

Single-cell transcriptomics in plants

Standard and advanced scRNA-seq data analysis workflows

Michael Schon

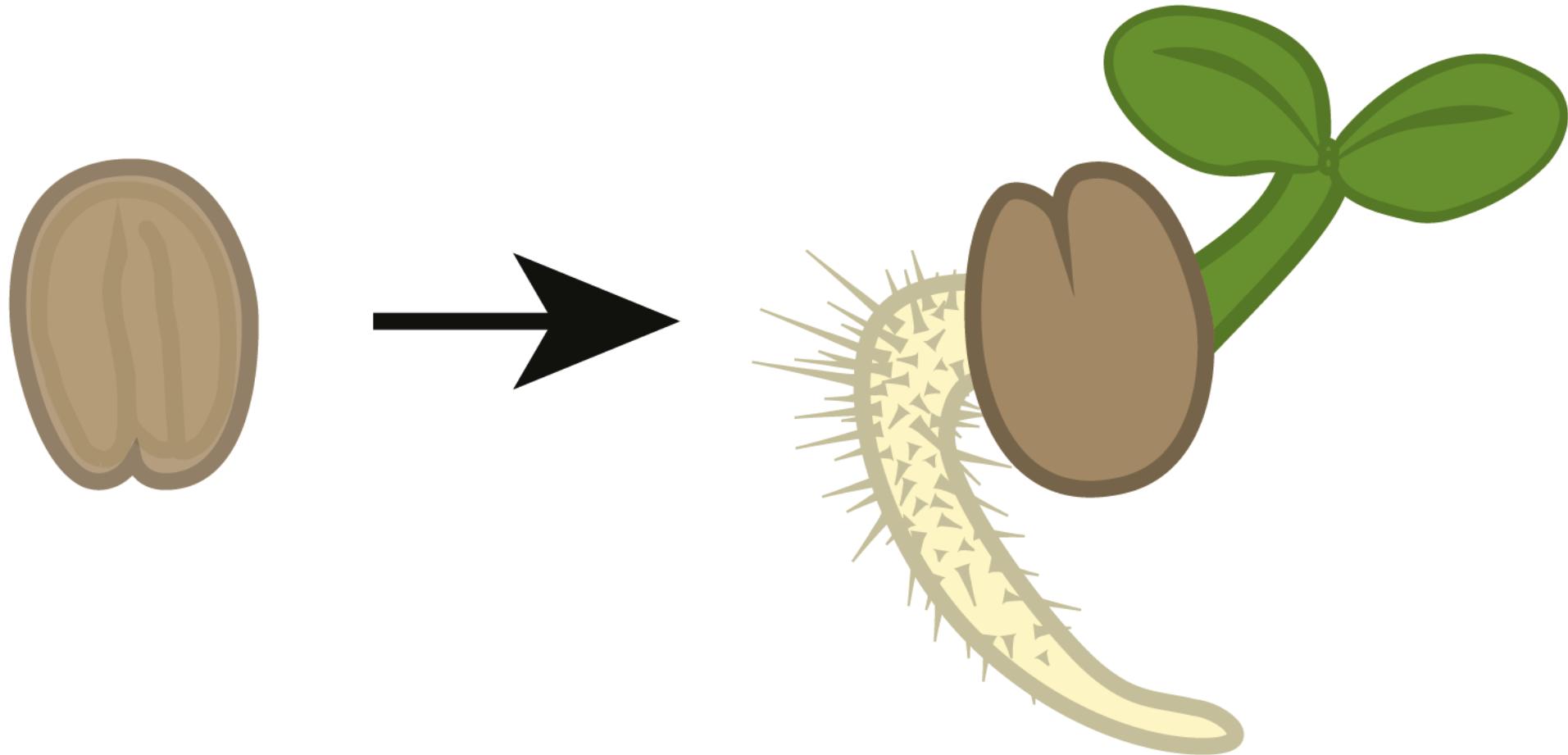
Postdoc, Nodine Group, Laboratory of Molecular Biology
Wageningen University & Research
23 March 2023

Overview

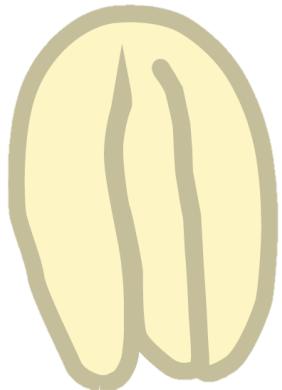
- Why single cells?
- How-to:
 - Biology → Hardware → Chemistry → Software → *Knowledge?*
- The standard workflow
 - Reads → Counts → Clusters → Markers
 - Dimensionality reduction
- Advanced techniques



Single cell RNA-seq: Why single cells?

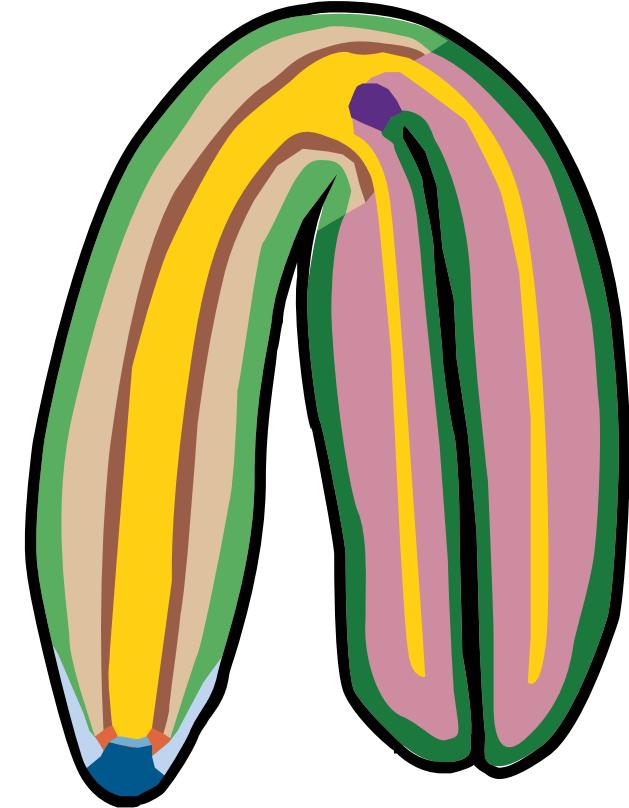


Single cell RNA-seq: Why single cells?

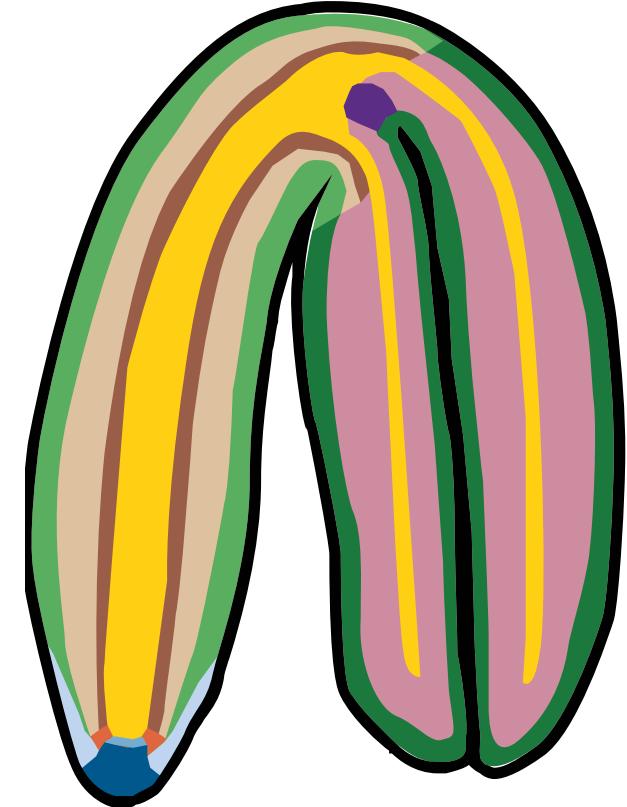
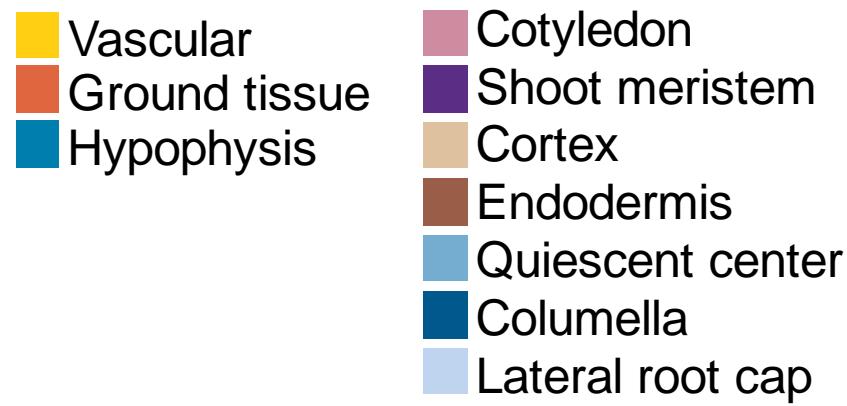


Vascular
Ground tissue
Hypophysis

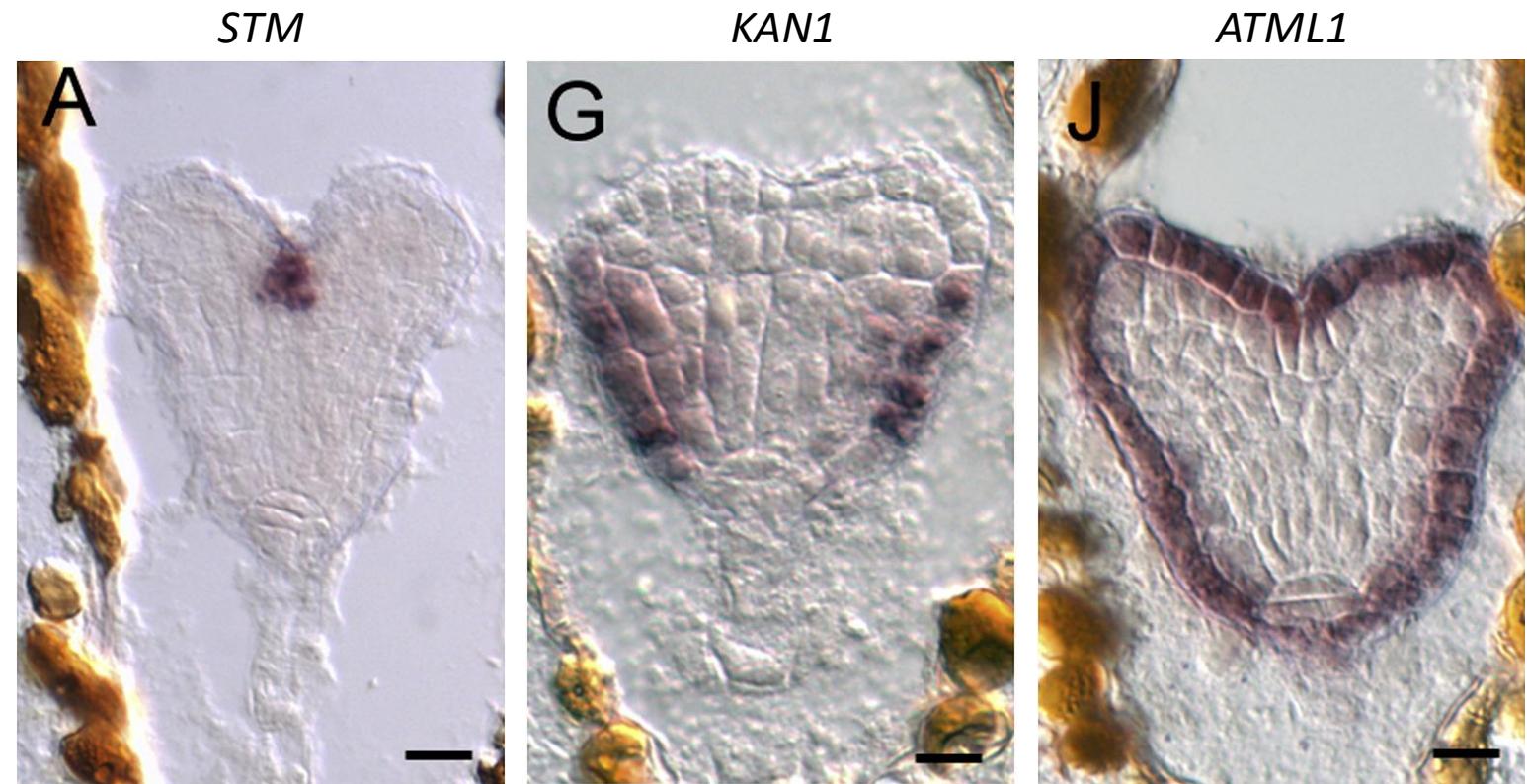
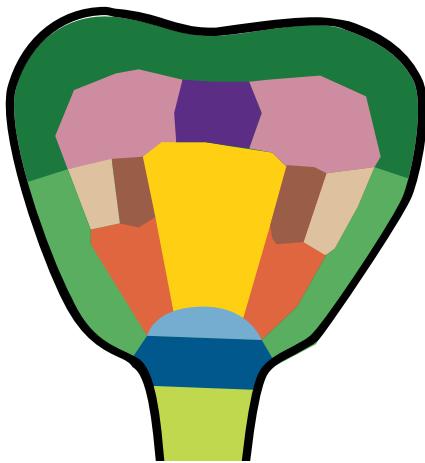
Cotyledon
Shoot meristem
Cortex
Endodermis
Quiescent center
Columella
Lateral root cap



Single cell RNA-seq: Why single cells?

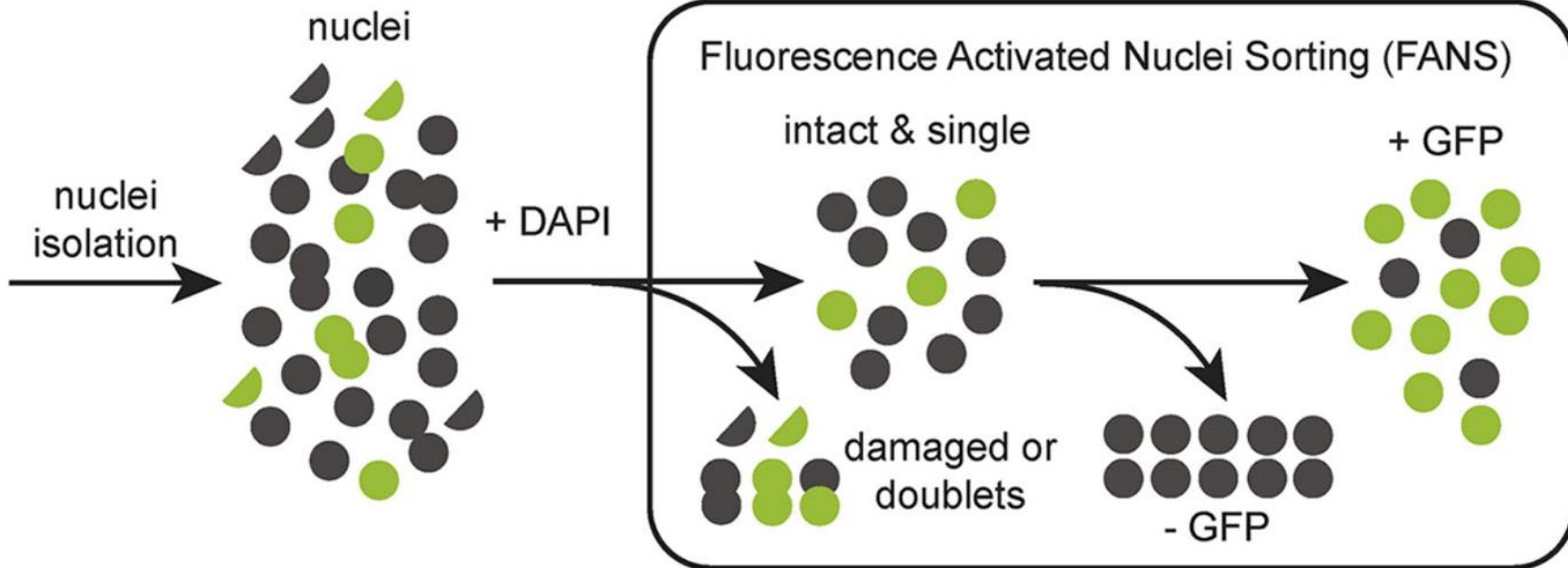
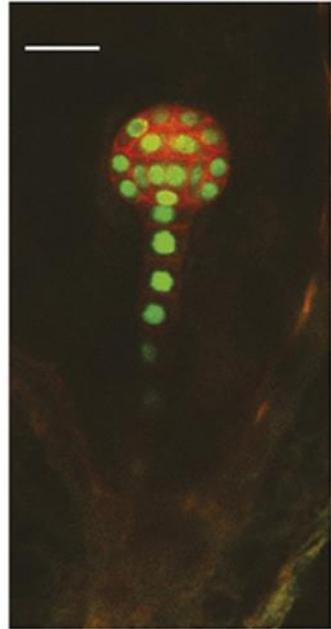


Single cell RNA seq: Why RNA?

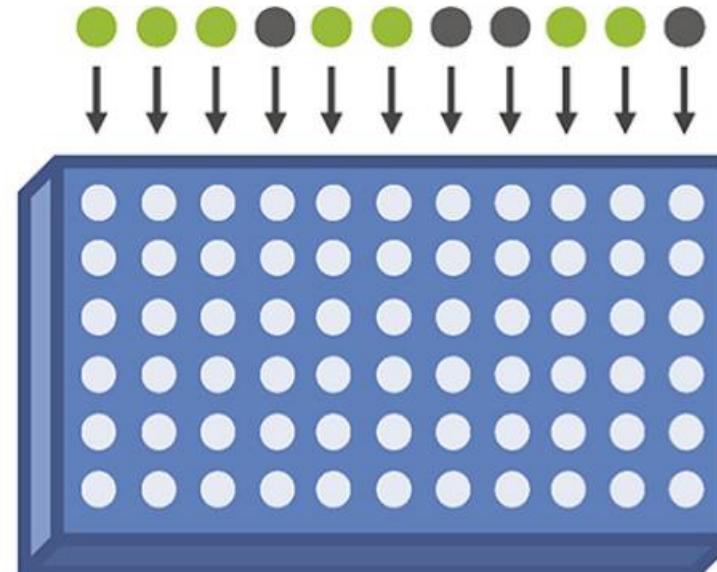


Grigg et al. 2009 *Current Biology*

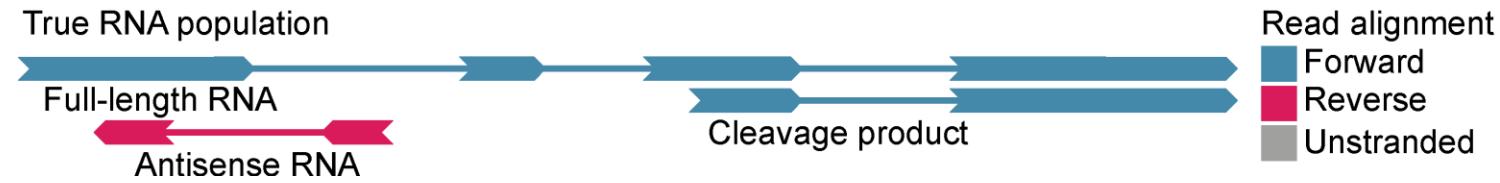
Capturing cells



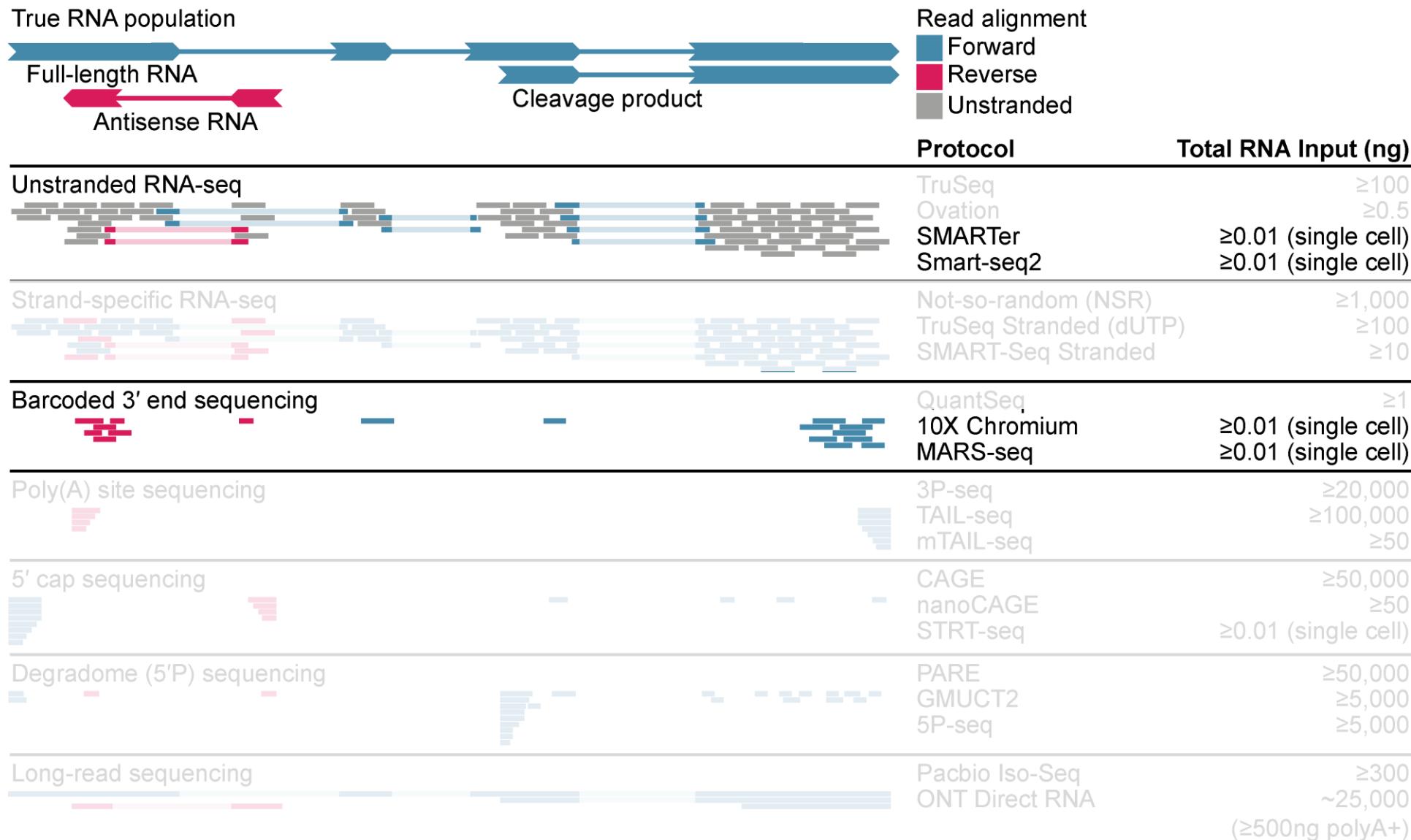
Kao et al. 2021 *Development*



Capturing the transcriptome

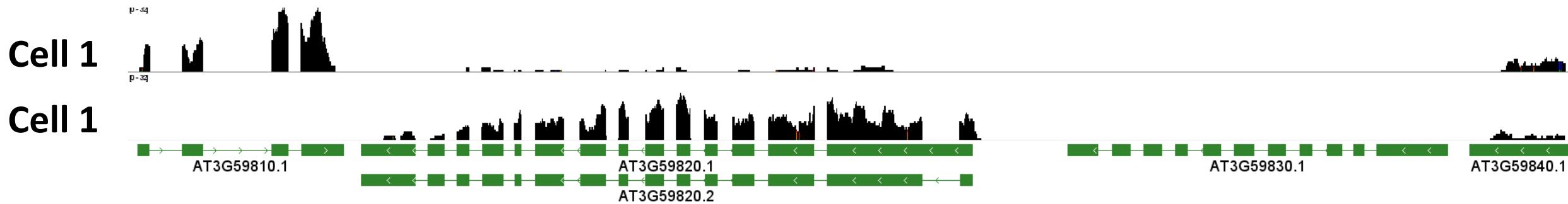


Capturing the transcriptome



RNA x Cells... Now what?

- Align sequenced reads with transcripts
 - Genome alignments (STAR, HiSat)
 - Pseudoalignment (Kallisto, Salmon)

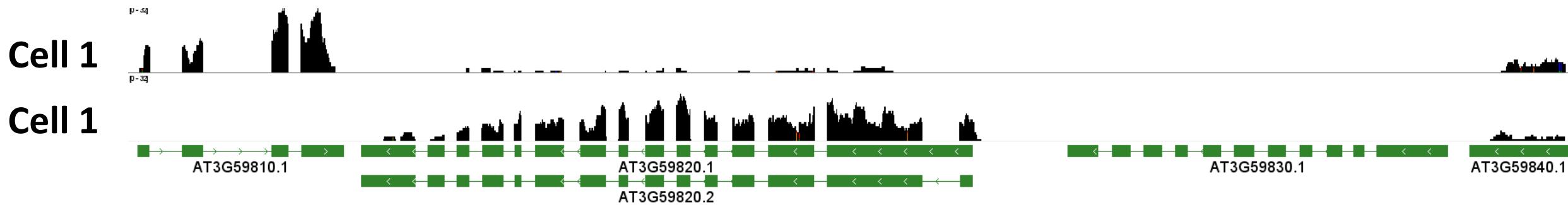


	Cell 1	Cell 2
AT3G59810	110	0
AT3G59820	20	304
AT3G59830	0	0
AT3G59840	22	13

More information in
Practical 1!

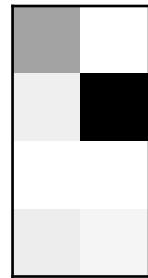
RNA x Cells... Now what?

- Align sequenced reads with transcripts
 - Genome alignments (STAR, HiSat)
 - Pseudoalignment (Kallisto, Salmon)

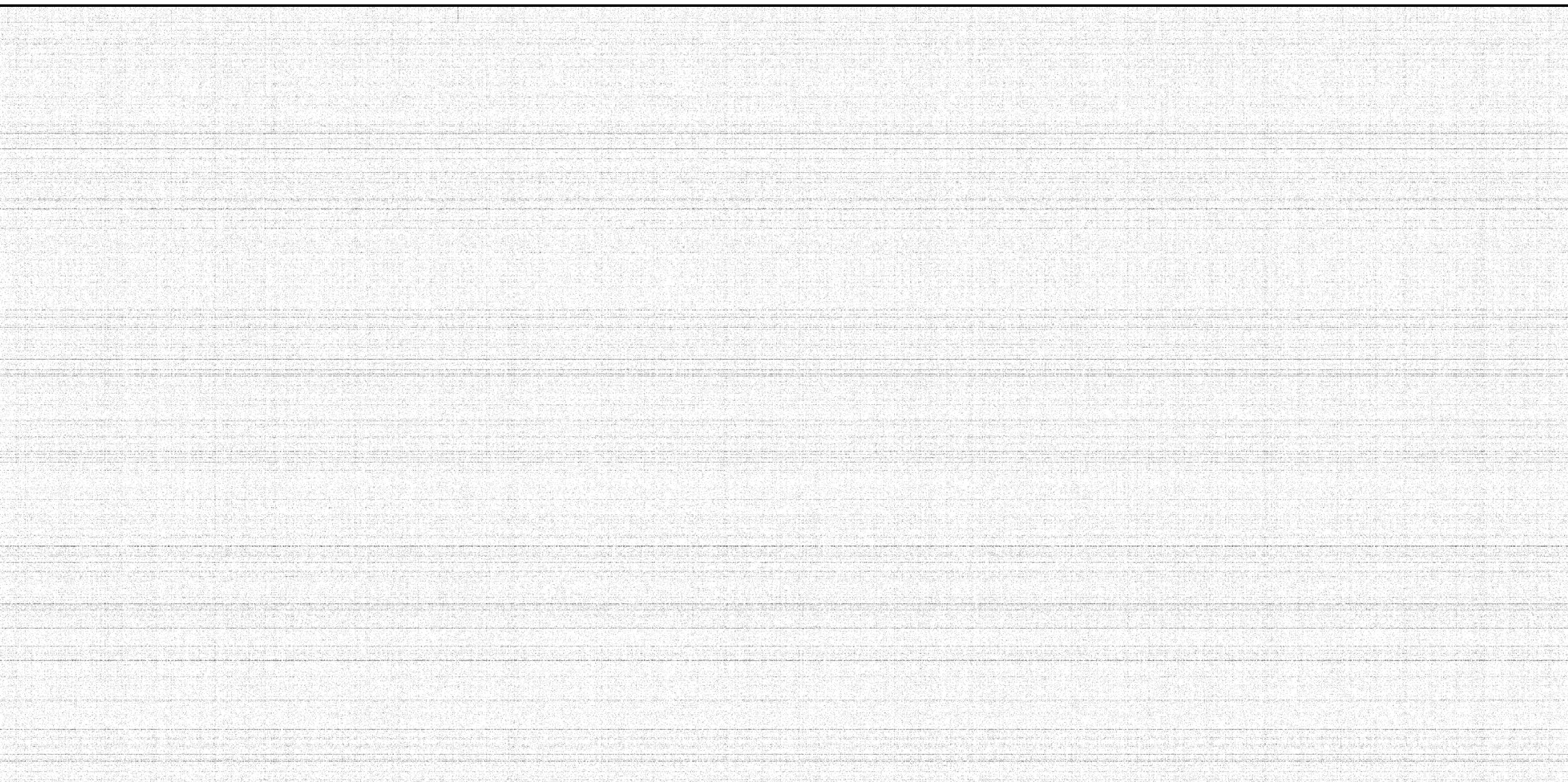


	Cell 1	Cell 2
AT3G59810	110	0
AT3G59820	20	
AT3G59830	0	0
AT3G59840	22	13

RNA x Cells... Now what?

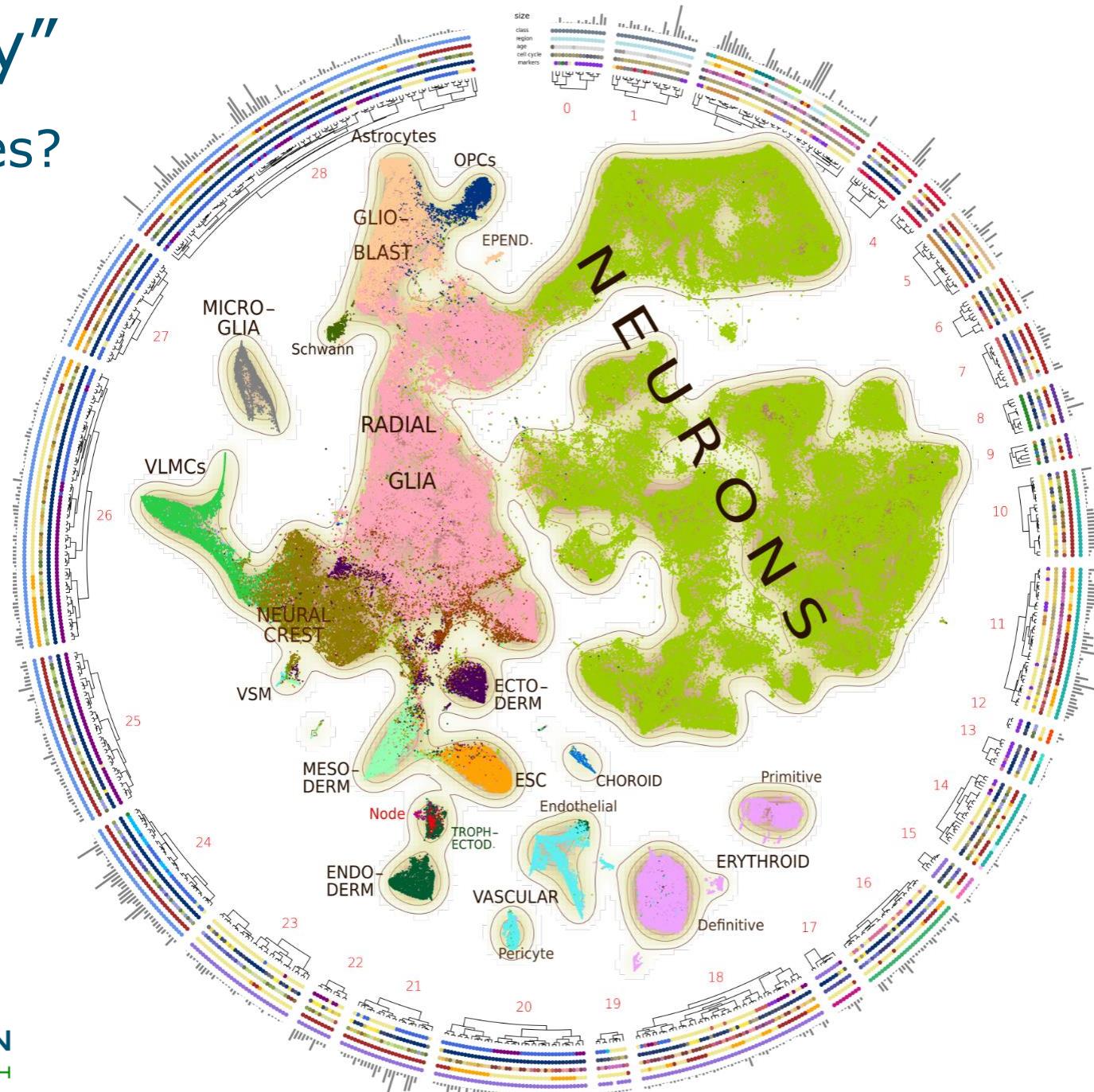


RNA x Cells... Now what?



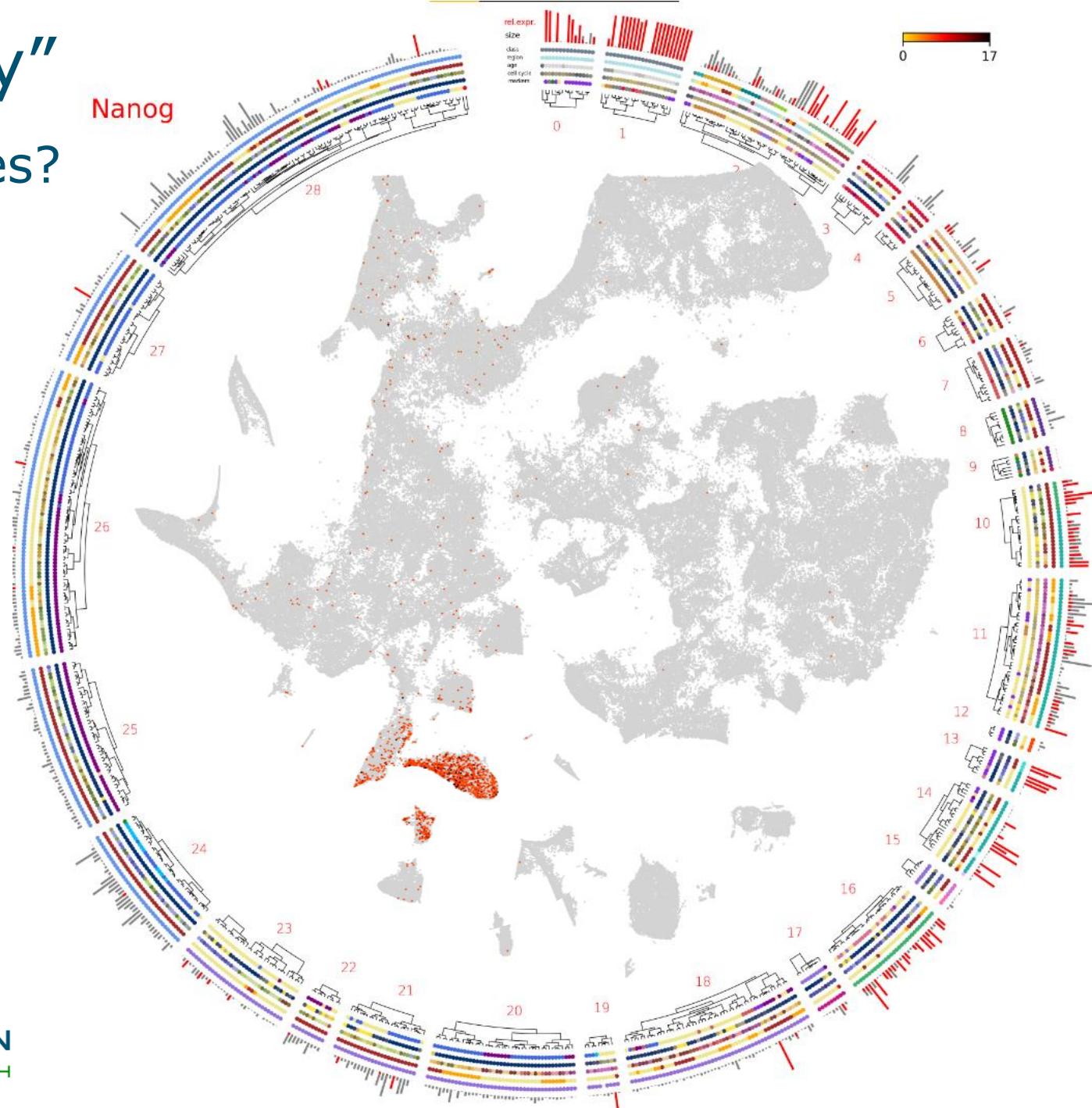
“Similarity”

- Marker genes?
 - *NANOG*



“Similarity”

- Marker genes?
 - *NANOG*



“Similarity”

- Marker genes?

- *NANOG*

⚠ Not all ESCs have *NANOG*

⚠ Not all *NANOG*+ cells are ESCs



Neighbors

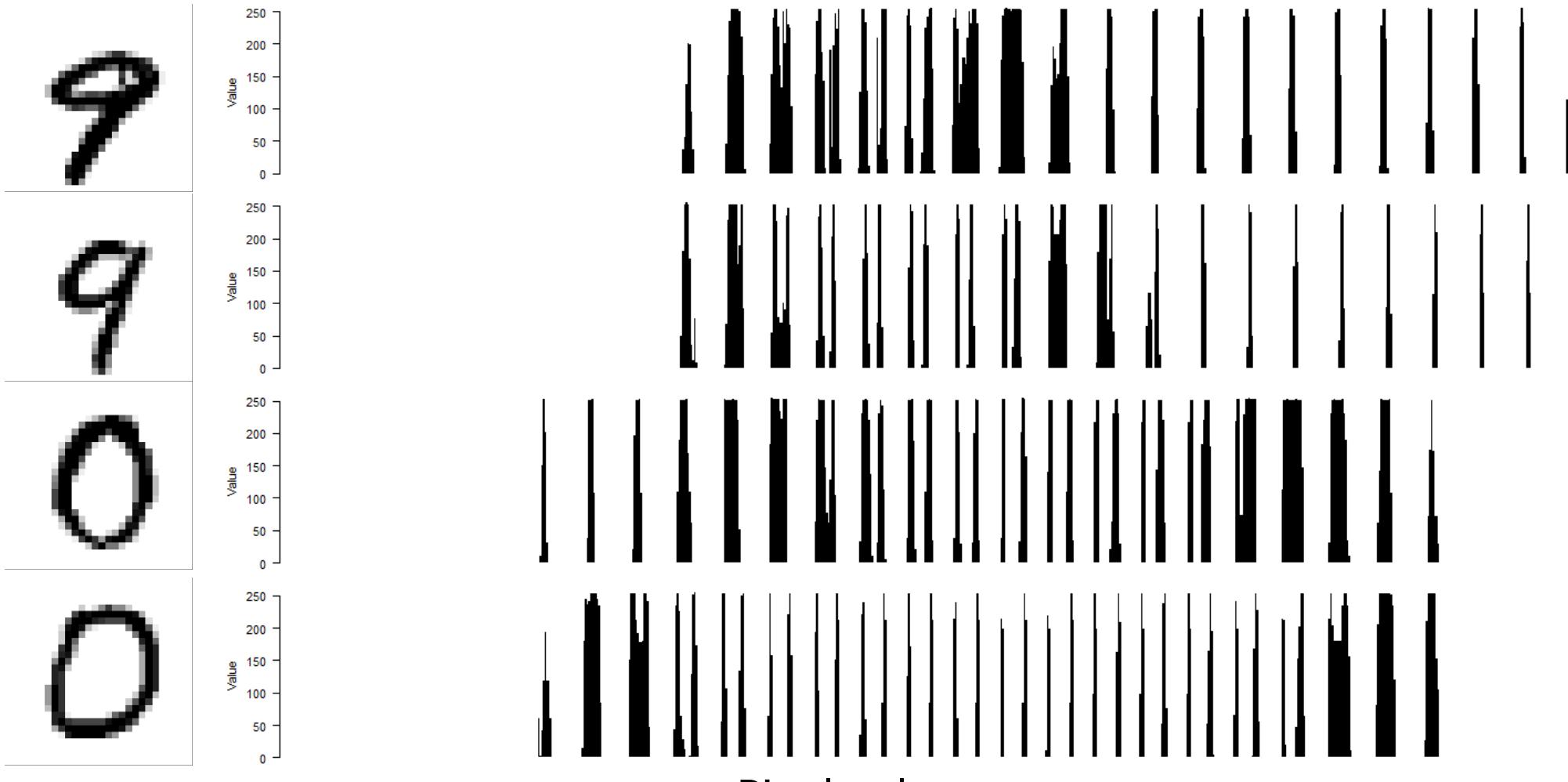
- t-SNE: stochastic neighbor embedding
 - Define neighbors in high-D space
 - Minimize a loss function
 - Many-body simulation
 - High initial learning rate (should scale with input)
- UMAP: uniform manifold approximation
 - Also minimizes neighbor distance, but by “folding space”

Neighbors in 784 dimensions

MNIST (70,000 hand-drawn digits, 28x28 pixels)

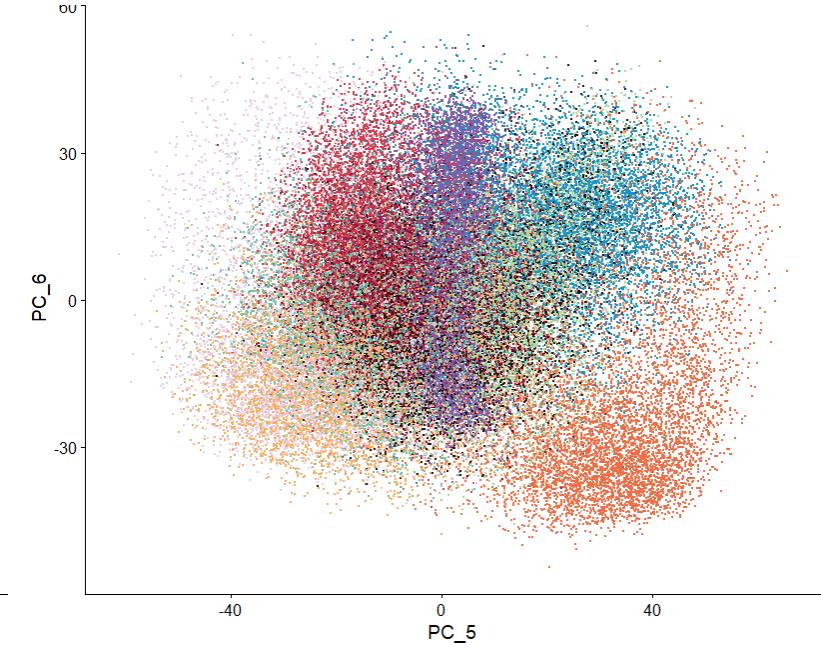
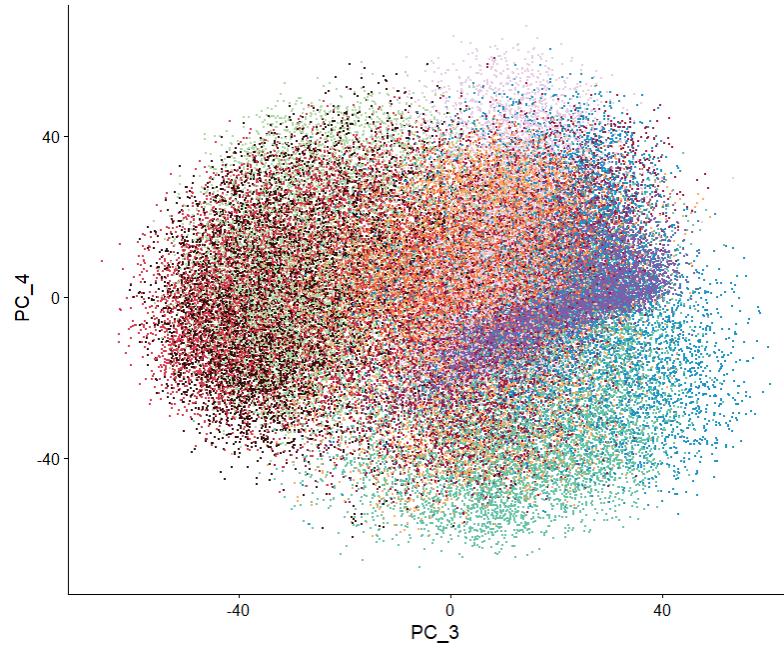
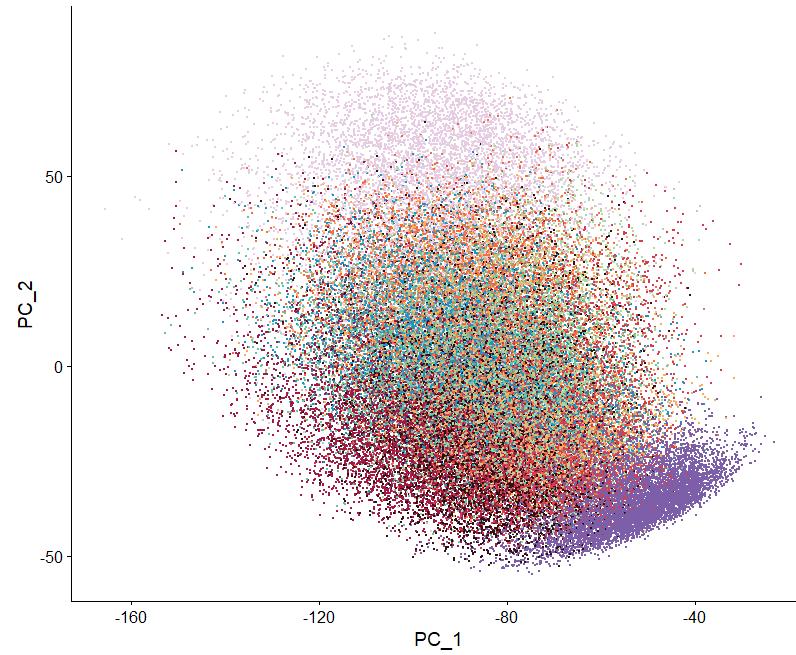
0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9 9 9

What dimensions are we reducing?



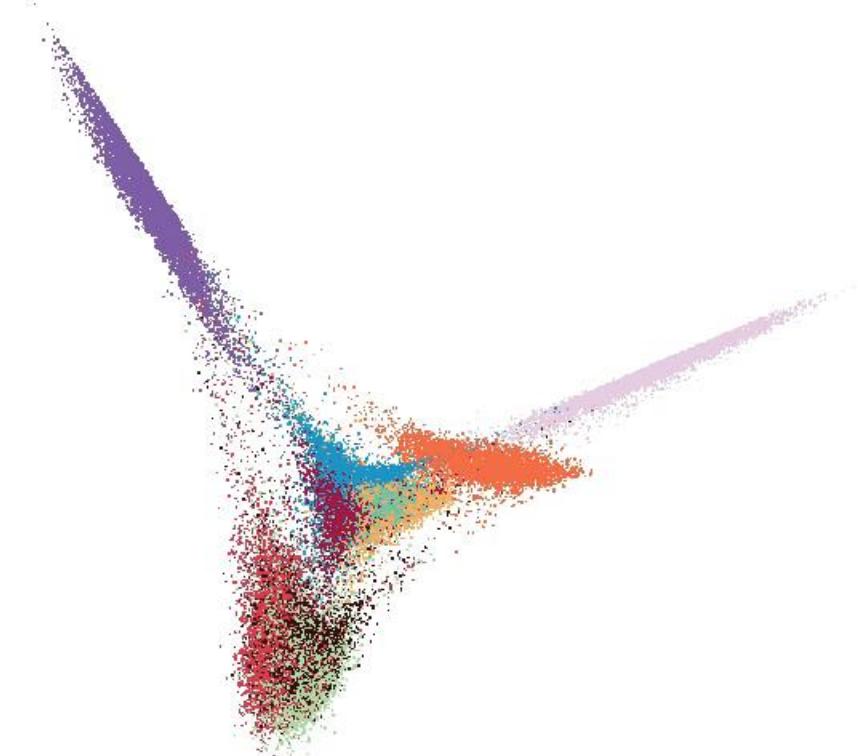
What dimensions are we reducing?

- “Traditional” DR technique: Principal component analysis (PCA)



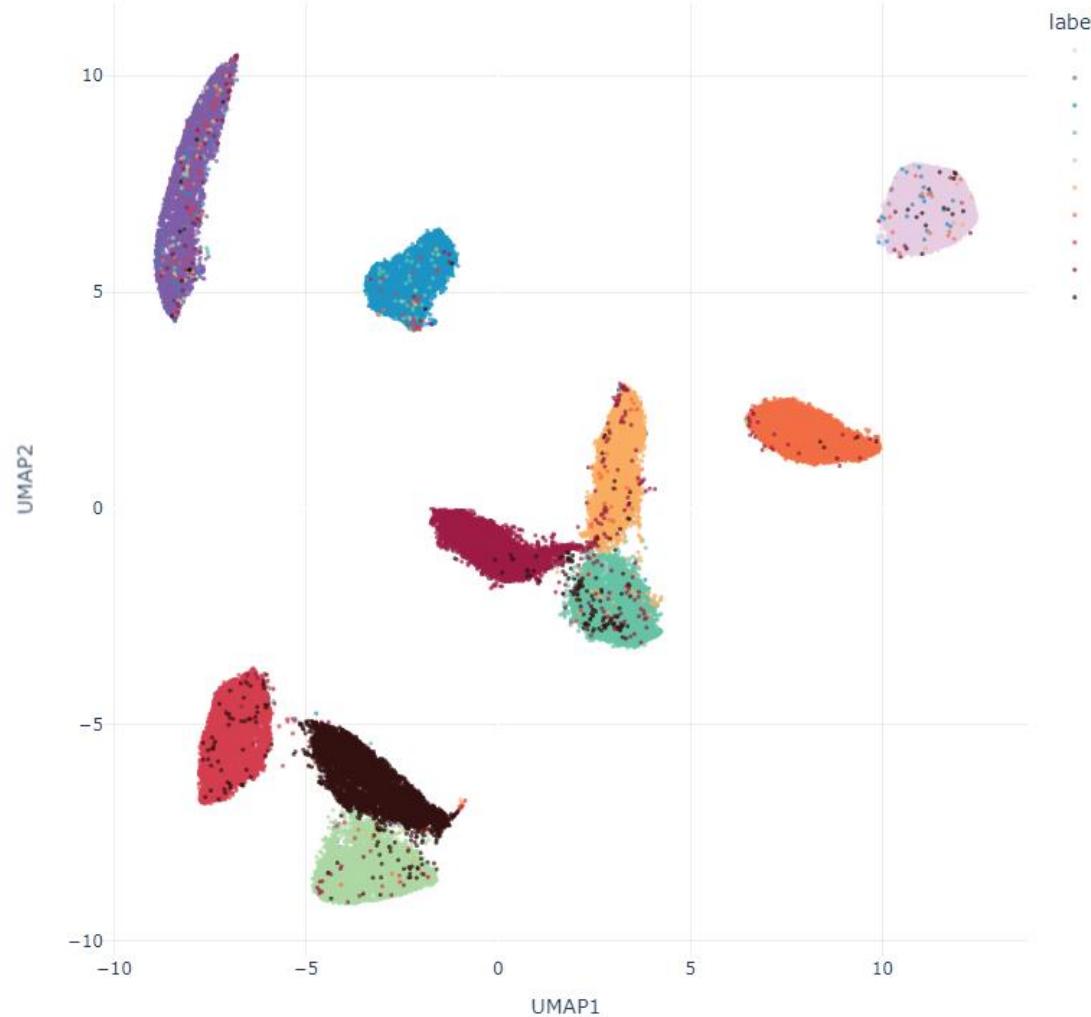
Dimensional reduction (UMAP)

- Begin with an embedding
 - PCA
 - “Spectral” embedding
(Laplacian eigenvectors)
- Minimize a loss function that penalizes high distances between neighbors
- “High-energy” early, followed by annealing



Dimensional reduction (UMAP)

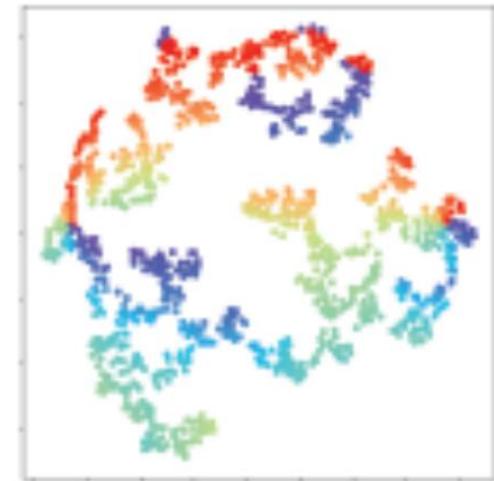
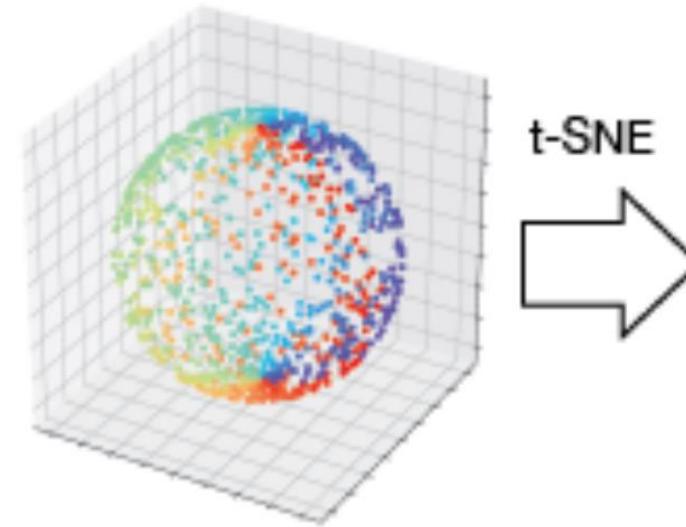
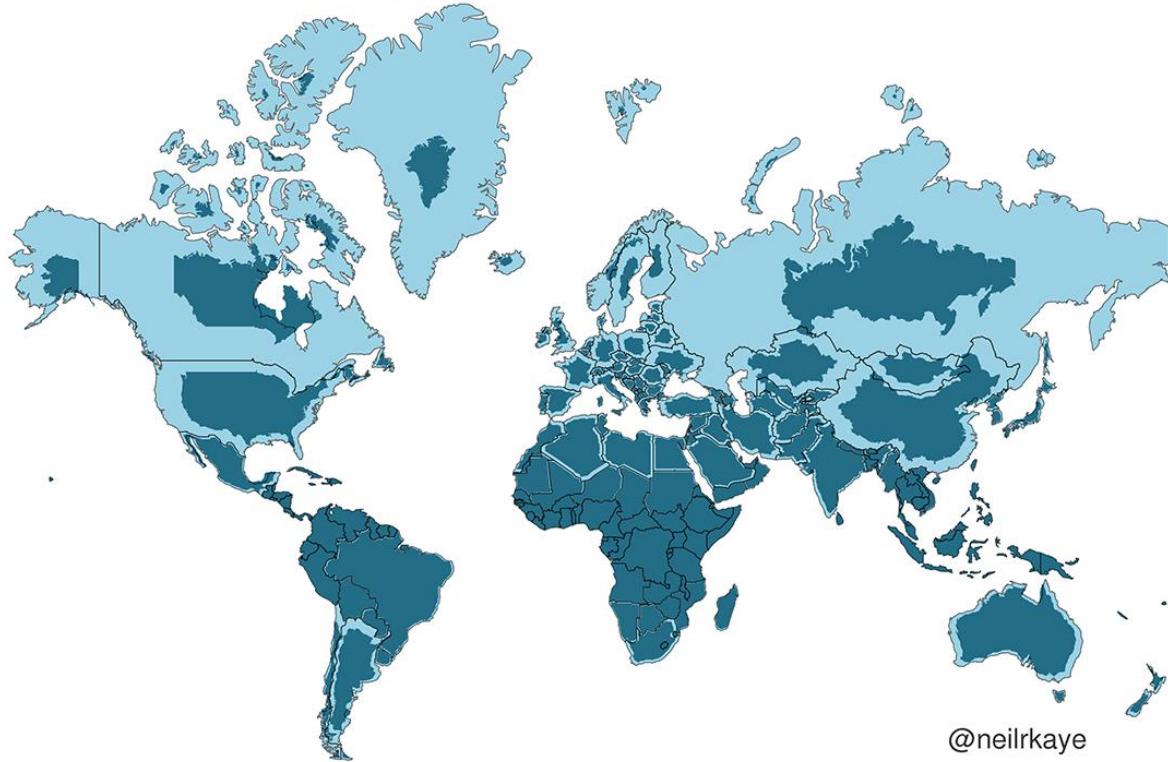
Interactive UMAP Plot of MNIST Dataset





Reducing dimensions distorts distance

MERCATOR PROJECTION VS THE TRUE SIZE OF COUNTRIES



Data



Initialization matters!

Finding causes of variation

- Differences will either be **discrete** or **continuous**
 - Two unrelated cell types = **discrete**
 - Clustering → marker genes
 - Cellular differentiation = **continuous**
 - Pseudotime trajectories

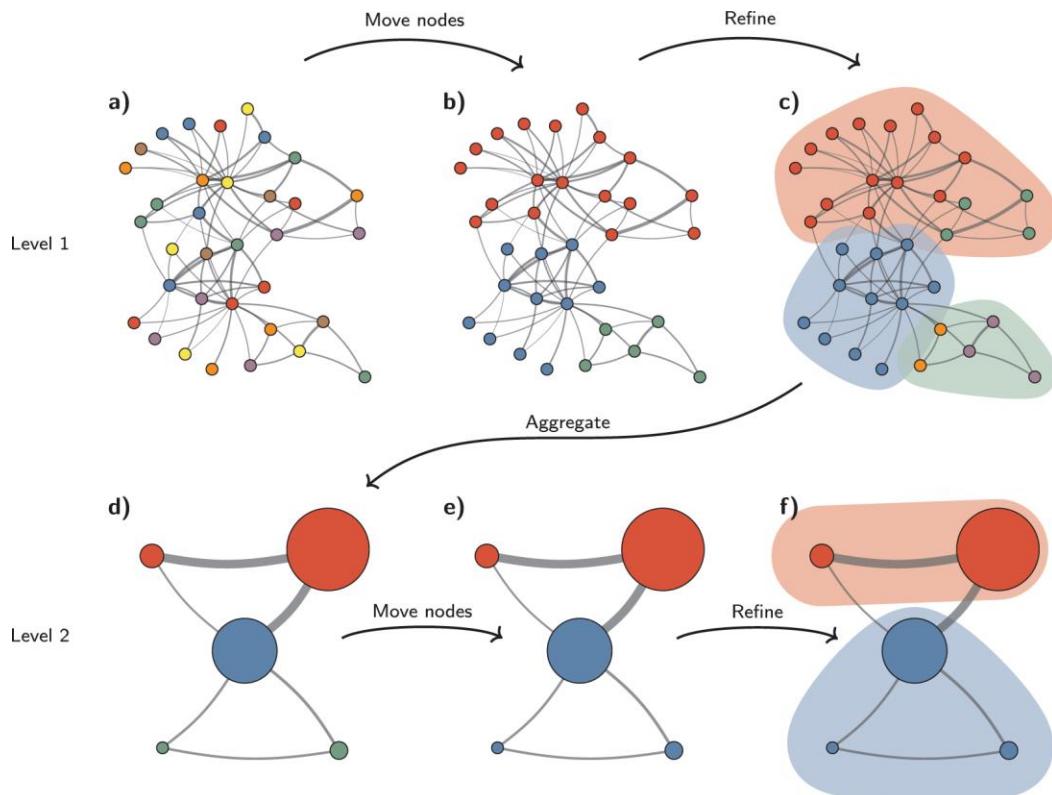
⚠ Not all variation is good variation!

- Mitochondrial / chloroplast RNA
- Doublets
- Dropouts
- Ambient RNA
- Cell cycle

Clustering

- Leiden algorithm

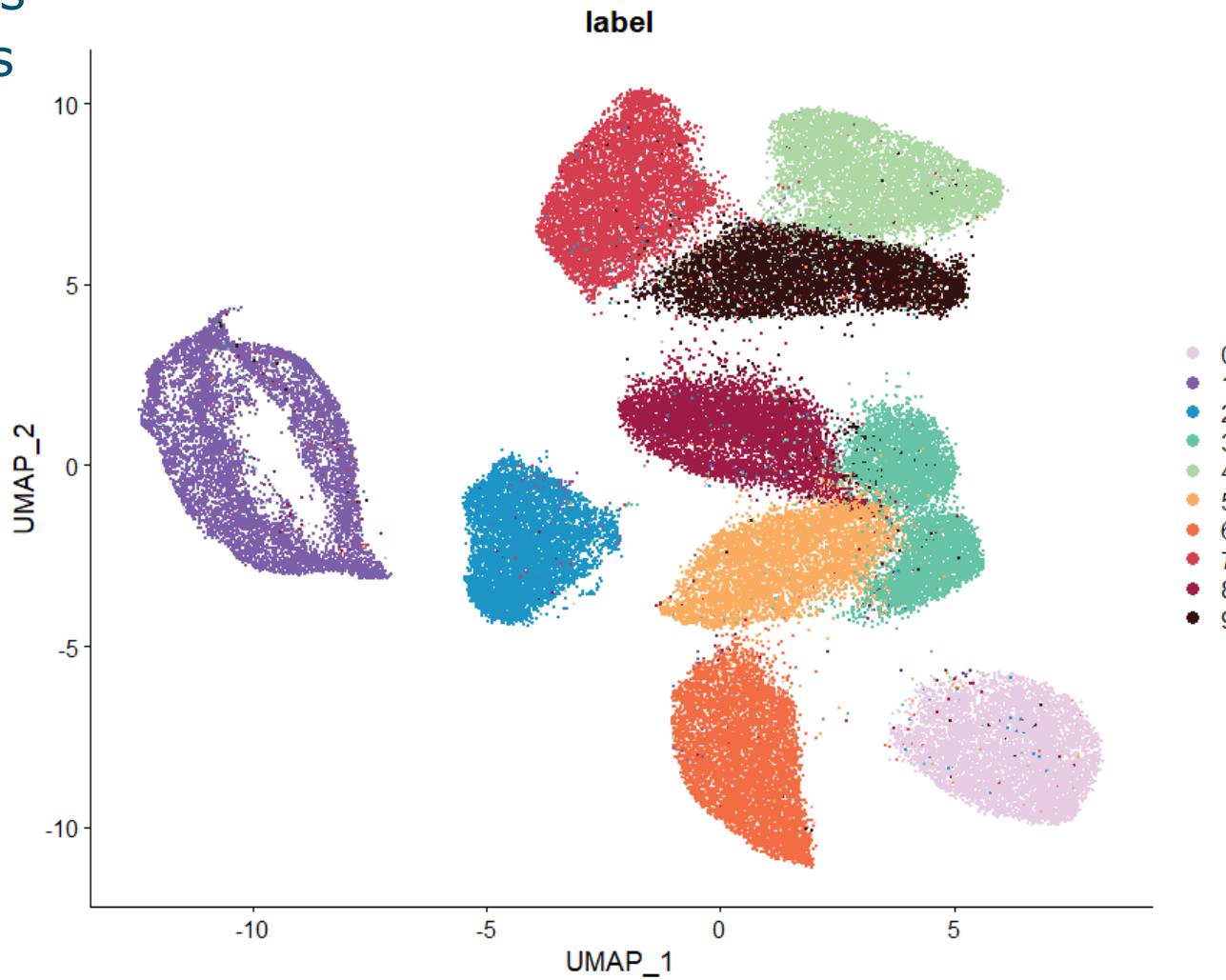
- Graph traversal of k-nearest neighbors
- Identifies well-connected communities
- Higher resolution → more clusters



Clustering

- Leiden algorithm

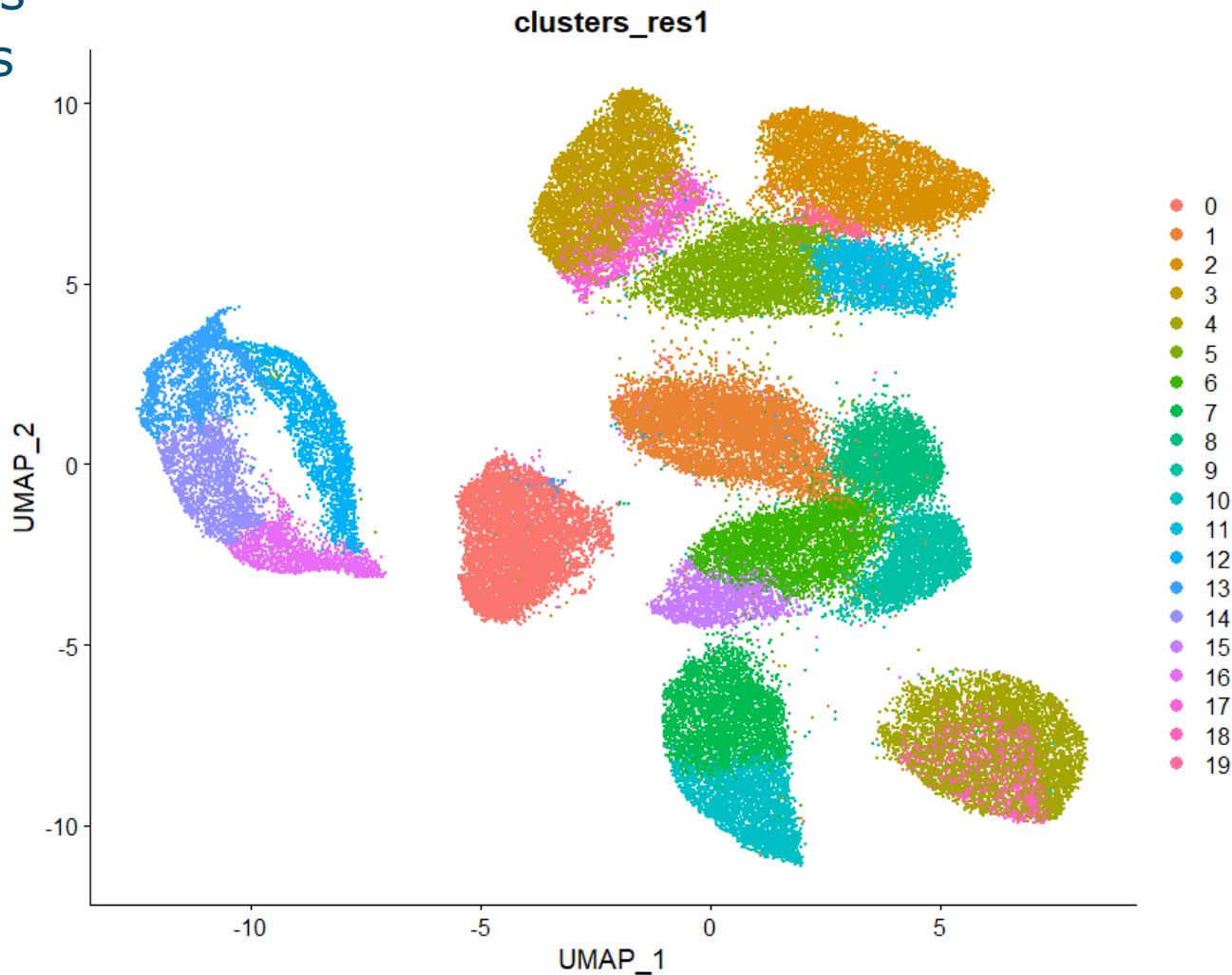
- Graph traversal of k-nearest neighbors
- Identifies well-connected communities
- Higher resolution → more clusters



Clustering

- Leiden algorithm

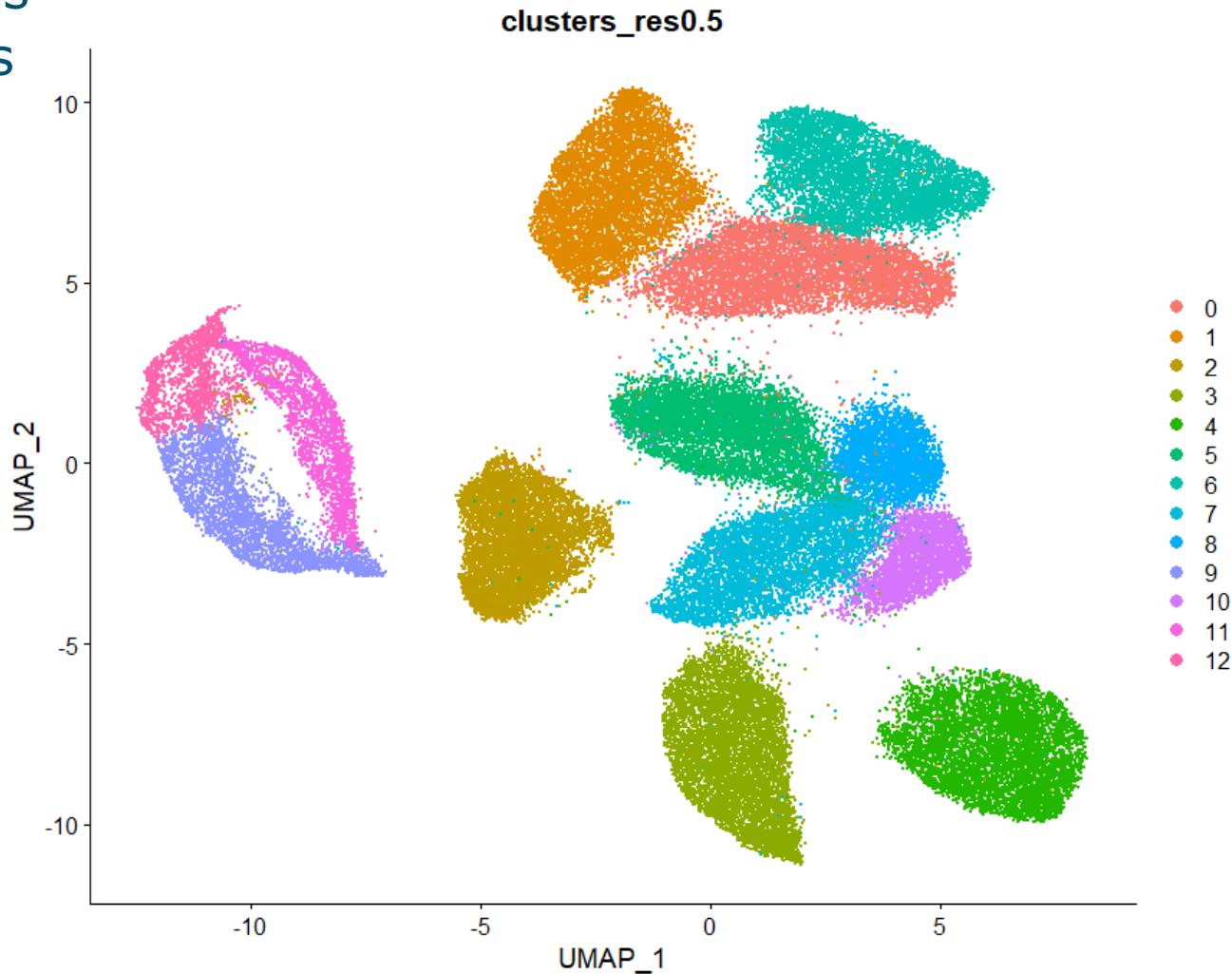
- Graph traversal of k-nearest neighbors
- Identifies well-connected communities
- Higher resolution → more clusters



Clustering

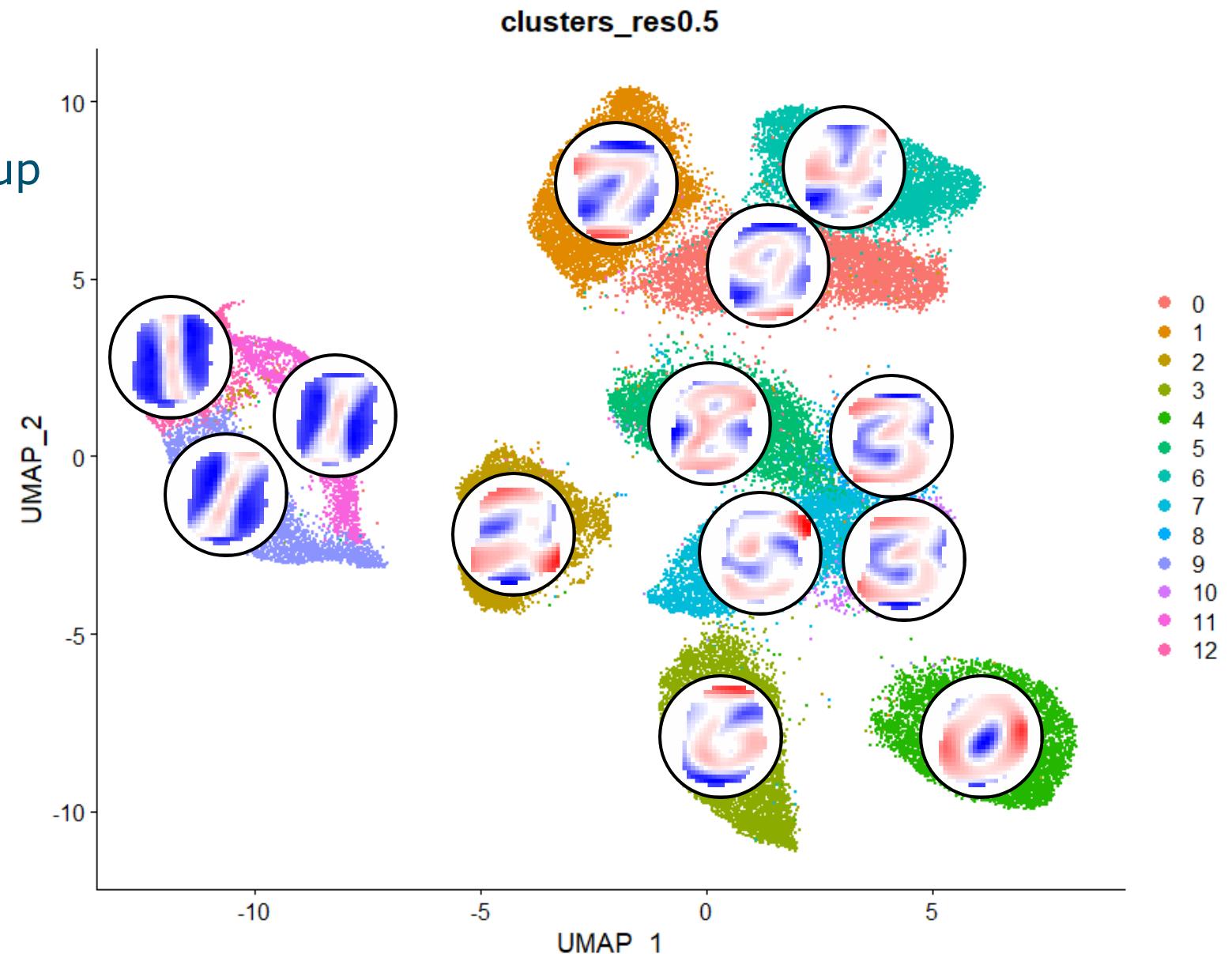
- Leiden algorithm

- Graph traversal of k-nearest neighbors
- Identifies well-connected communities
- Higher resolution → more clusters



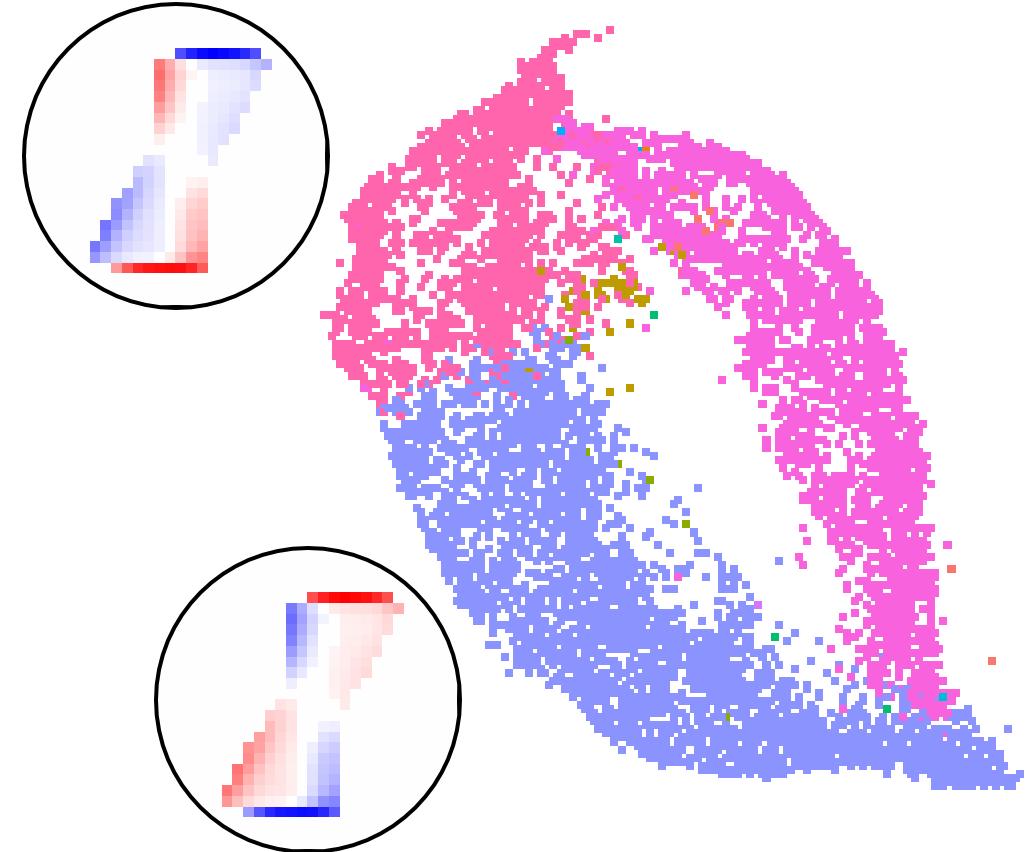
Marker genes

- One vs. all
 - Red – higher in cluster
 - Blue – higher in outgroup



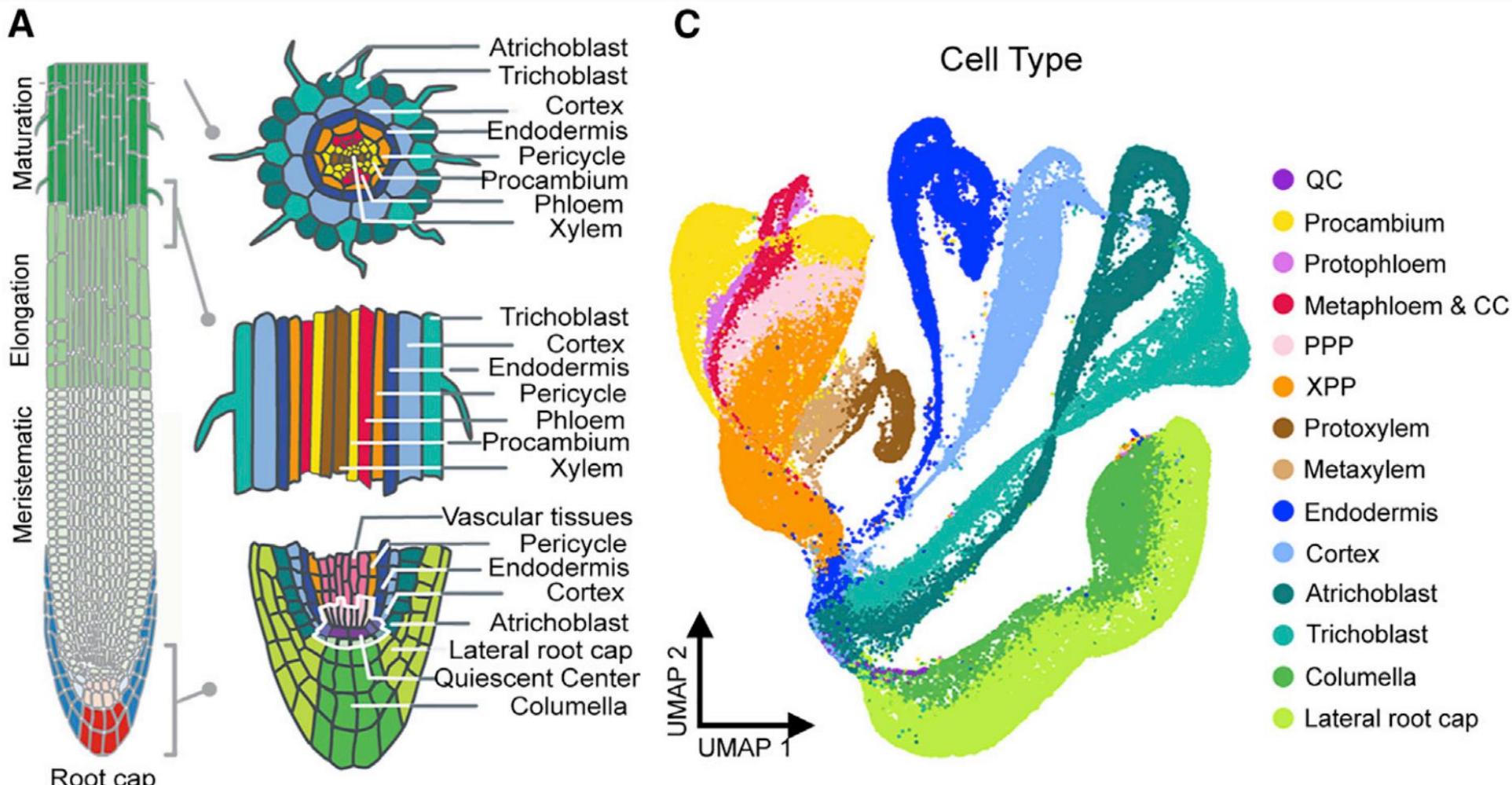
Marker genes

- One cluster vs. another
 - All other clusters are irrelevant
 - $A \text{ vs. } B = -(B \text{ vs. } A)$



Pseudotime trajectories

- Developing tissues show more **continuous** variation



Pseudotime trajectories

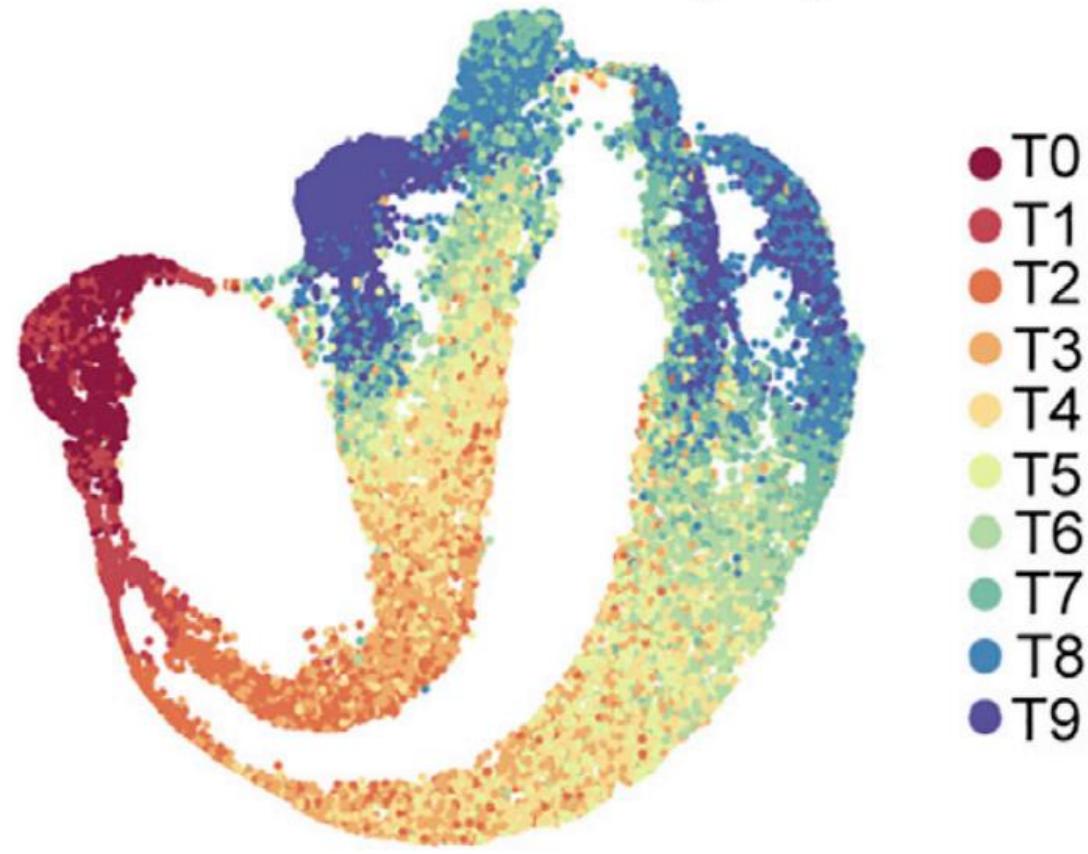
B

Developmental stage

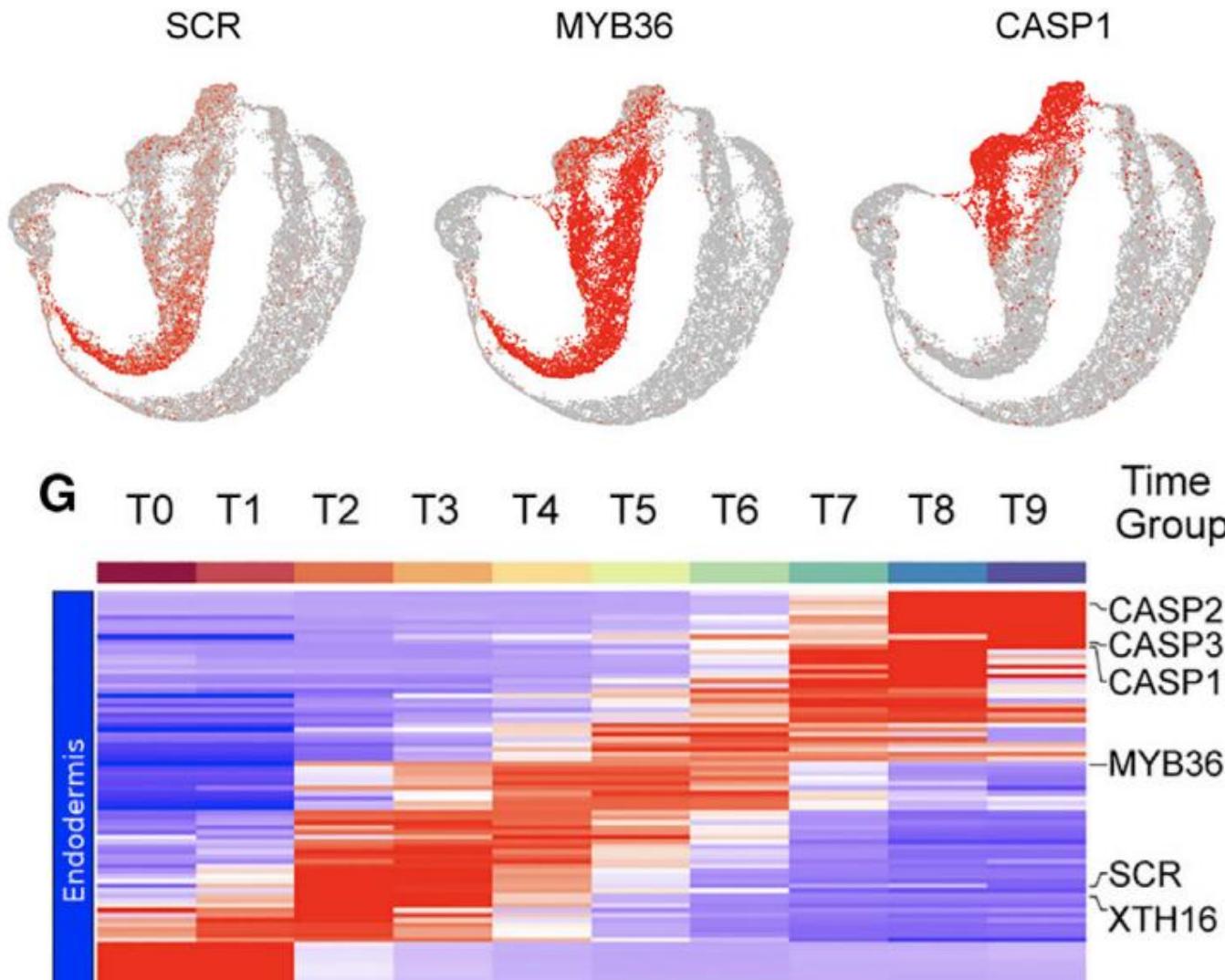


C

Pseudotime estimation
consensus time group

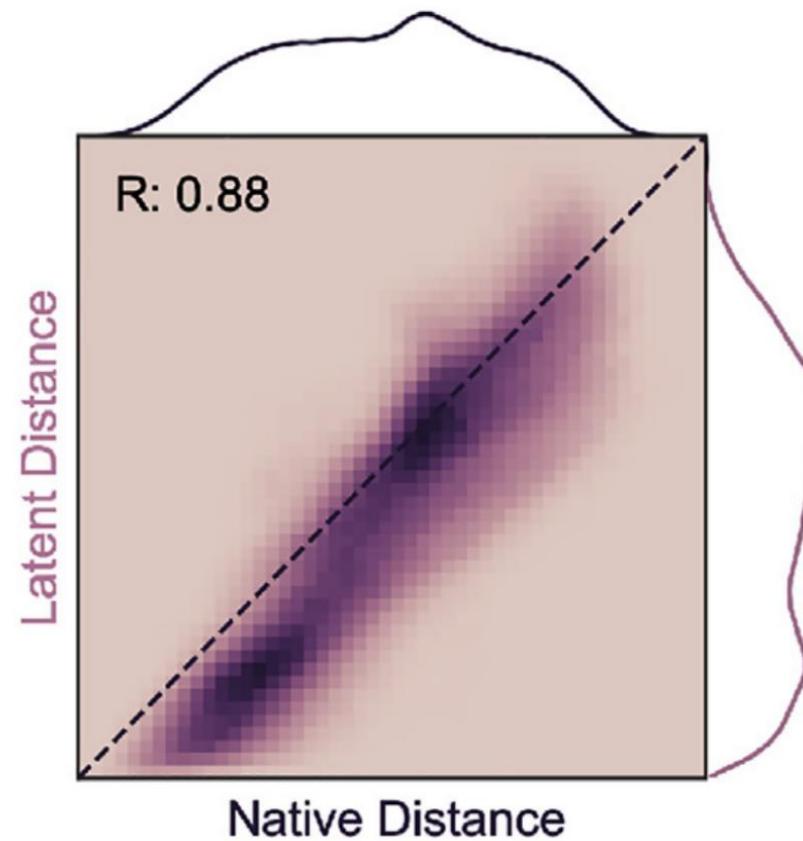
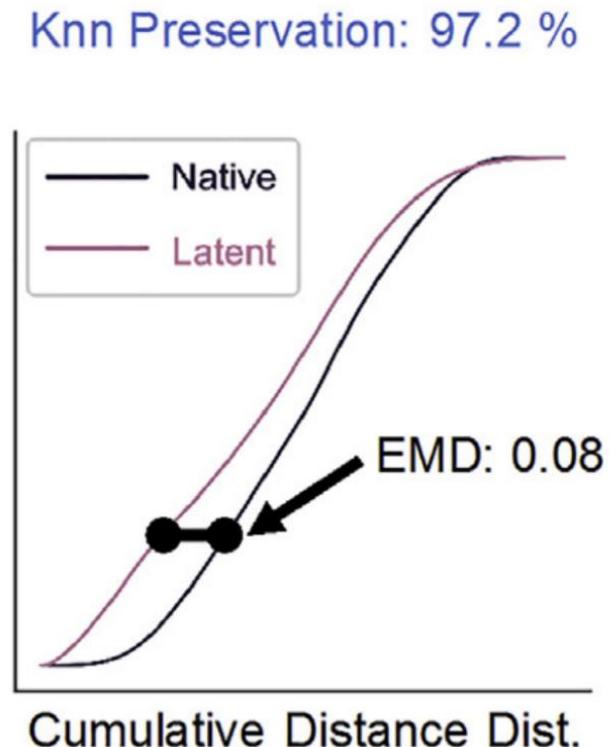


Pseudotime trajectories



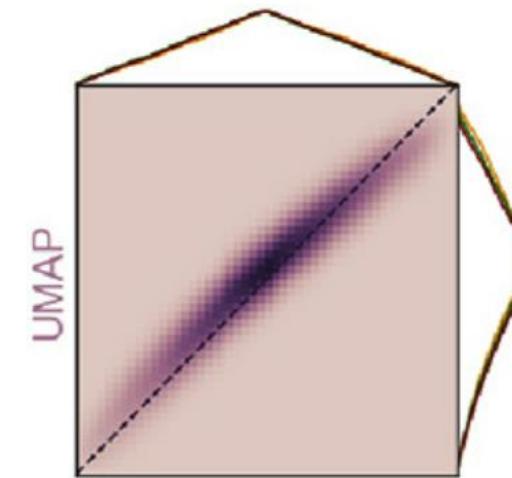
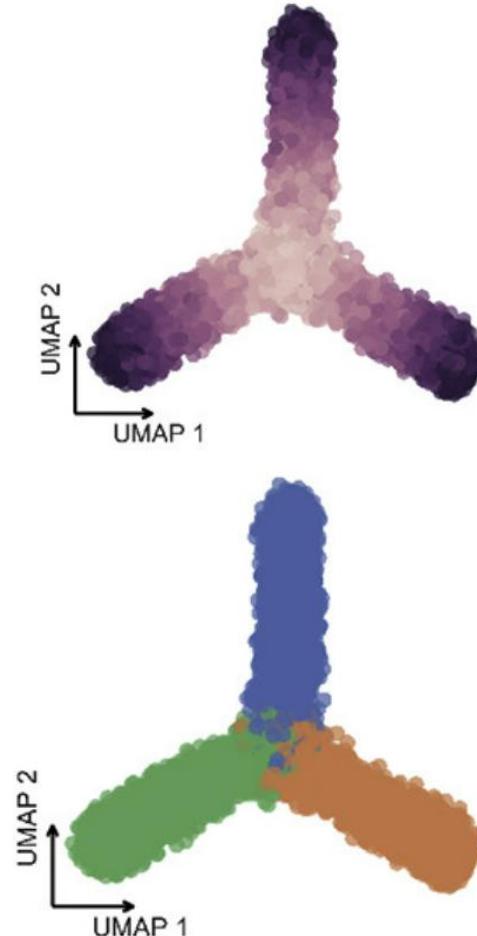
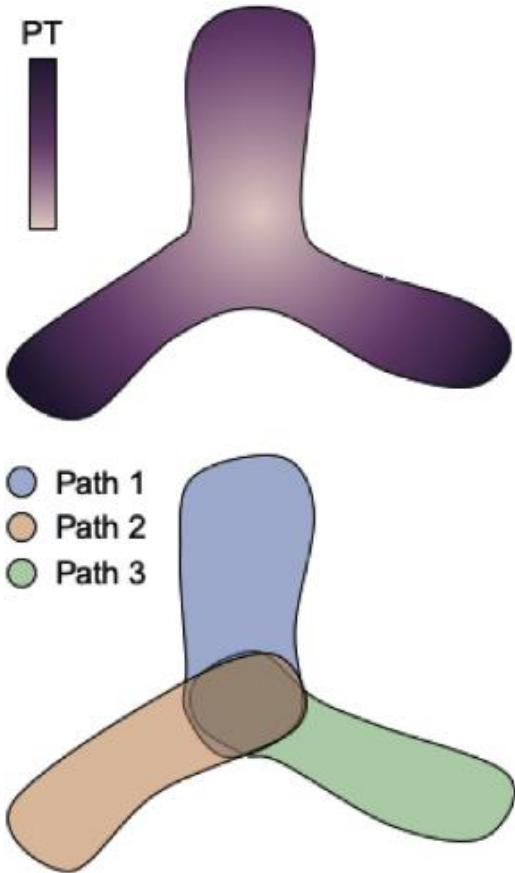
Pseudotime trajectories

- But how can we know if a trajectory is accurate?
 - Calibrate against known marker genes
 - Measure distortions of “native” distances in the reduced space



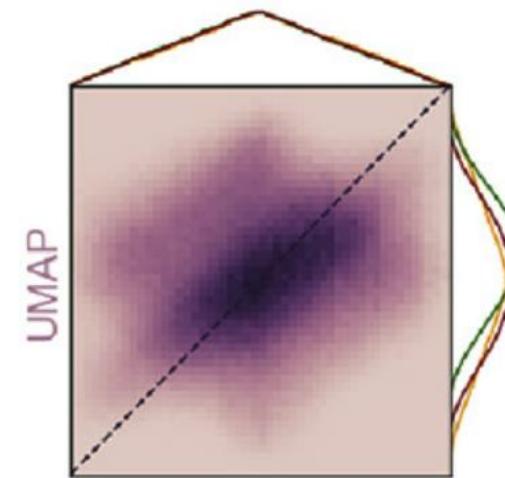
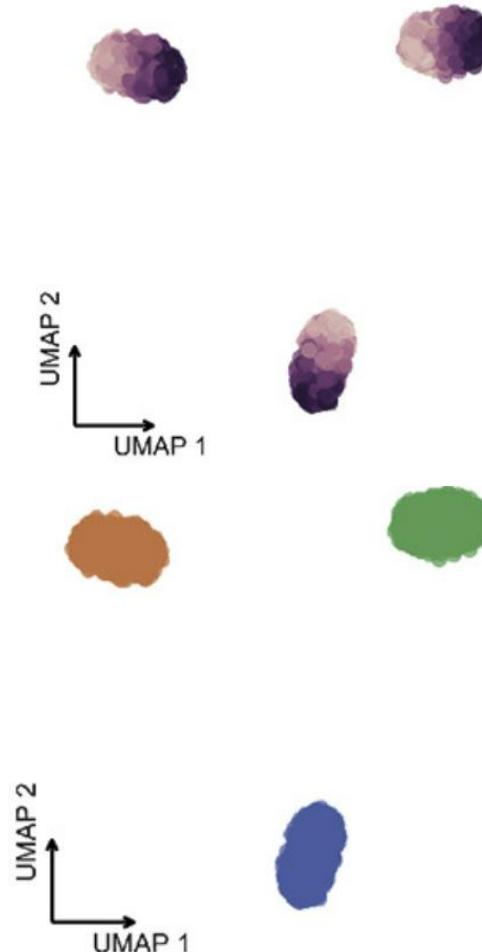
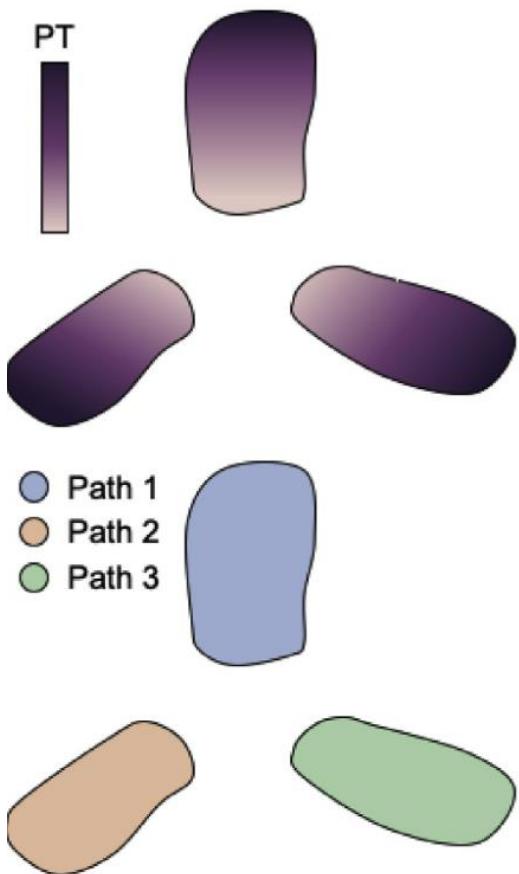
Pseudotime trajectories

- UMAP excels when data are continuous and uniform



Pseudotime trajectories

⚠️ UMAP fails to globally orient discontinuous data



Pseudotime trajectories

⚠ Assumes a continuous path between cell types

- Incomplete data creates bad embeddings
- Not all heterogeneity is developmental!

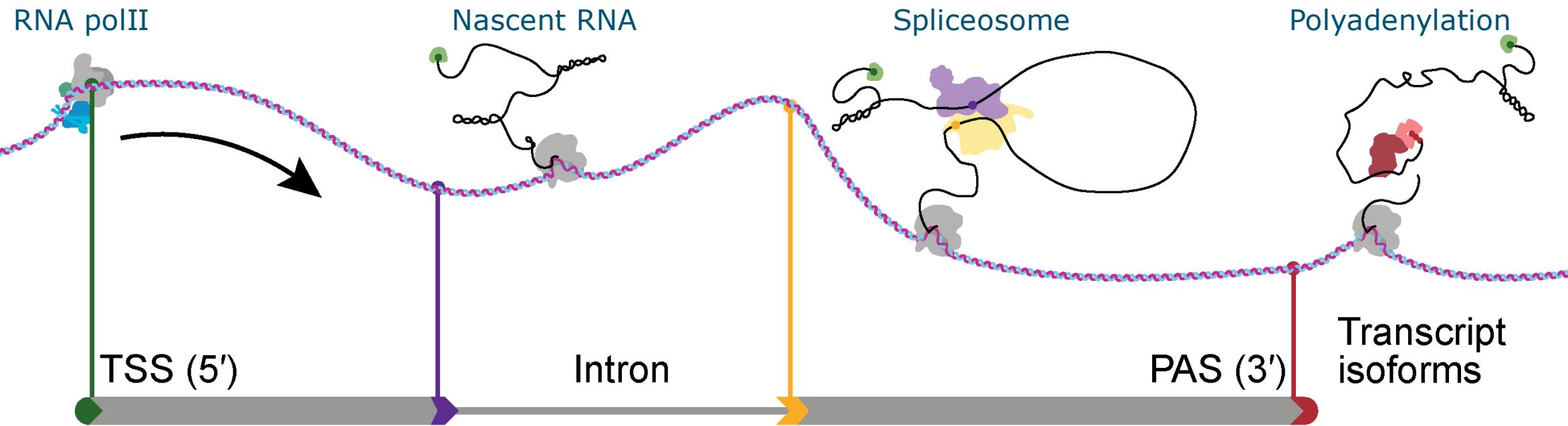
⚠ Operates in reduced dimensions

- 2D projections can have massive distortion!
- Need to be “anchored” in reality
- DR methods (UMAP, t-SNE) are not globally coherent

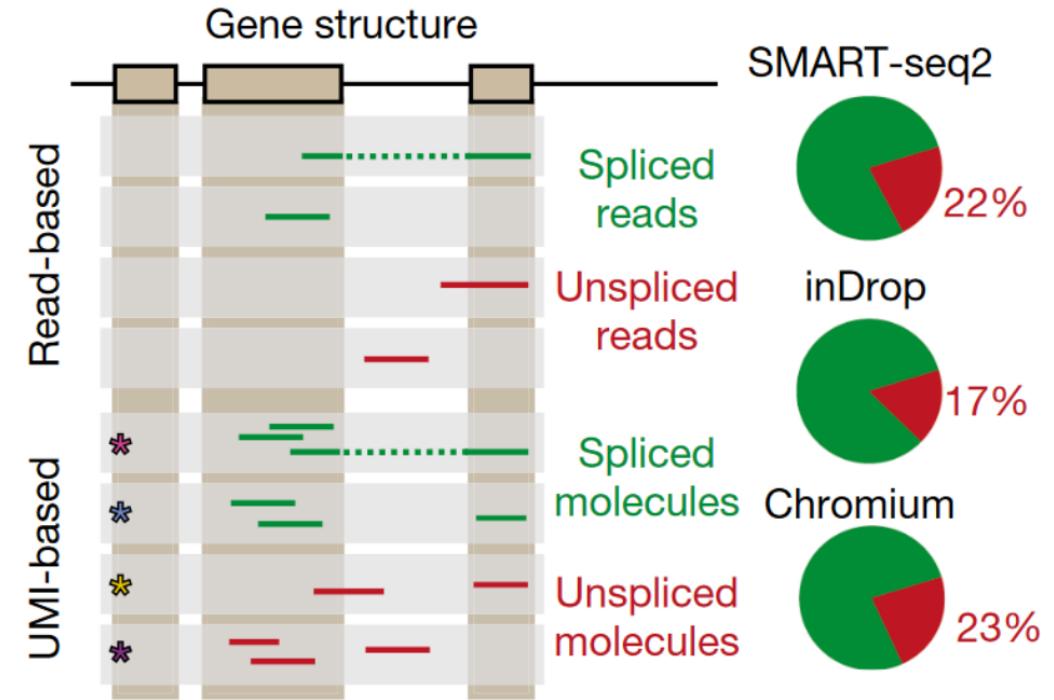
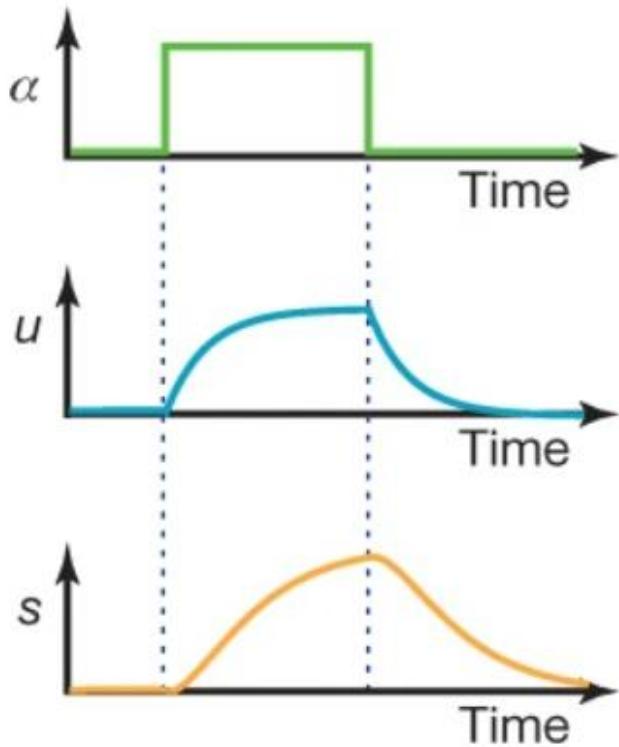
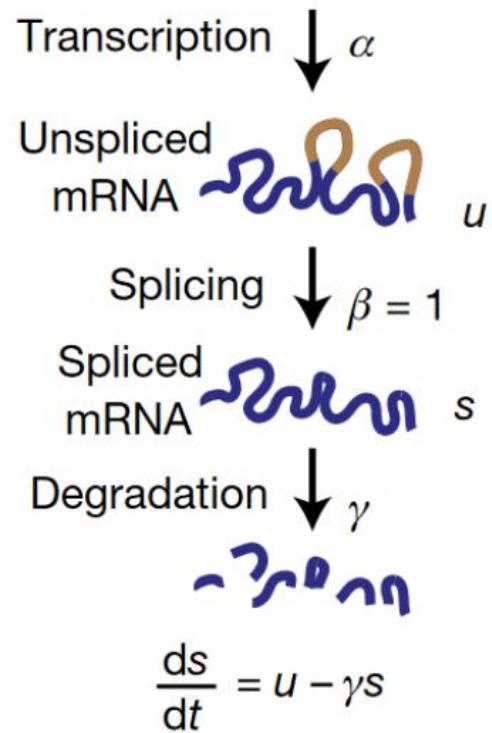
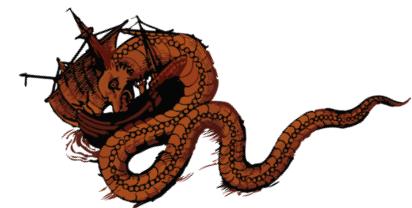
What else can we do with single-cell RNA?



Measuring RNA birth and death



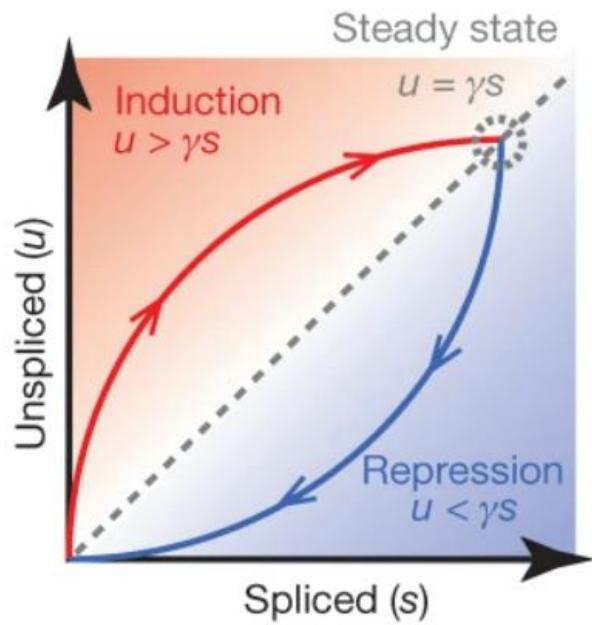
RNA “velocity”: predicting the future?



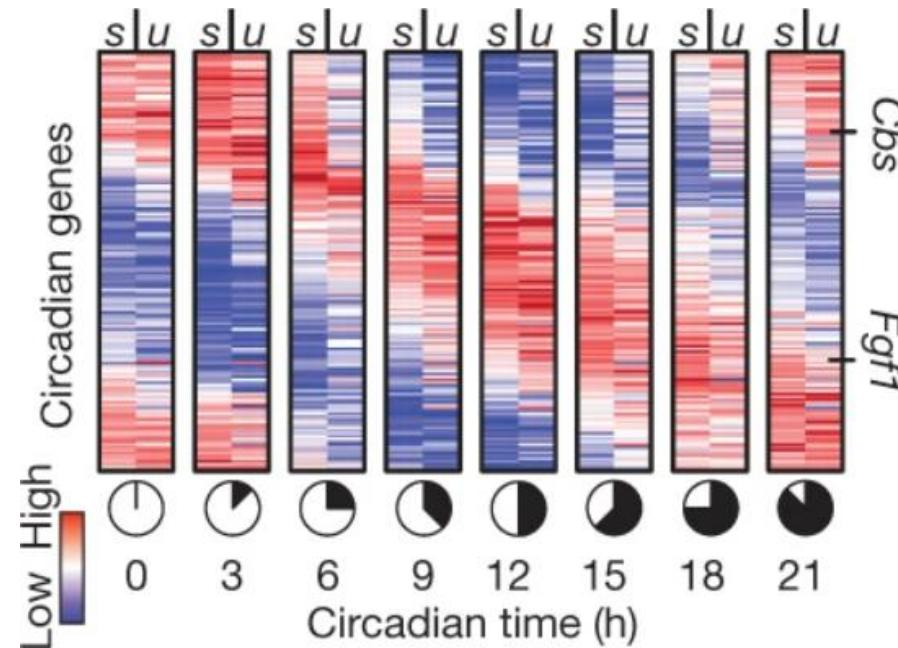
RNA “velocity”: predicting the future?



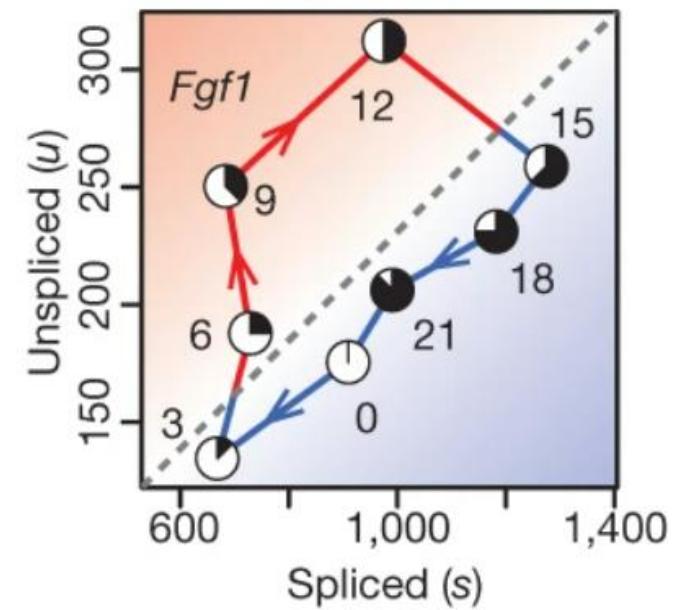
Model RNA kinetics



Unspliced RNA precedes spliced RNA



Kinetics of a single circadian gene

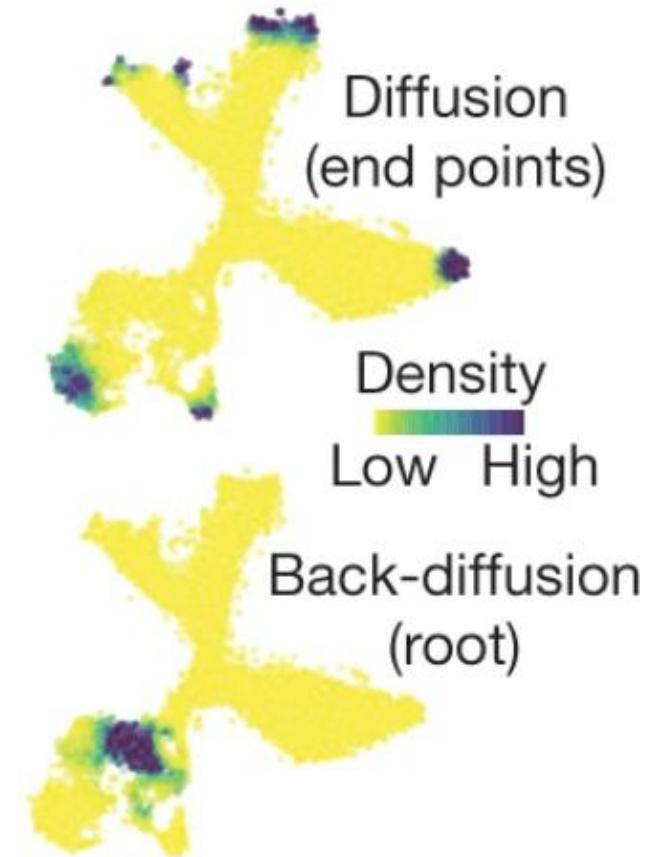
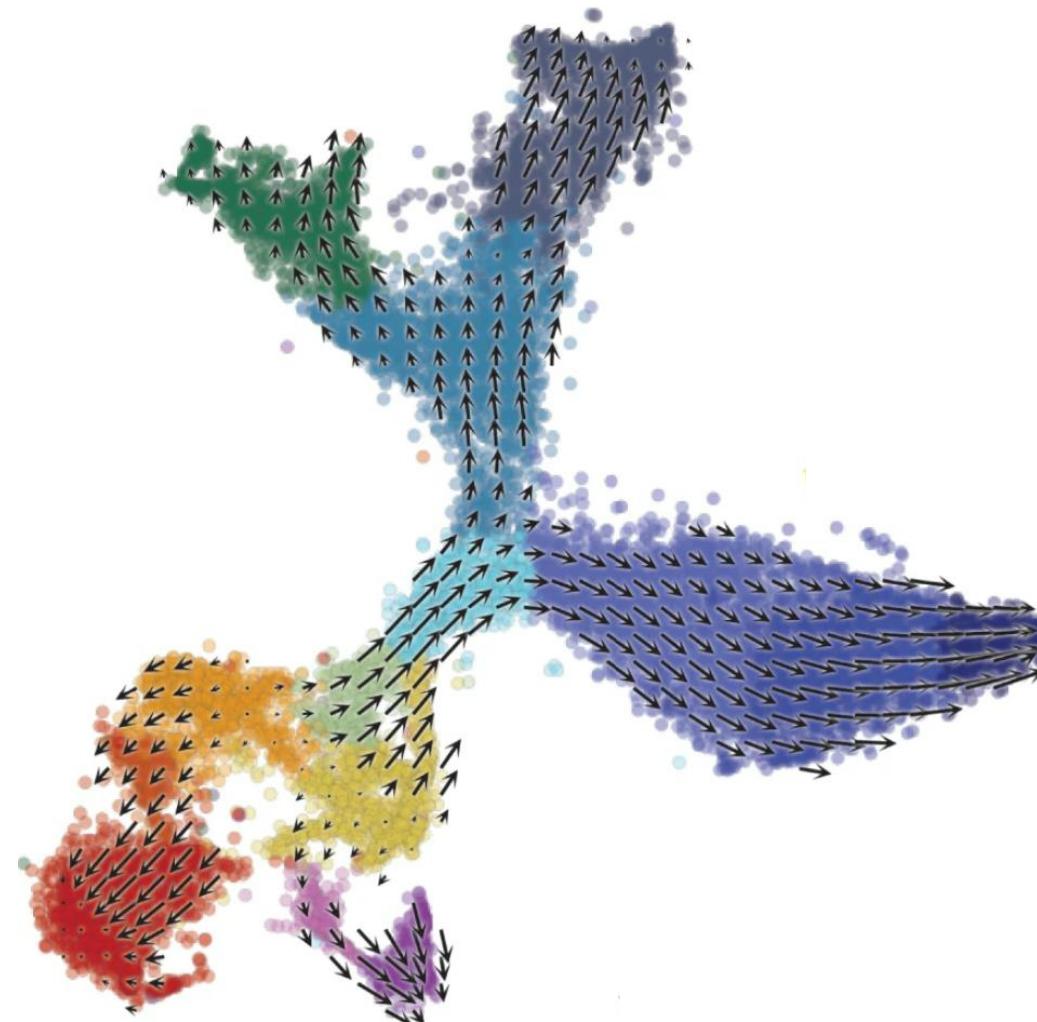
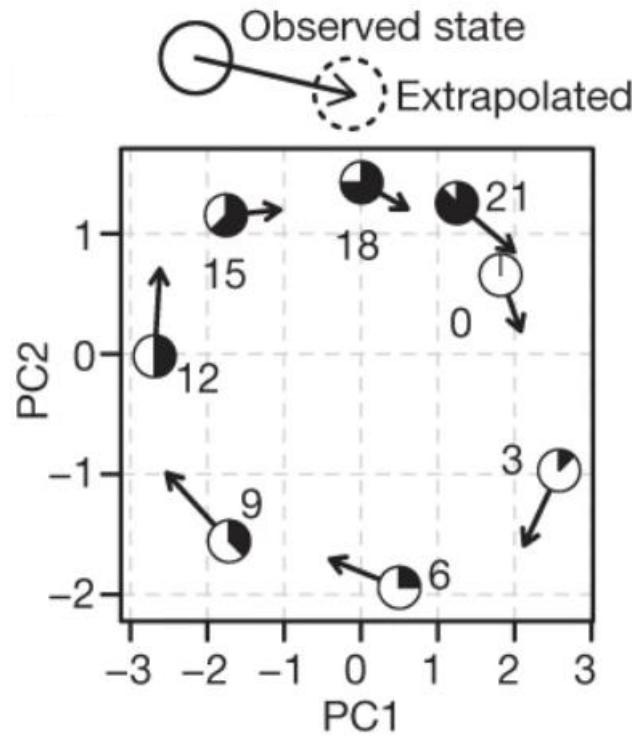


RNA “velocity”: predicting the future?



Developing hippocampus cells

Circadian cycle



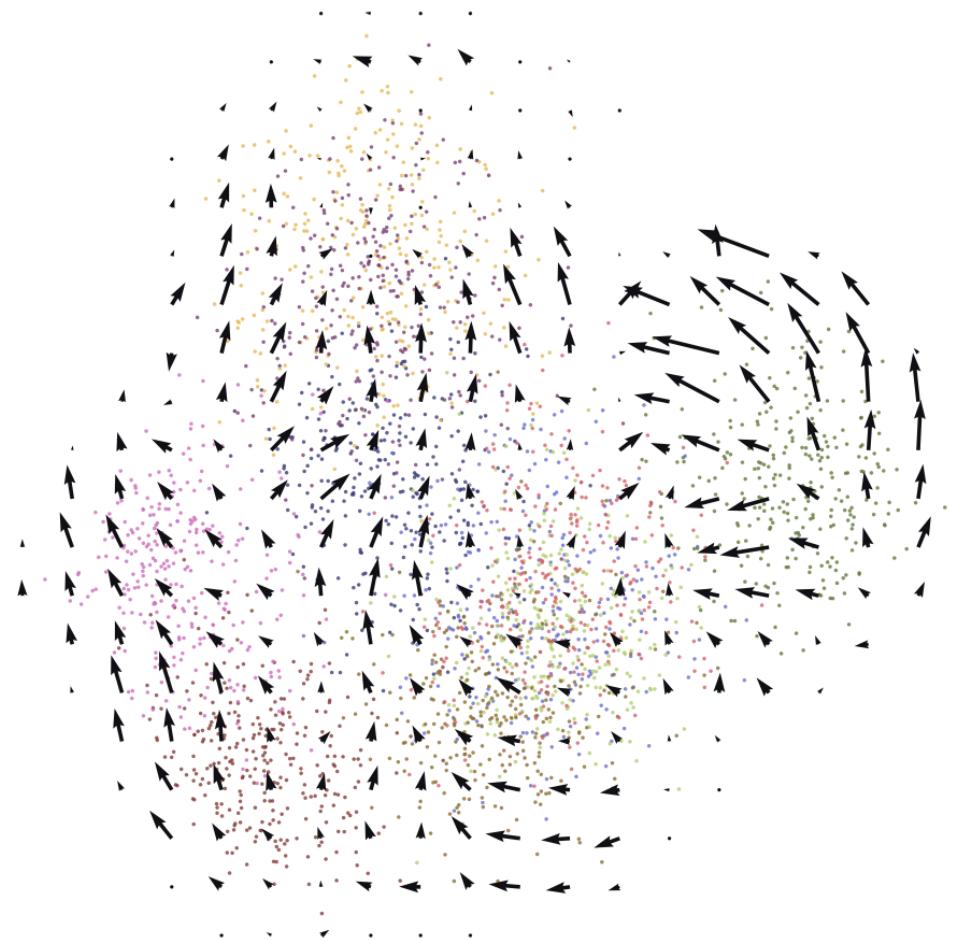
La Manno et al. 2018 *Nature*

RNA “velocity”: predicting the future?

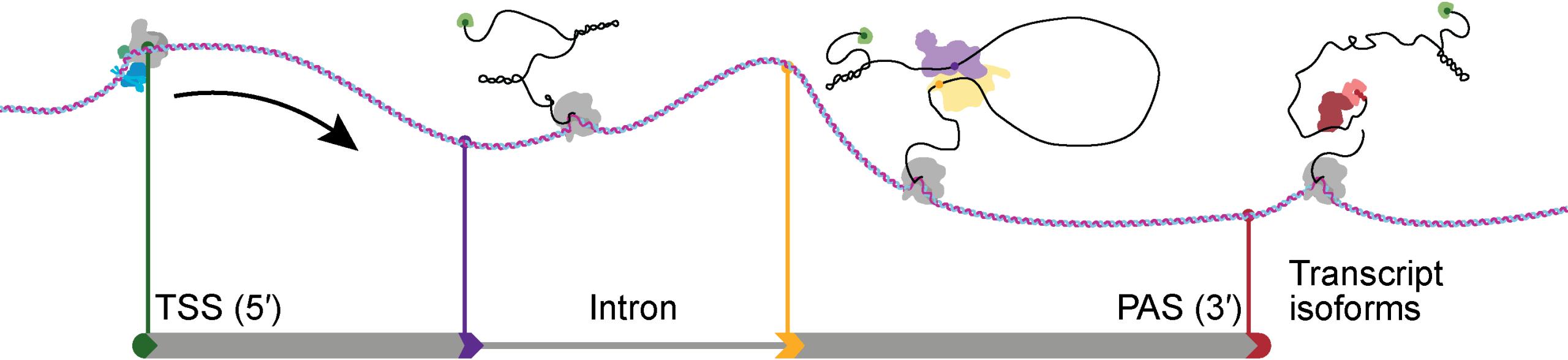


⚠️ Vectors will always be found, even when none exist

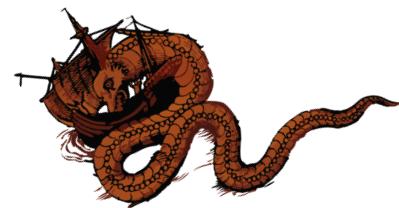
Noisy simulation of 10 discrete cell types



Genes → isoforms?



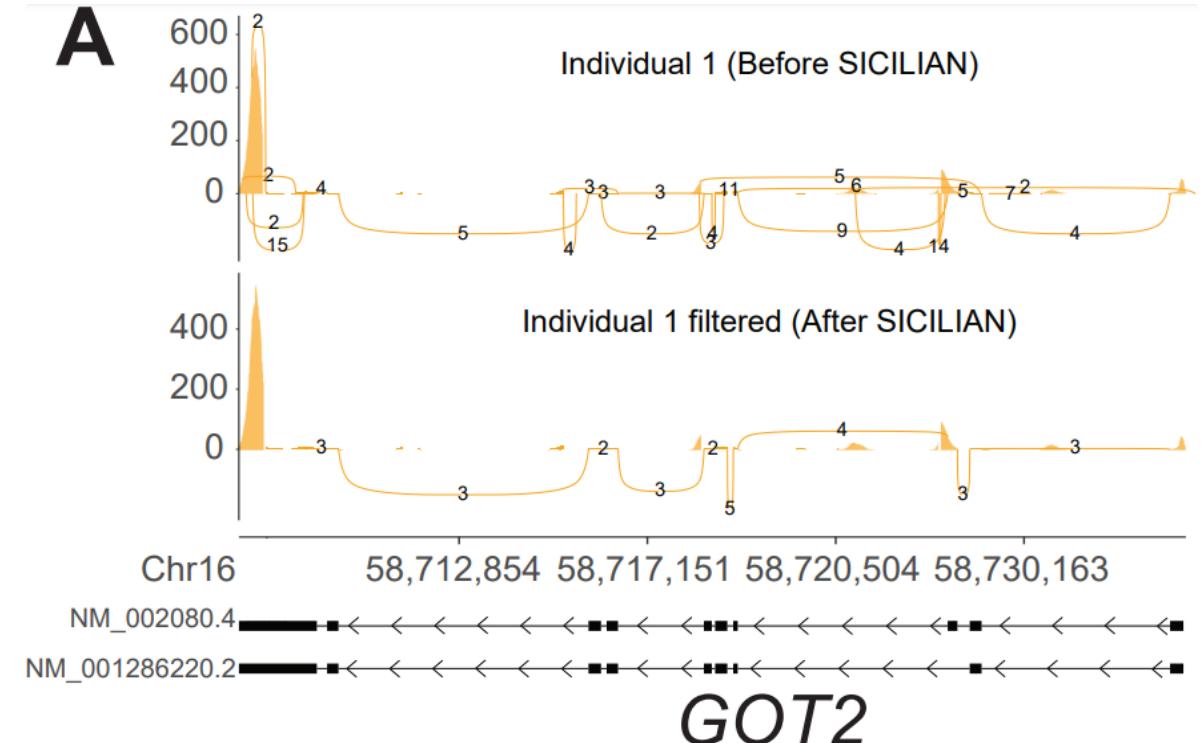
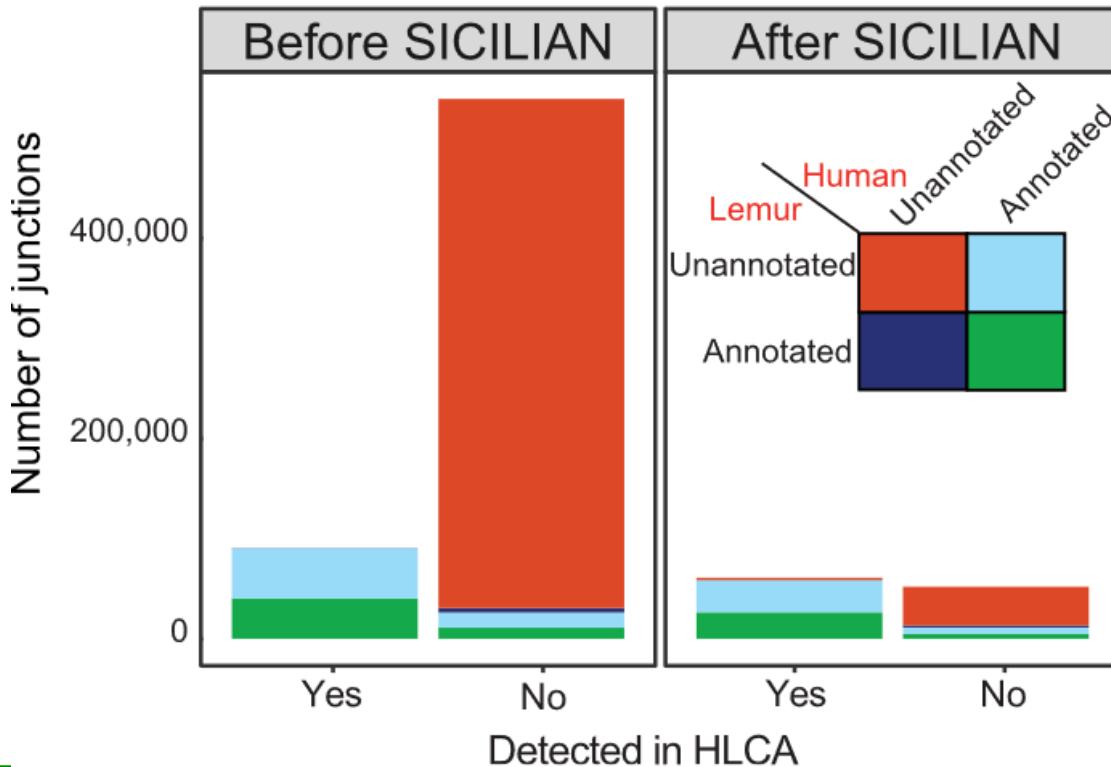
Limits to isoform identification



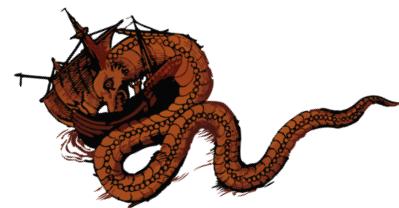
⚠ Single-cell sequencing is noisy

- Gapped alignments yield more false splice junctions than true ones
- SICILIAN filters splice junctions by sequence entropy

⚠ Most scRNA-seq data is heavily biased to 3' ends (10X Genomics)

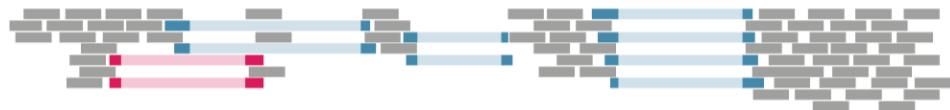


Single-cell isoform assembly?

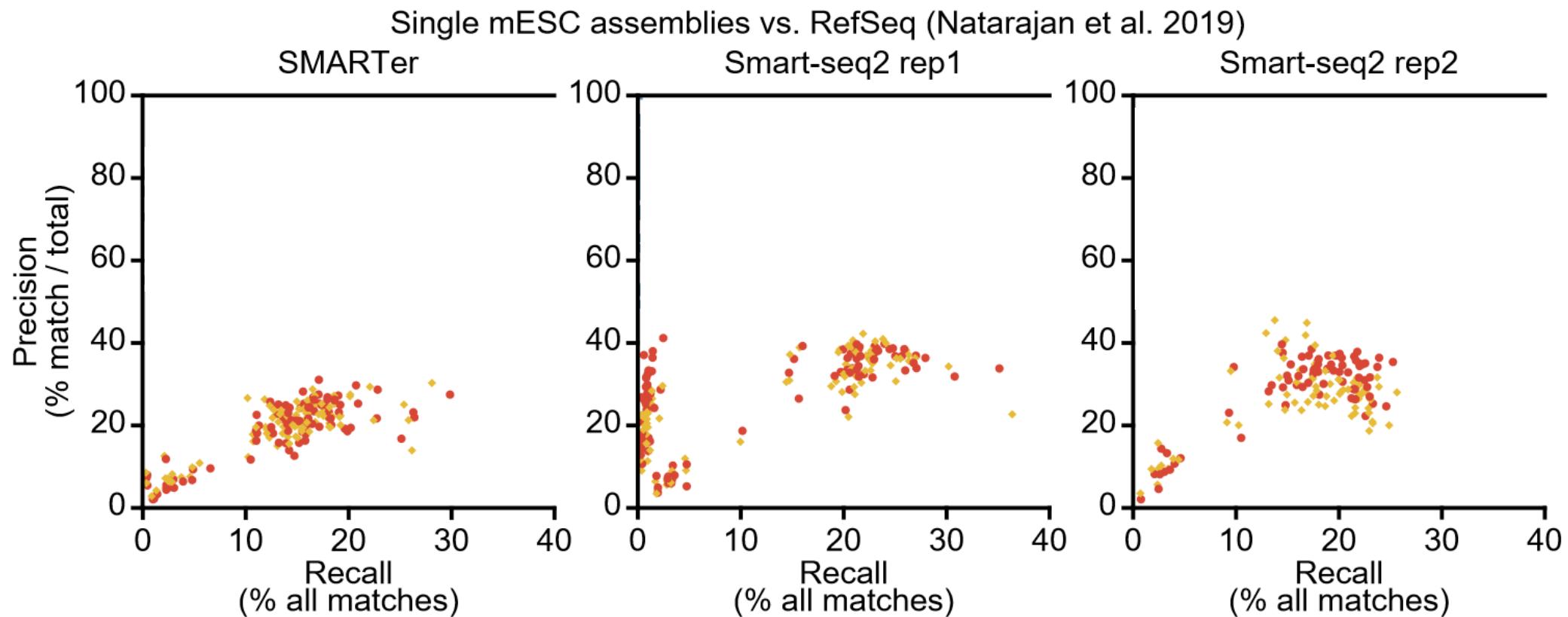


- Full-length single cell sequencing? (Smart-seq)

- Transcript assemblers have abysmal precision
(4 false transcripts per true transcript)



StringTie2
Scallop



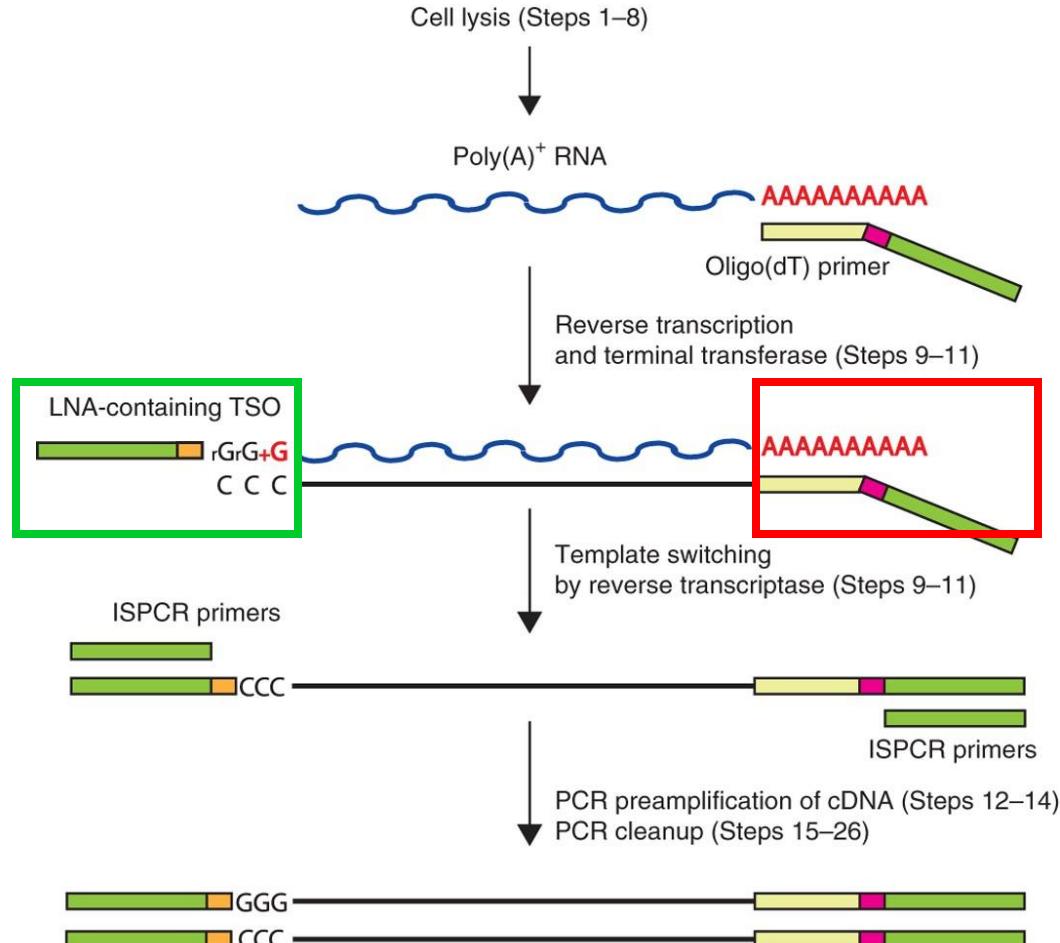


Something is missing...

Smart-seq2 protocol

Template-switching
oligo (TSO)

Poly(A) tail

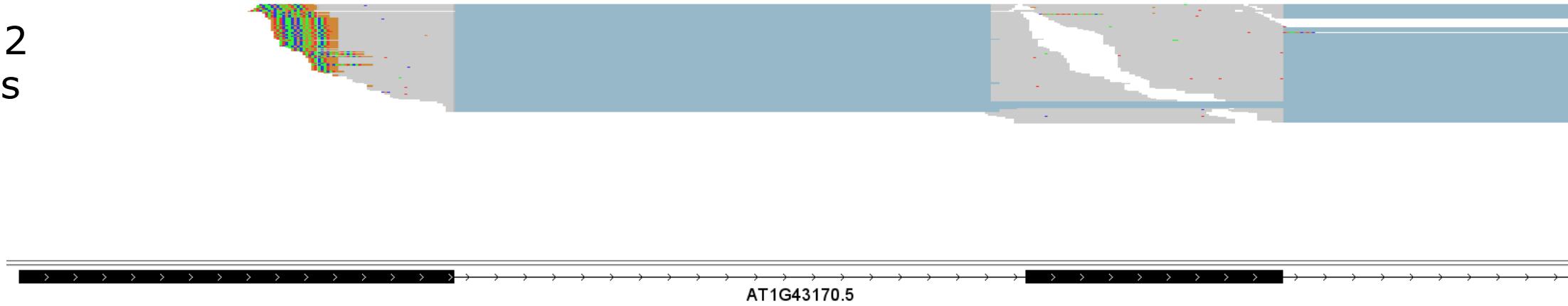


Something is missing...



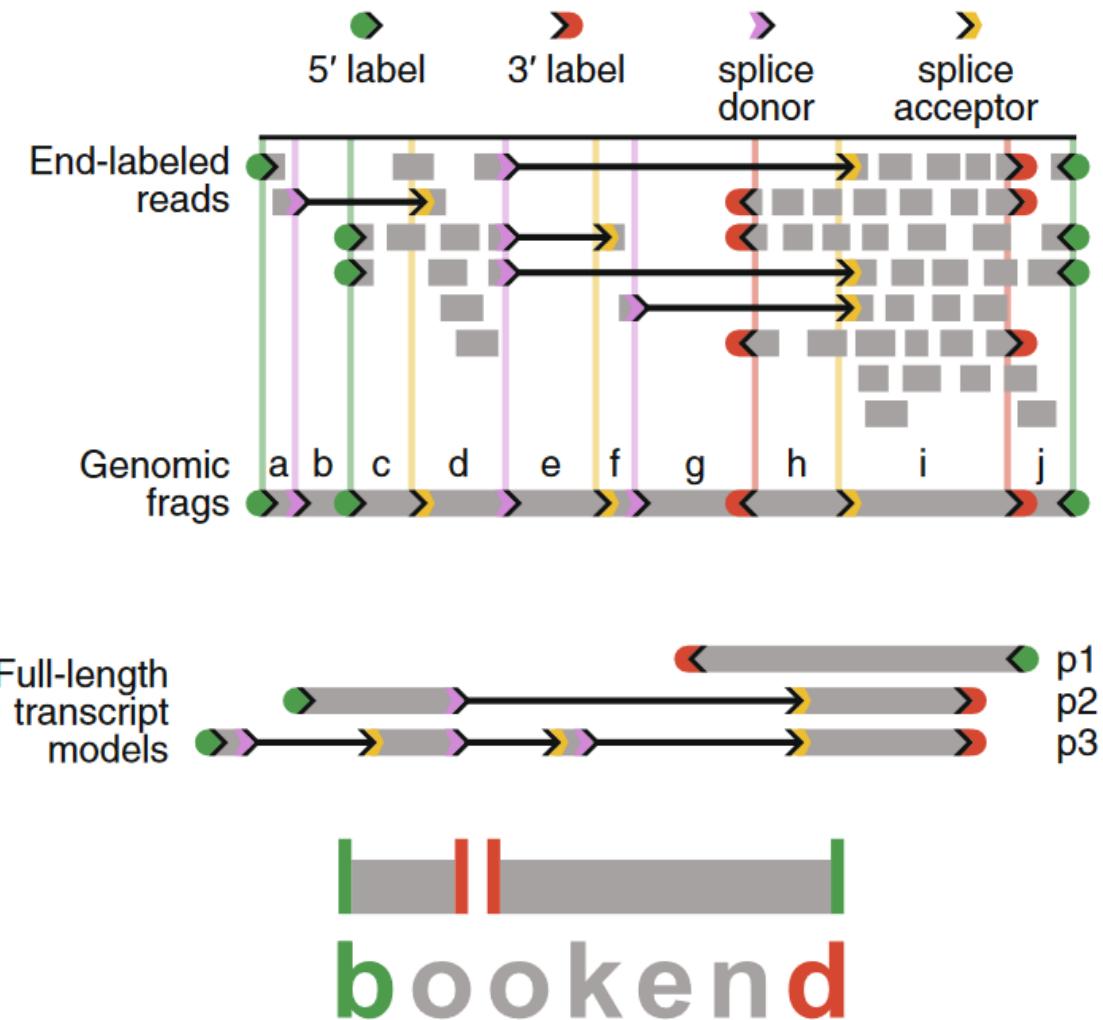
TSO positions

Smart-seq2
alignments

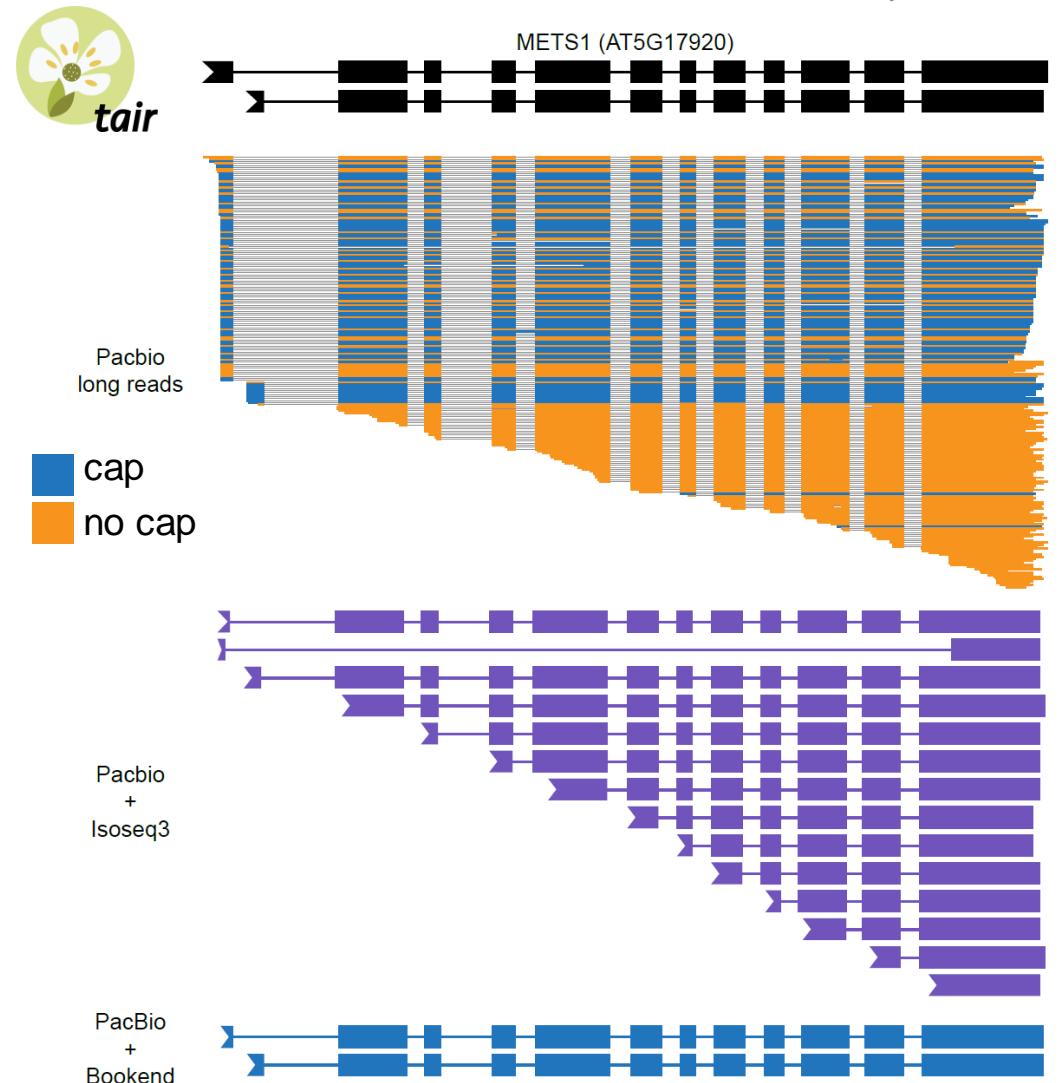
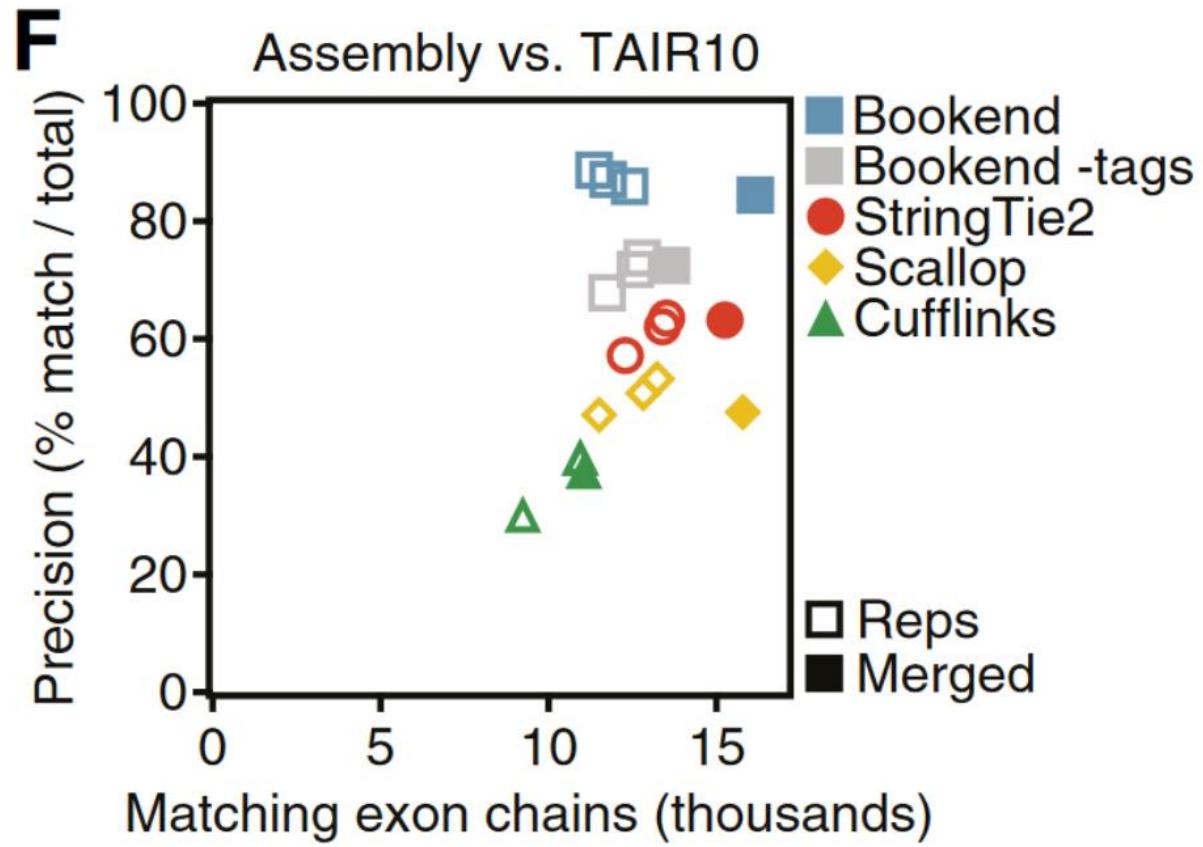


- Could RNA end labels improve transcript assembly?

End-guided assembly



End-guided assembly

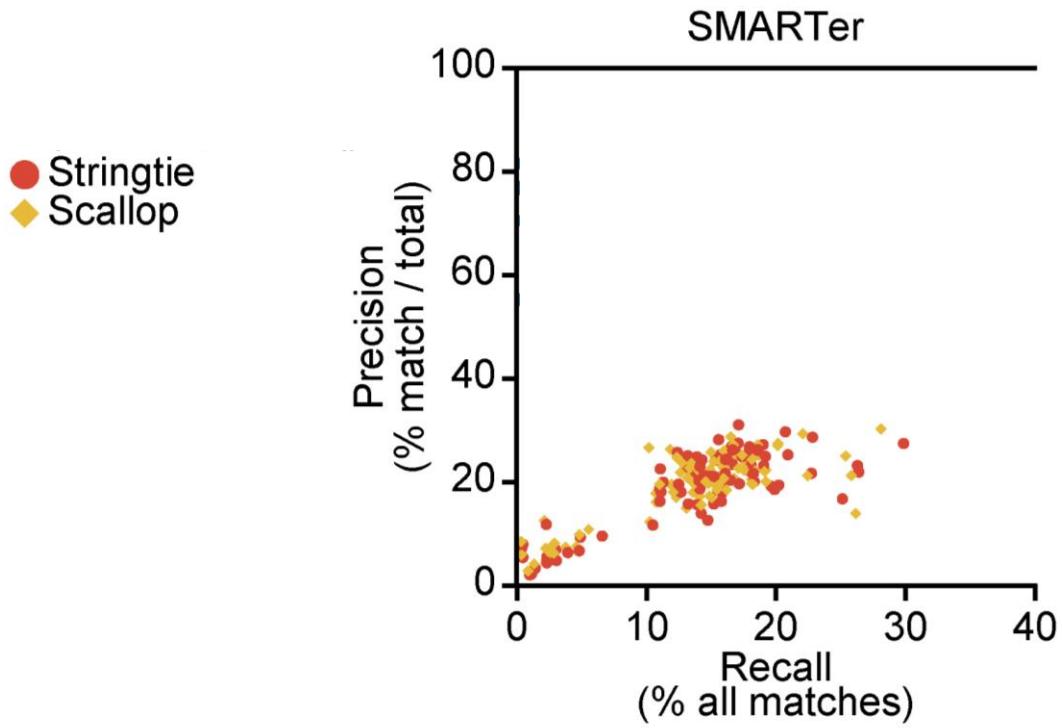


Schon et al. 2022 *Genome Biology*

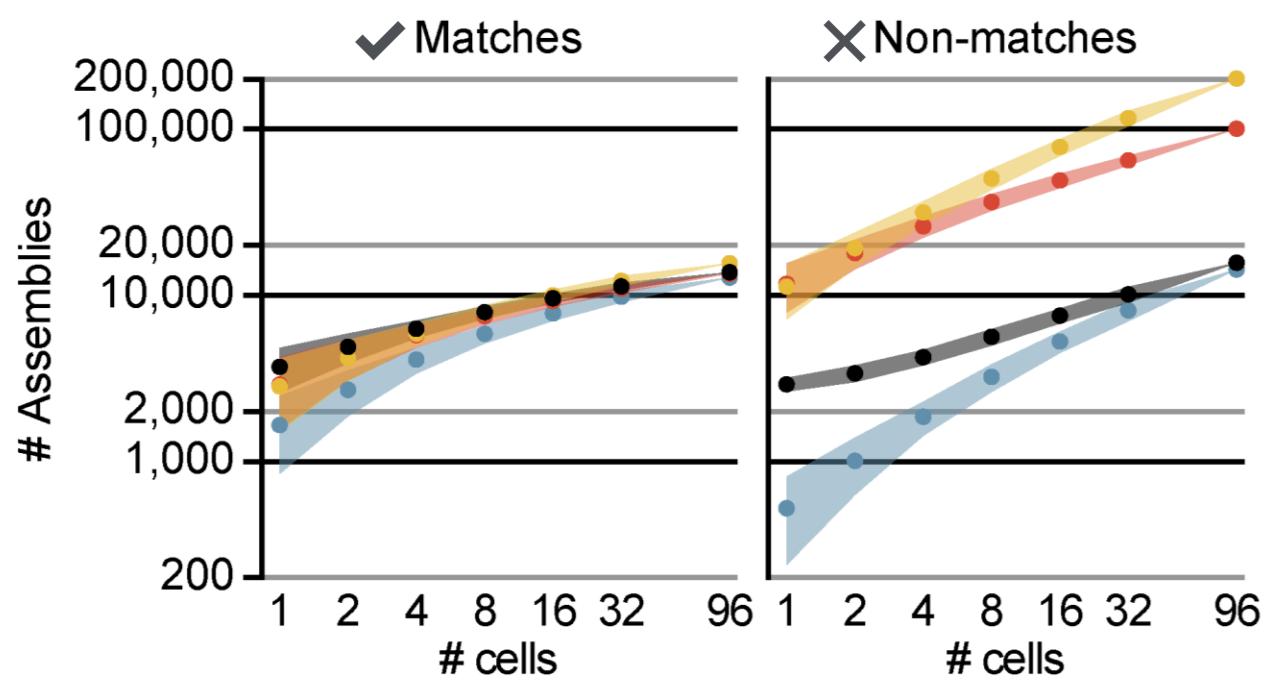
End-guided assembly



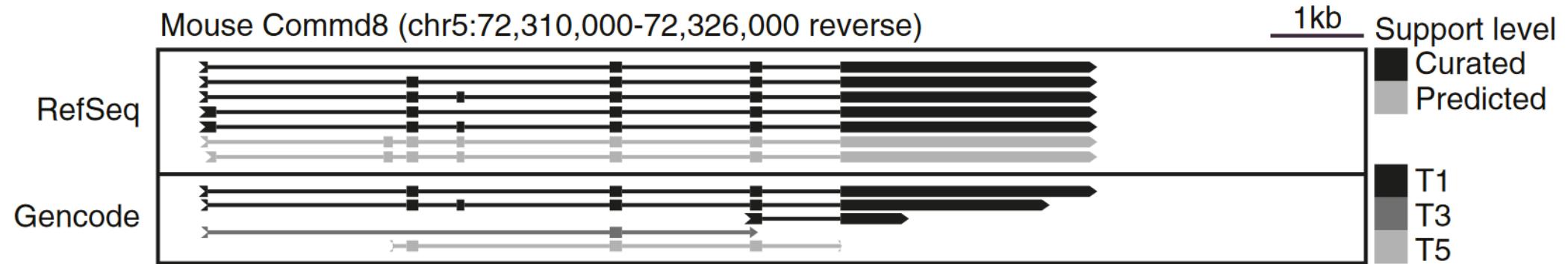
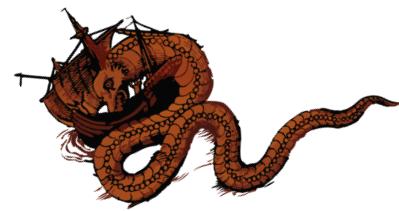
Mouse embryonic stem cell assemblies
(96 single cells)



Multi-cell assemblies
vs. reference (RefSeq)



Meta-assembly: each cell is unique



mESC SLIC-CAGE

mESC 3P-seq

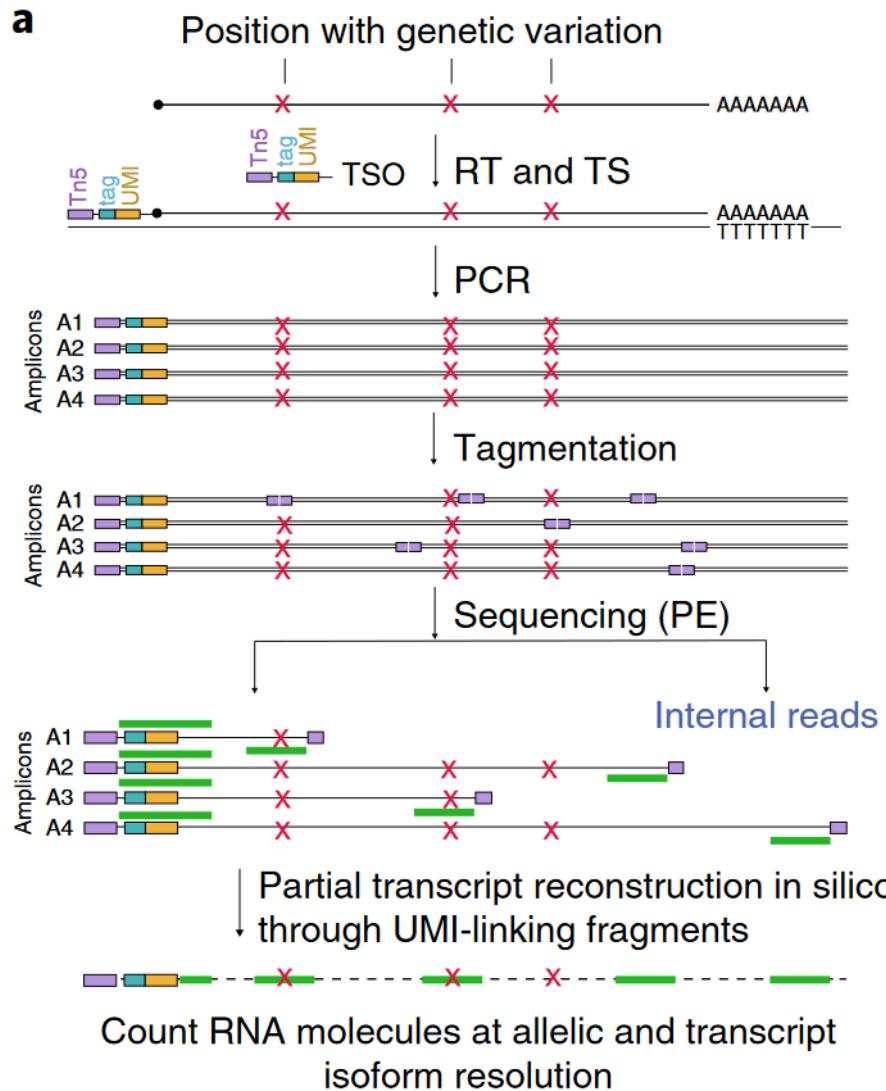
mESC
single-cell
SMARTer
condensed
(96 cells)

Max depth
141

242

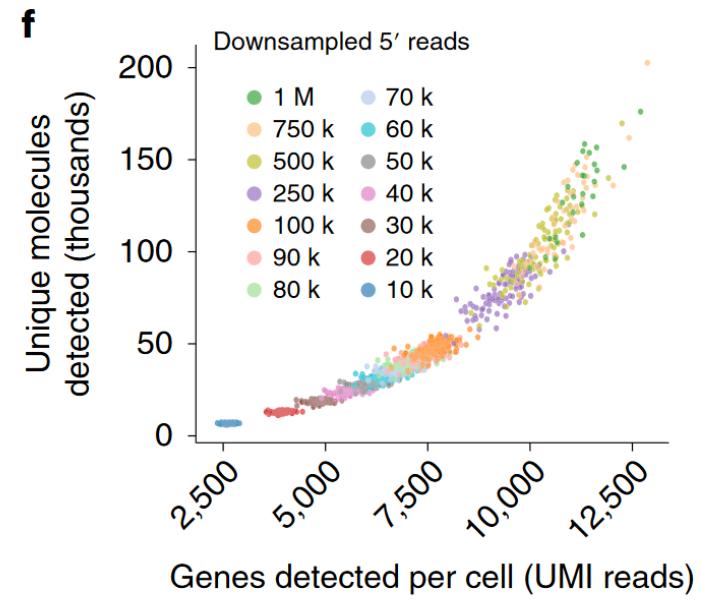
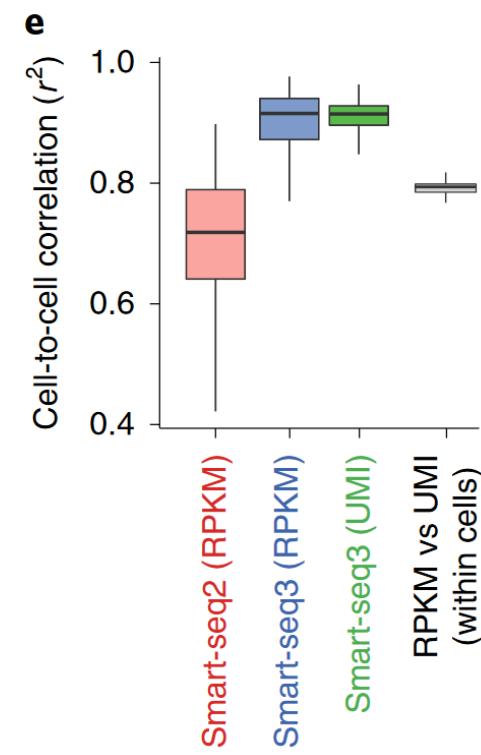
- Start Tag
- Cap Tag
- Start + End
- Cap + End
- End Tag
- Unlabeled

Partial molecule reconstruction

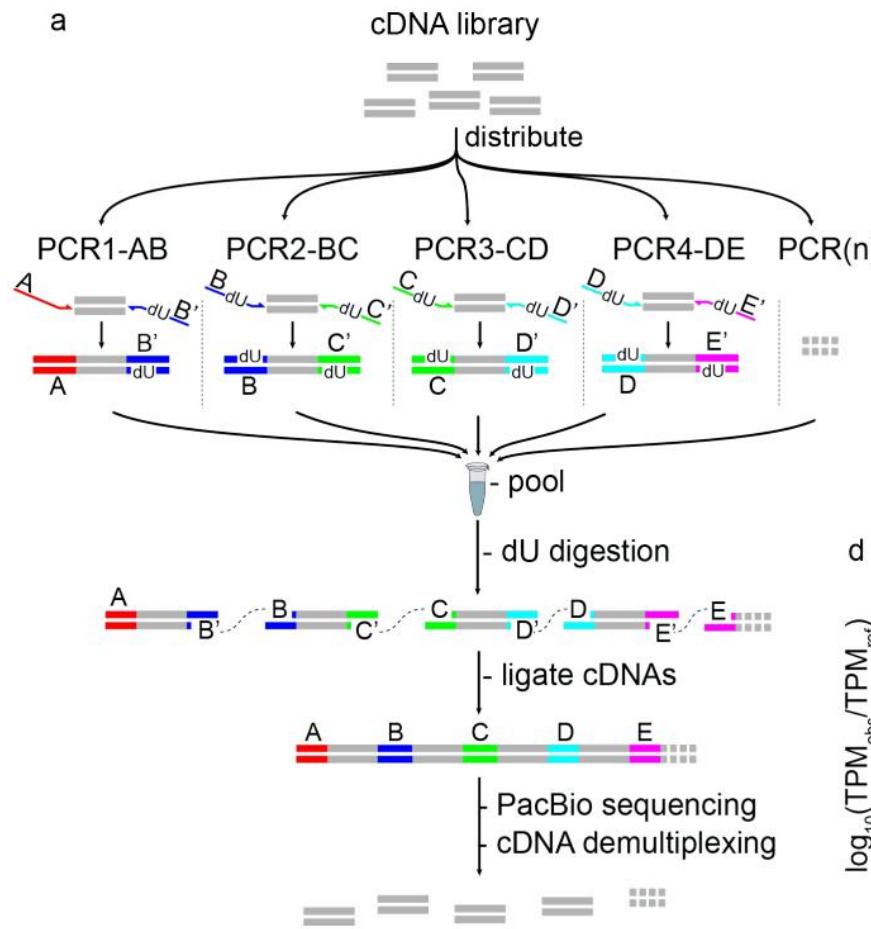


Smart-seq3

- Tn5-UMI-TSO → 20-30% 5' PE reads
- Can capture >100,000 UMIs/cell
~25% of smFISH!



Long reads in single cells



short-read



d
long-read
(MAS-ISO-seq)



Conclusions



- Development is a cellular process
- We now have many methods for sequencing many RNAs x many cells
- Downstream analysis is both art and science
 - Well-grounded in reality:
marker genes
timepoints
validation
 - “Number of clusters” is (partly) subjective
 - Continuous methods are possible, but prone to distortions
- Not all data is created equal!
 - Make sure the protocol is sufficient to answer your question
 - Isoform analysis requires full-length sequencing:
Smart-seq (2,3,xpress)
MAS-ISO-SEQ

Acknowledgements



(former) Nodine Lab

Max Kellner
Alexandra Plotnikova
Falko Hofmann
Stefan Lutzmayer
Ping Kao
Ranjith Papareddy
Katalin Paldi
Magda Mosiolek
Balaji Enugutti

Thesis Advisory Committee

Stefan Ameres
Mike Axtell
Andrea Pauli

Thesis Reviewers

Kelly Swarts
Dolf Weijers

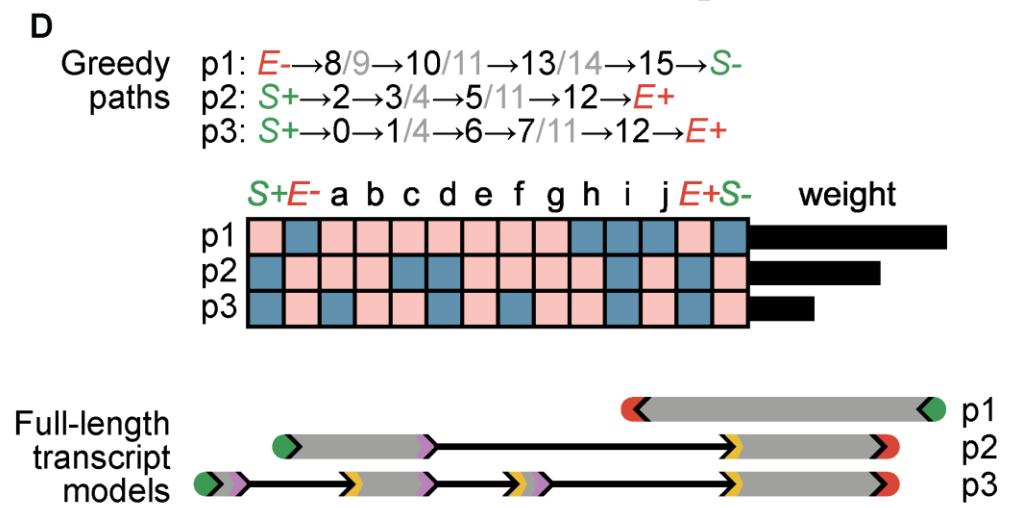
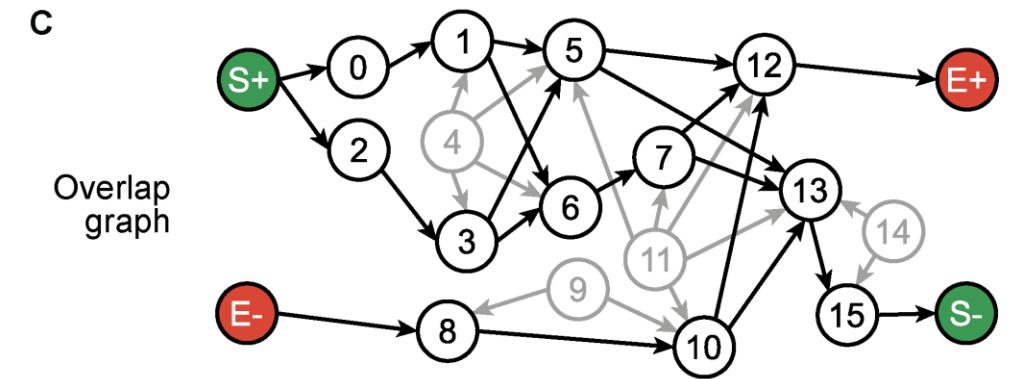
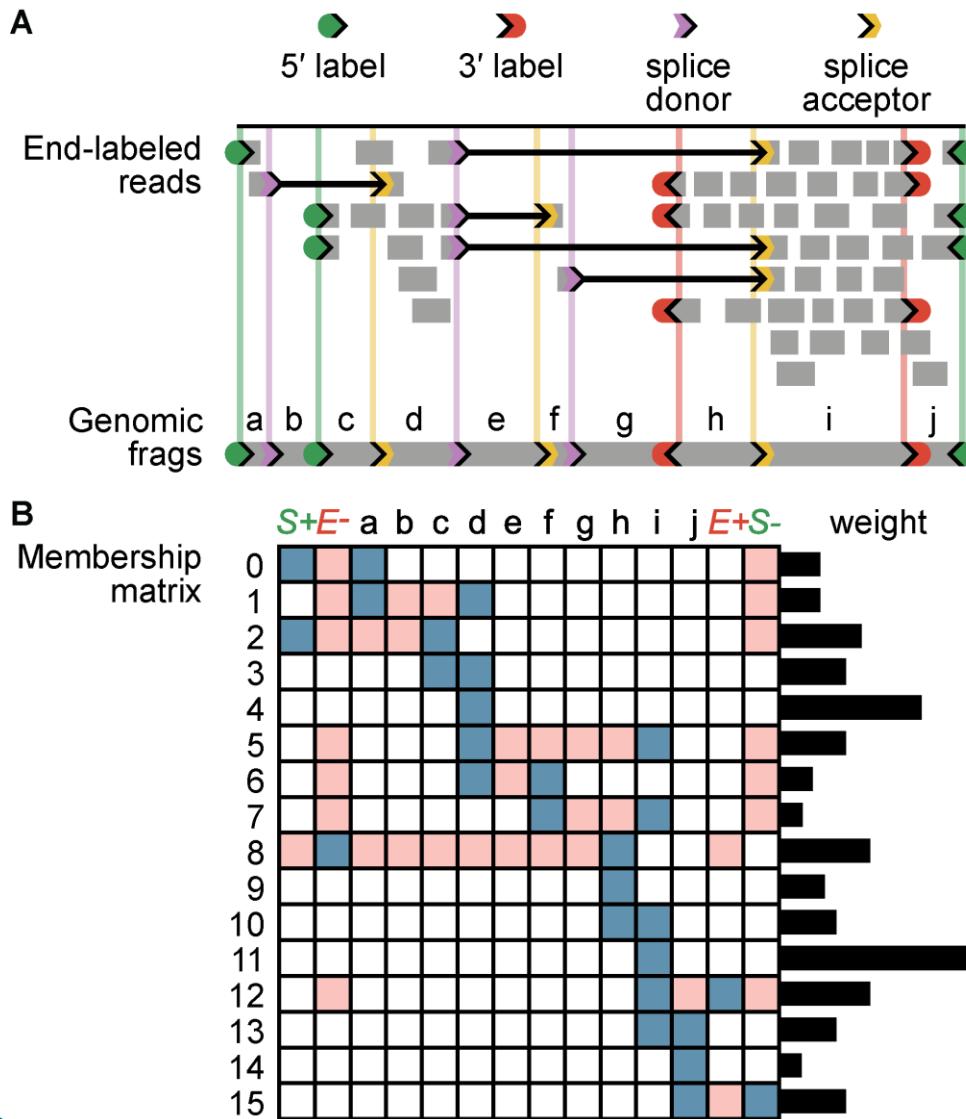
VBCF

Next Generation Sequencing
Andreas Sommer
Ido Tamir



European Research Council
Established by the European Commission

The Bookend framework



bookend

A

scRNA-seq SIRV spike-ins

