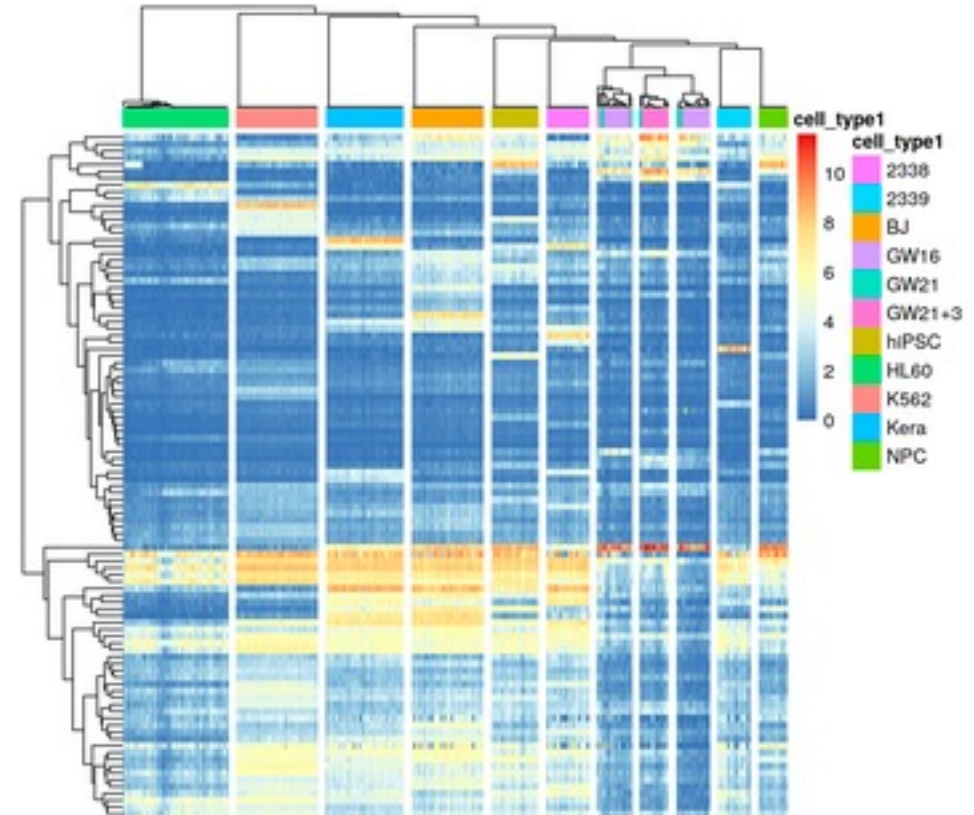
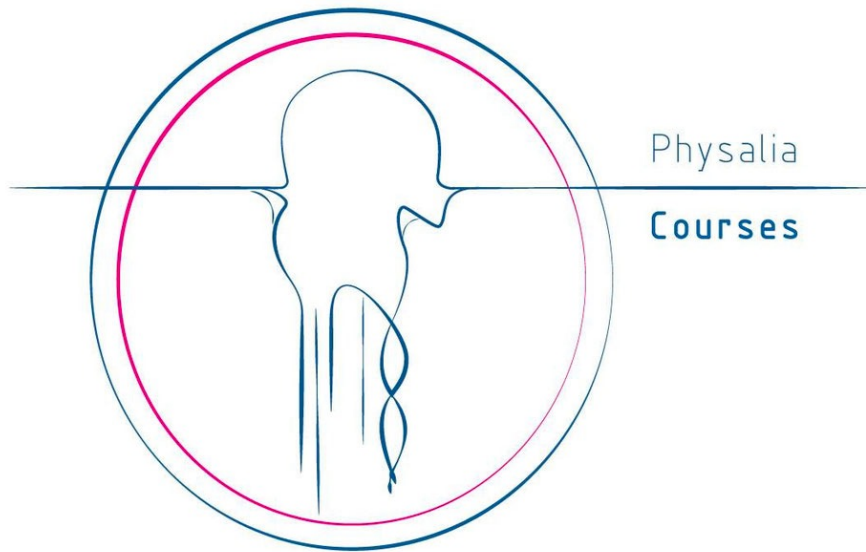


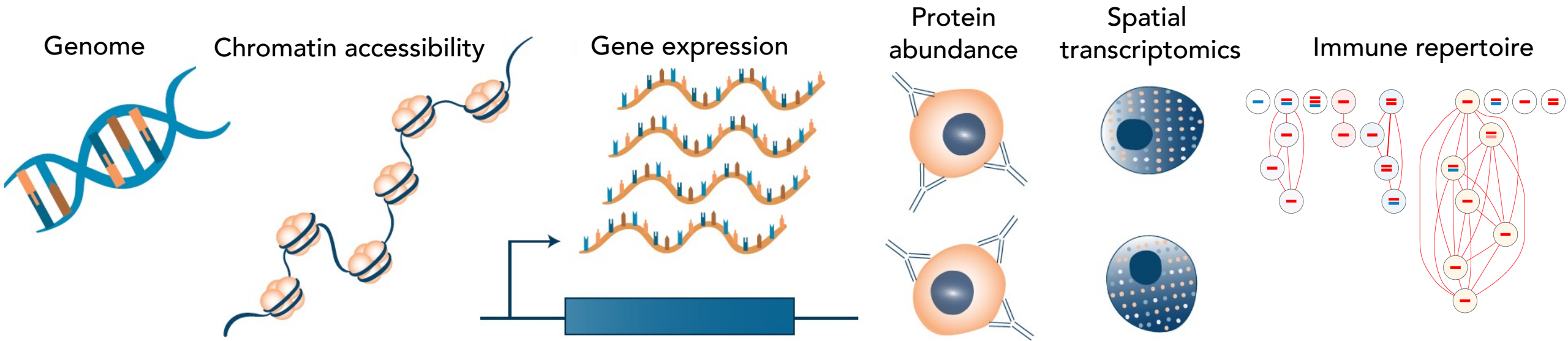
Analysis of single-cell ATAC-seq data

Orr Ashenberg, Jacques Serizay

June, 2023



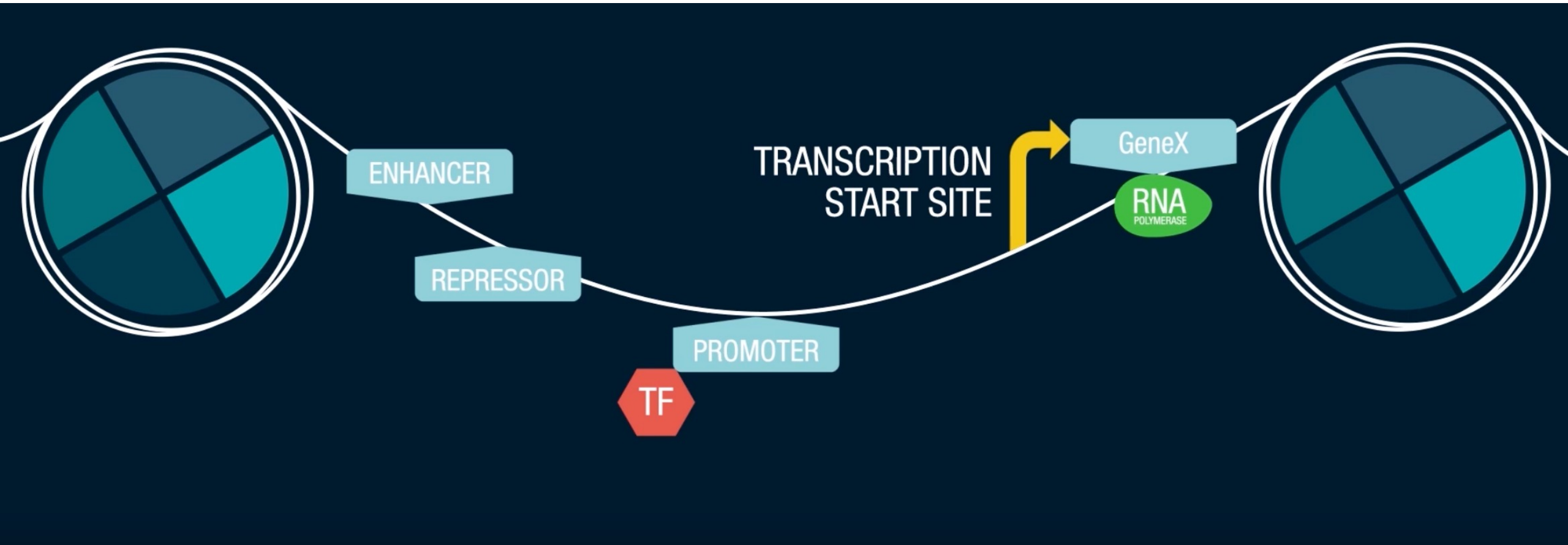
Multimodal measurements



Combining single-cell transcriptomic measurements with other data modalities can reveal gene function and gene regulation.

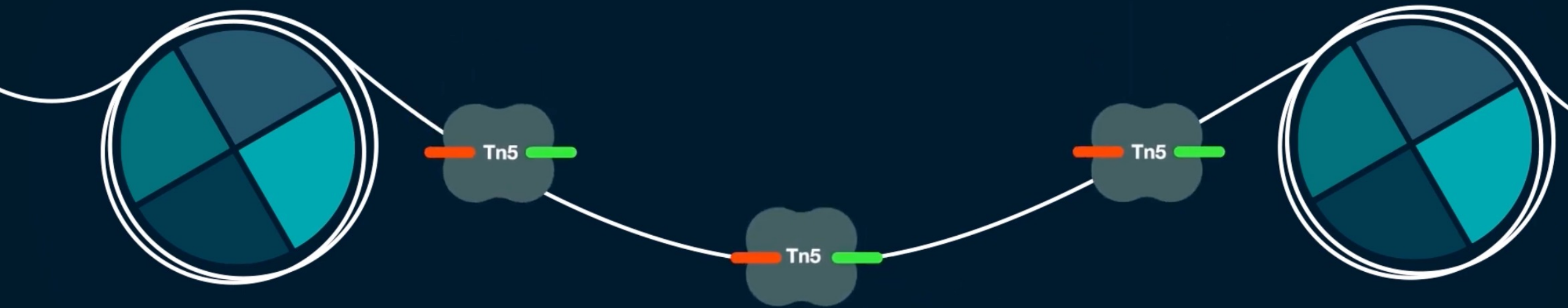
ATAC-Seq detects accessible chromatin regions

In the cell nucleus, the chromosomes contain tightly packed chromatin material. Part of the chromatin is open and accessible to many regulatory factors who control the expression and suppression of a variety of genes.

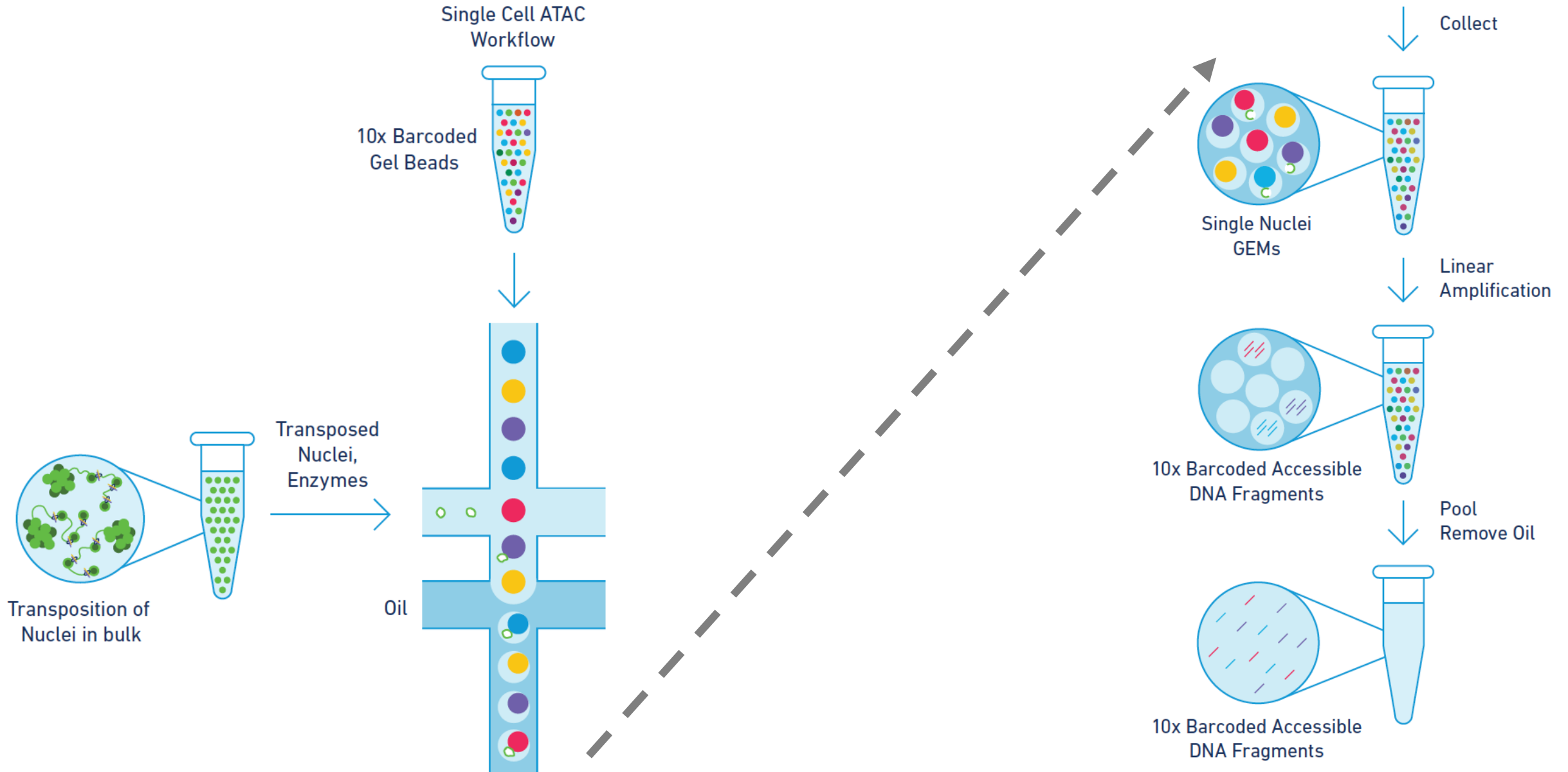


ATAC-Seq detects accessible chromatin regions

ATACseq (as well as scATACseq) measures how open this piece of DNA is. This openness is a proxy of how easily a transcription factor can bind these parts of the genome. ATACseq measures by using an enzyme called Tn5 transposase which binds open chromatin and inserts DNA sequencing adapters.



Chromium Single Cell ATAC-Seq (10x)



Single cell resolution reveals cell-type specific regulatory elements

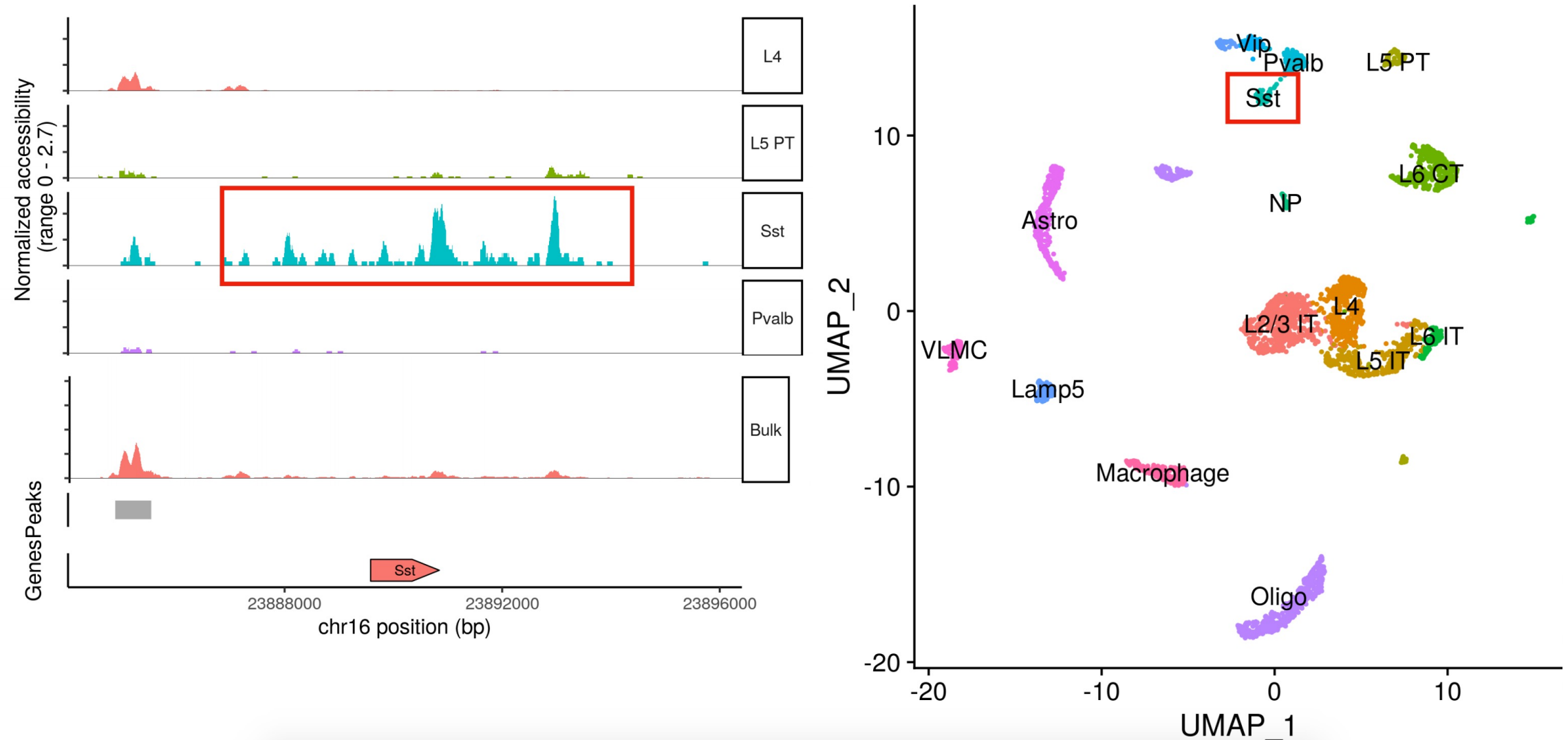


Figure adopted from "Analyzing single-cell ATAC-seq datasets" lecture by Tim Stuart

Pre-processing generates a fragment file and a peak/cell matrix

A full list of **all** unique fragments across all single cells, as opposed to only reads that map to peaks.

1. Indexed fragment file

chrom	start	stop	barcode	reads
chr1	3000141	3000517	GGTTGCGAGCCGCAAA-1	3
chr1	3000159	3000373	CTCAGCTAGTGTCCT-1	1
chr1	3000431	3000621	GAAGTCTGTAACACTC-1	1

2. Large sparse matrix

	AAACGAAAGAGTTTGA-1	AAACGAAAGCGAGCTA-1
chr1:565107-565550	.	.
chr1:569174-569639	.	.
chr1:713460-714823	.	2
chr1:752422-753038	.	.
chr1:762106-763359	.	4

Each value in the matrix represents the number of Tn5 cut sites for each single barcode (i.e. cell) that map within each peak

scATAC-Seq data is highly sparse

1. Indexed fragment file

chrom	start	stop	barcode	reads
chr1	3000141	3000517	GGTTGCGAGCCGCAAA-1	3
chr1	3000159	3000373	CTCAGCTAGTGTCACT-1	1
chr1	3000431	3000621	GAAGTCTGTAACACTC-1	1

2. Large sparse matrix

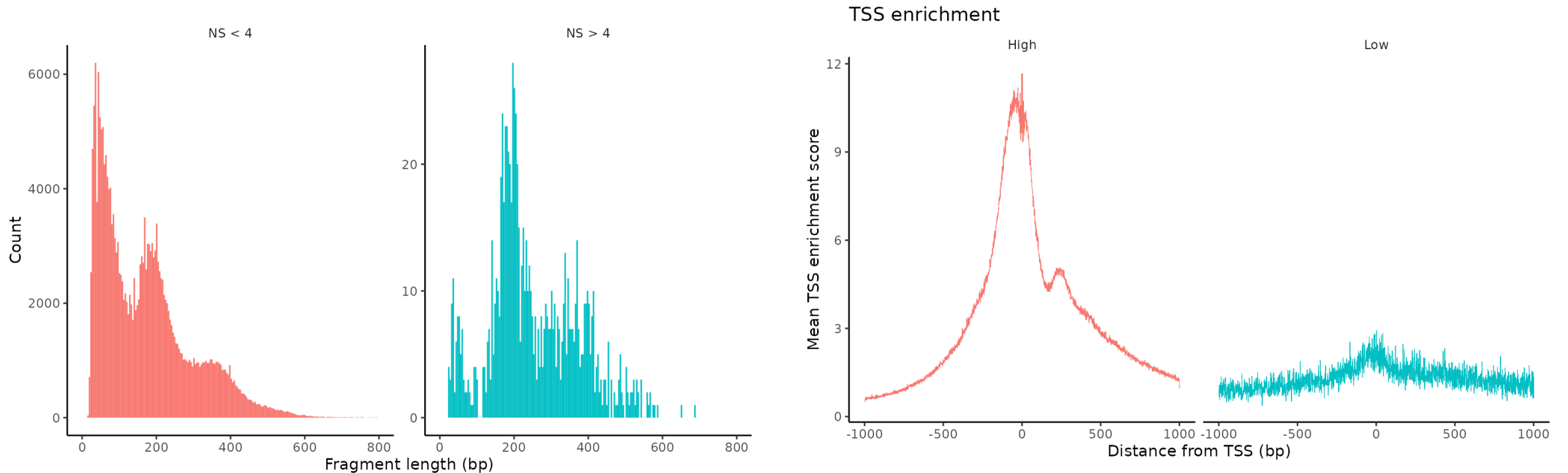
	AAACGAAAGAGTTTGA-1	AAACGAAAGCGAGCTA-1
chr1:565107-565550	.	.
chr1:569174-569639	.	.
chr1:713460-714823	.	2
chr1:752422-753038	.	.
chr1:762106-763359	.	4

Challenges in comparison to scRNA:

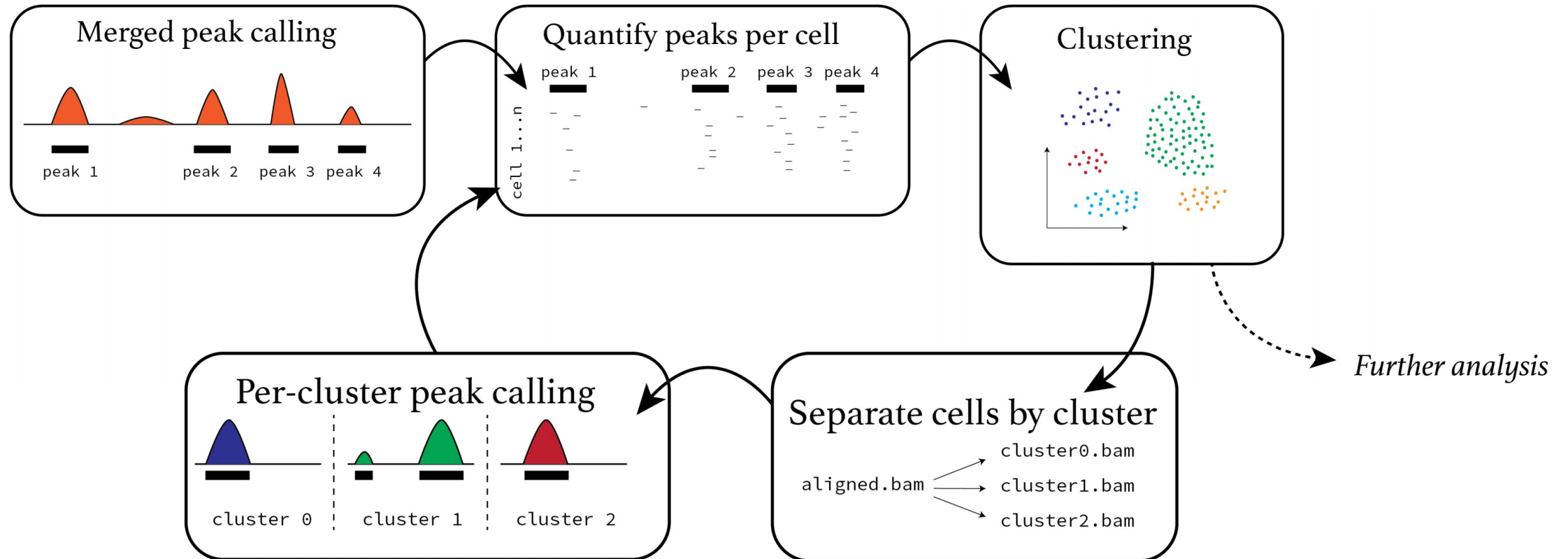
1. More sparse
2. Near-binary data
3. Non-fixed feature set
4. Order of magnitude more features

Quality control metrics for scATAC-seq data

1. Nucleosome banding pattern
2. Transcriptional start site (TSS) enrichment
3. Total number of fragments in peaks
4. Fraction of fragments in peaks



Overview of scATAC-Seq analysis



scATAC-Seq often uses latent semantic indexing (LSI) for dimensionality reduction

- Originally developed for topic modeling / natural language processing (Deerwester et al. 1990) and first applied to scATAC-Seq in 2015 (Cusanovich et al. *Science*)

- term frequency-inverse document frequency (TF-IDF) normalization followed by singular value decomposition (SVD)

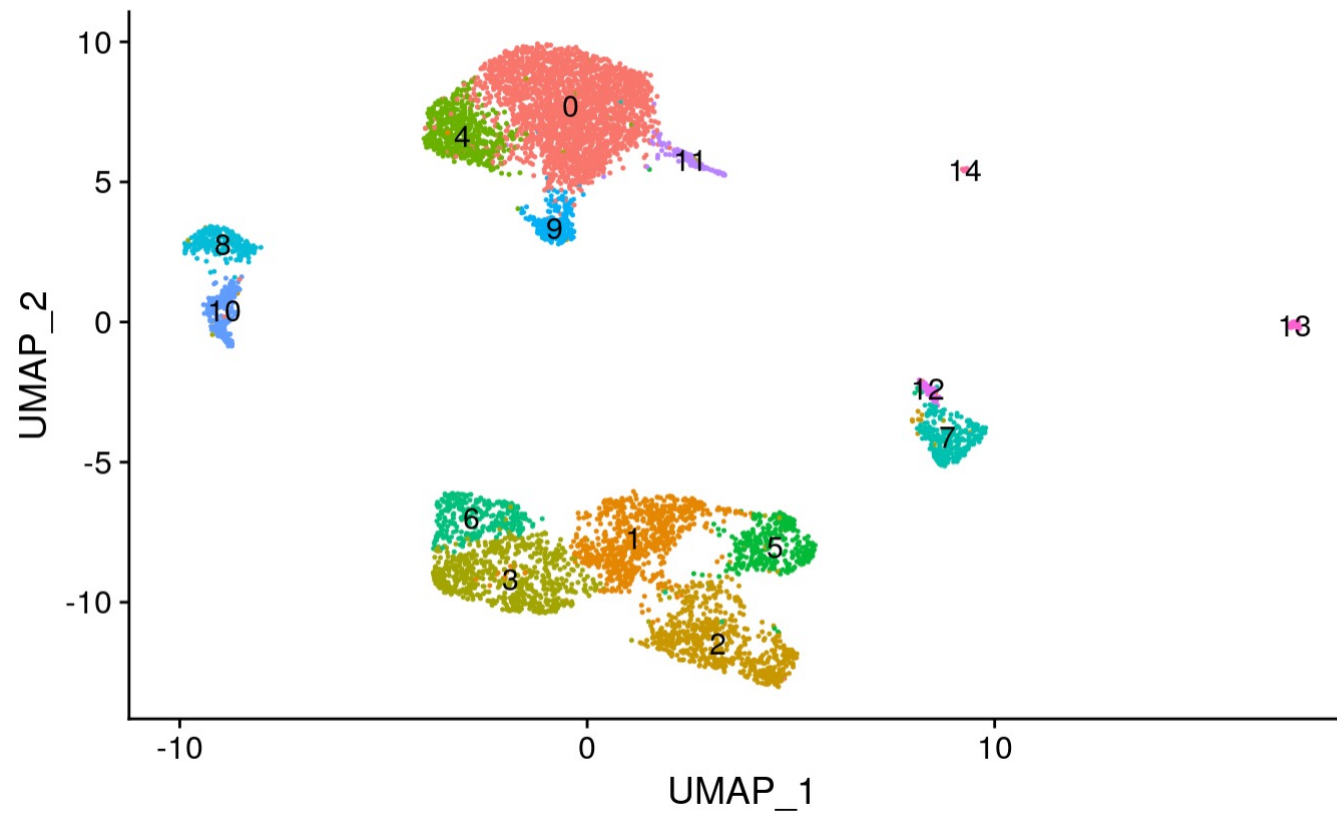
- Cell = document and peak = term

```
pbmc <- RunTFIDF(pbmc)
pbmc <- FindTopFeatures(pbmc, min.cutoff = 'q0')
pbmc <- RunSVD(pbmc)
```

- Term frequency: normalize across cells to correct for differences in sequencing depth
- Inverse document frequency: give higher values to more rare peaks

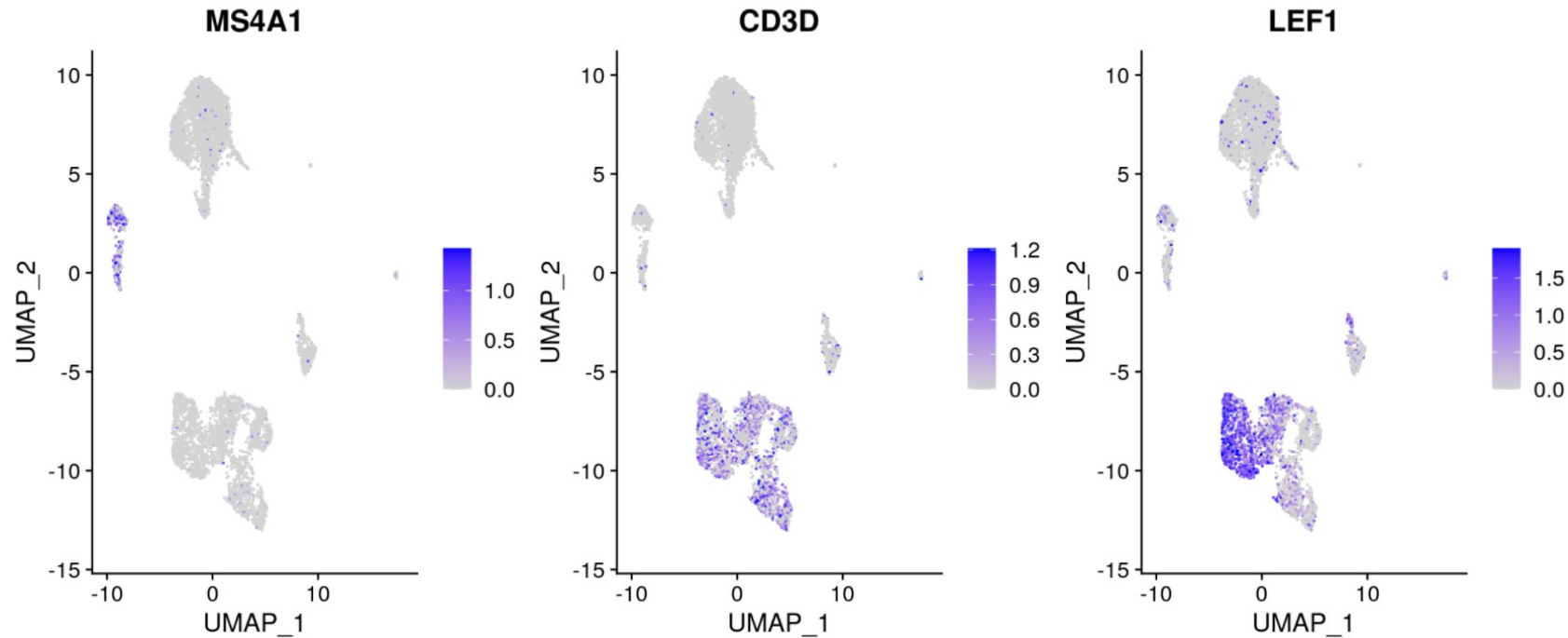
Visualize clusters in scATAC-Seq

```
pbmc <- RunUMAP(object = pbmc, reduction = 'lsi', dims = 2:30)
pbmc <- FindNeighbors(object = pbmc, reduction = 'lsi', dims = 2:30)
pbmc <- FindClusters(object = pbmc, verbose = FALSE, algorithm = 3)
DimPlot(object = pbmc, label = TRUE) + NoLegend()
```

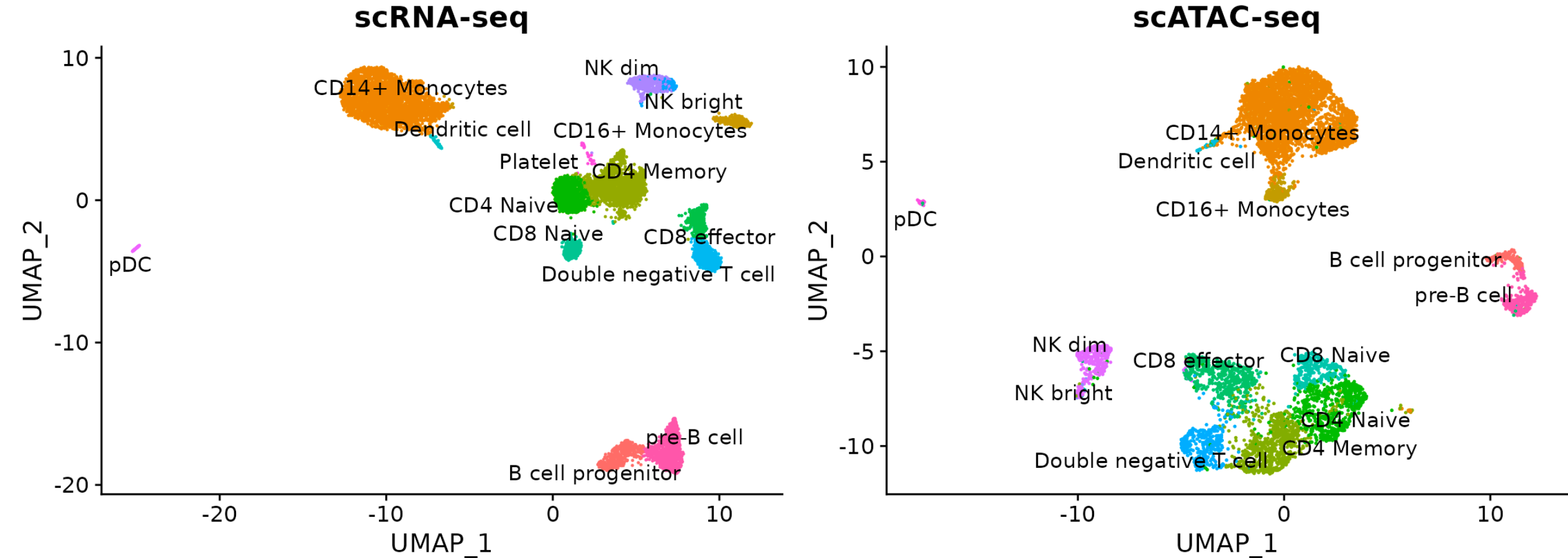


Fragment file helps infer gene “activity” and annotate clusters

Quantify the activity of each gene in the genome by assessing the chromatin accessibility associated with each gene: count the number of fragments for each cell that map to the promoter + gene body



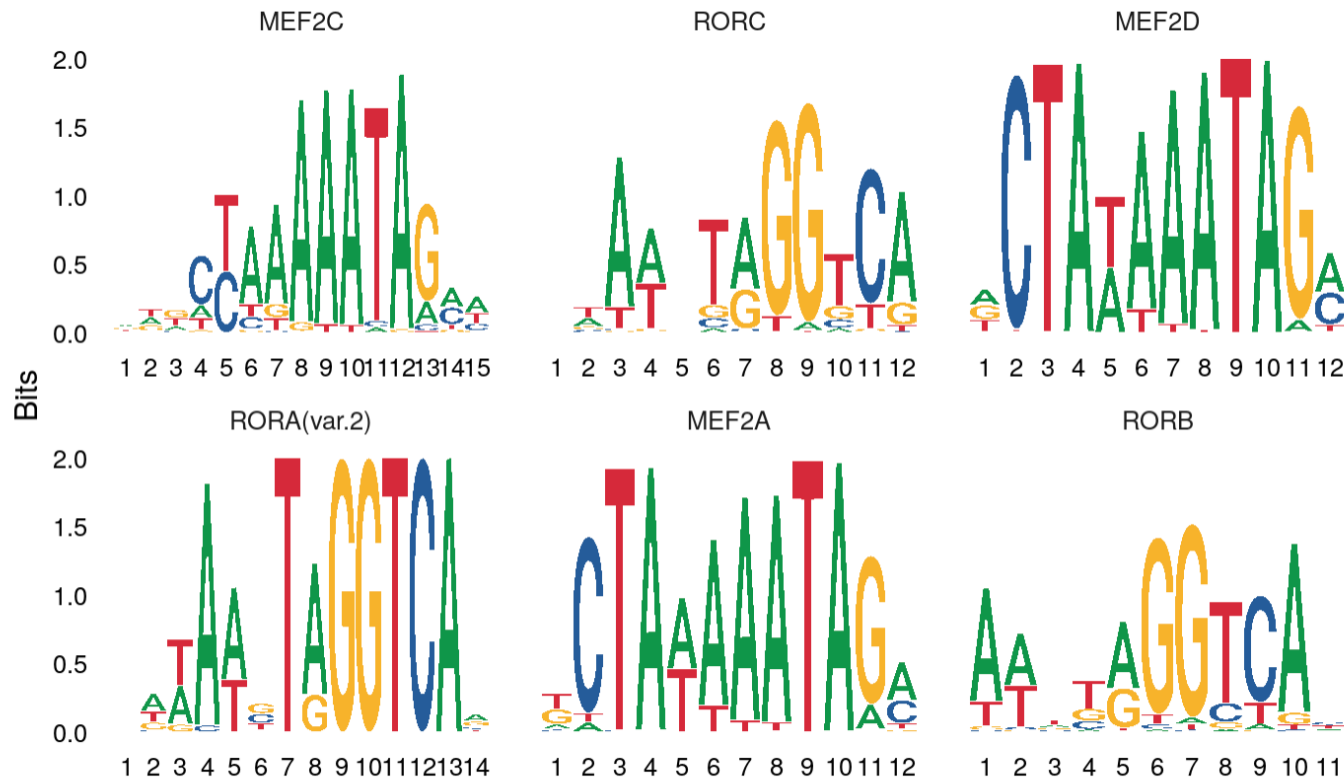
Integrating with scRNA-seq data using CCA + MNN (Seurat v4)



Critical assumption: there is generally a positive correlation between chromatin accessibility and gene expression!!!!

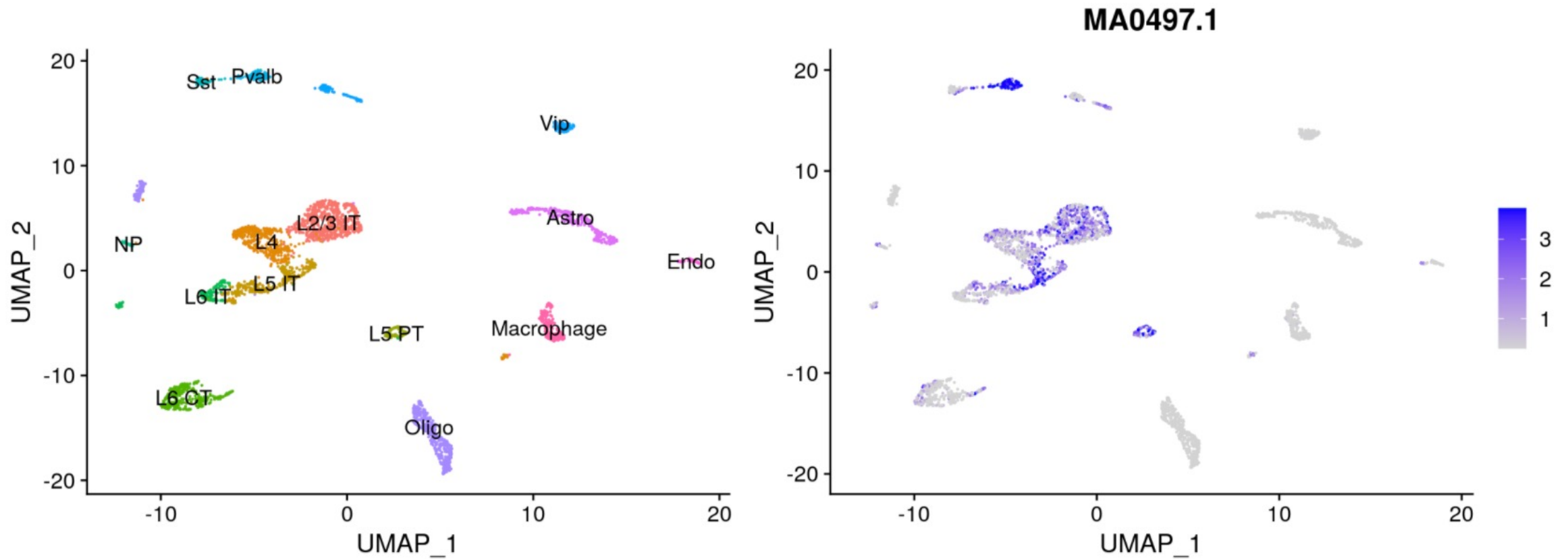
Finding overrepresented motifs

To identify potentially important cell-type-specific regulatory sequences, signac searches for DNA motifs that are overrepresented in a set of peaks that are differentially accessible between cell types.



Computing motif activities

ChromVAR identifies motifs associated with variability in chromatin accessibility between cells.



Software for scATAC-Seq analysis

<i>Signac</i>	https://satijalab.org/signac/index.html
<i>SnapATAC</i>	https://github.com/r3fang/SnapATAC
<i>ArchR</i>	https://www.archrproject.com/
<i>cisTopic</i>	https://github.com/aertslab/cisTopic
<i>chromVAR</i>	https://github.com/GreenleafLab/chromVAR
<i>CICERO</i>	https://cole-trapnell-lab.github.io/cicero-release/

Software for scATAC-Seq analysis

	ArchR [*]	Signac	SnapATAC	
Pre-processing	NR	NA	✓	Data Import
Data import / base file type creation	✓	NA	✓	
QC filter cells	✓	✓	✓	
Matrix creation	✓ (Tile)	✓ (Peak)	✓ (Tile)	
Doublet removal	✓	NP	NP	Doublet Removal
Data imputation with MAGIC	✓	NP	✓	Gene Scores
Genome-wide gene score matrix	✓	✓	✓	
Dimensionality reduction and clustering	✓	✓	✓	Clustering
UMAP and tSNE plotting	✓	✓	✓	
Cluster peak calling	✓	NP	✓	Standard ATAC-seq Analyses
Cluster-based peak matrix creation	✓	NP	✓	
Motif enrichment	✓	✓	✓	
chromVAR motif deviations	✓	✓	✓	
Footprinting	✓	NP	NP	
Feature set annotation	✓	NP	NP	Data Visualization
Track plotting	✓	✓	NP	
Co-accessibility	✓	NP	NP	
Interactive genome browser	✓	NP	NP	
Cellular trajectory analysis	✓	NP	NP	Advanced ATAC-seq Analyses
Project bulk data into scATAC embedding	✓	NP	NP	
Integration of RNA-seq and ATAC-seq	✓	✓	✓	Integration of RNA-seq and ATAC-seq
Genome-wide peak-to-gene links	✓	NP	NP	

NR = Not Required NA = Not Applicable NP = Not Possible

Some approaches to multiome data (scRNA-seq and scATAC-seq)

ArchR: https://greenleaflab.github.io/ArchR_2020/Ex-Analyze-Multiome.html

Signac (same people who built Seurat): https://satijalab.org/signac/articles/pbmc_multiomic.html

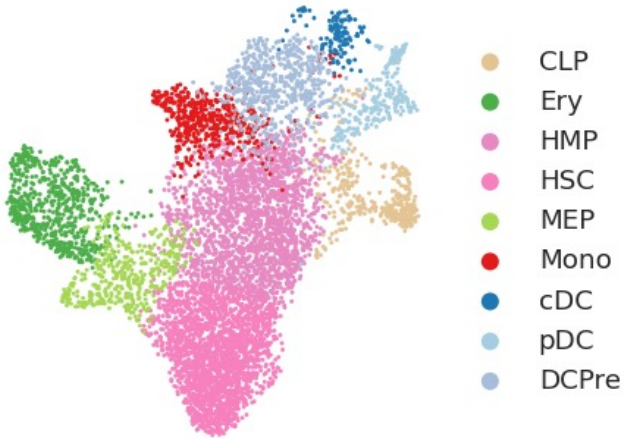
Methods to overcome sparsity in ATAC-seq data

Computing meta-cells is one methodology used to overcome sparsity in scATAC-seq data

Computing meta-cells (e.g. SEACells algorithm) can improve computation of peak-gene associations

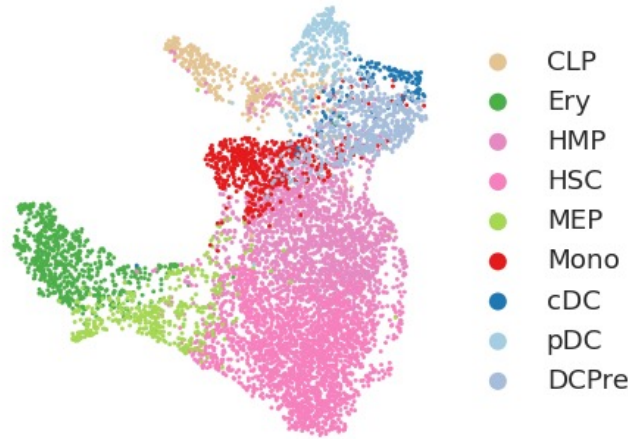
scRNA-seq

celltype



scATAC-seq

celltype



Metacell Assignments



Methods to overcome sparsity in ATAC-seq data

Generate summarized ATAC and RNA metacells

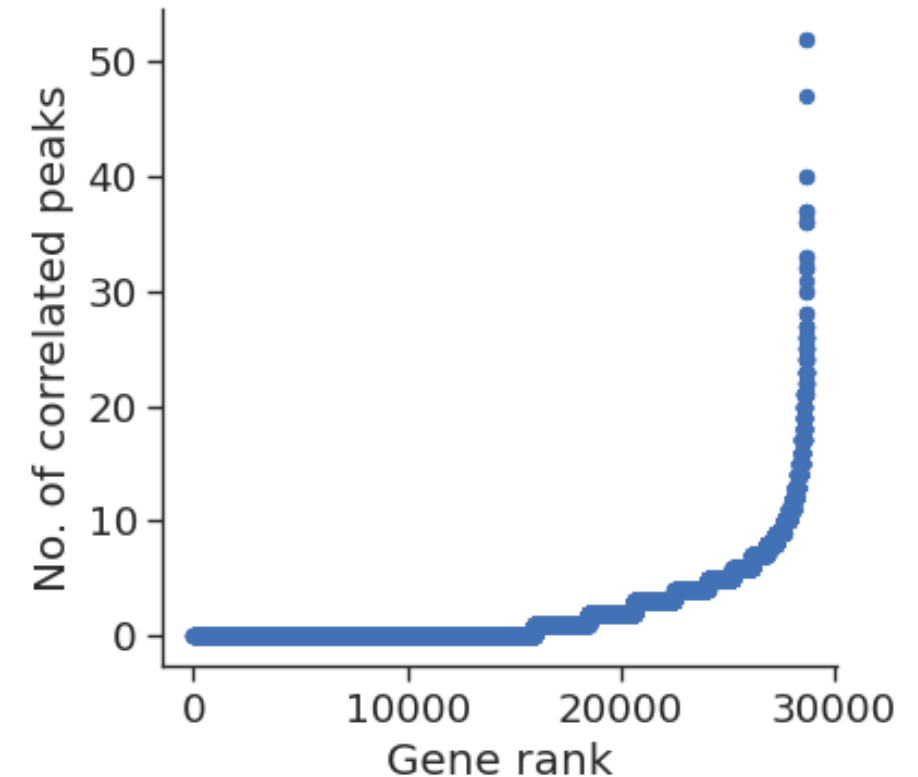
Used the paired multiome metacells to compute the correlation of gene expression and accessibility of peaks within the vicinity of the gene

Pick out highly regulated genes that are correlated with multiple peaks

Calculate gene scores as the weighted sum of the accessibility of correlated peaks

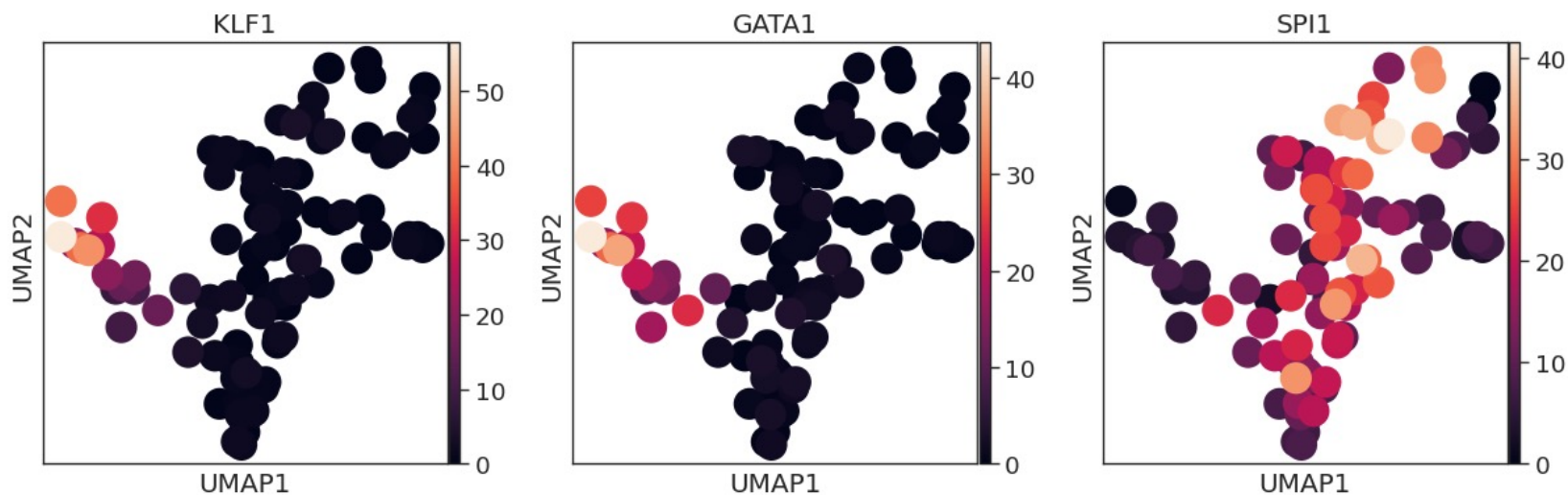
To calculate gene accessibility metrics identify the subset of peaks that are open in each metacell

Open peaks are used to compute a gene accessibility metric which represents the fraction of correlated open peaks

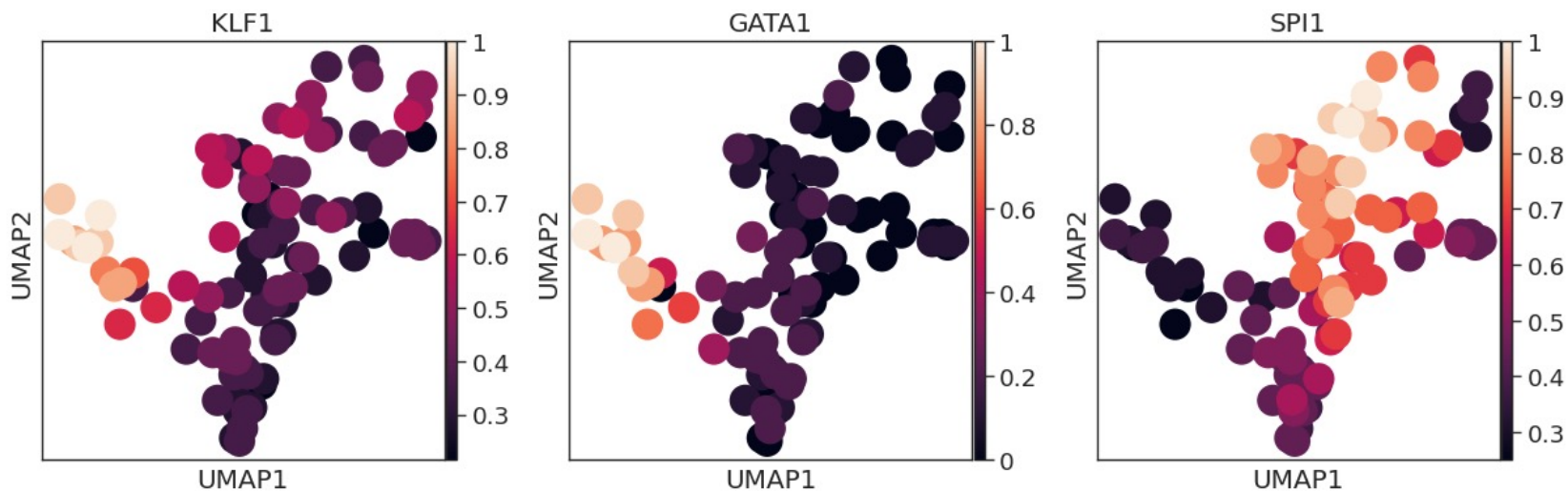


Methods to overcome sparsity in ATAC-seq data

Gene expression

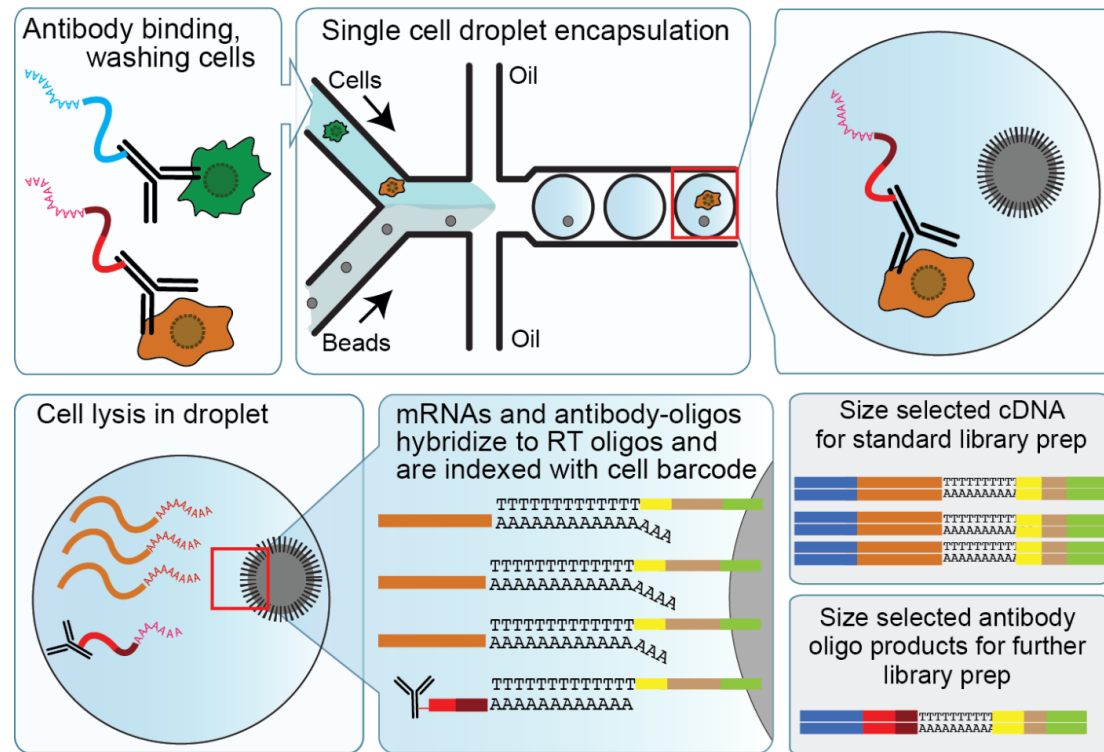


Accessibility



APPENDIX: OTHER MULTIMODAL METHODOLOGIES

CITE-Seq allows simultaneous measurement of RNA expression and surface protein expression in same cell

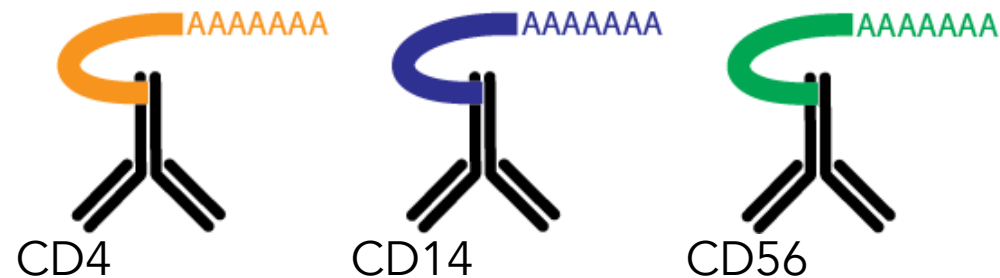


Count matrix1: scRNAseq

Count matrix2: CITEseq

Antibodies are conjugated with DNA-oligos

- DNA-oligo indicates antibody identity and contains poly-A tail



Pre-processing generates RNA count and protein count matrices

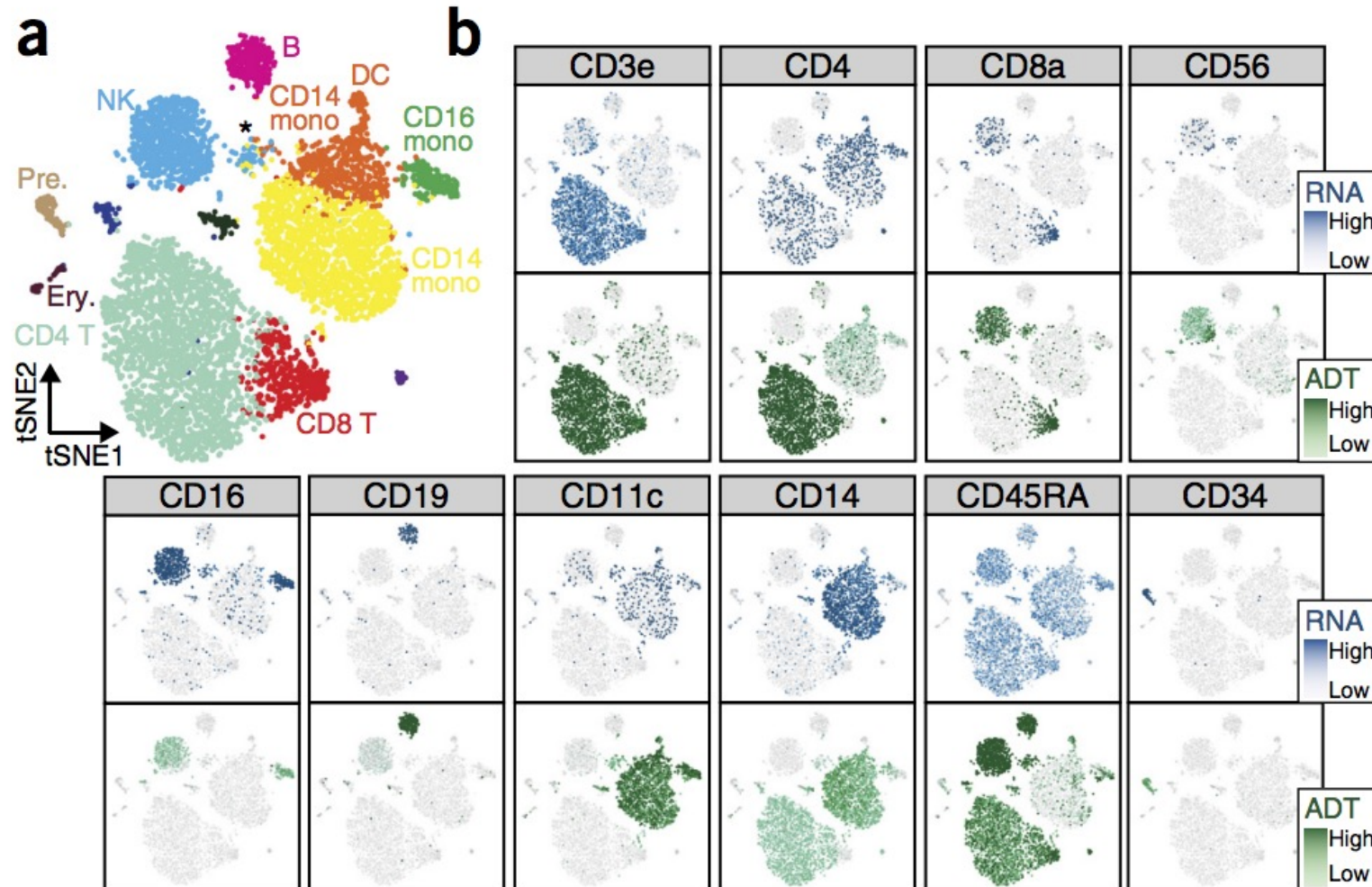
scRNA-Seq

	cell1	cell2	cell3	cell4
Gene1	0	0	0	4
Gene2	0	2	0	0
Gene3	0	0	0	0

CITE-Seq

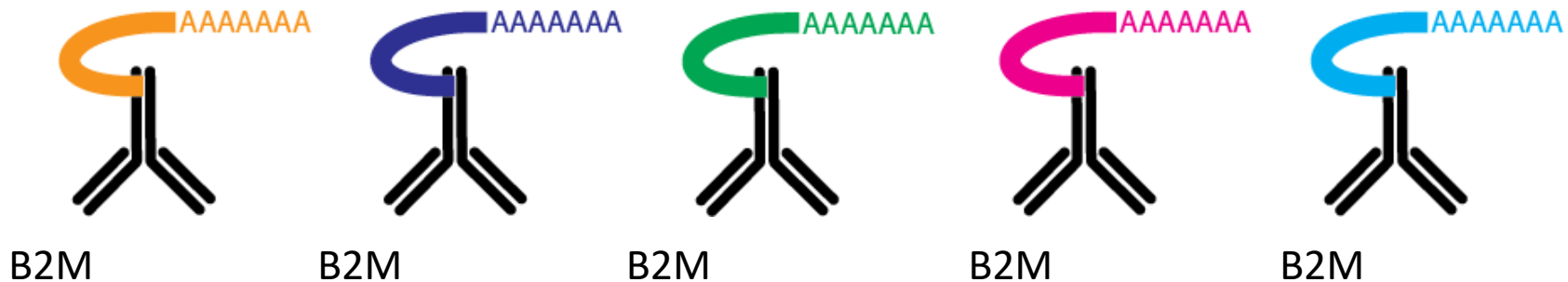
	cell1	cell2	cell3	cell4
Antibody for CD4	100	0	0	38
Antibody for CD14	0	0	45	0
Antibody for CD56	2	0	0	0

Simultaneous protein and RNA measurement in single cells



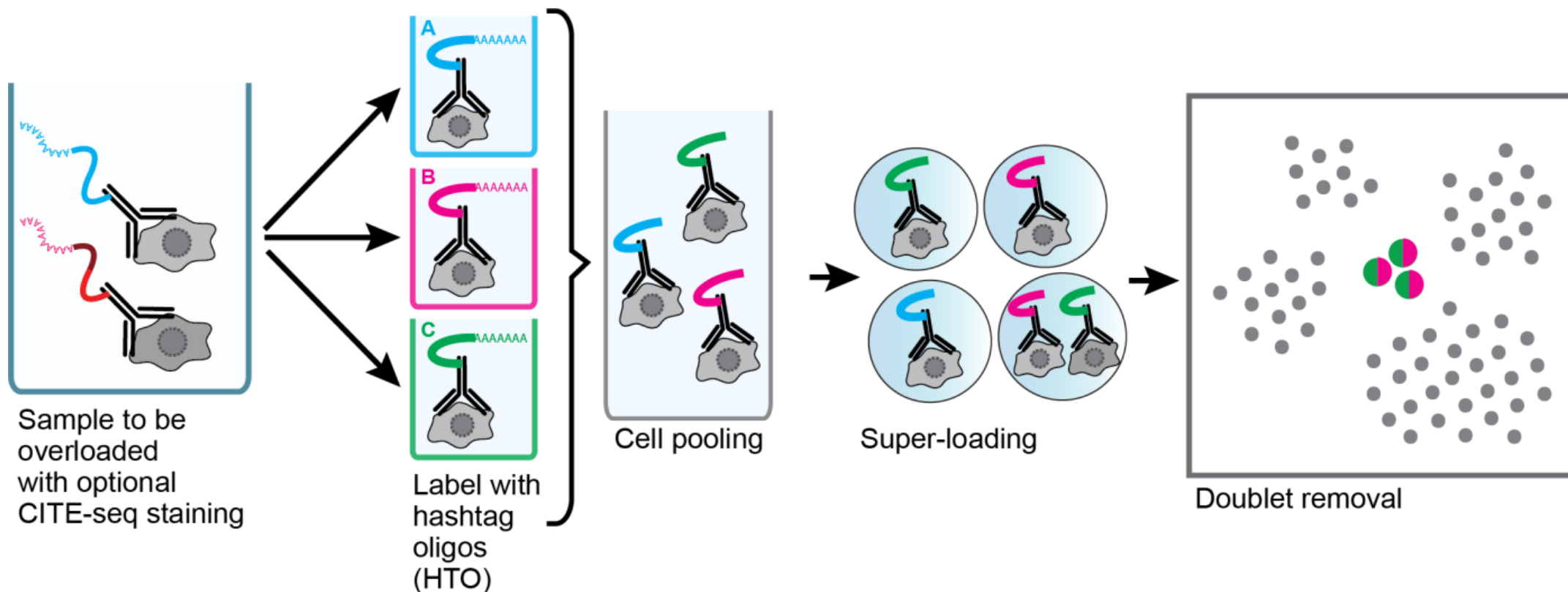
Cell hashing using CITE-Seq

- Cell Hashing enables sample multiplexing and super-loading on single cell RNA-sequencing platforms.
- Cell Hashing uses a series of oligo-tagged antibodies against ubiquitously expressed surface proteins with different barcodes to uniquely label cells from distinct samples, which can be subsequently pooled in one scRNA-seq run.



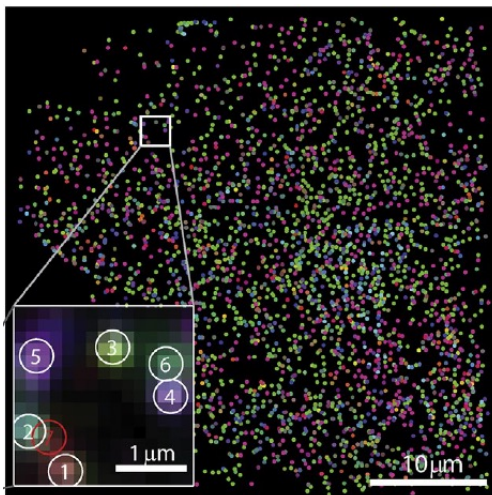
With cell hashing, we can “overload” the 10x channel

By sequencing these tags alongside the cellular transcriptome, we can assign each cell to its sample of origin, and robustly identify doublets originating from multiple samples.



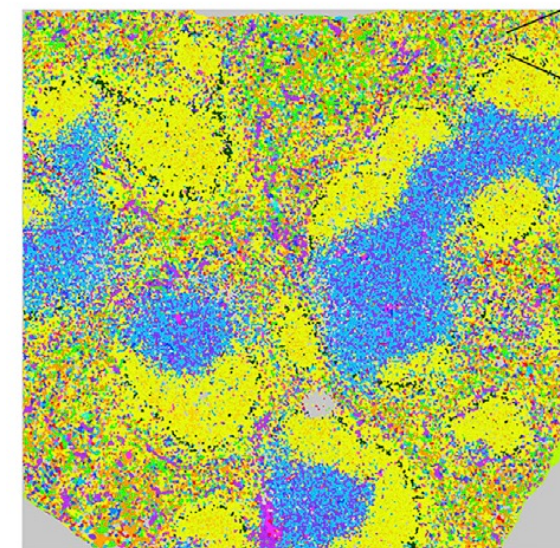
Ongoing developments in single-cell genomics: Growing toolbox for spatial genomics

MERFISH (RNA)



breast cancer

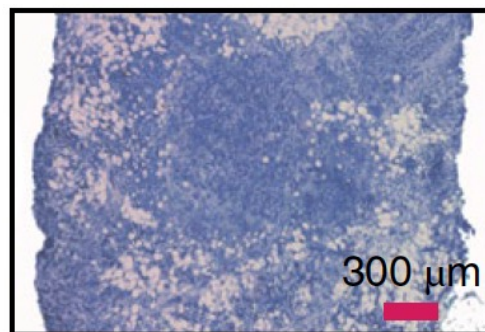
CODEx (protein)



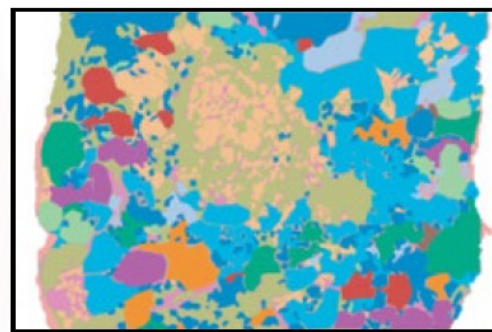
mouse spleen

High Density Spatial Transcriptomics (RNA)

H&E



Annotations



- Fatty tissue, immune/lymphoid ● Fatty tissue, invasive cancer
 ● Fibrous tissue, invasive cancer ● Fibrous tissue, immune/lymphoid
 ● Invasive cancer, immune/lymphoid ● Immune/lymphoid
 ● Fatty tissue, fibrous tissue, invasive cancer ● Fibrous tissue
 ● Fibrous tissue, invasive cancer, immune/lymphoid ● Fatty tissue
 ● Fatty tissue, fibrous tissue, invasive cancer, immune/lymphoid
 ● Fatty tissue, invasive cancer, immune/lymphoid ● Invasive cancer

- | | |
|---------|--|
| capsule | CD31(thi) vascular |
| | plasma cells |
| | CD106(-)/CD16/32(-)/Ly6C(+)/CD31(+)/stroma |
| | CD4(+)/CD8(-)/cDC |
| | ERTR7(+)/stroma |
| | CD4(+)/MHCII(+) |
| | CD4(-)/CD8(-)/cDC |
| | CD4(-)/CD8(+)/cDC |
| | CD106(+)/CD16/32(-)/Ly6C(+)/CD31(+) |
| | megakaryocytes |
| | CD106(-)/D16/32(+)/Ly6C(+)/CD31(-) |
| | CD106(+)/CD16/32(+)/Ly6C(-)/CD31(+)/stroma |
| | granulocytes |
| | CD3(+)/other markers (-) |
| | NK cells |
| | CD3(+) |
| | F4/80(+)/mphs |
| | erythroblasts |
| | CD11c(+)/B cells |
| | CD106(+)/CD16/32(+)/Ly6C(-)/CD31(-)/stroma |
| | FDCs |
| | marginal zone mphs |
| | B cells |
| | CD8(+)/T cells |
| | CD4(+)/T cells |
| | no id |
| | B220(+)/DN T cells |

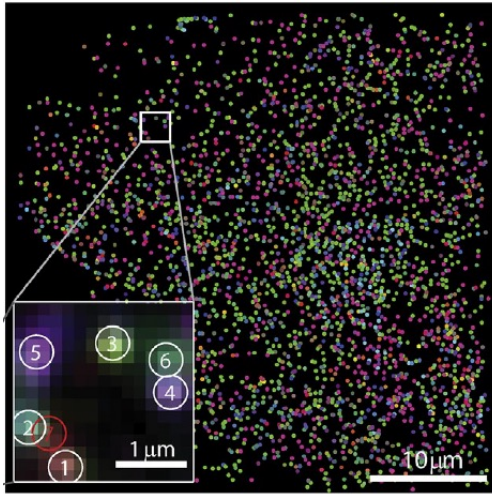
Moffitt J.R. et al. (2016) *PNAS* 113: 11046-11051.

Goltsev Y. et al. (2018) *Cell*. 174: 968-981.

Vickovic S. et al. (2019) *Nature Methods*. 16: 987–990.

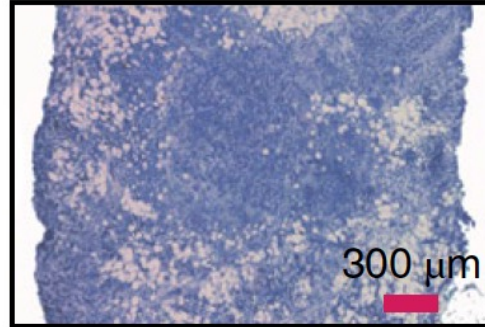
Ongoing developments in single-cell genomics: Growing toolbox for spatial genomics

MERFISH (RNA)



High Density Spatial Transcriptomics (RNA)

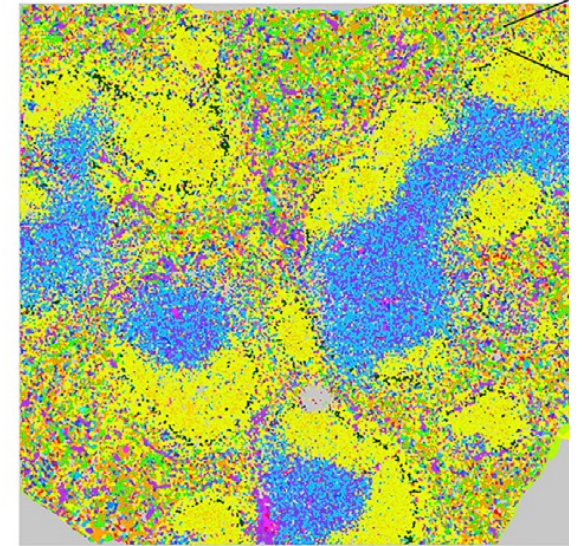
H&E



Annotations



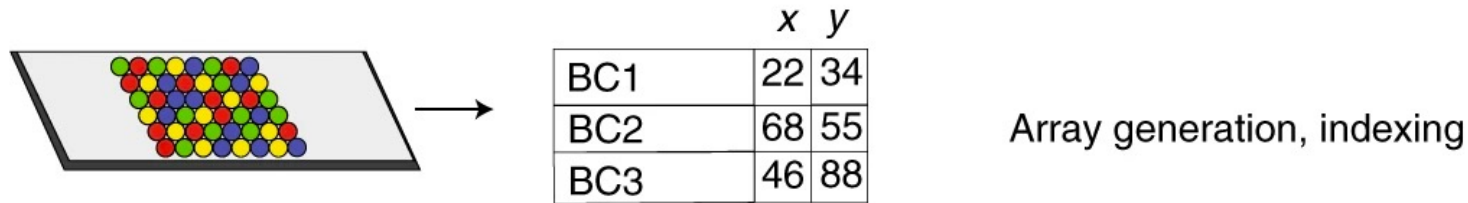
CODEX (protein)



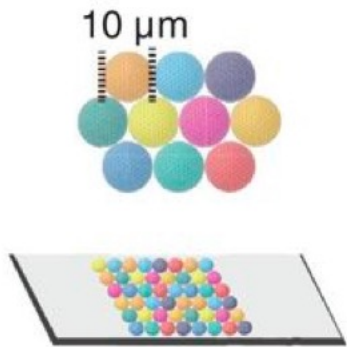
Major method differences

- detecting RNA or protein
- choosing target genes/proteins or full transcriptome
- number of cells profiled
- single-cell or almost single-cell resolution
- *in situ* vs traditional sequencing
- ...

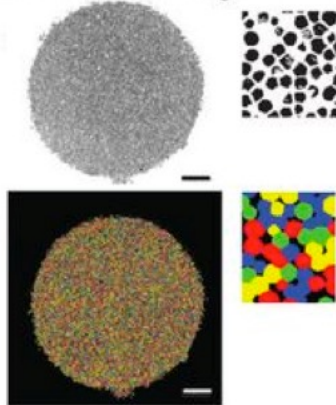
Slide-Seq overview (full transcriptome)



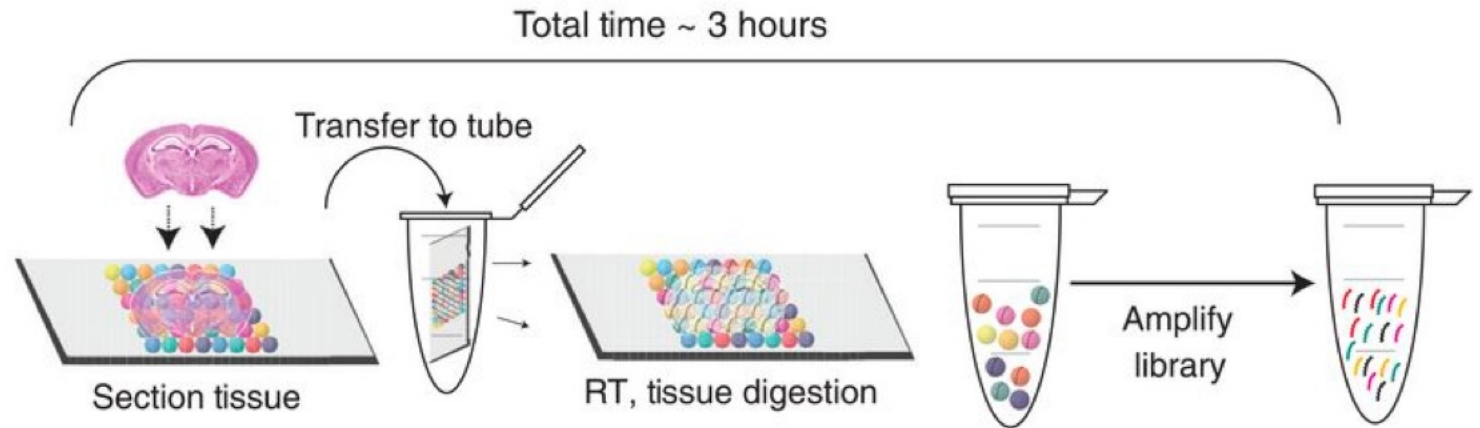
A Bead deposition



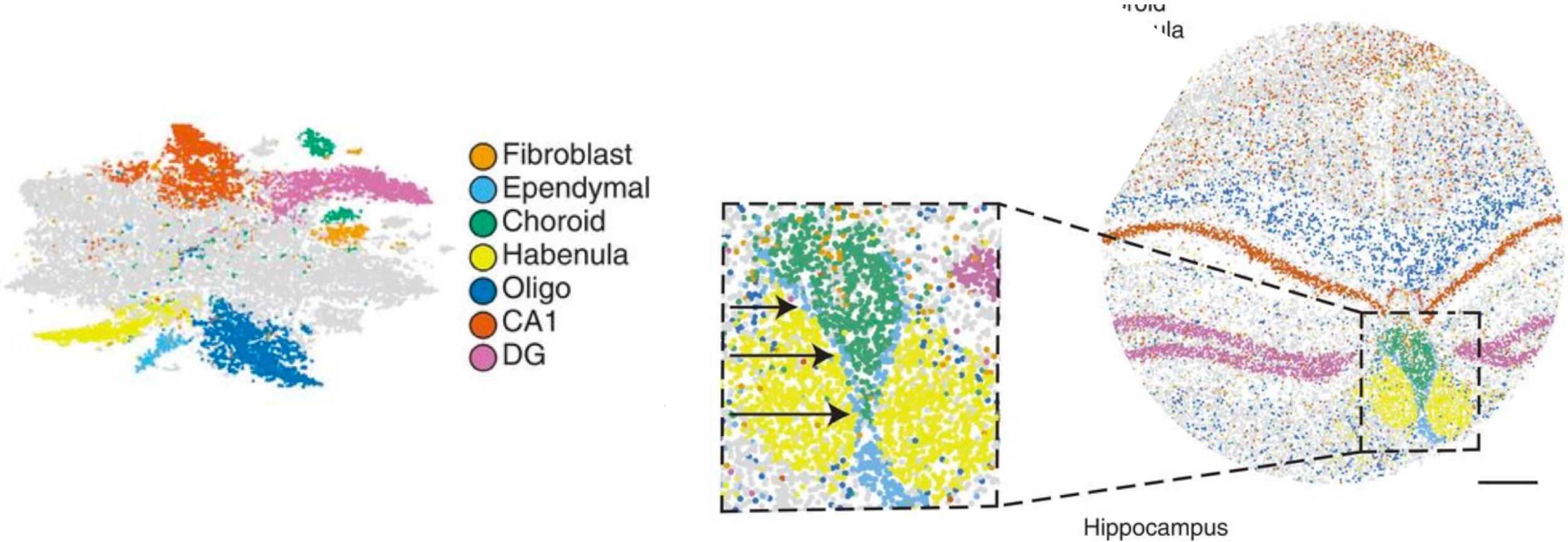
In situ indexing



B

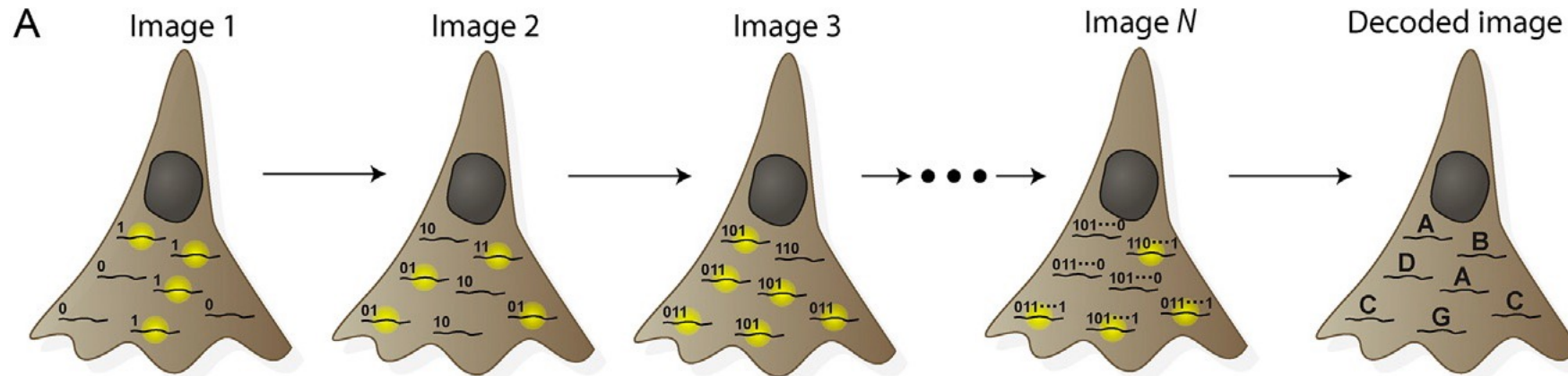


Slide-Seq overview



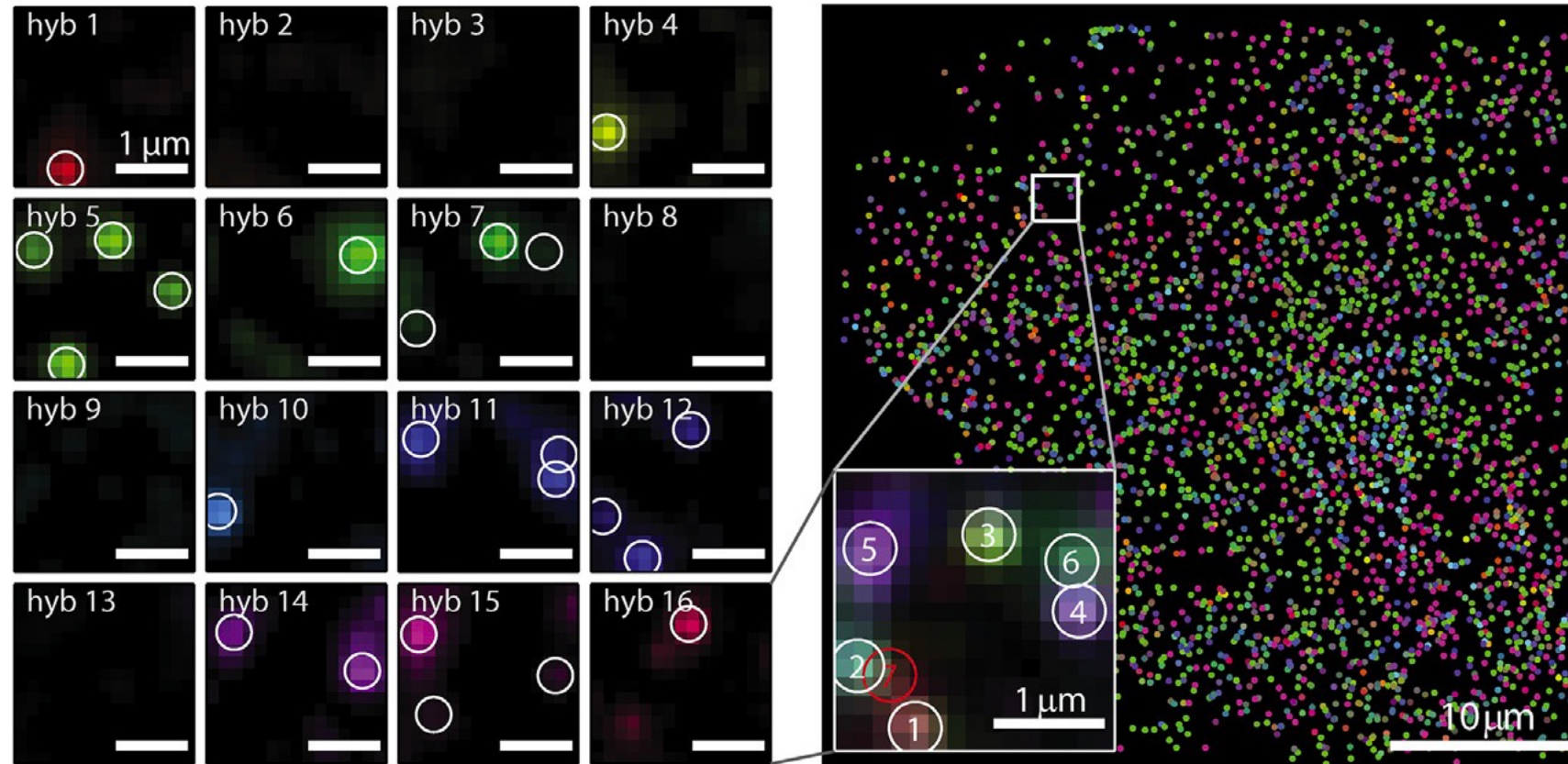
- *Slide-Seq is close to single-cell resolution (10 μm bead) but usually requires deconvolution of the beads to identify the contributing single cells.*
- *Visium (10x) is a related spatial transcriptomic technology.*

MERFISH overview (choose target genes)



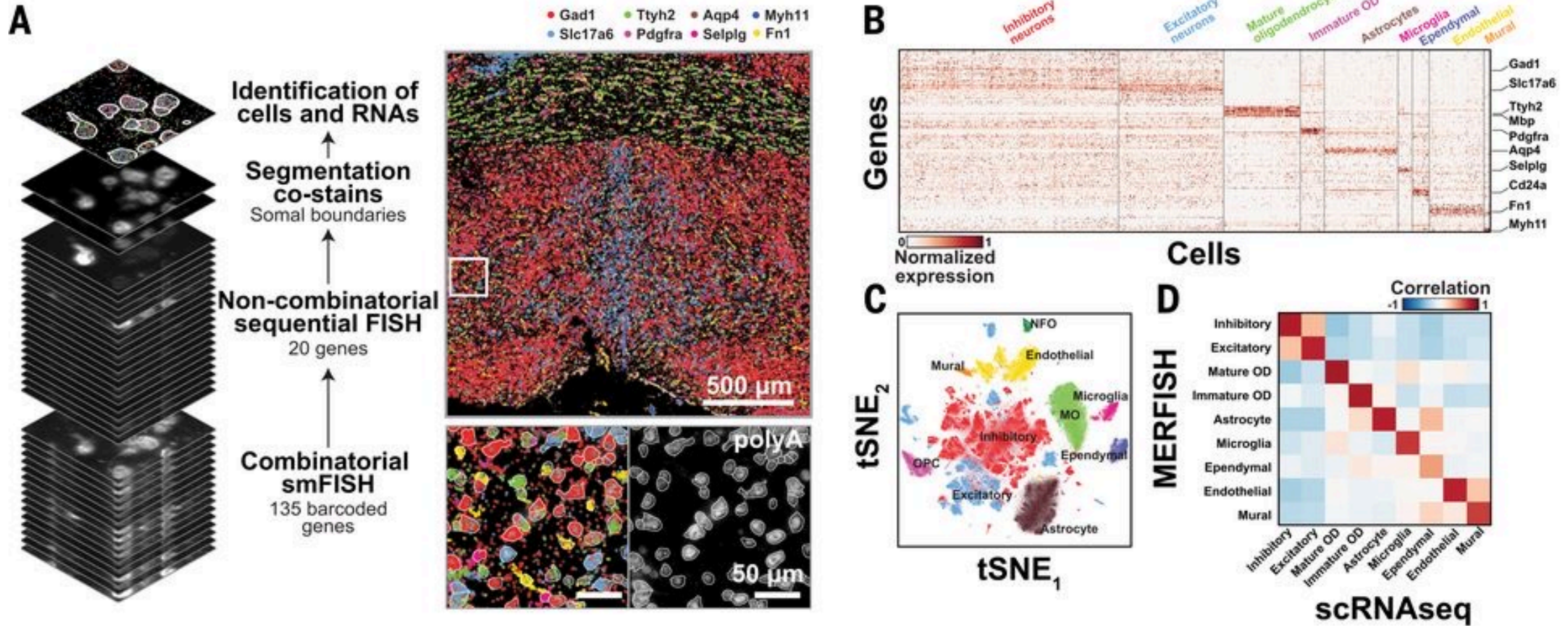
"multiplexed error-robust fluorescence in situ hybridization (MERFISH), a massively parallelized form of smFISH that can image and identify hundreds to thousands of different RNA species simultaneously with high accuracy in individual cells in their native spatial context."

MERFISH overview

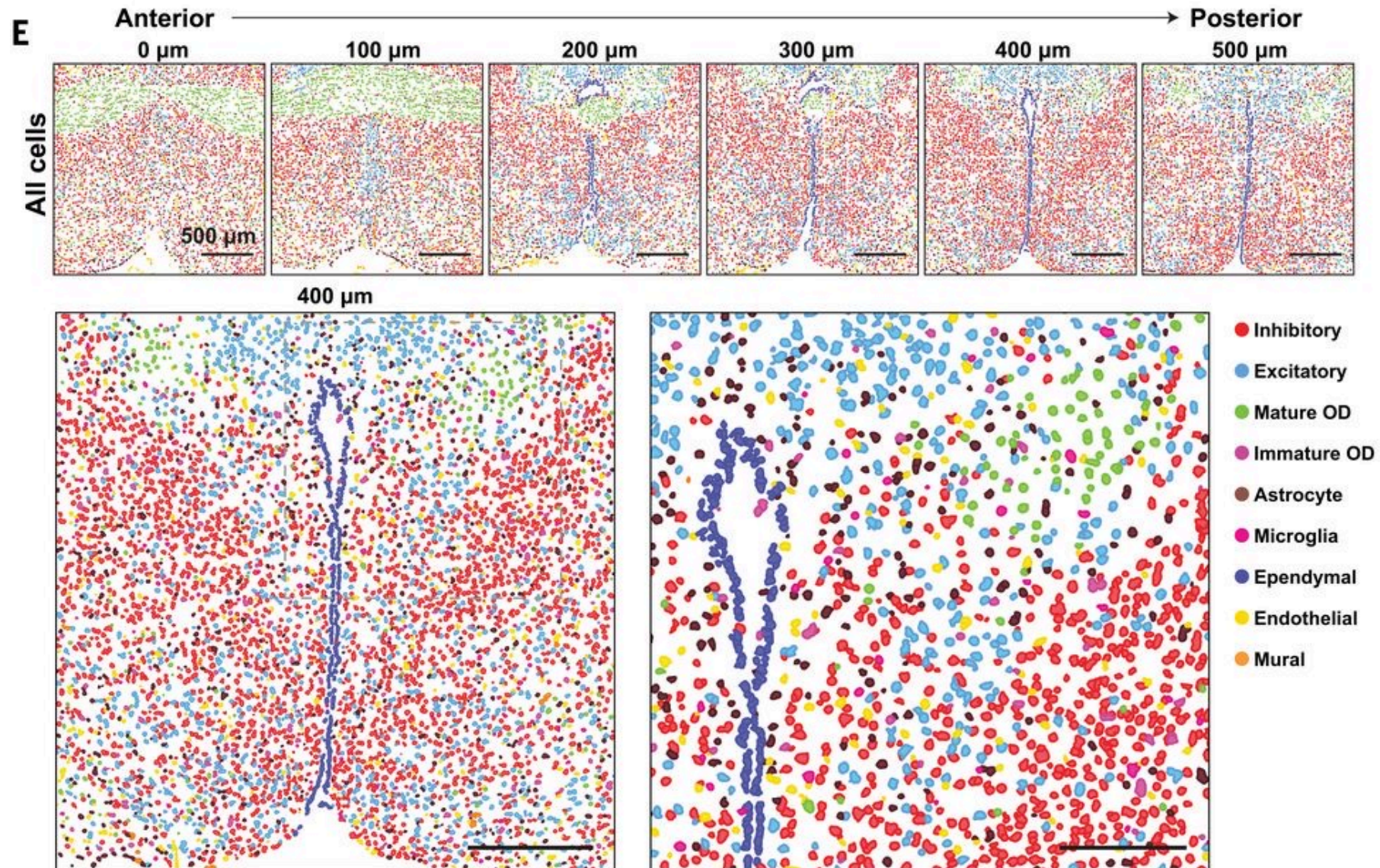


Imaging 100s-1000s RNA species in their native cellular and tissue environment

MERFISH overview

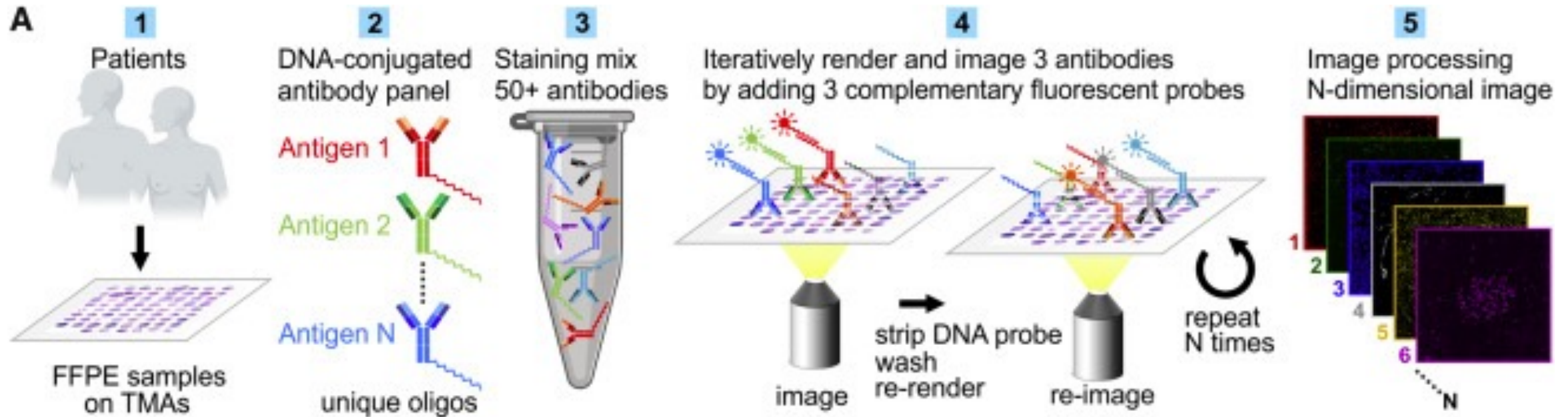


MERFISH overview



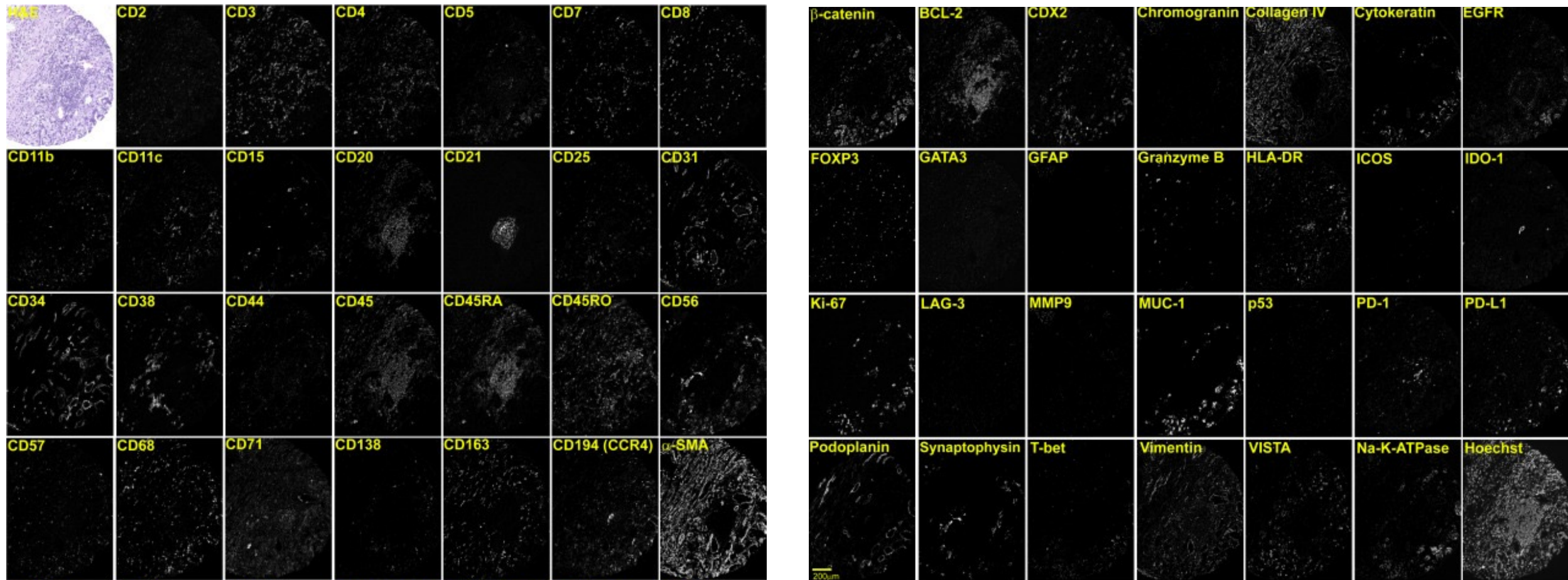
CODEX overview (choose target proteins)

Stain and iteratively image antibodies applied to a tissue sample

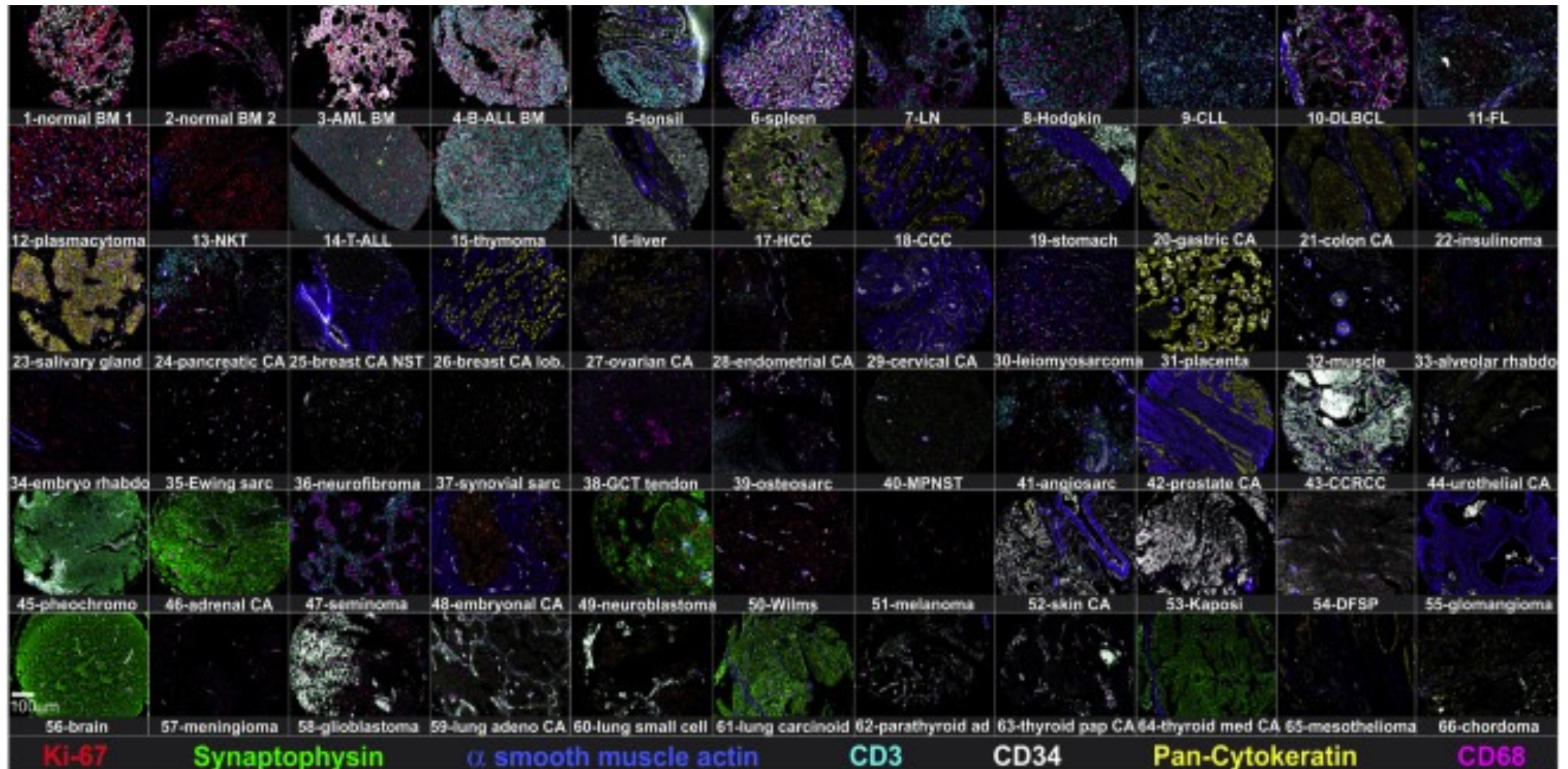


CODEX overview: 56 antibody marker panel

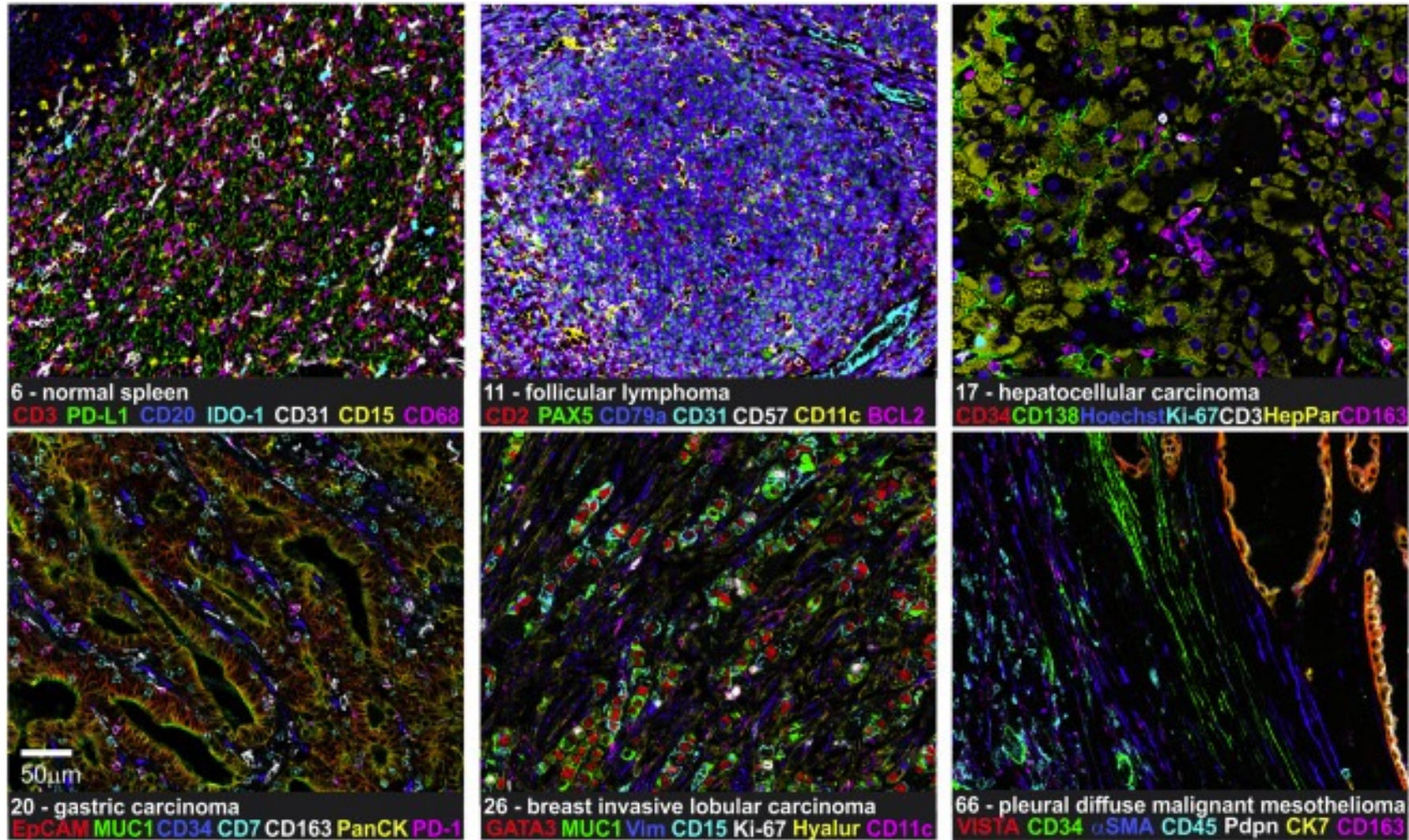
The antibody panel must be validated, and the panel typically includes many immune targets



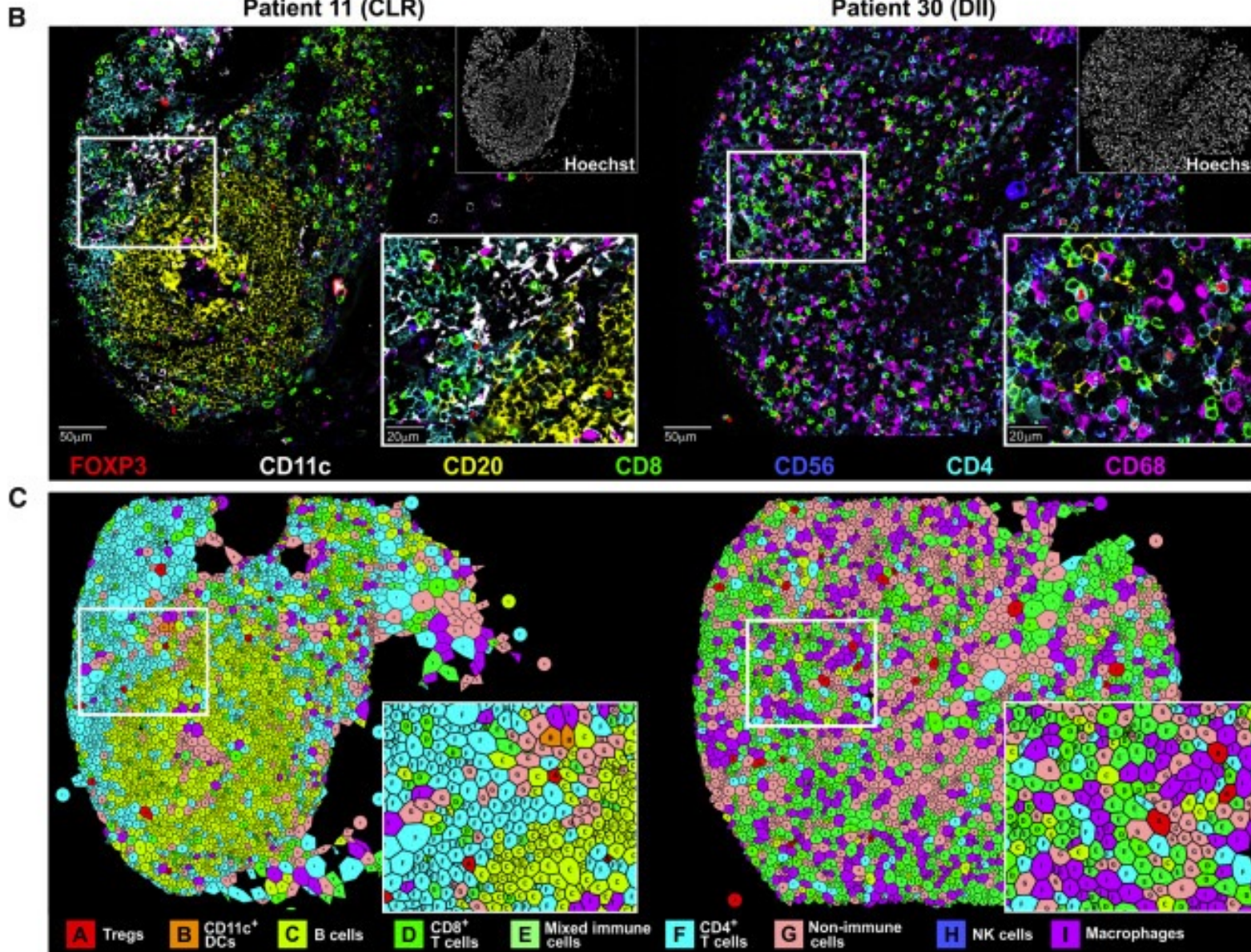
CODEX overview: multi-tumor tissue microarray



CODEX overview: multi-tumor tissue microarray



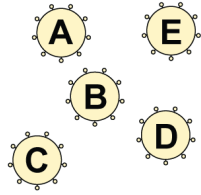
CODEX overview: 56 antibody marker panel



Big challenges

- Imaging
- Segmentation
- Analysis

Perturb-Seq overview



I. Pooled sgRNA library,
polydenylated guide barcode

II. Transduce cells
(low or high MOI)

III. Capture cells,
barcode mRNA

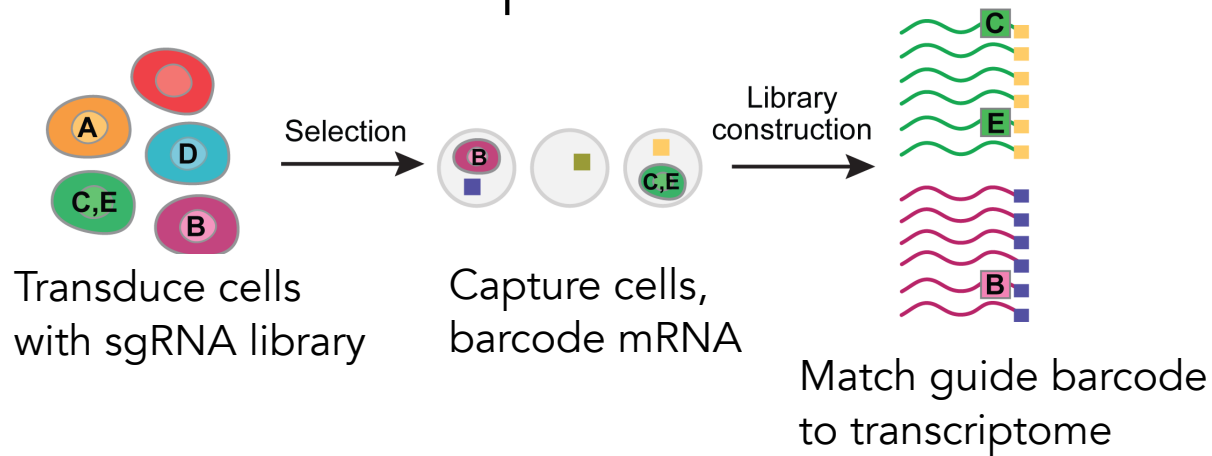
IV. Match guide barcode
to transcriptome

Perturb-Seq advantages

- no need to define a phenotype beforehand
- single cell readout
- pooled format
- higher order interactions

Perturb-Seq overview

Perturb-Seq links perturbations to their effects on the transcriptome



$$\log(\text{Cells}) \begin{matrix} \text{Y} \\ \text{Expression} \\ \text{matrix} \\ \text{Genes} \end{matrix} = \begin{matrix} \text{X} \\ \text{Design} \\ \text{matrix} \\ \text{Covariates} \end{matrix} \begin{matrix} \beta \\ \text{Coefficient} \\ \text{matrix} \\ \text{Genes} \end{matrix} + 1$$

Infer perturbation regulatory effects

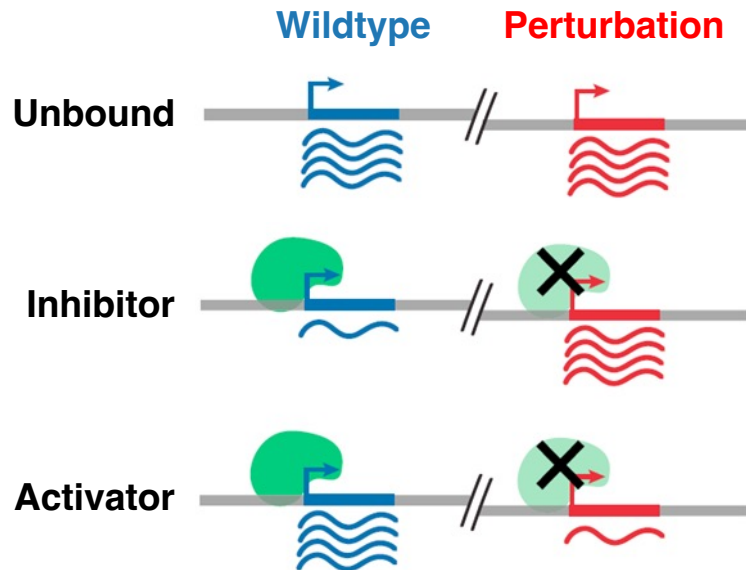
The equation shows the relationship between the log-transformed expression matrix (Y), the design matrix (X) which includes sgRNAs and confounders, and the coefficient matrix (β) representing the regulatory effects of perturbations on genes.

Ongoing developments in single-cell genomics:

Perturb-Seq

$$\log(\text{Cells}) \begin{matrix} \text{Y} \\ \text{Expression} \\ \text{matrix} \\ \begin{bmatrix} 0 & 8 & 0 & 1 & 0 & 2 & \dots & 0 \\ 1 & 2 & 0 & 0 & 1 & 1 & \dots & 0 \\ 0 & 7 & 0 & 1 & 0 & 0 & \dots & 1 \\ \dots & & & & & & & \\ 0 & 5 & 0 & 1 & 0 & 0 & \dots & 0 \end{bmatrix} \\ \text{Genes} \\ \text{Expression} \\ \text{(observed)} \end{matrix} + 1 = \begin{matrix} \text{X} \\ \text{Design} \\ \text{matrix} \\ \begin{bmatrix} 1 & 0 & 0 & \dots & -0.1 & \dots \\ 0 & 1 & 0 & \dots & 0.3 & \dots \\ 0 & 0 & 1 & \dots & 0.2 & \dots \\ \dots & & & & & \\ 0 & 0 & 0 & \dots & -0.2 & \dots \end{bmatrix} \\ \text{sgRNAs} \quad \text{Confounders} \\ \text{Covariates} \\ \text{GBCs} \\ \text{(observed)} \end{matrix} \begin{matrix} \beta \\ \text{Coefficient} \\ \text{matrix} \\ \begin{bmatrix} \beta_{1,1} & \beta_{1,2} & \beta_{1,3} & \dots & \beta_{1,G} \\ \beta_{2,1} & \beta_{2,2} & \beta_{2,3} & \dots & \beta_{2,G} \\ \beta_{3,1} & \beta_{3,2} & \beta_{3,3} & \dots & \beta_{3,G} \\ \dots & & & & \dots \\ \beta_{C,1} & \beta_{C,2} & \beta_{C,3} & \dots & \beta_{C,G} \end{bmatrix} \\ \text{Genes} \\ \text{Regulation} \\ \text{(fit)} \end{matrix} \begin{matrix} \text{Covariates} \end{matrix}$$

β is the regulatory effect a given guide has on a given gene.



$\beta = 0$, no change

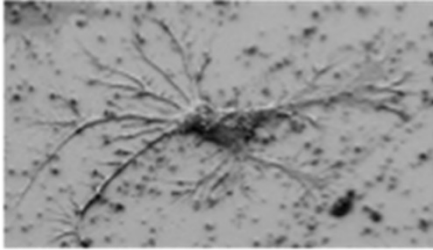
$\beta > 0$, increased expression

$\beta < 0$, decreased expression

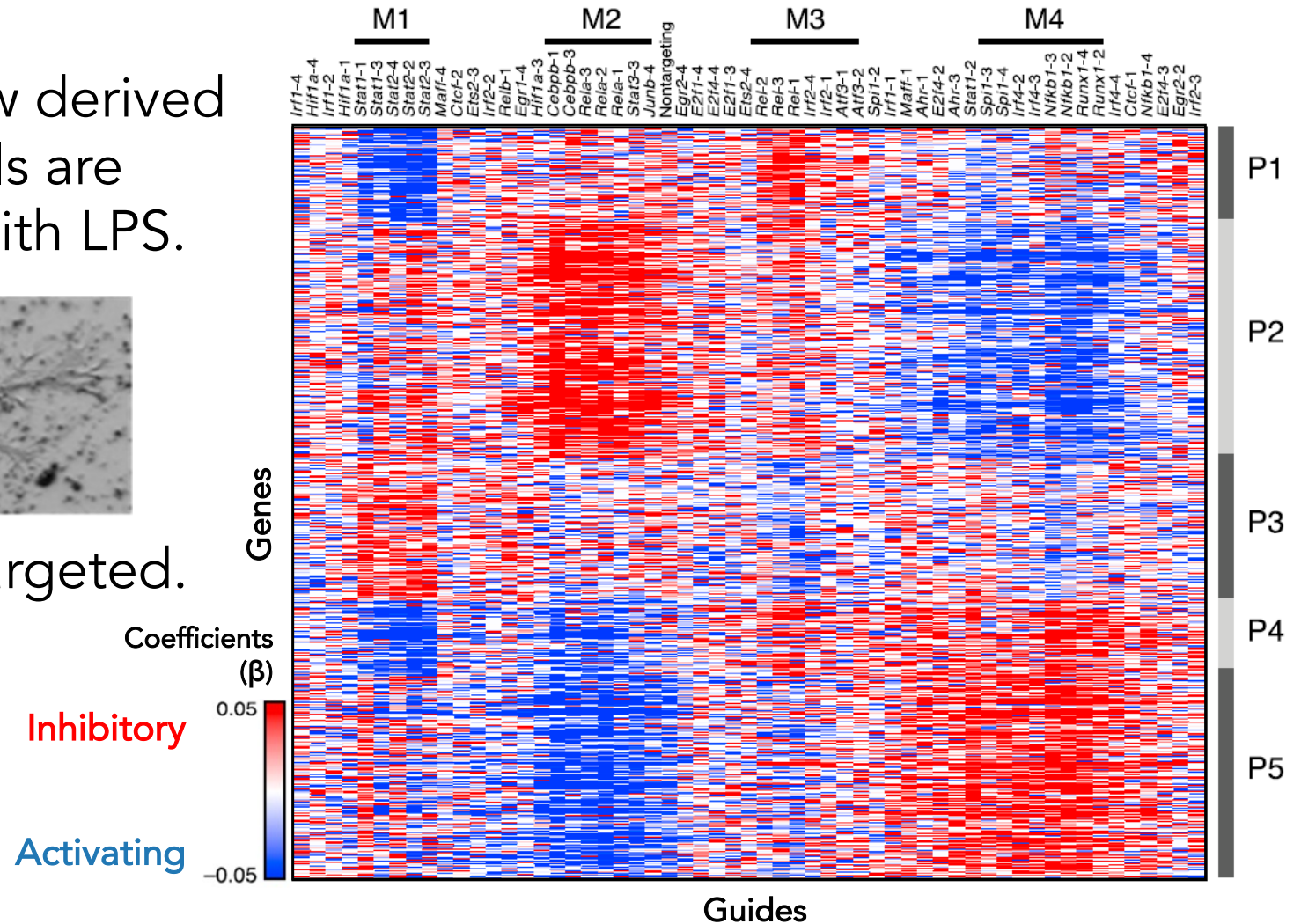
Perturbations reveal gene regulation logic

Guides and genes cluster into transcription factor modules (M) and regulated expression programs (P) respectively.

- Bone-marrow derived dendritic cells are stimulated with LPS.



- 24 TFs are targeted.



Perturbations reveal gene regulation logic

