



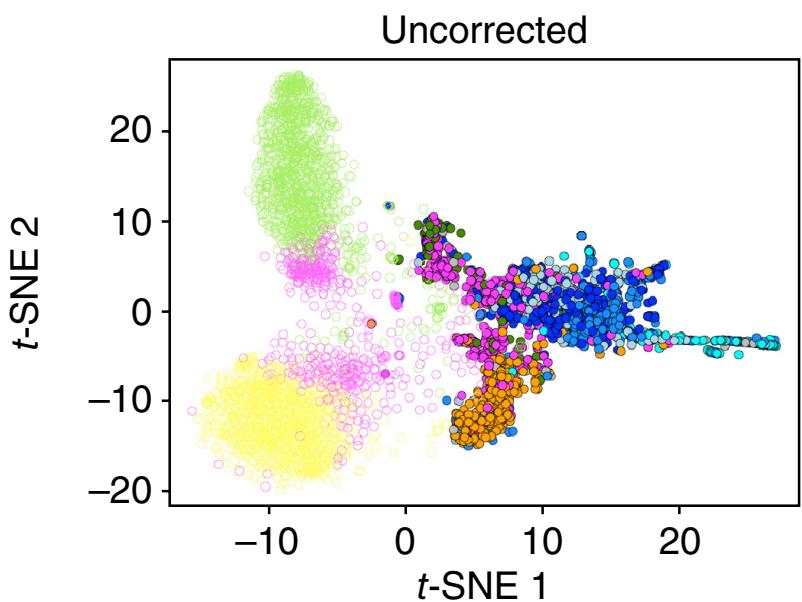
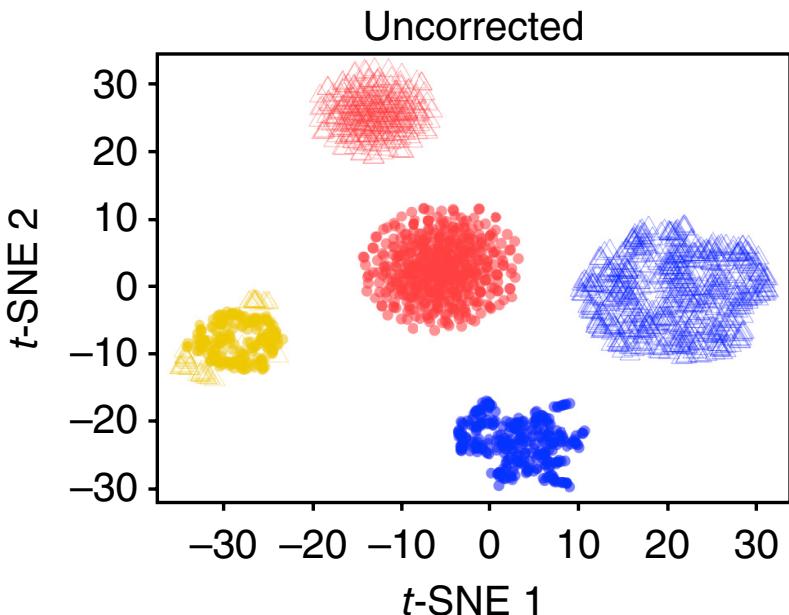
# Theory II

Single cell RNAseq tutorial

ISMB 2018, Chicago



# Correcting for confounders (batch correction)



# Technical/experimental sources of batch difference

- Sequencing related differences
  - Sequencing depth
  - Sequencing saturation
  - Sequencing device (e.g. Nextseq vs HiSeq)

# Technical/Experimental sources of batch difference

- Poor choice of experimental design
  - E.g. batch tracks biology
- Differences in capture technology
  - SmartSeq2 vs 10x

# Biological confounders

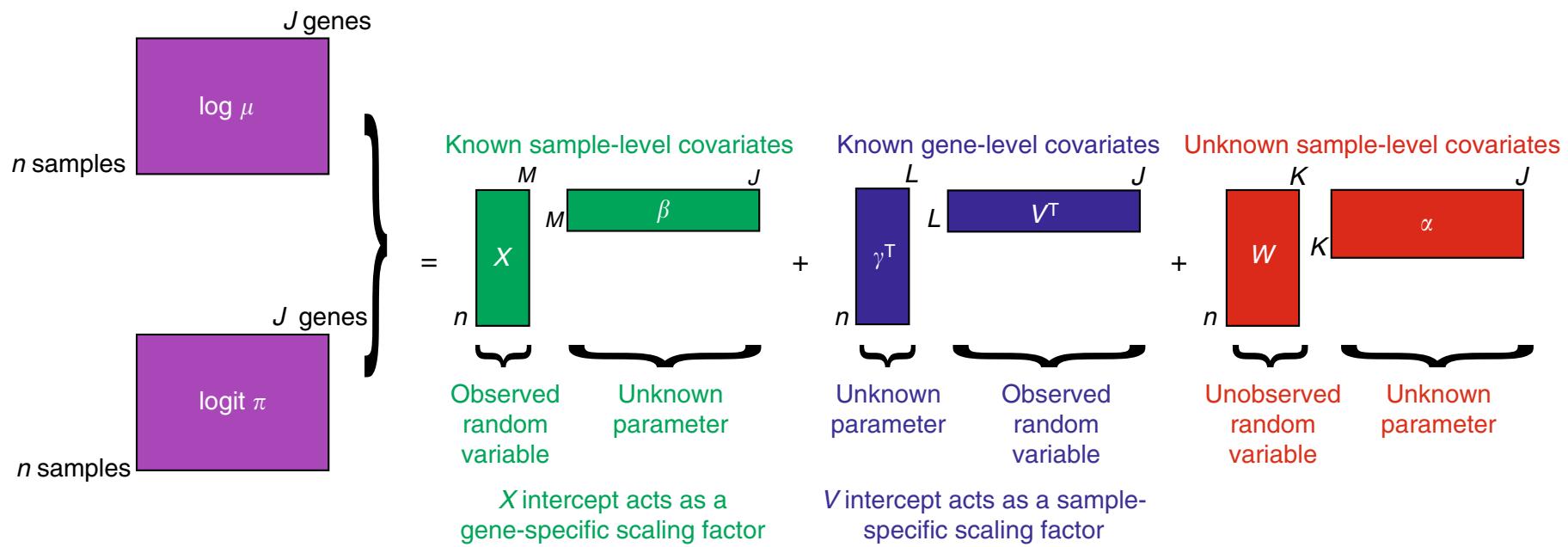
- Sample degradation
- Cell cycle
- Mitochondrial expression (or DNA content)
- Ribosomal fraction
- Copy number differences

# Classic approaches to correction

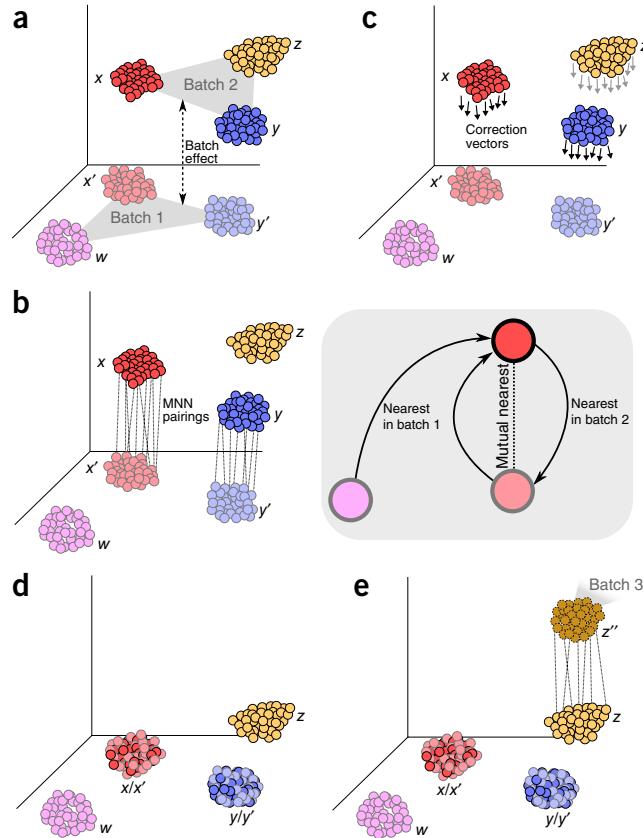
- Down sampling
- Normalization  
(scaling, centering, variance, log)
- Removing genes associated with batch
  - Removing PCs associated with batch
- ComBat
- Regression of residuals (linear or glm)
  - Experimental factor
  - Number of genes
  - Total UMI
  - % mitochondrial

# Modeling based approaches

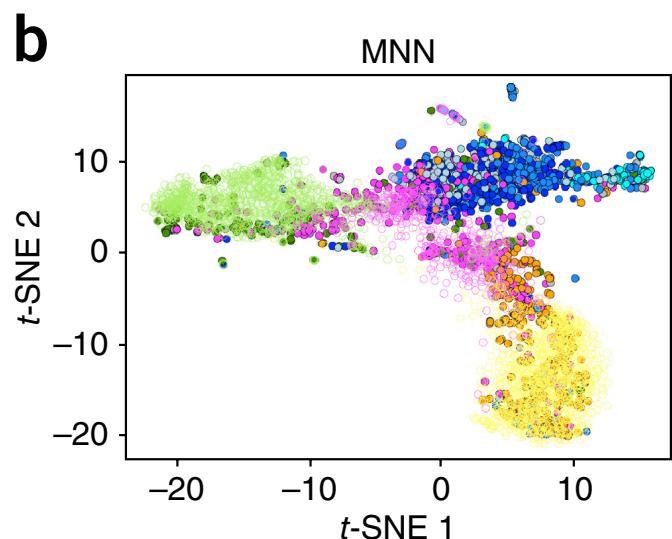
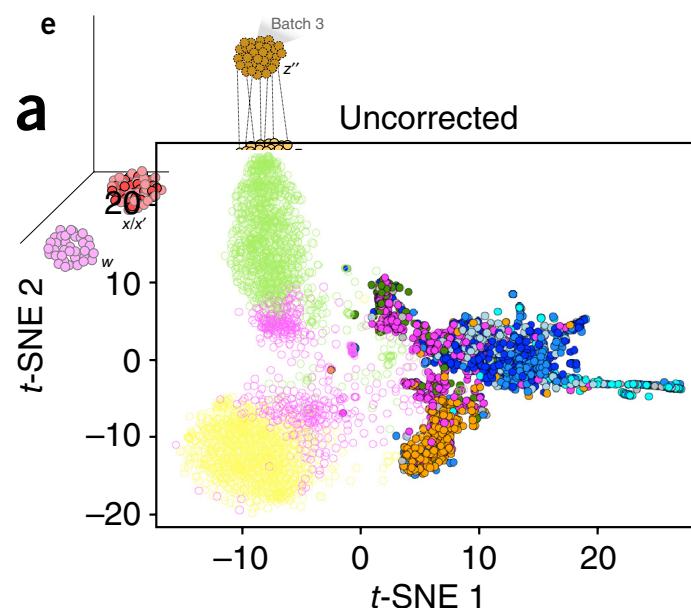
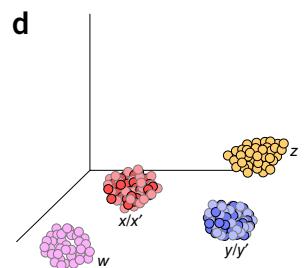
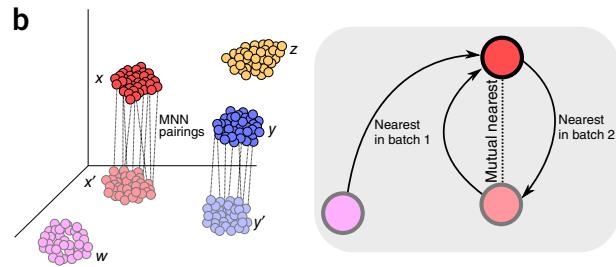
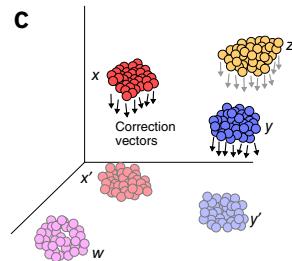
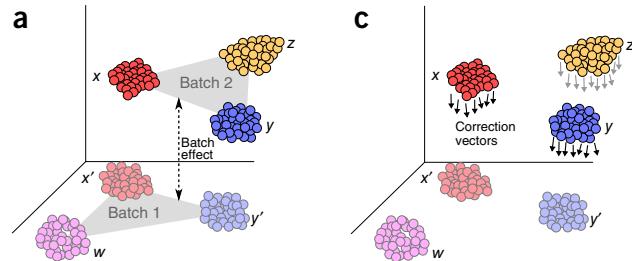
- ZINB-wave



# NN graph correction – mutually nearest neighbors

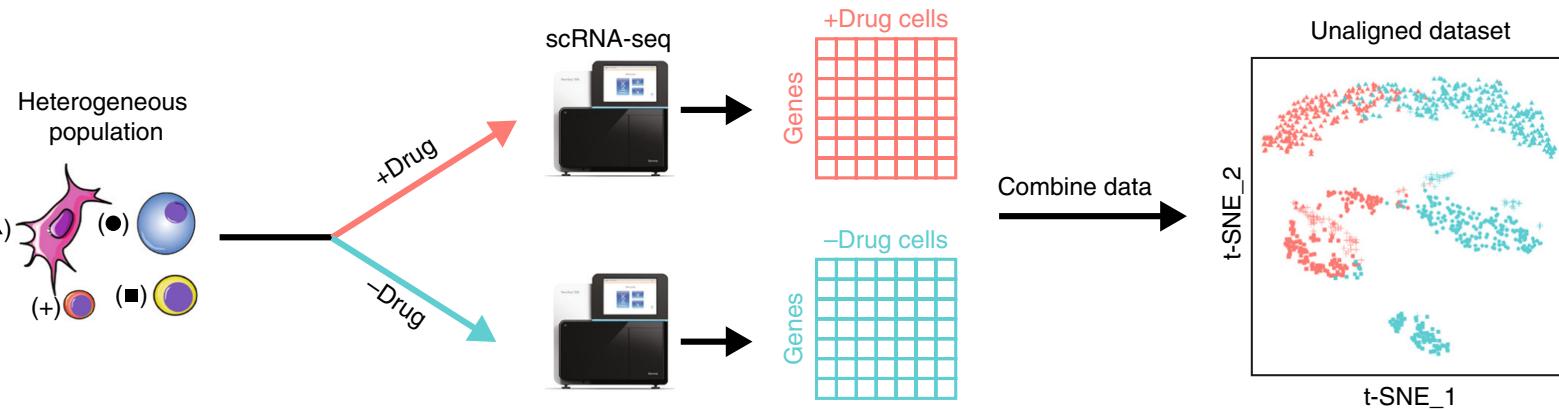


# Mutually nearest neighbors

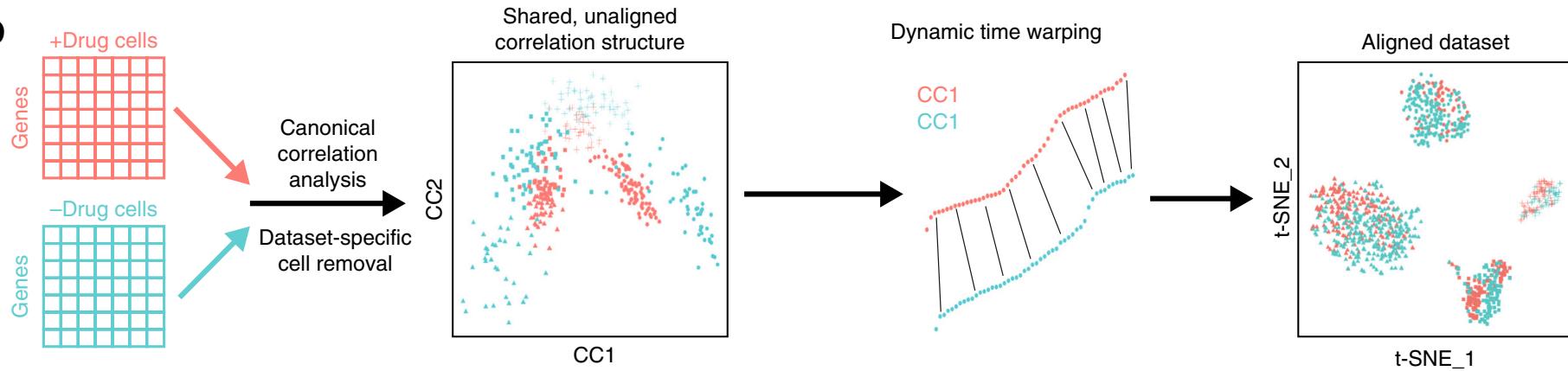


# Canonical correlation analysis

a



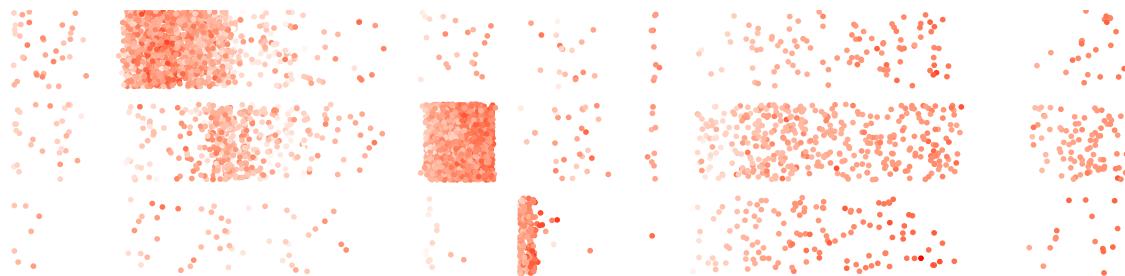
b



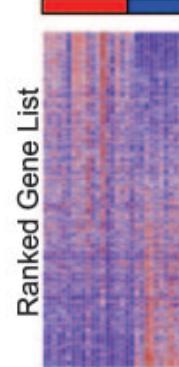
# Geneset and pathway enrichment

# Gene set enrichment analysis and scRNAseq

Sparse gene x cell matrix



A Phenotype  
Classes  
A B

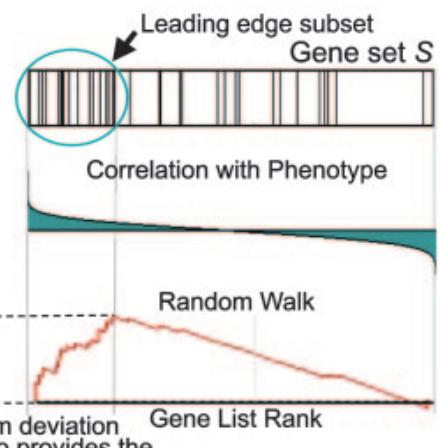


B

Gene set S



Leading edge subset  
Gene set S  
Correlation with Phenotype  
Random Walk  
 $ES(S)$   
Maximum deviation  
from zero provides the  
enrichment score  $ES(S)$



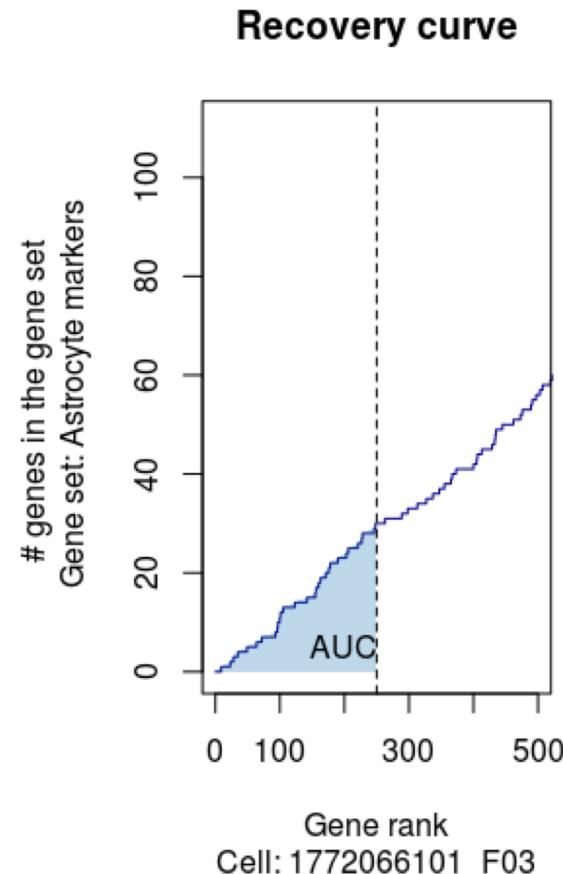
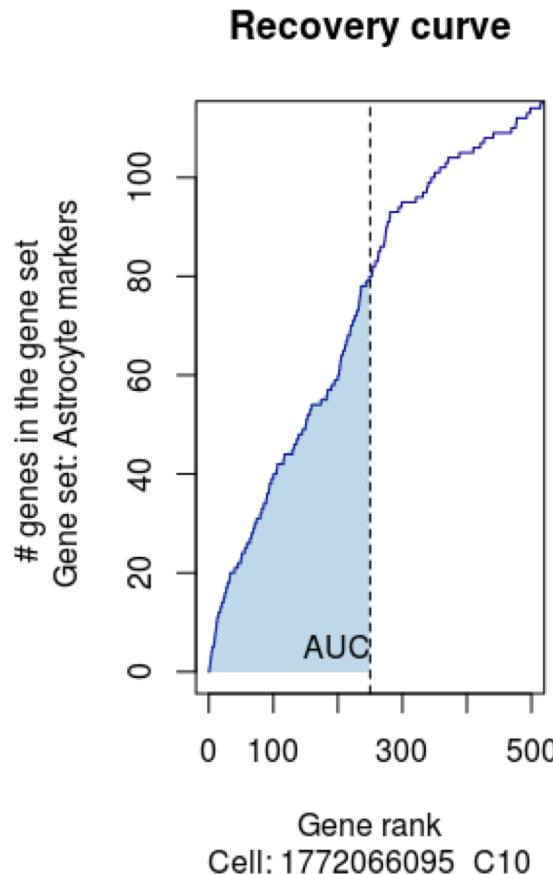
# Useful sources for gene set information

- GO/KEGG/BioCarta (old school cool)
- MsigDB (Broad)
- Enricher (Ma'ayan lab, Mount Sinai)
- GSKB – multi species (Ge lab, South Dakota State)
- G:profiler (Vilo lab, University of Tartu, Estonia)

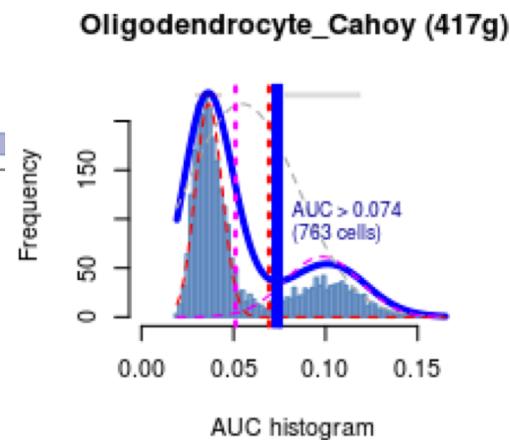
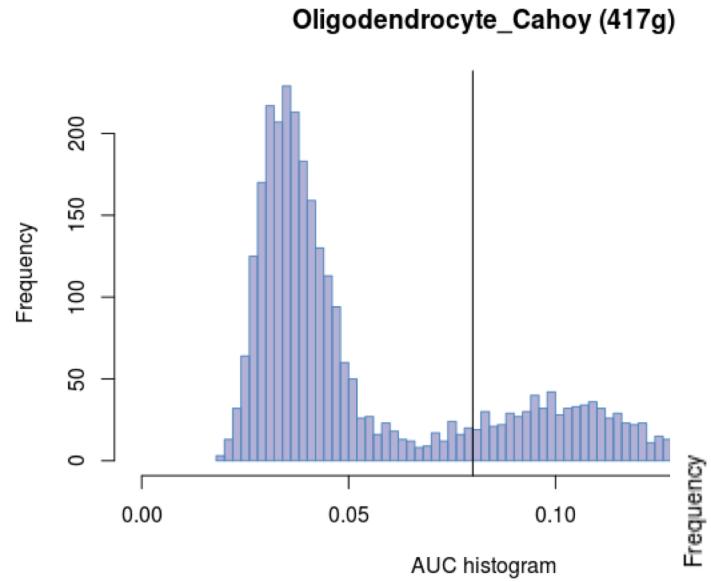
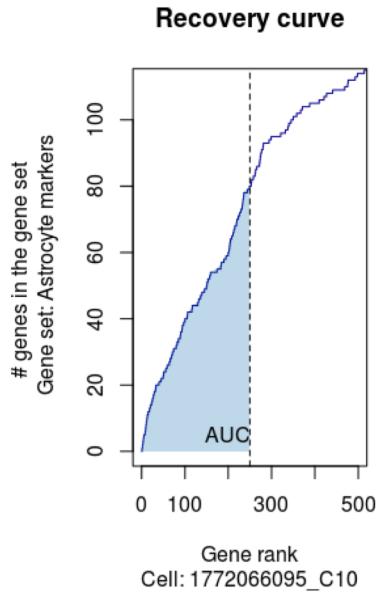
# Additional bulk RNA enrichment tools

- GSA (R)
- GAGE (R)
- fgsea (R)
- HOMER (Perl)

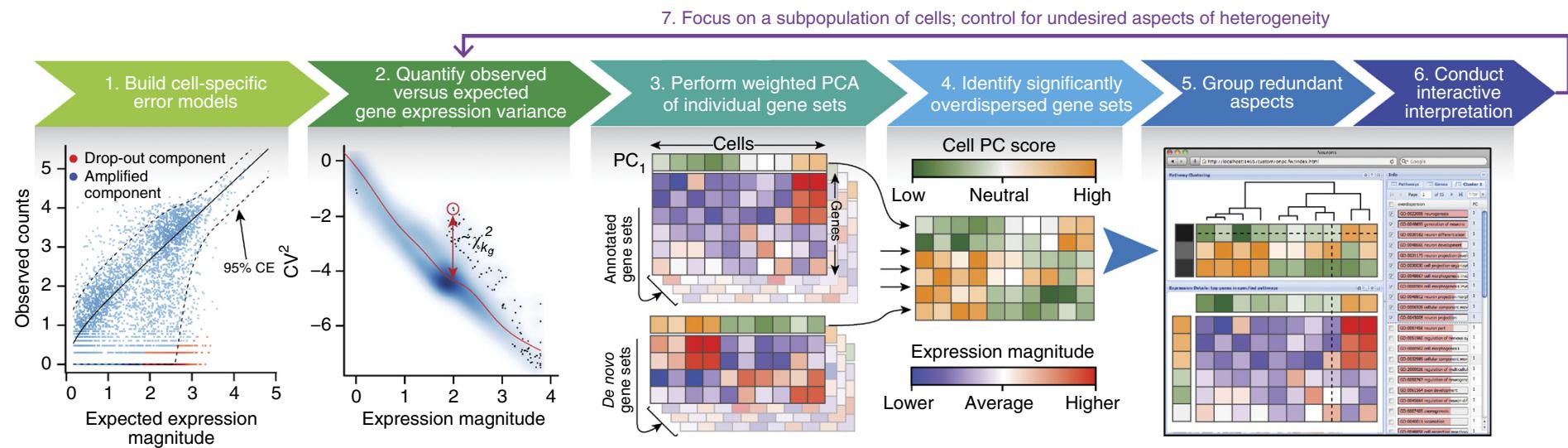
# Single cell enrichment score - AUCell



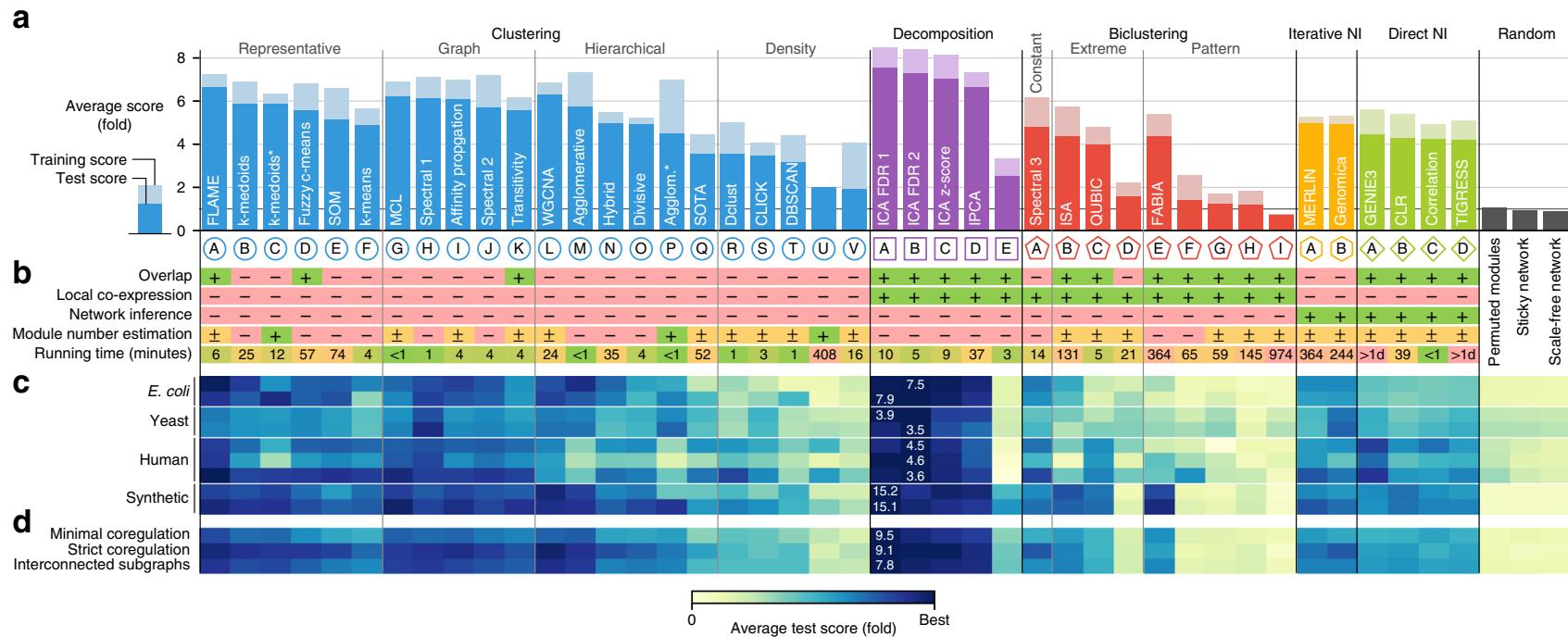
# Single cell enrichment score - AUCell



# Enriched variance within a gene-set - Pagoda



# Survey of module discovery methods



## **COMMUNICATIONS**

## ARTICLE

DOI: 10.1038/s41467-018-03424-4

OPEN

# A comprehensive evaluation of module detection methods for gene expression data

Wouter Saelens<sup>1,2</sup>, Robrecht Cannoodt<sup>1,3</sup> & Yvan Saeys<sup>1,2</sup>

Trajectory/pseudo-time  
analysis

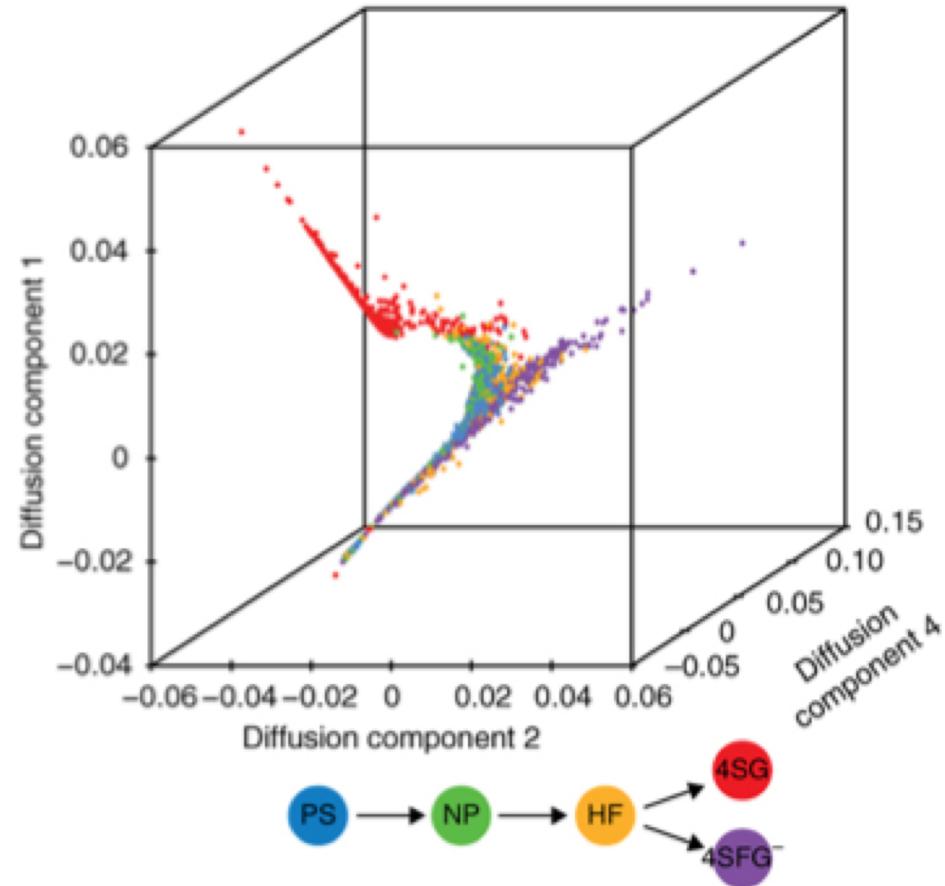
# Visualizing Continuous Variation using Diffusion Maps

For many systems we can not assume clear states.

Assumes that there exists a “continuous manifold” of states that cells are diffusing through (e.g. during differentiation).

Diffusion maps construct a random walk (Markov process) on the data, and infer a low dimensional representation (DCs) that describes these paths.

Of the analysis types, trajectory inference is modeling data in a way almost impossible using bulk



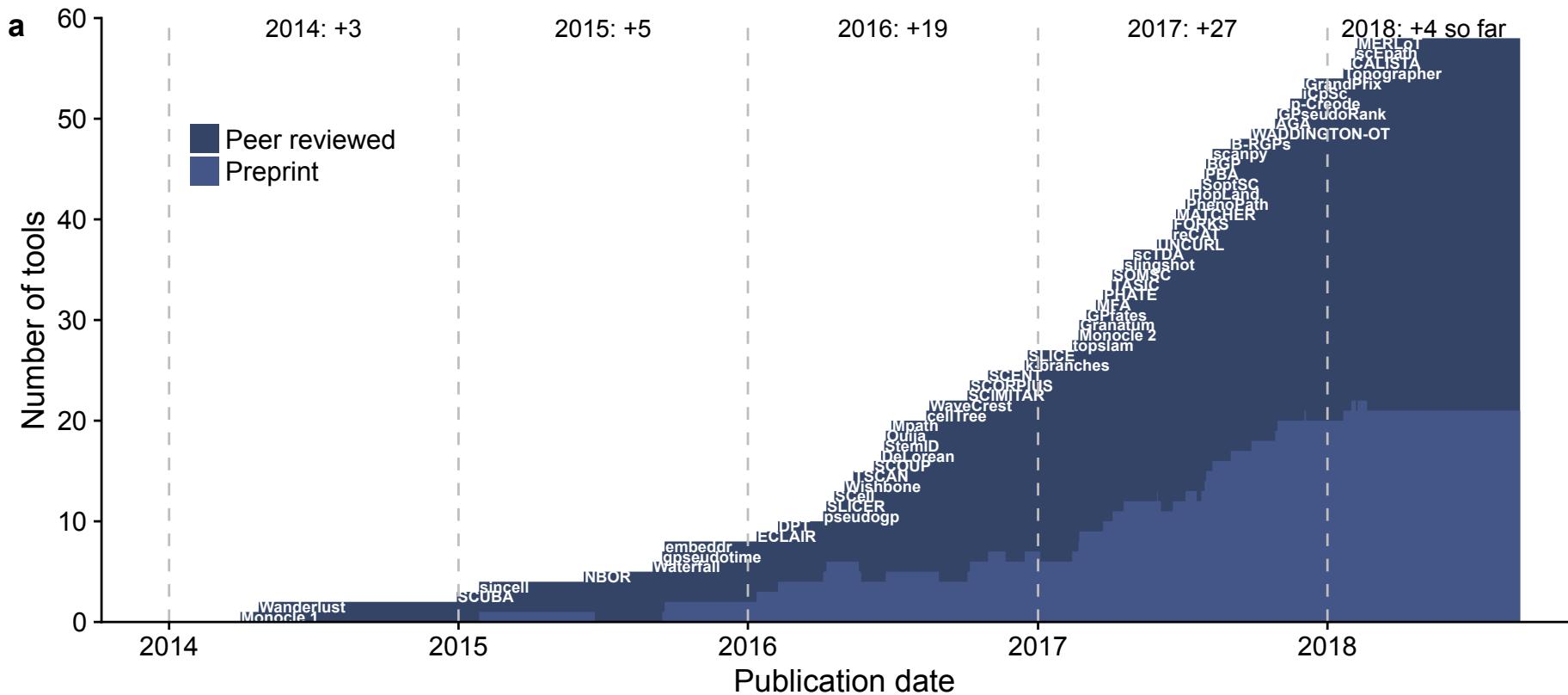
# The pseudotime/trajectory estimation problem

**Given:** single-cell expression measurements for a heterogeneous collection of cells that is transitioning among states.

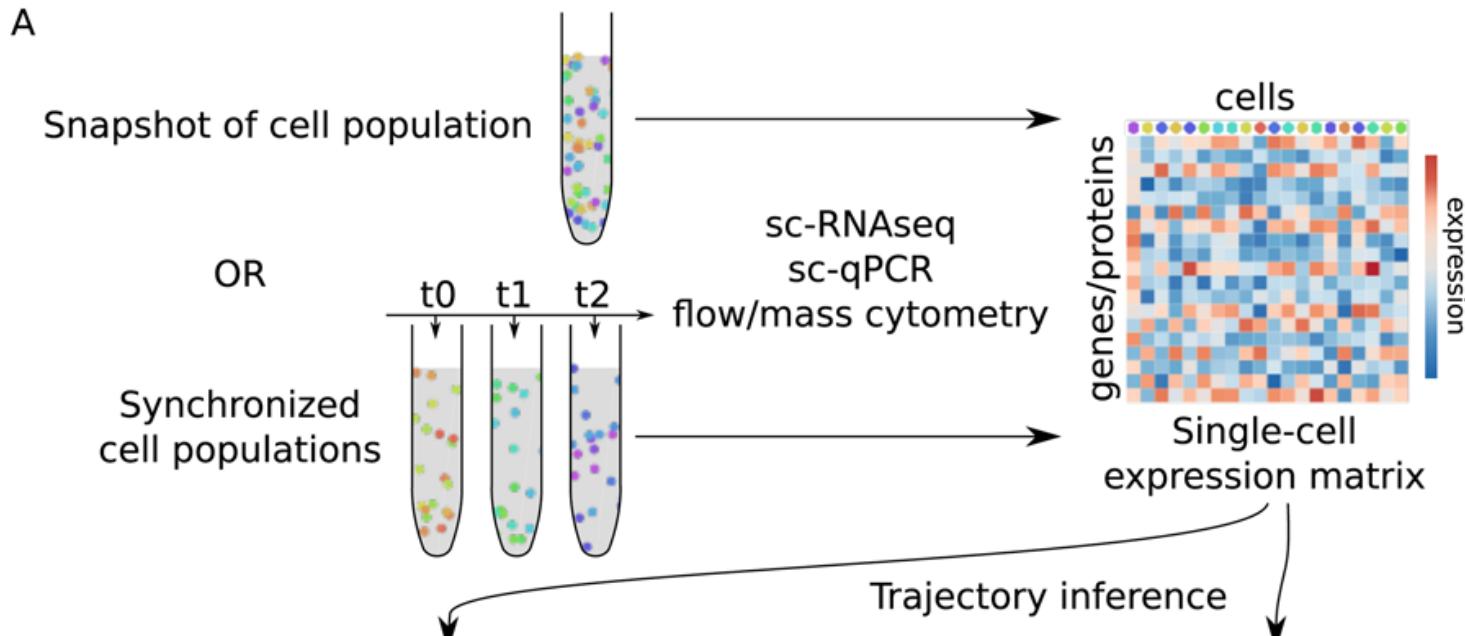
**Return:**

1. A quantitative value for each cell that represents its progress in the transition.
2. A path describing the likely transition between biological states.

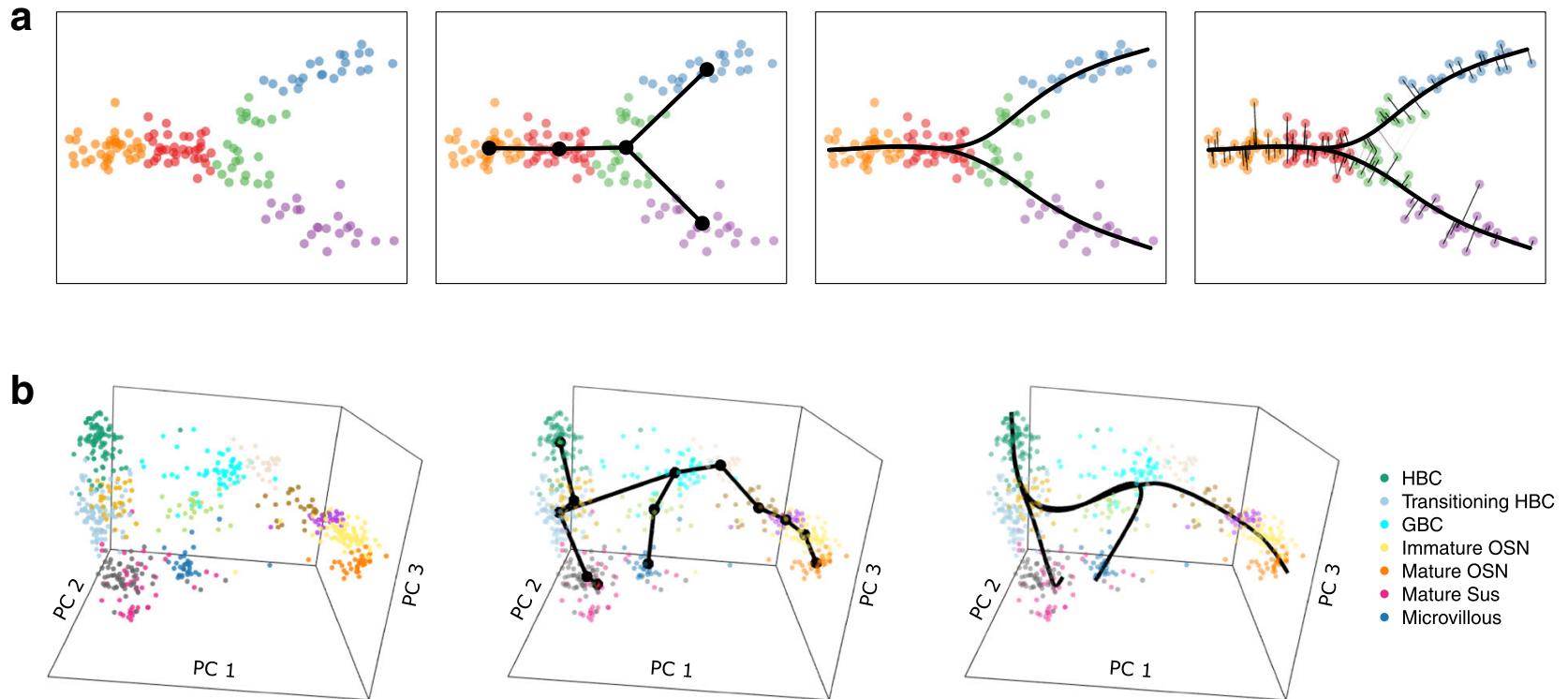
# The trajectory inference gold rush



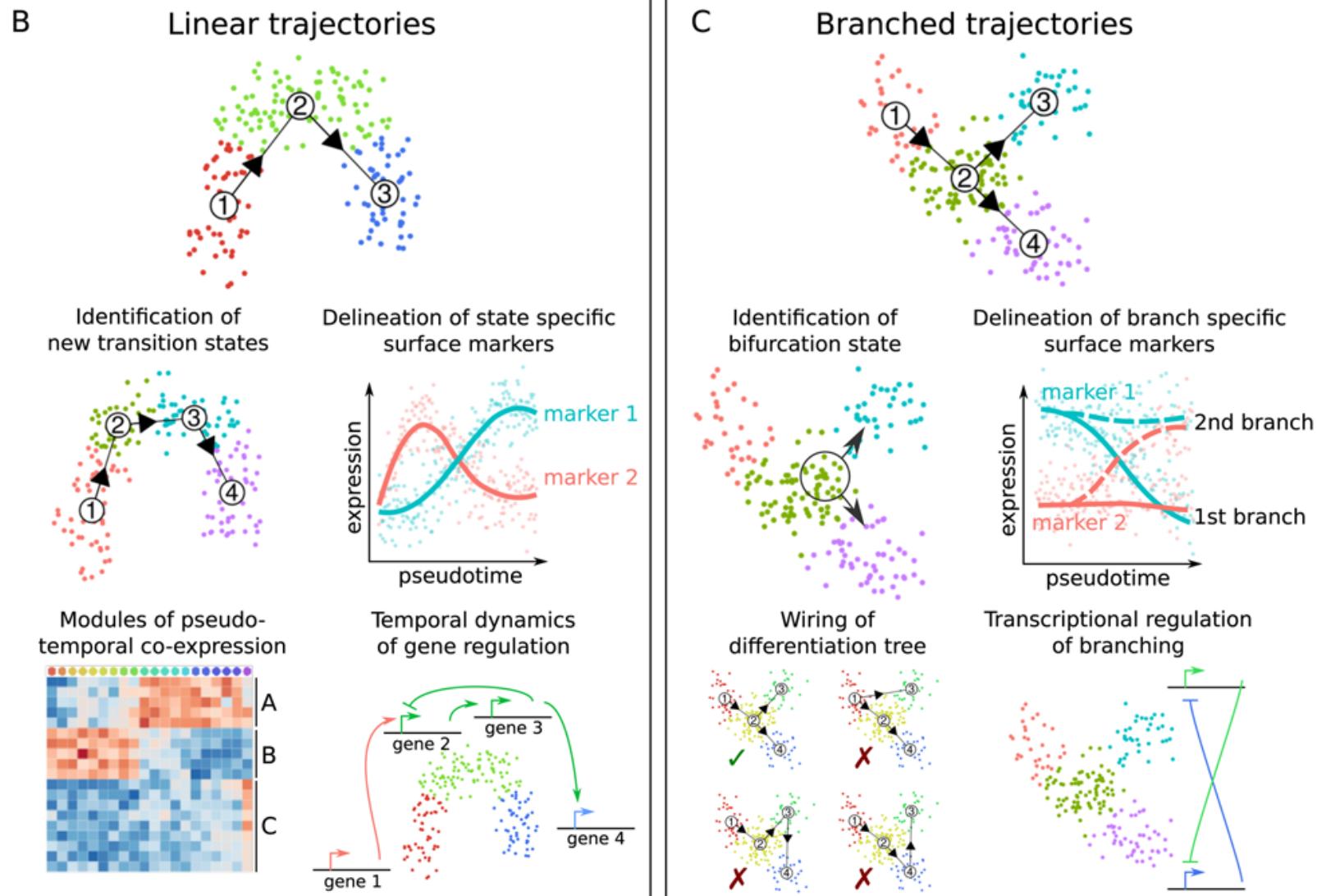
# Inferring expression dynamics from a snapshot



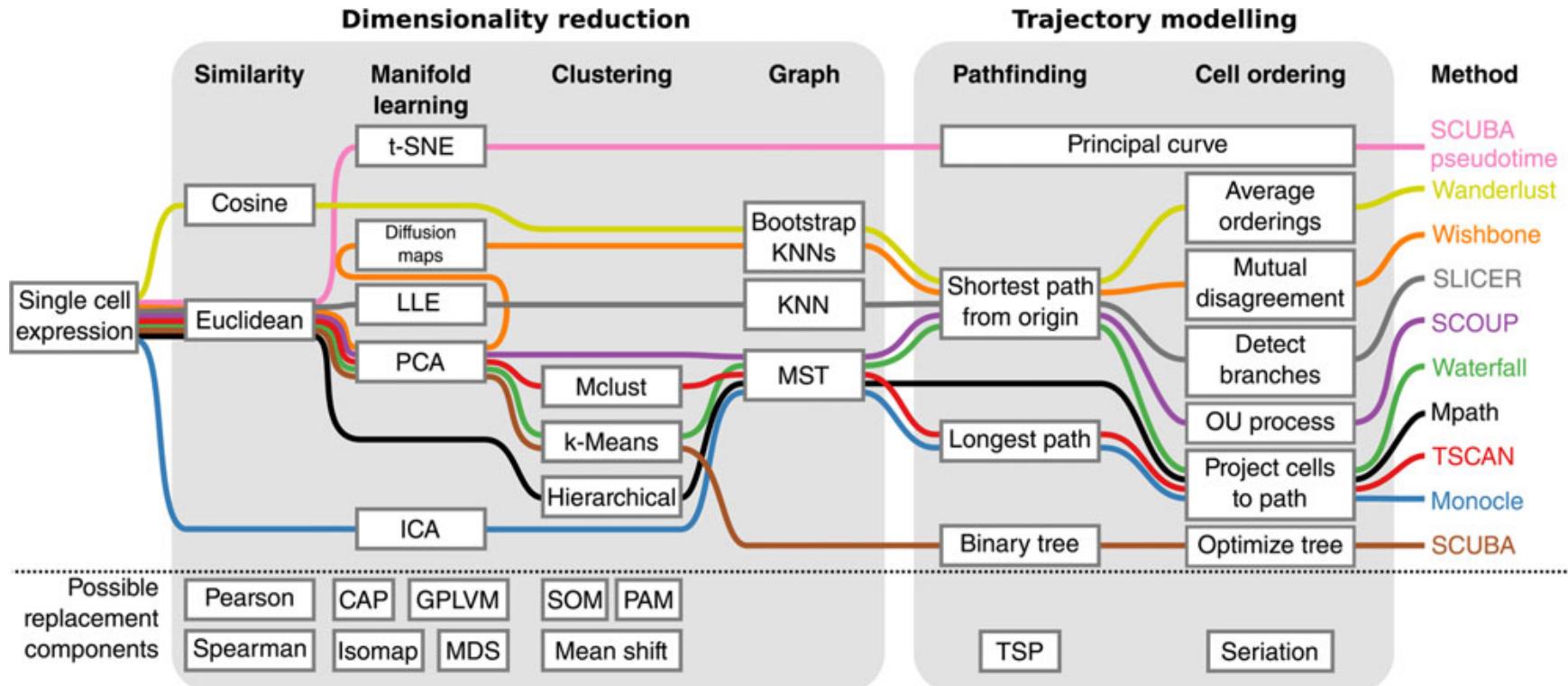
# Inferring expression dynamics from a snapshot



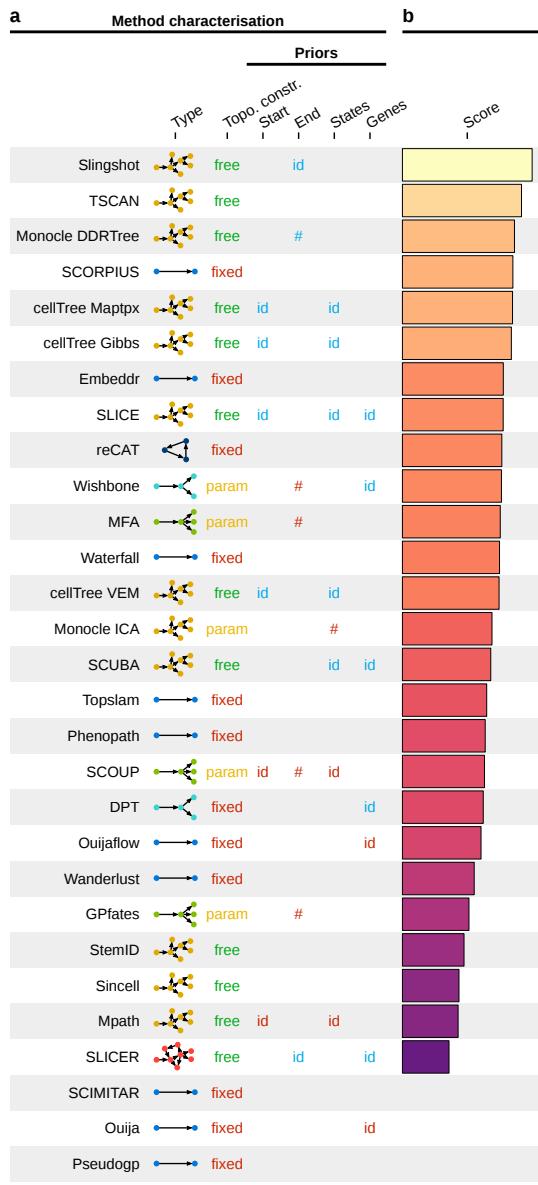
# Inferring expression dynamics from a snapshot



# An honest attempt to chart the landscape (circa 2016)

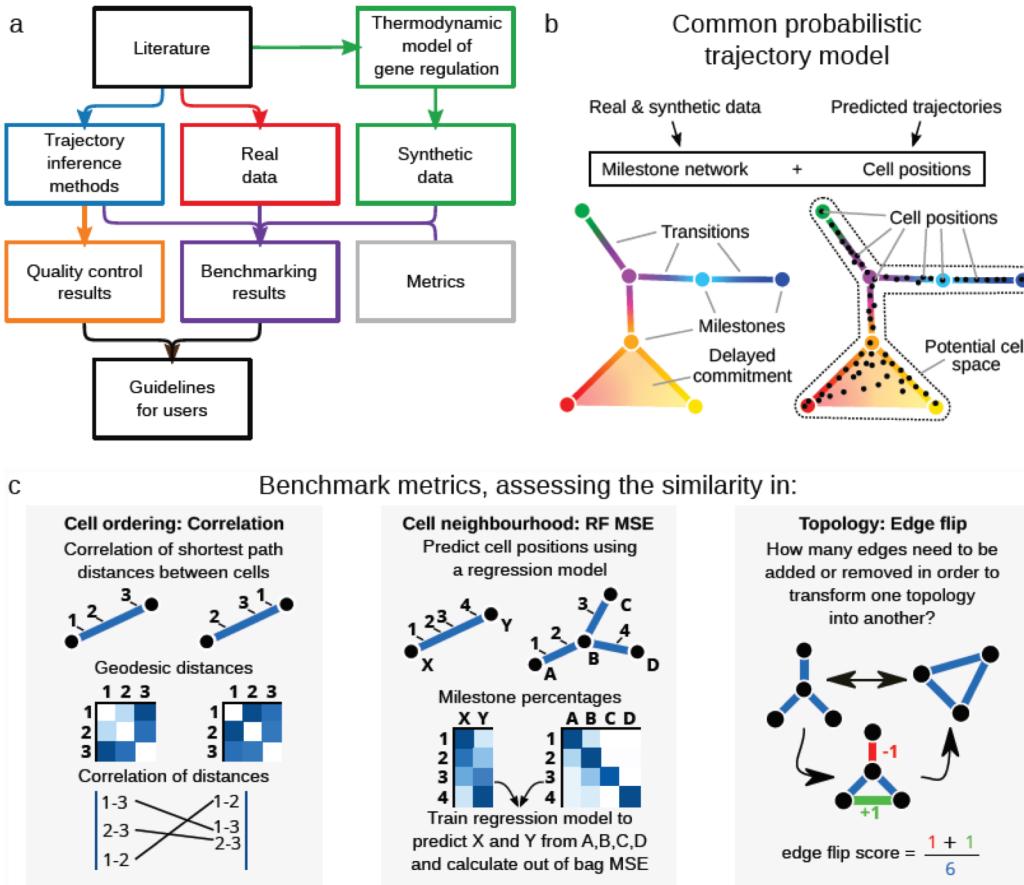


# In 2018 ...



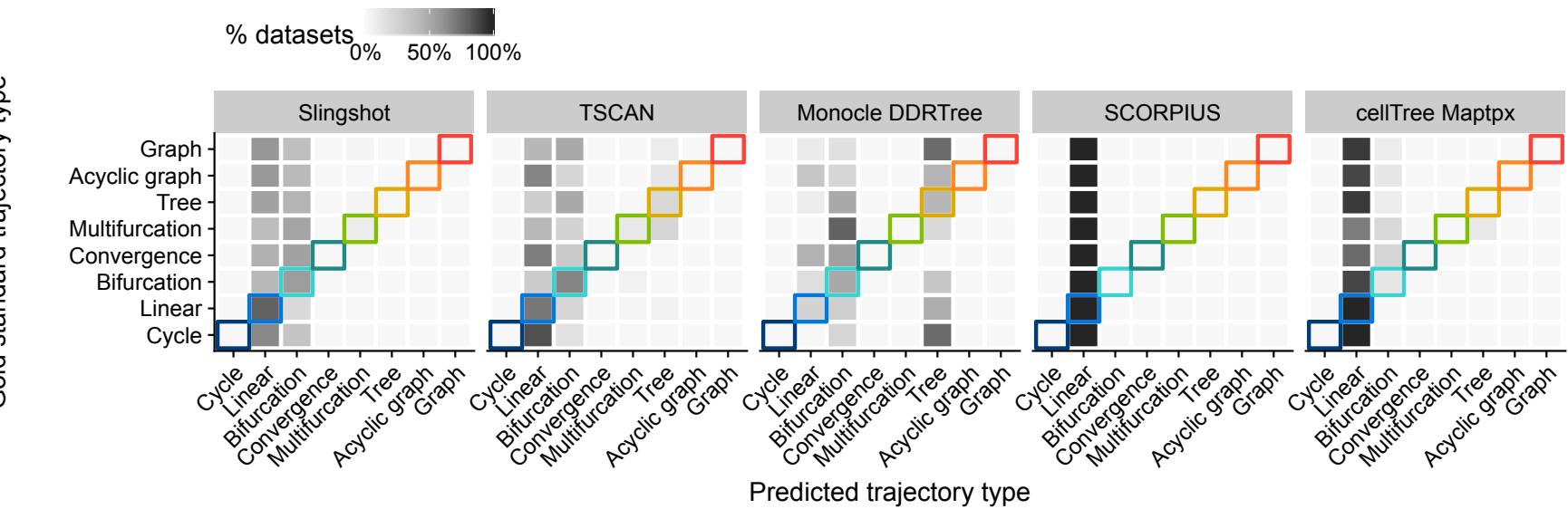
Saelens et al. 2018

# TI method evaluation



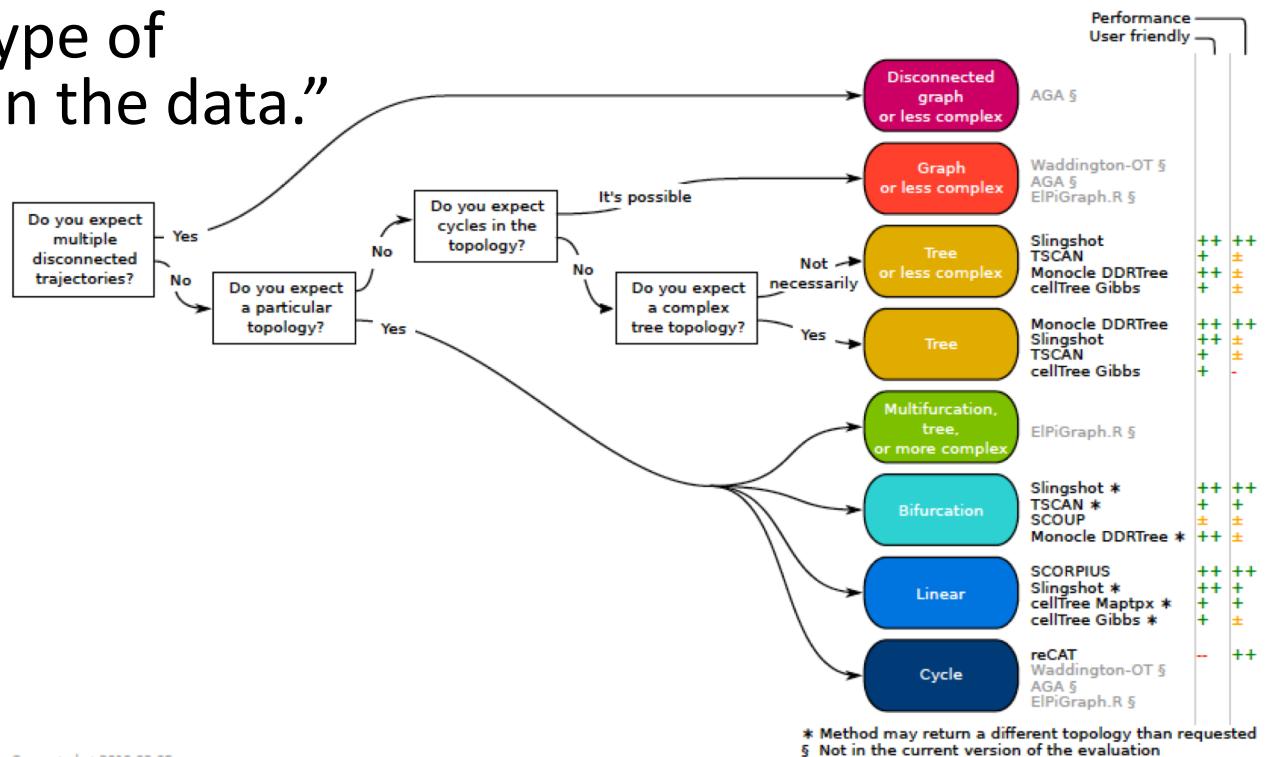
# Not quite a solved problem?

a

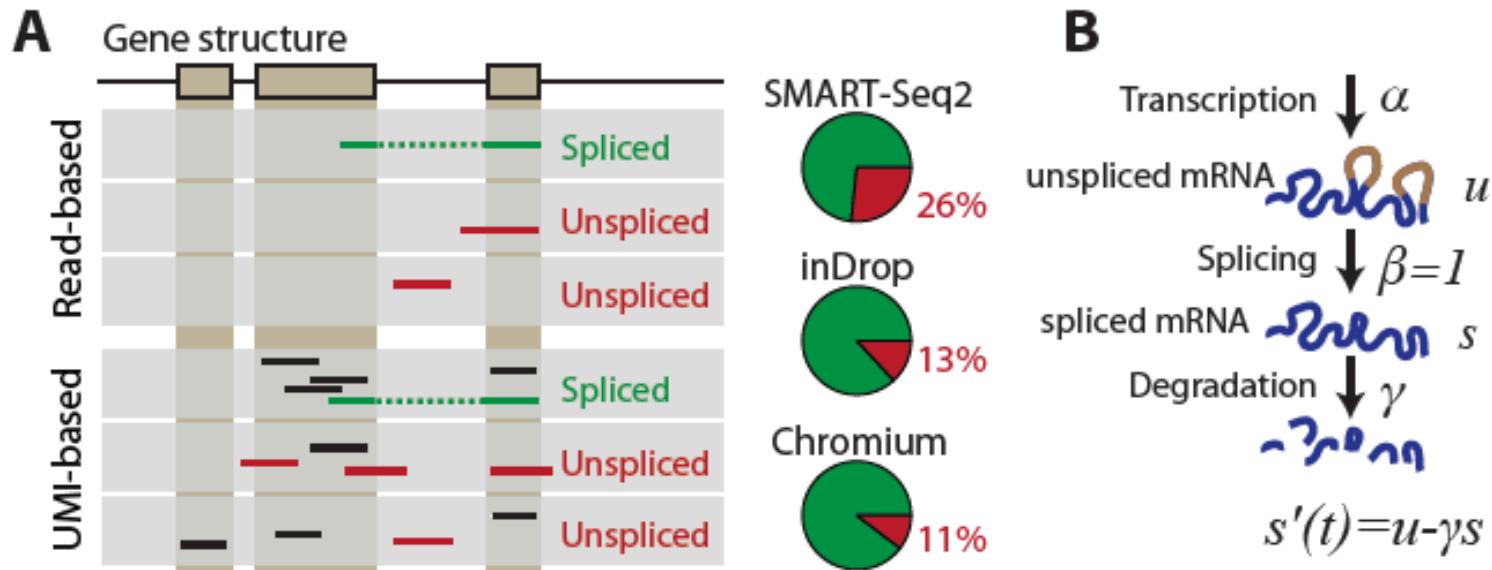


# Decision Tree for Method Use

“We found that some methods, including Slingshot, TSCAN and Monocle DDRTree, clearly outperform other methods, although their performance depended on the type of trajectory present in the data.”

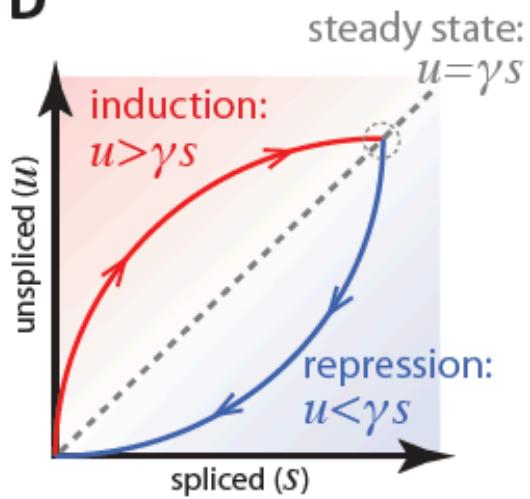


# RNA velocity - Inferring dynamics from RNA degradation

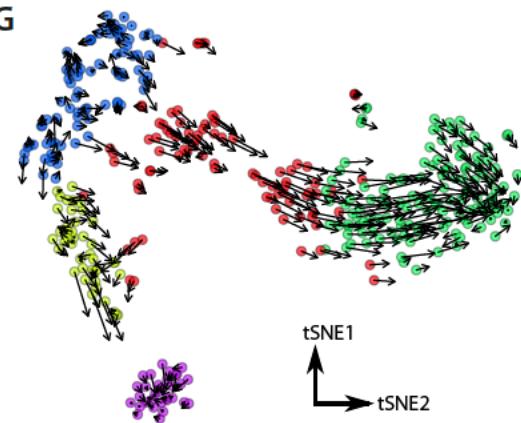


# Inferring dynamics -

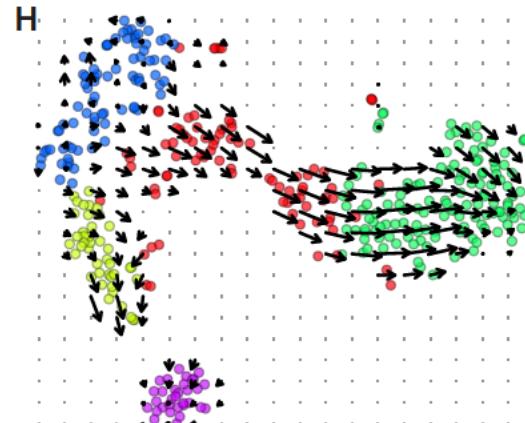
D



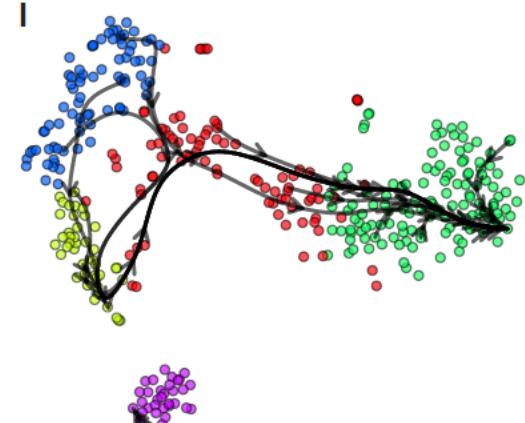
G



H



I



# Interactive tools

# Single Cell Portal

Single Cell Portal BETA

Help ▾

Sign In

## Single Cell Portal BETA

Visualization portal for single cell RNA-seq data.

Now featuring **41** studies with **456,607** cells.

### Browse Studies

Search Studies...



Most Recent

Most Popular

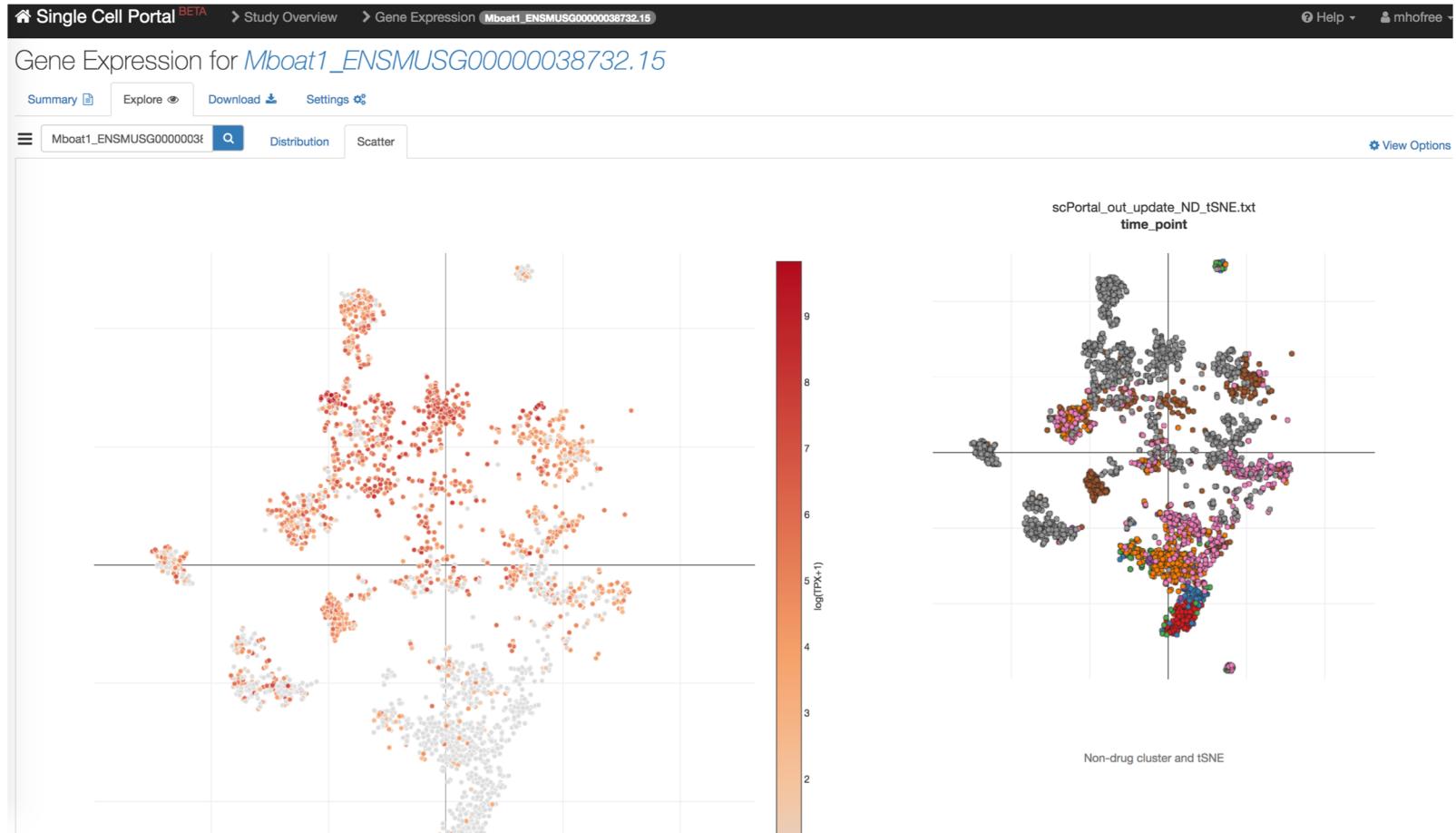
Reset Filters

### Single nucleus RNA-seq of cell diversity in the adult mouse hippocampus (sNuc-Seq) ▾

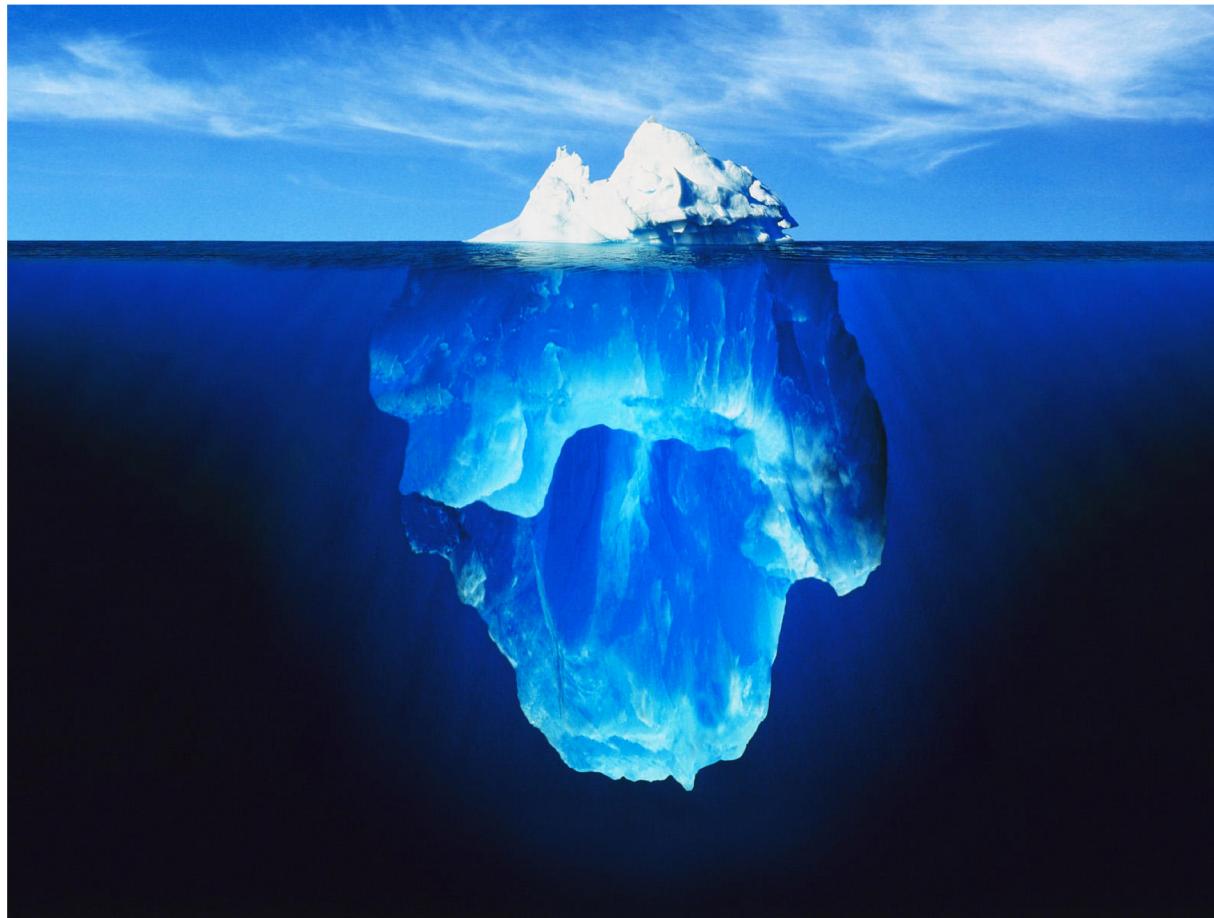
**View Study**

Single nucleus RNA-seq of cell diversity in the adult mouse hippocampus. Habib N, Li Y, Heidenreich M, Swiech L, Avraham-David I, Trombetta J, Hession C, Zhang F, Regev A. *Div-Seq: Single-nucleus RNA-Seq reveals dynamics of rare adult newborn neurons*. *Science* 28 Jul 2016 DOI: 10.1126/science.aad7038 Contact: naomi@broadinstitute.org Single cell RNA-Seq provides rich information about cell types and states. However, it is difficult to capture rare dynamic processes, such as adult neurogenesis, because isolation of rare neurons from adult tissue is challenging and markers for each phase are limited. Here, we develop Div-Seq, which combines scalable single-nucleus RNA-Seq (sNuc-Seq) with pulse labeling of proliferating cells by EdU... (continued)

# Single cell portal - interactive



# What did we cover?



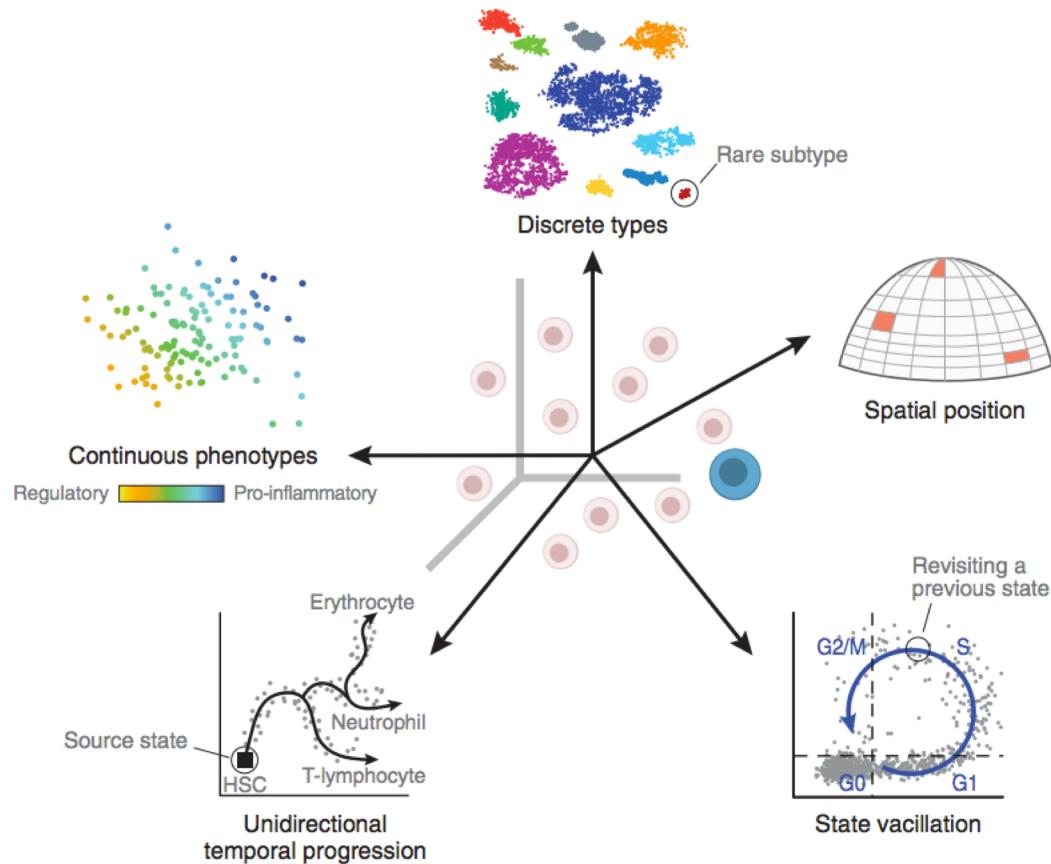
# What did we cover...

## REVIEW

nature  
biotechnology

Revealing the vectors of cellular identity with single-cell genomics

Allon Wagner<sup>1</sup>, Aviv Regev<sup>2,3,5</sup> & Nir Yosef<sup>1,4,5</sup>



# Acknowledgements



Special thanks:

Brian Haas

Timothy Tickle

Karthick Shekar

Aviv Regev

