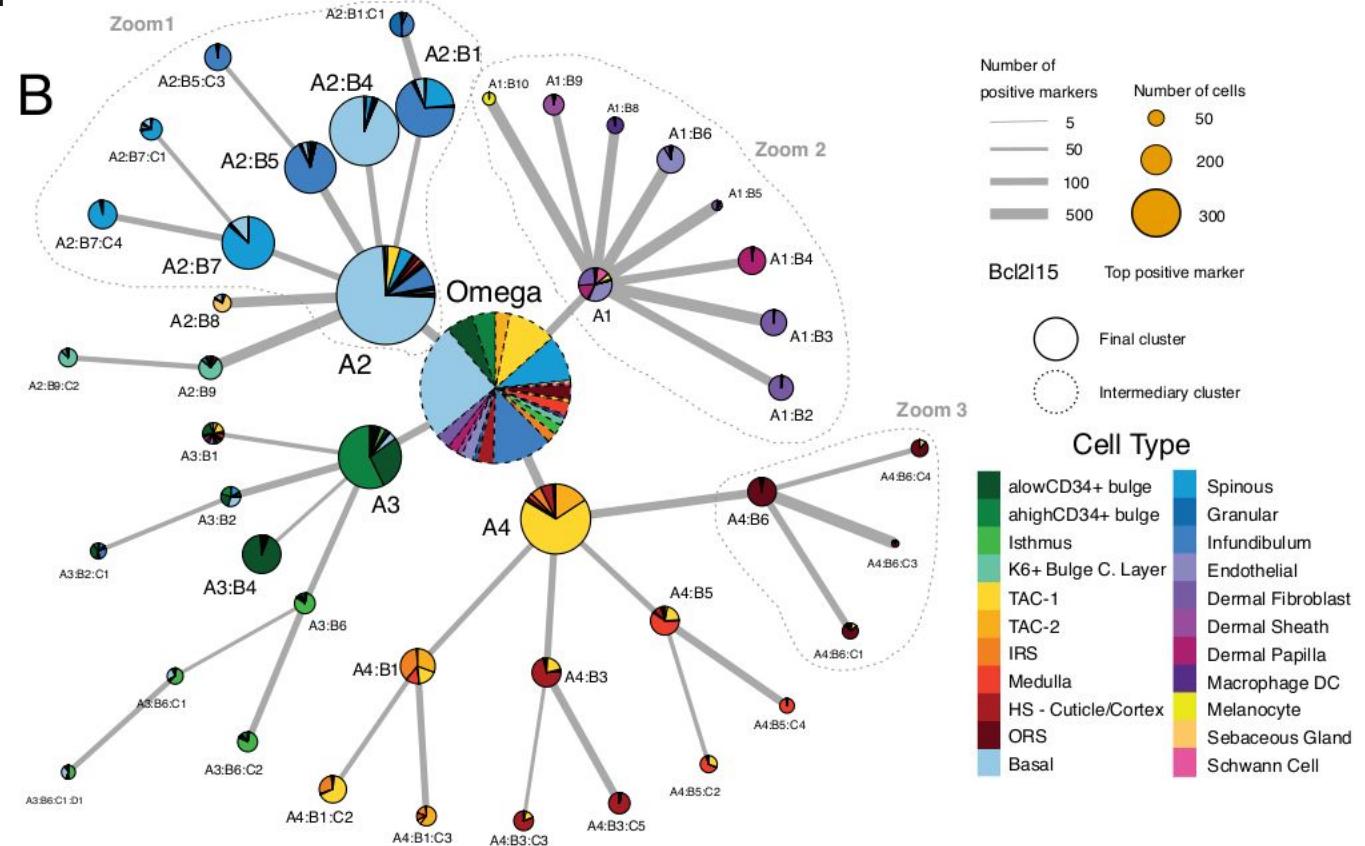
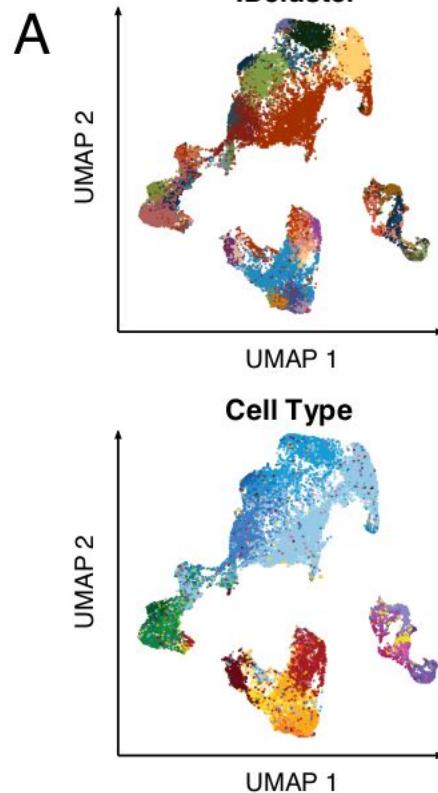
---

- IDclust is a method that finds cluster based on their differences
- IDclust can find different kind of hierarchies in the clusters

# Limitations

- The framework is ‘as good’ as the differential analysis (e.g. limits in small clusters)
  - The inner clustering parameters can still have an impact on the final results
  - Effects such as cell cycle can influence the latest partition of the data

# SHARE-seq - ATAC



# PairedTag - H3K4me1

