

CCMI data driven modelling week

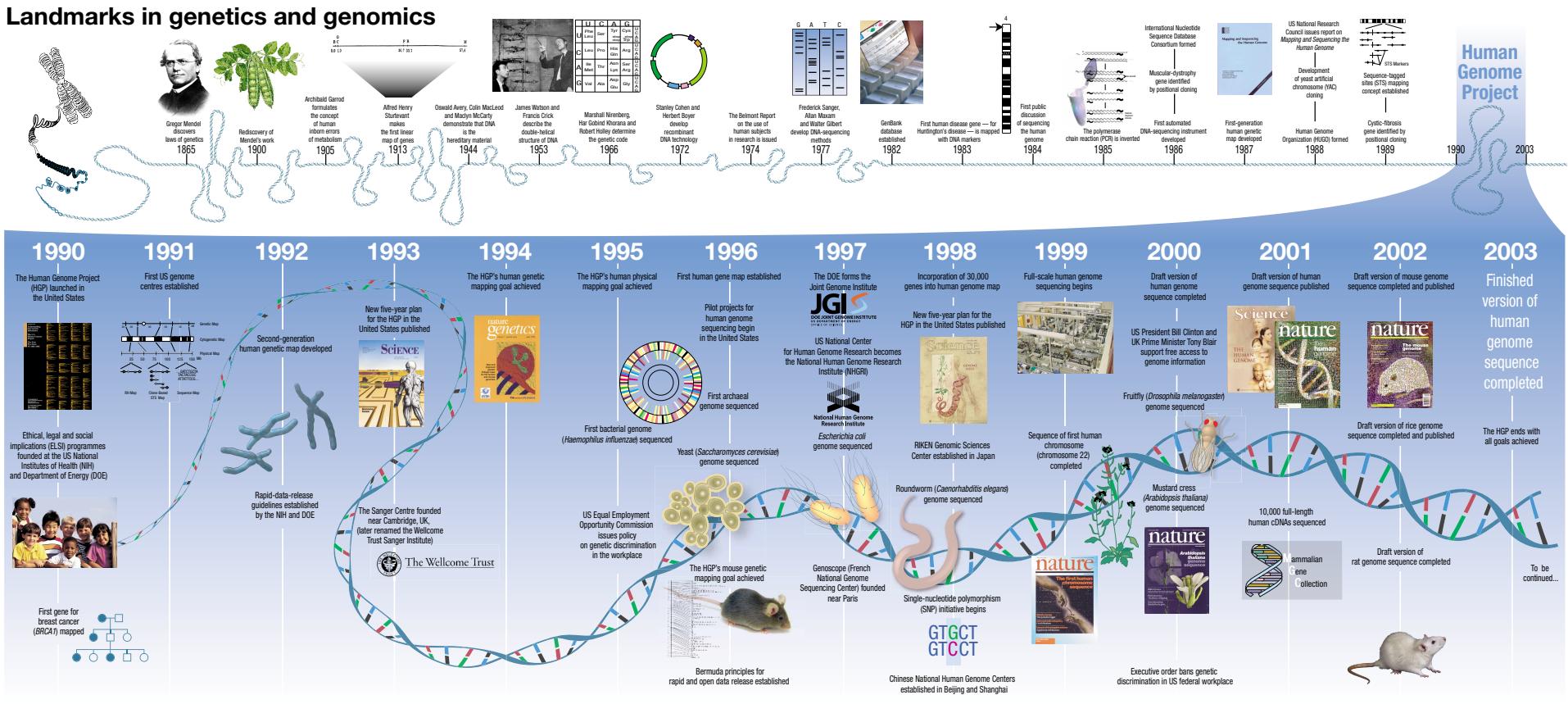
Data analysis and modelling:
Case study of single-cell RNA-sequencing
(transcription) data

Vahid Shahrezaei
Department of Mathematics

IMPERIAL

Development in genetics: Human Genome Project

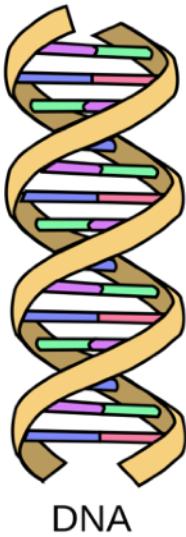
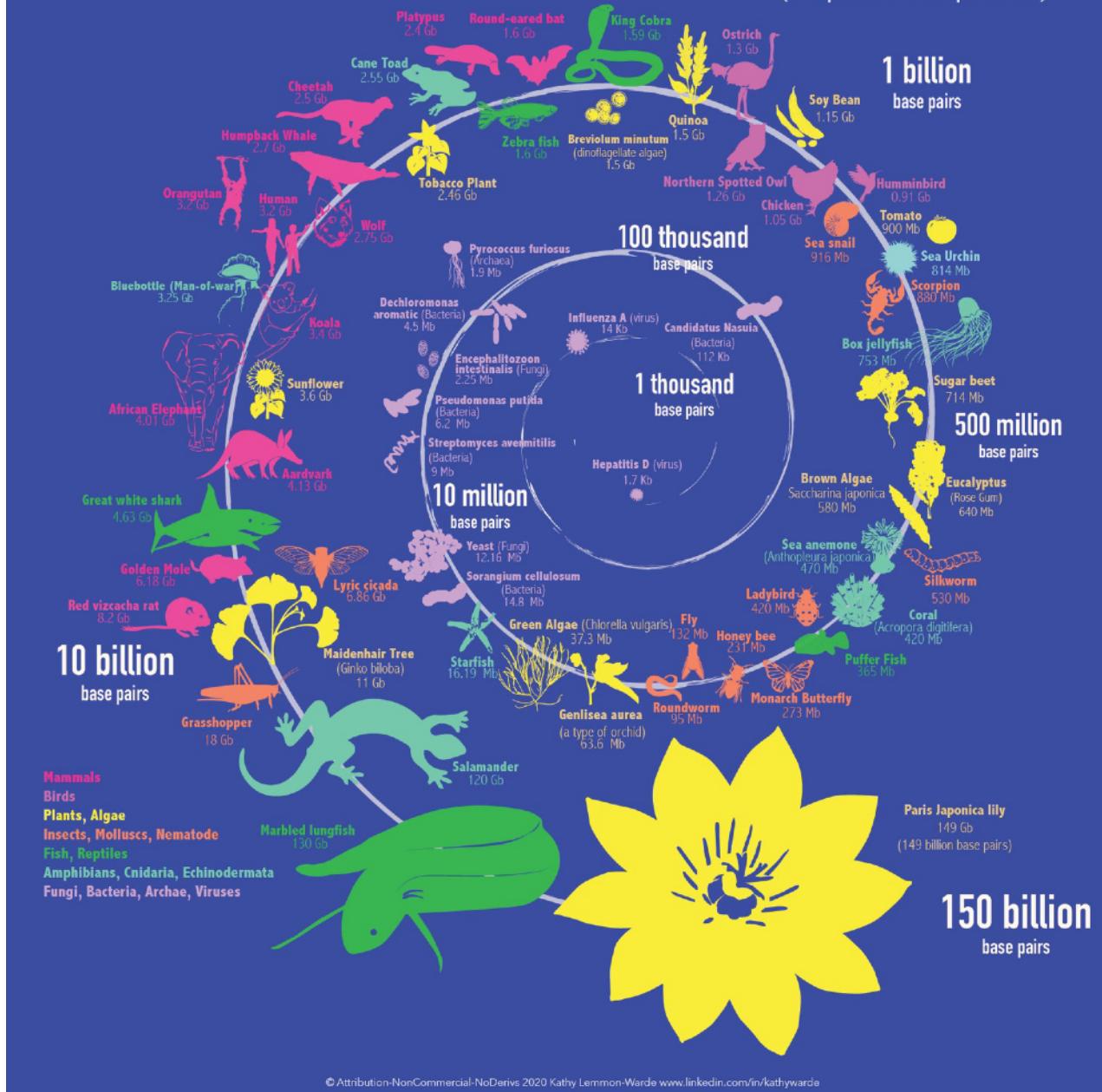
Relating genotype to phenotype is not easy!



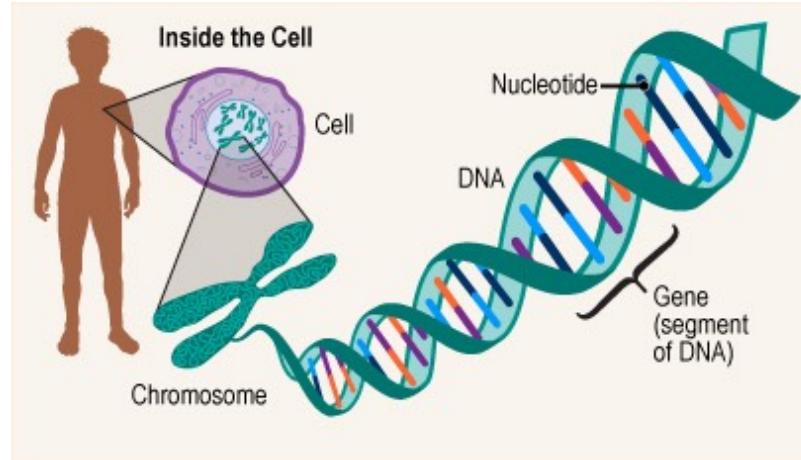
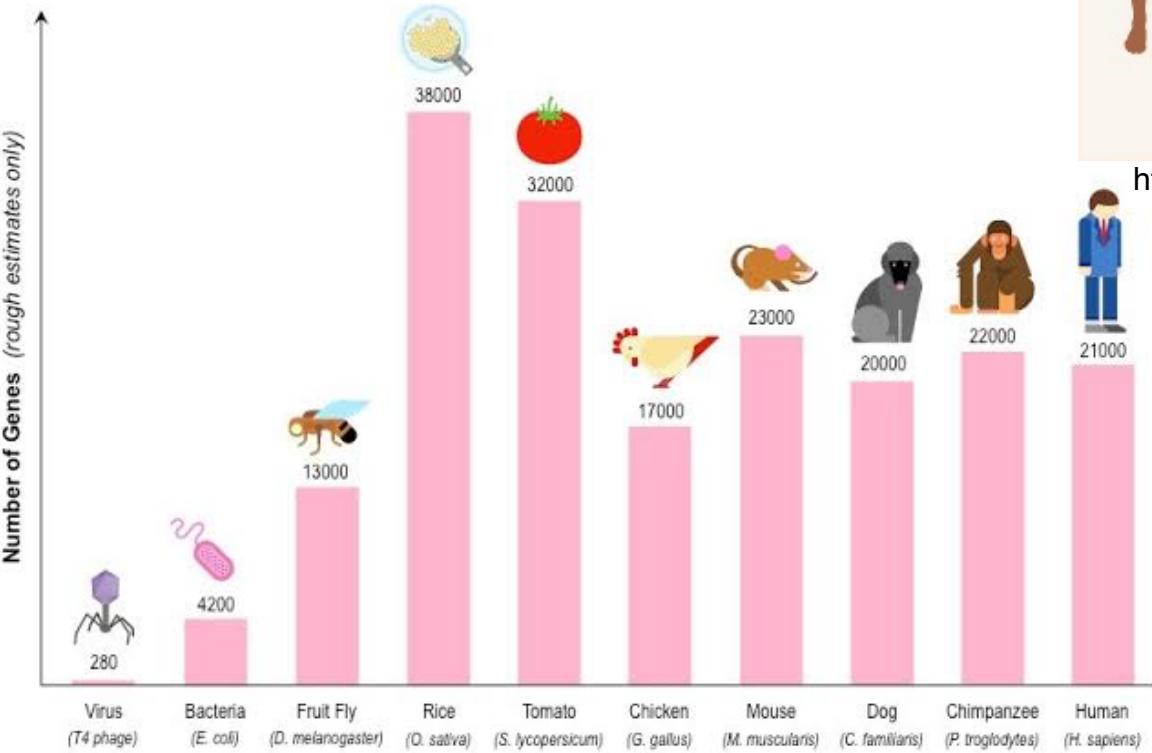
<https://www.genome.gov/human-genome-project>

Genome size

(A species comparison)



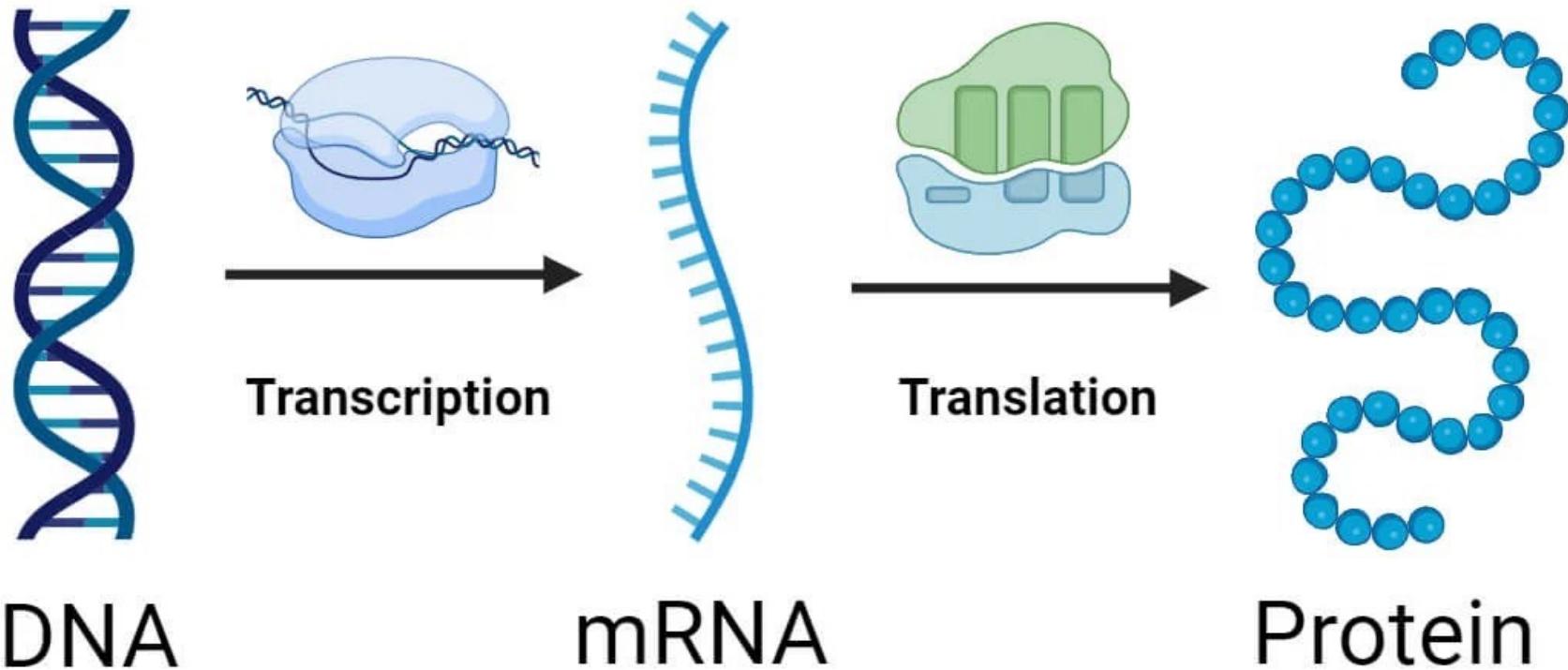
An organism's complexity, its genome size and its total number of (protein-coding) genes are only loosely correlated!



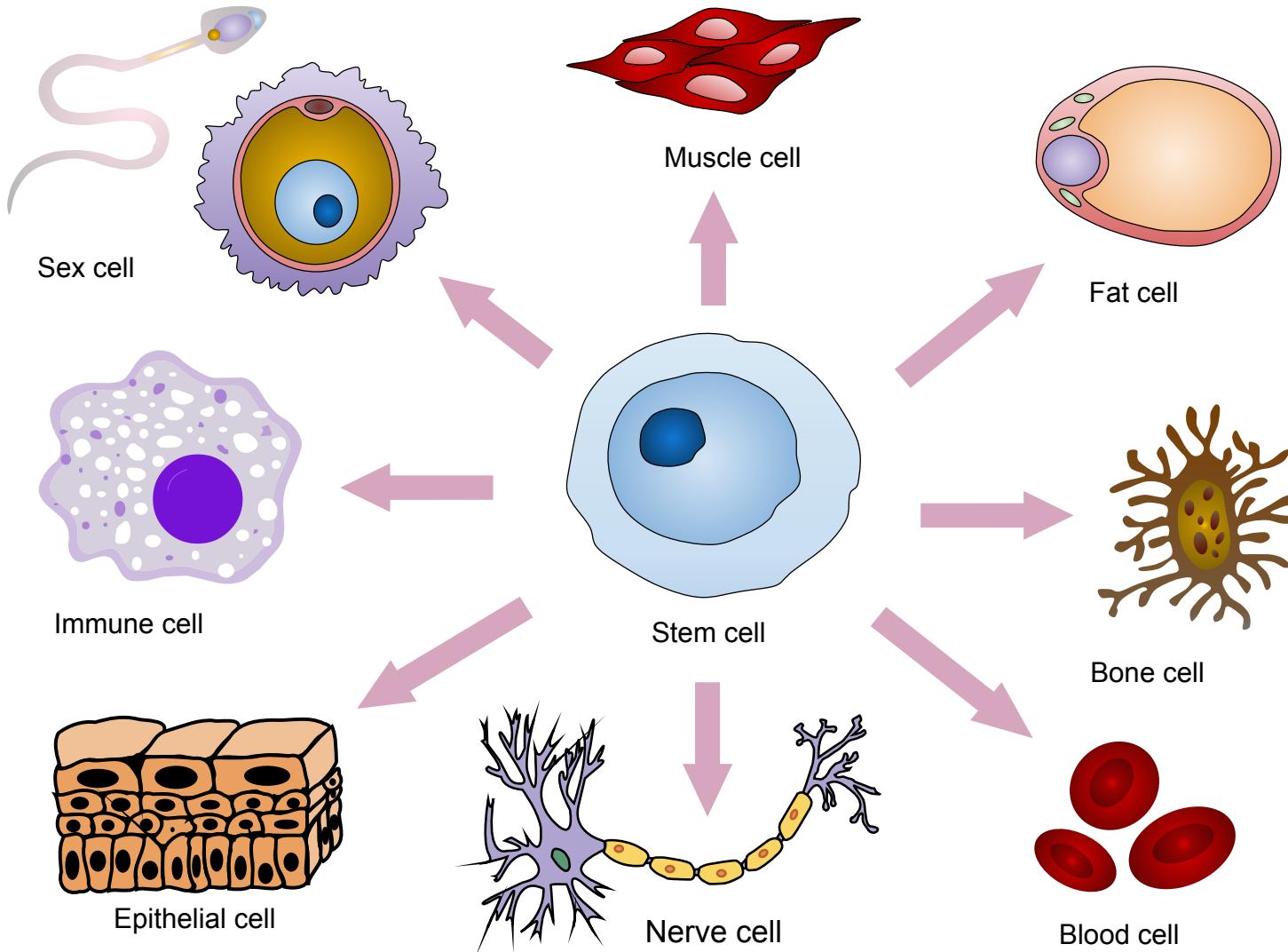
<https://kidshealth.org/en/parents/about-genetics.html>

Central dogma of Molecular Biology:

Gene Expression



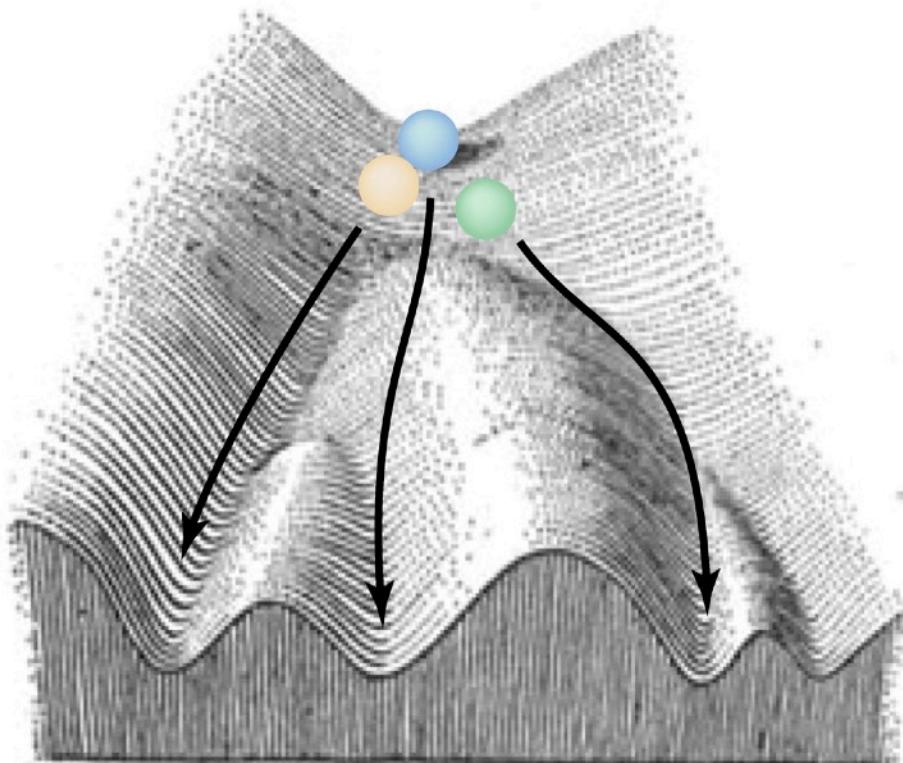
All the cell-types contain the same genome: Each cell-type expresses a different subset of genes



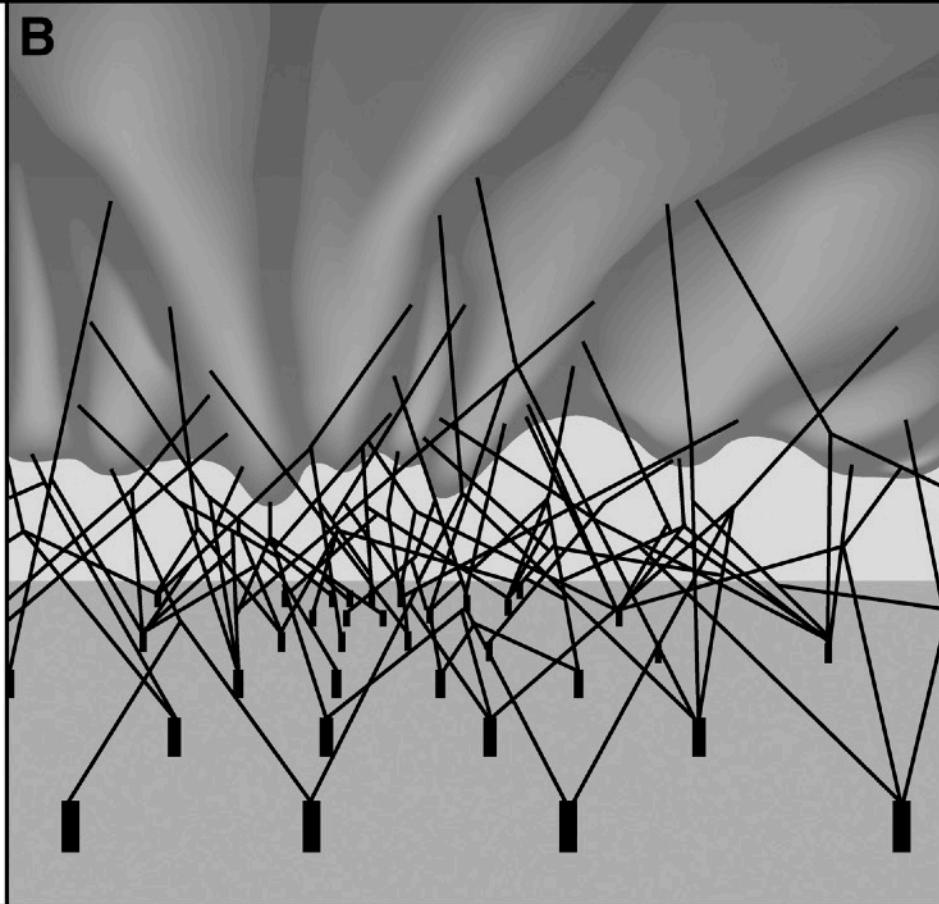
Human genome has about 20K-30K genes. Image from Wikipedia

Waddington development landscape and cellular differentiation

A

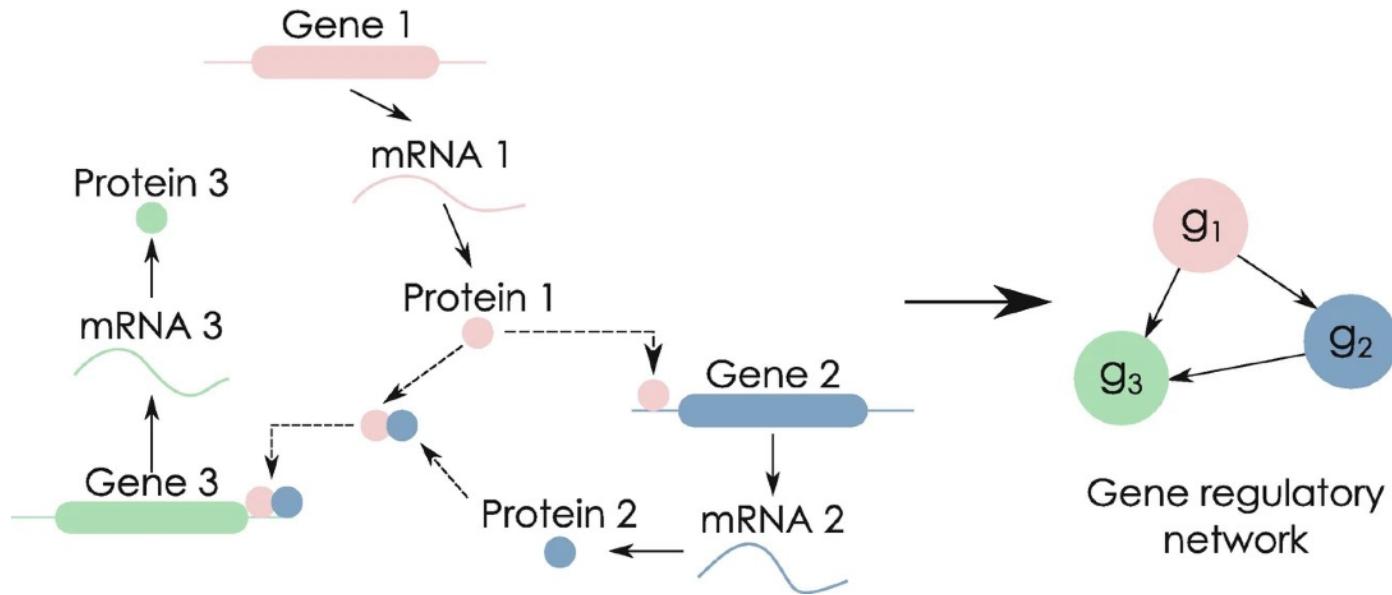


B



Adapted from *The Strategy of the Genes* (Waddington, 1957).
From Rajapogal & Stanger Developmental Cell 2016

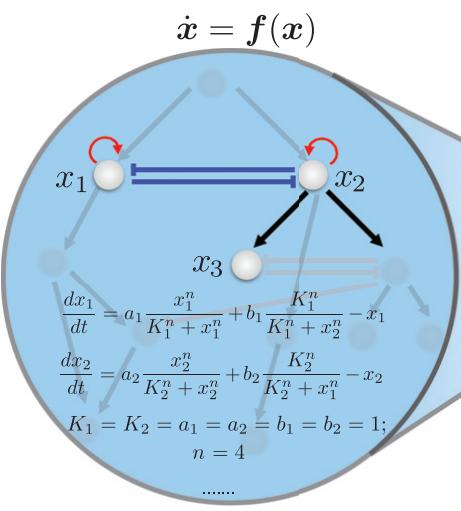
Gene Regulatory Networks (GRNs)



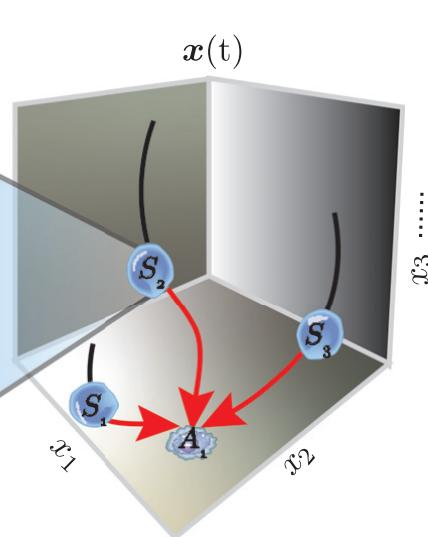
Huynh-Thu & Sanguinetti Gene Regulatory Network Inference, Methods in Molecular Biology Springer Protocols 2019

GRNs can be modelled as (stochastic, nonlinear) dynamical systems and cell states can be seen as (stable) fixed points of these systems

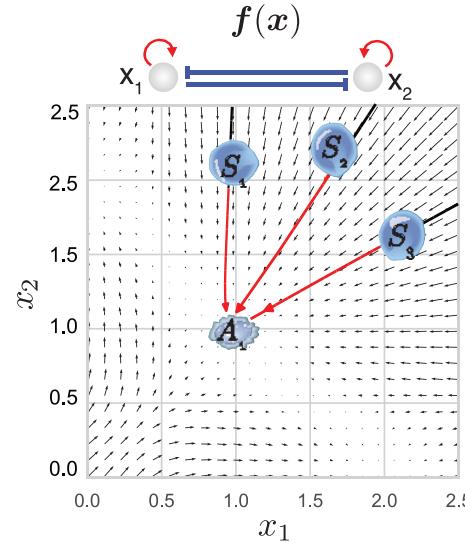
1) Regulatory network in single cells



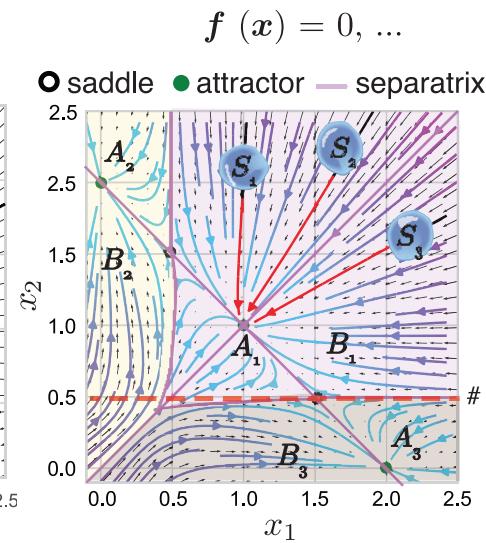
2) Cell dynamics in high dimension



3) Velocity vector field of cell dynamics



4) Topology of velocity vector field



Stochastic influences on phenotype

A



Fingerprints of
identical twins

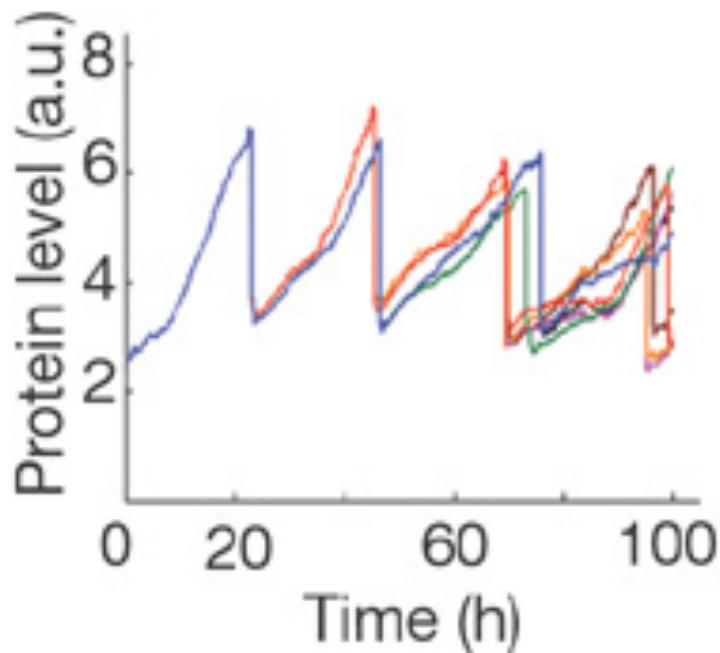
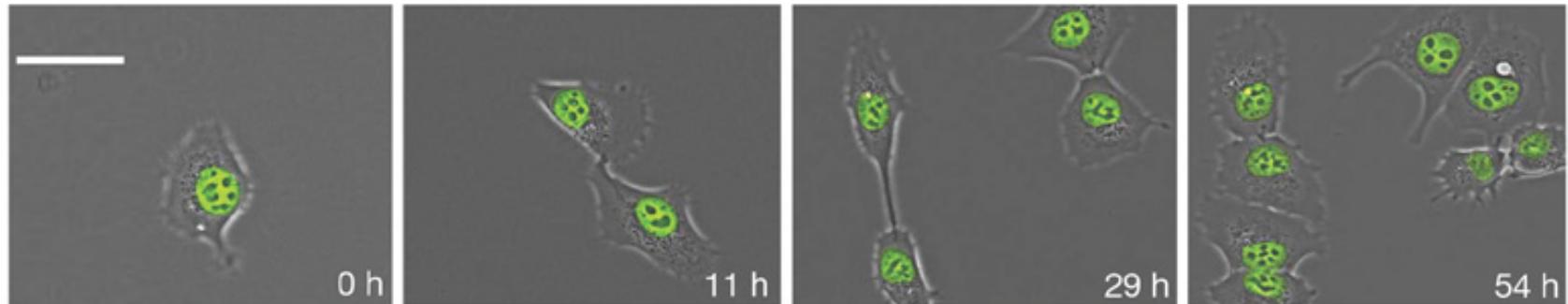
B



Cc the first
clone cat &
Rainbow, Cc's
genetic mother

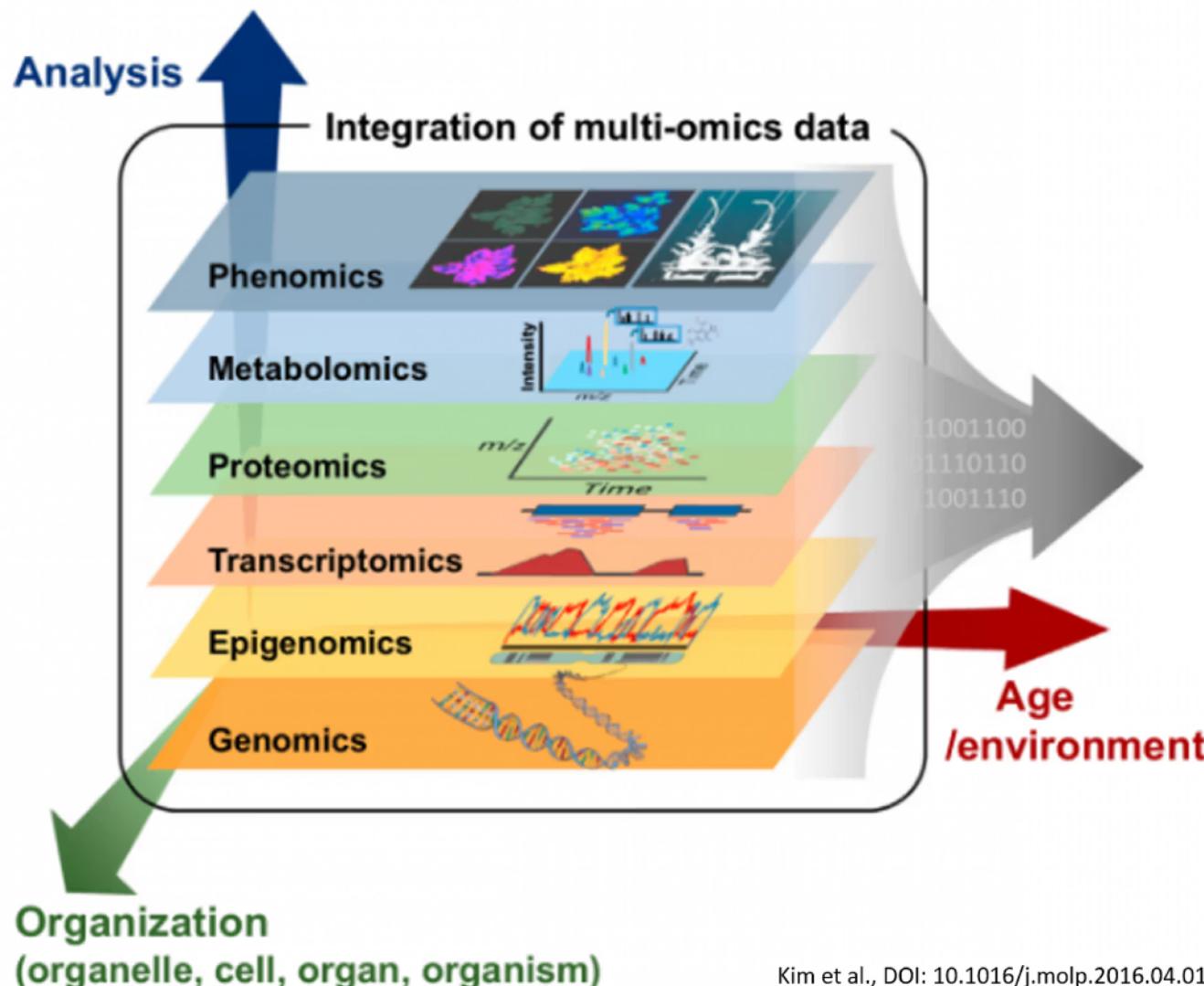
Raser and O'shea, Science 309, 5743 (2006)

Gene expression variability in single cells Fluorescent protein reporters



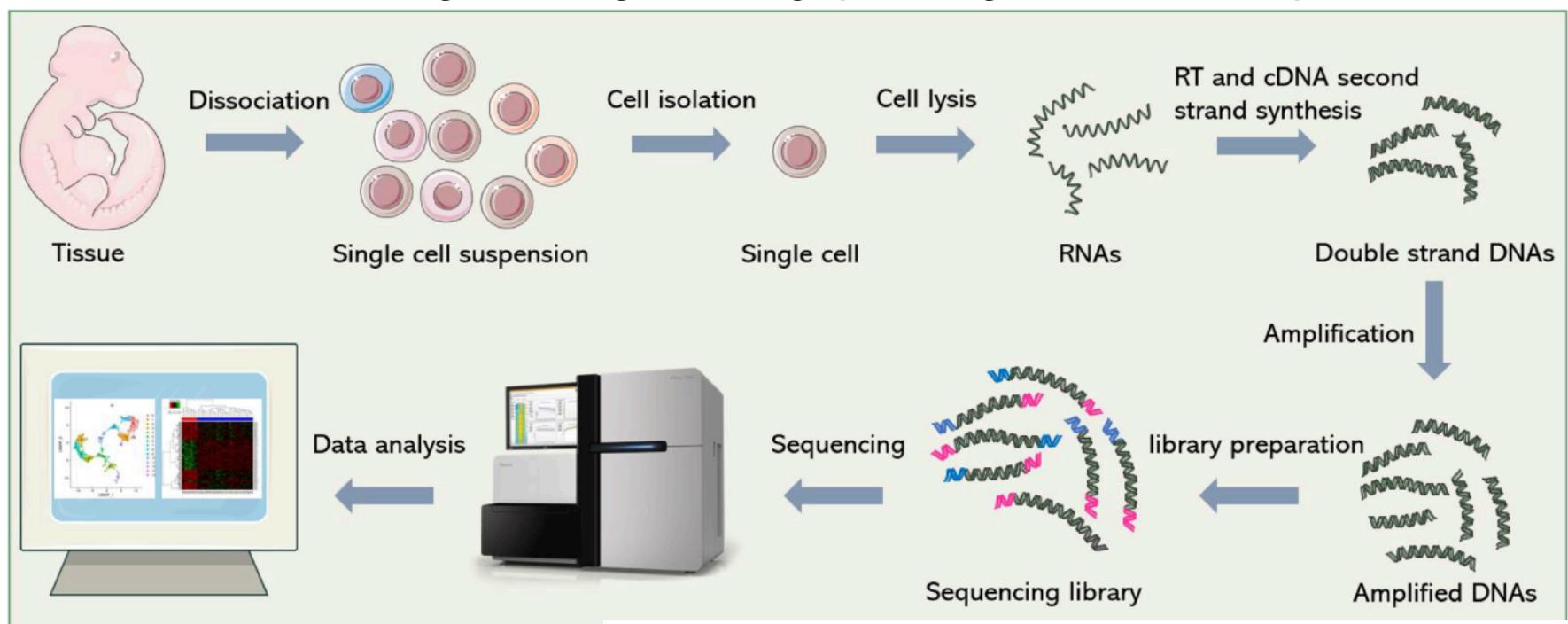
Sigal *et al.* Nature 444, 643 (2006)

High-throughput global biological data: Different types of omics data

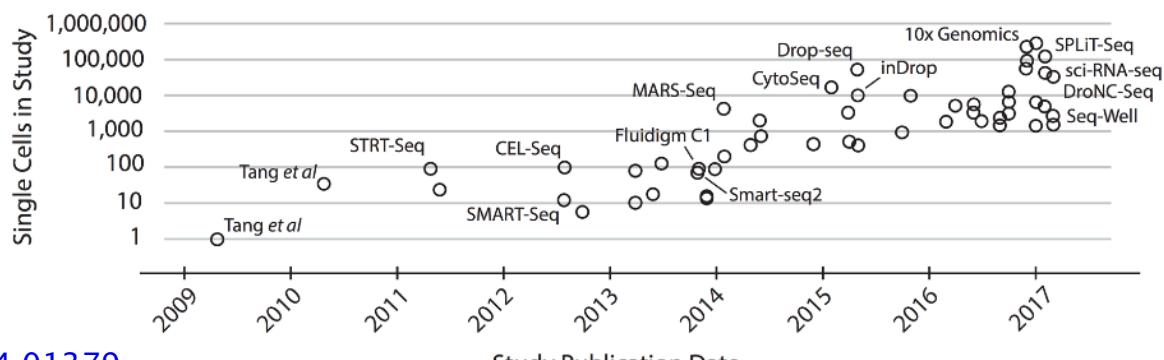


Development of single-cell RNA-sequencing

A method for global high-throughput single cell transcriptomics



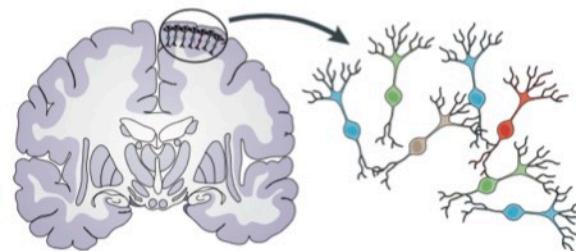
<https://encyclopedia.pub/entry/24618>



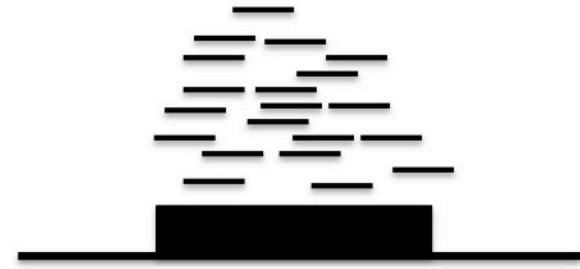
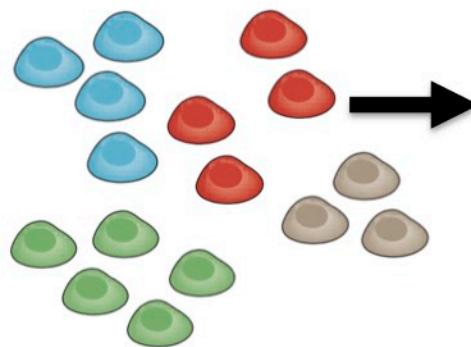
<https://doi.org/10.48550/arXiv.1704.01379>

Single-cell RNA-Seq (scRNA-Seq)

Tissue (e.g. tumor)



Isolate and sequence individual cells

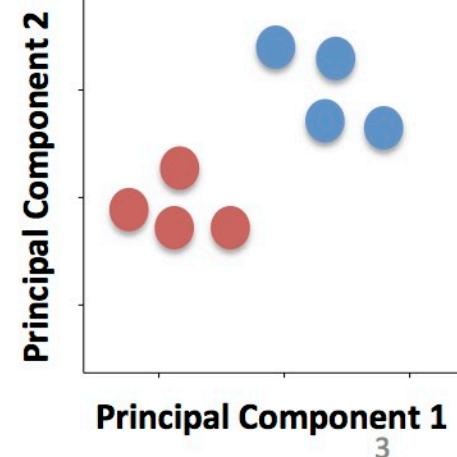
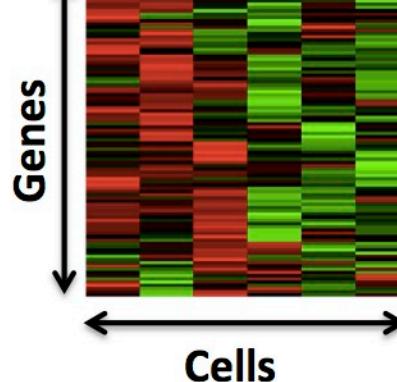


Gene 1
Cell 1

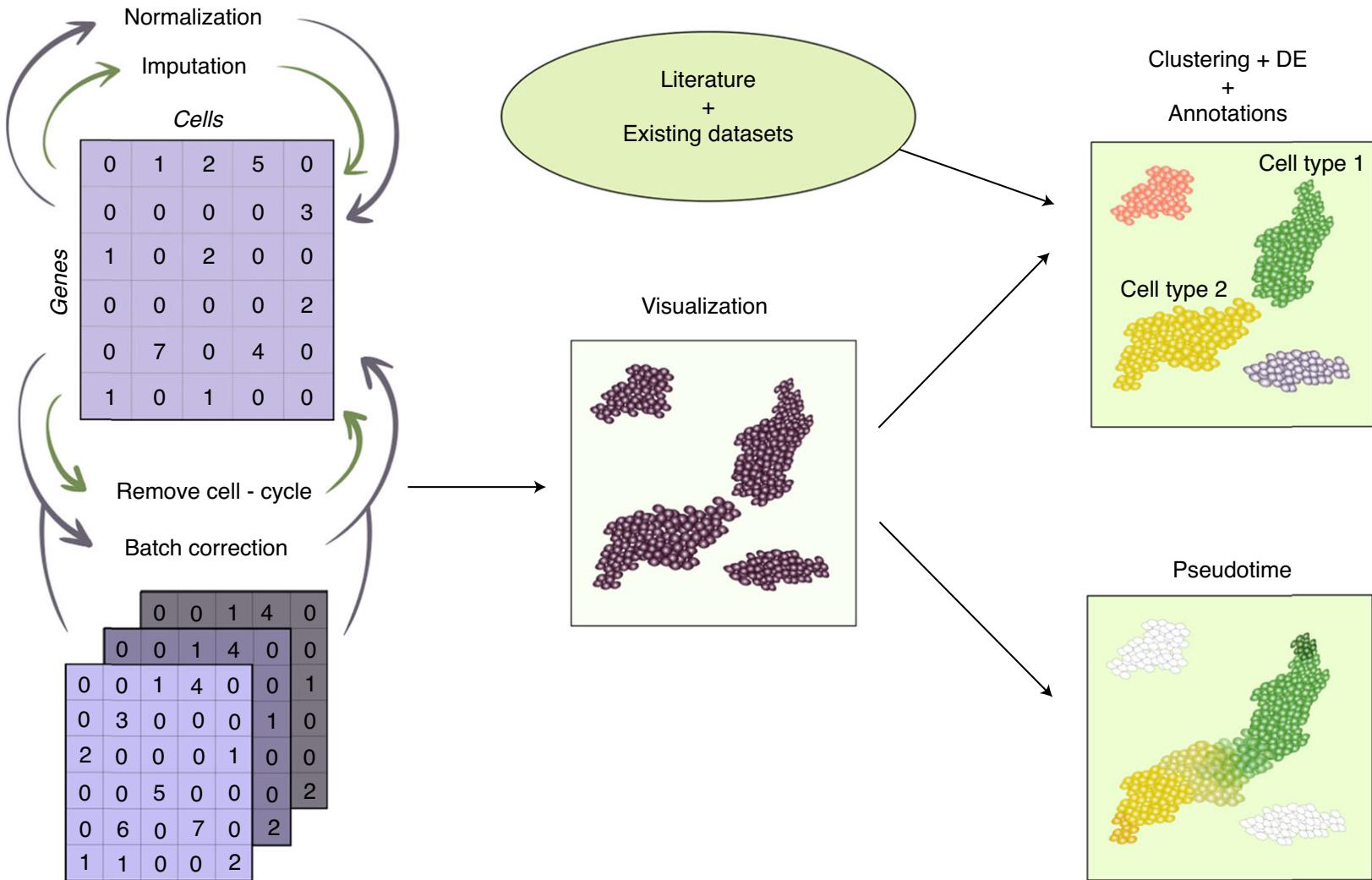
Read Counts

	Cell 1	Cell 2	...
Gene 1	18	0	
Gene 2	1010	506	
Gene 3	0	49	
Gene 4	22	0	
...			

Compare gene expression profiles of single cells



Mathematical, statistical and computational analysis of scRNA-Seq data

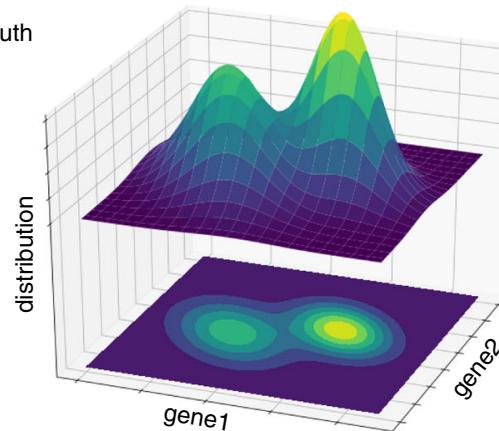


1. Experimental design

a

Sequencing budget allocation problem

Ground truth



Sequencing budget

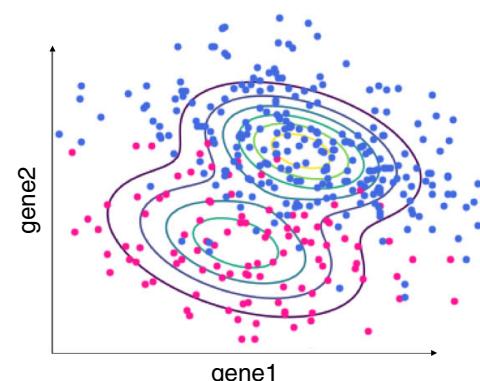
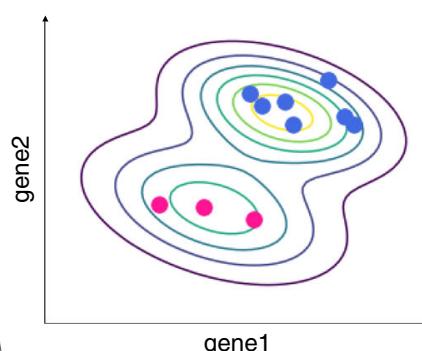


Deep sequencing
of a few cells

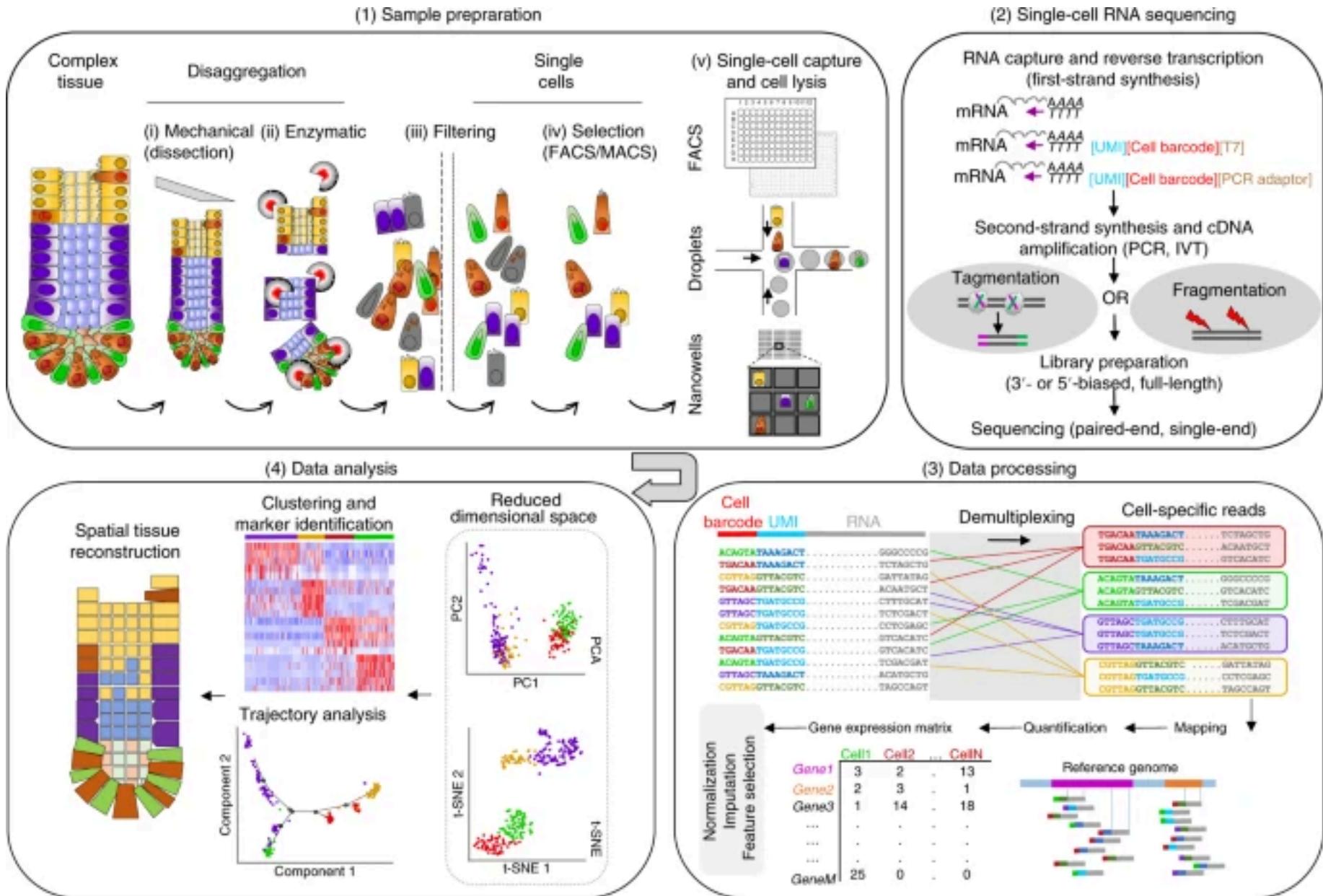


v.s.

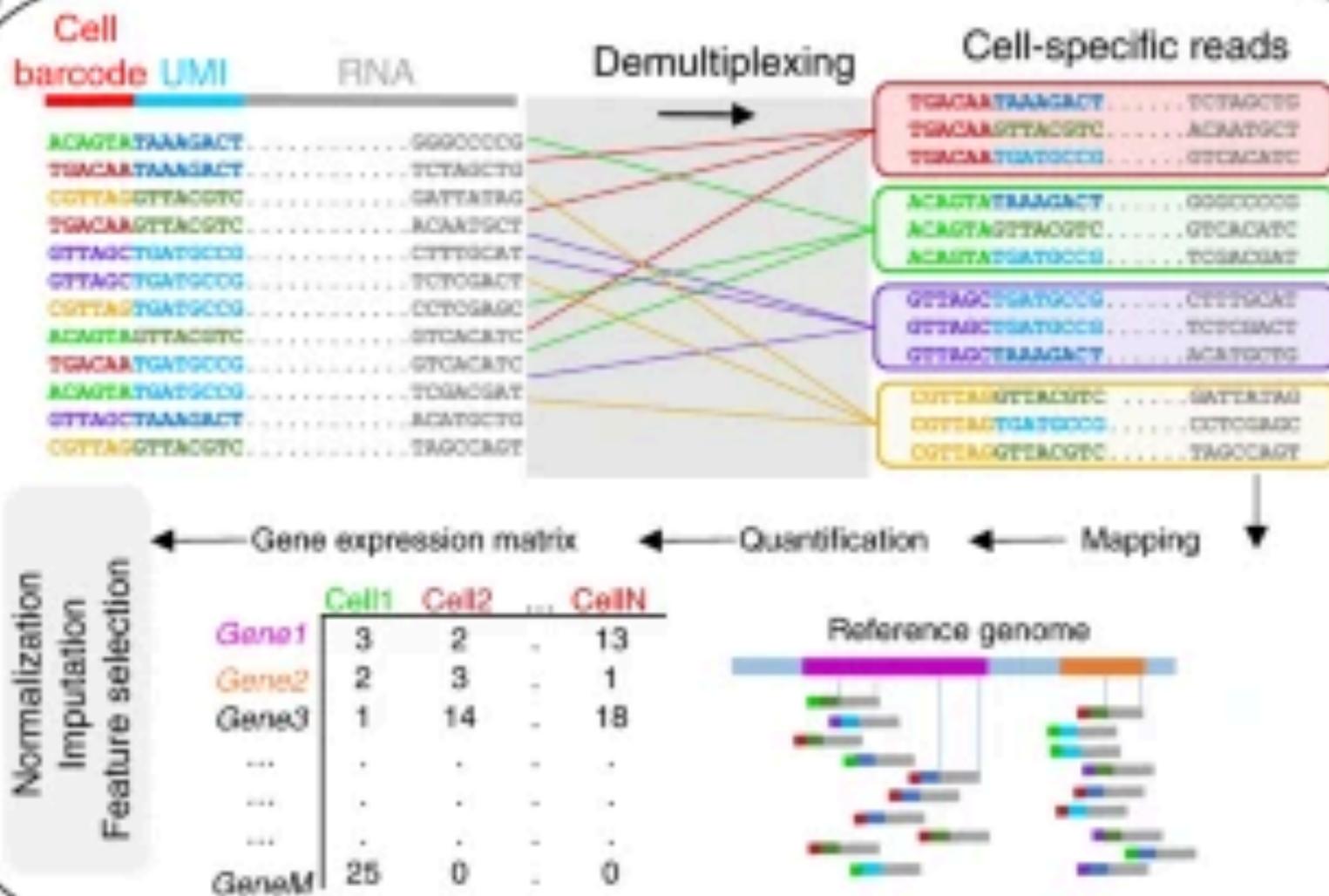
Shallow sequencing
of many cells



2. From raw reads to a count matrix



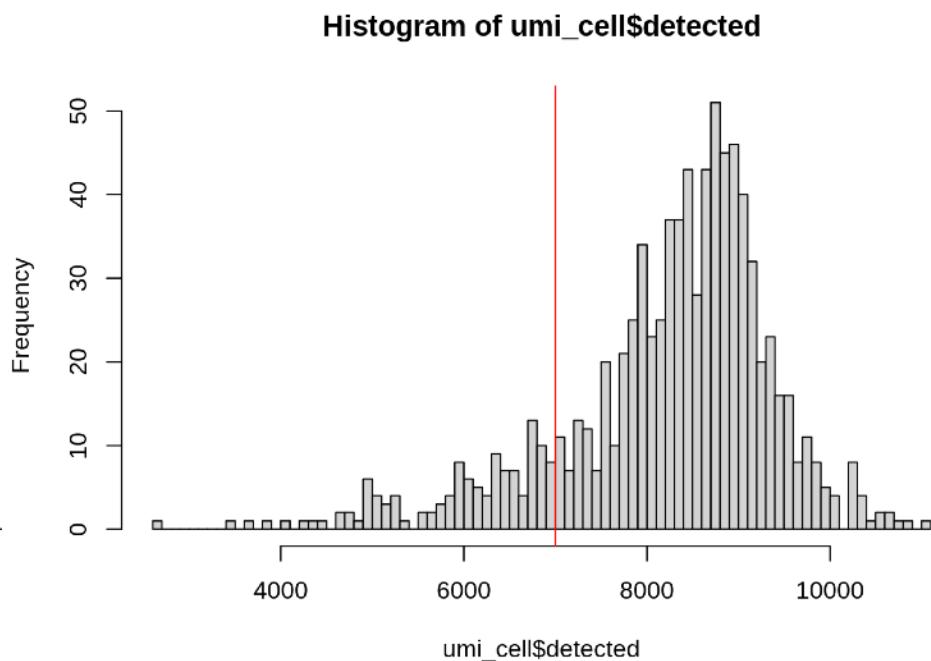
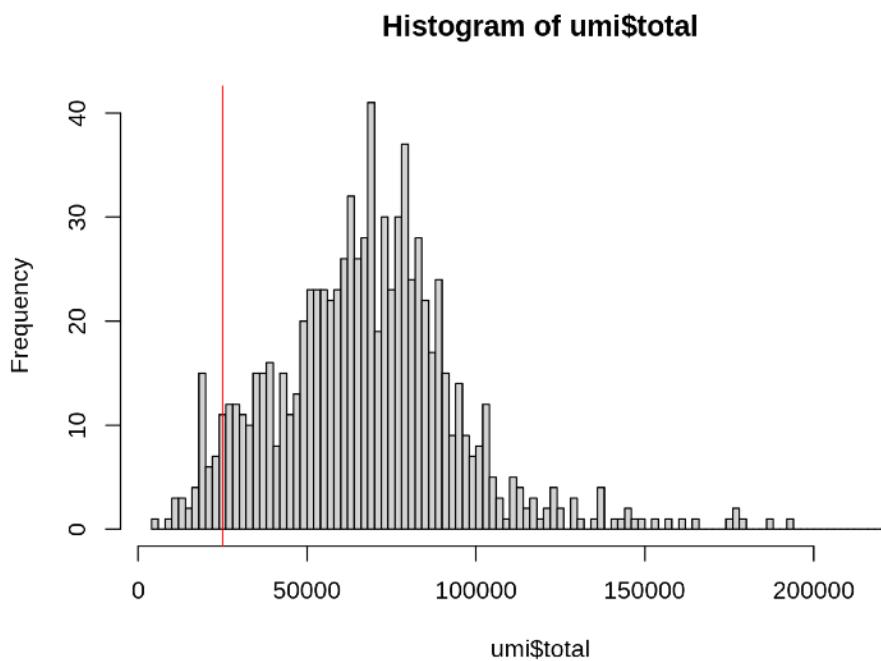
(3) Data processing



Sequence alignment to reference genome and demultiplexing unique molecular identifiers (UMIs) and cell barcodes to remove amplification biases and map reads to specific genes and cells.

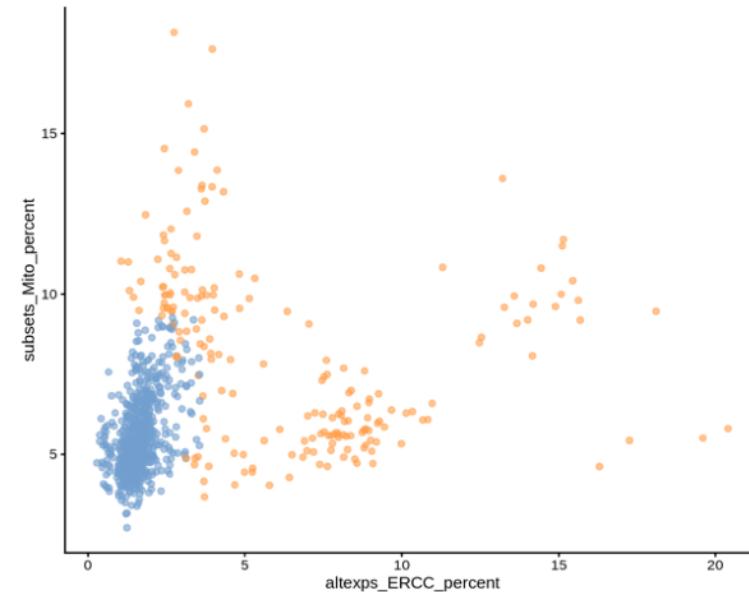
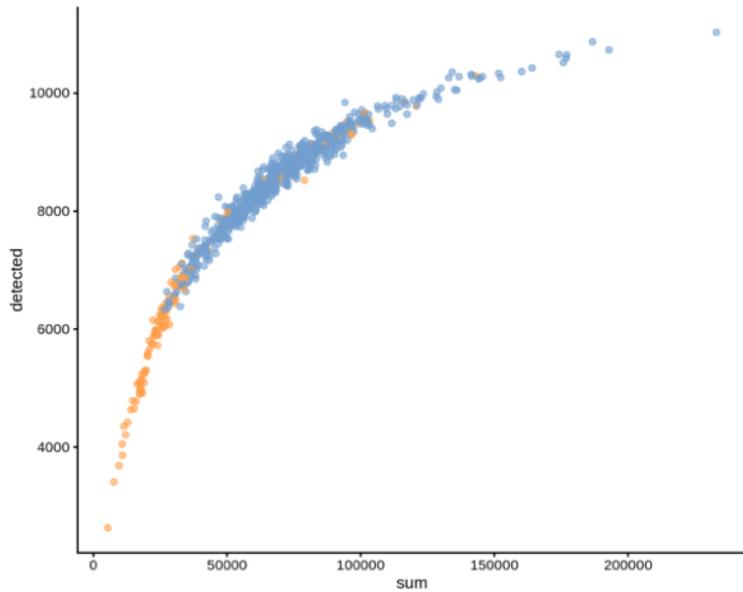
3. Basic quality control

Filtering cells based on total number of transcripts per cell and detected genes

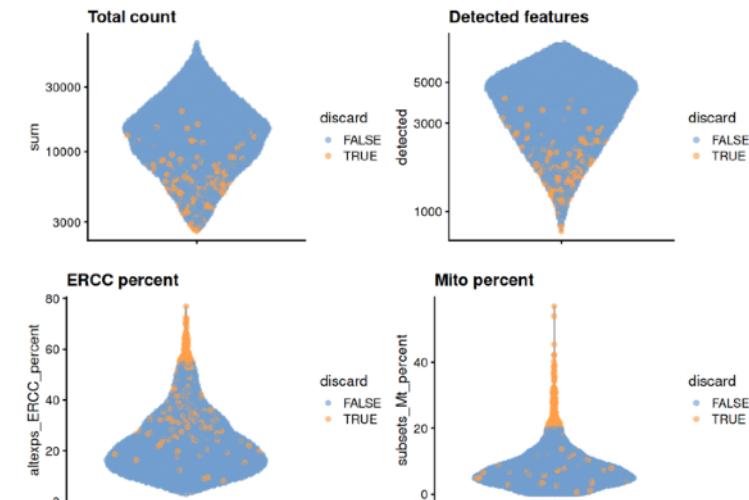


3. Basic quality control

Additional information such as percentage of mitochondrial genes and use of control spike-in (ERCCs) are common for filtering low-quality cells

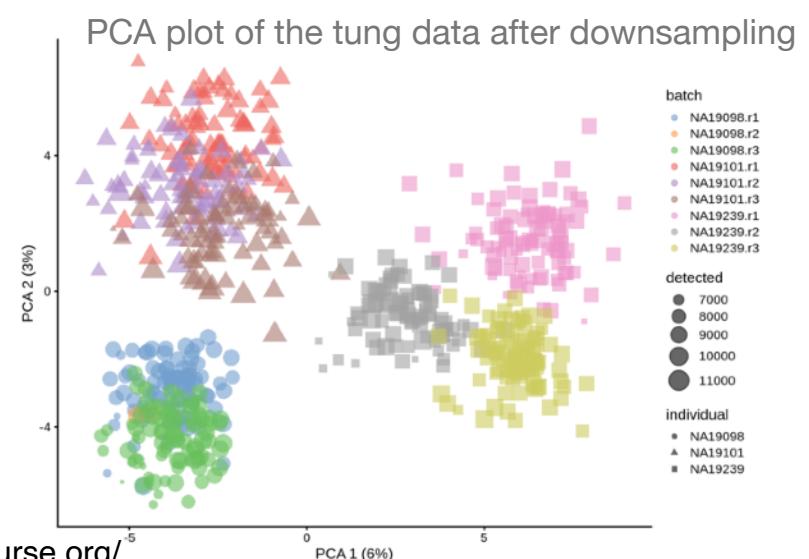
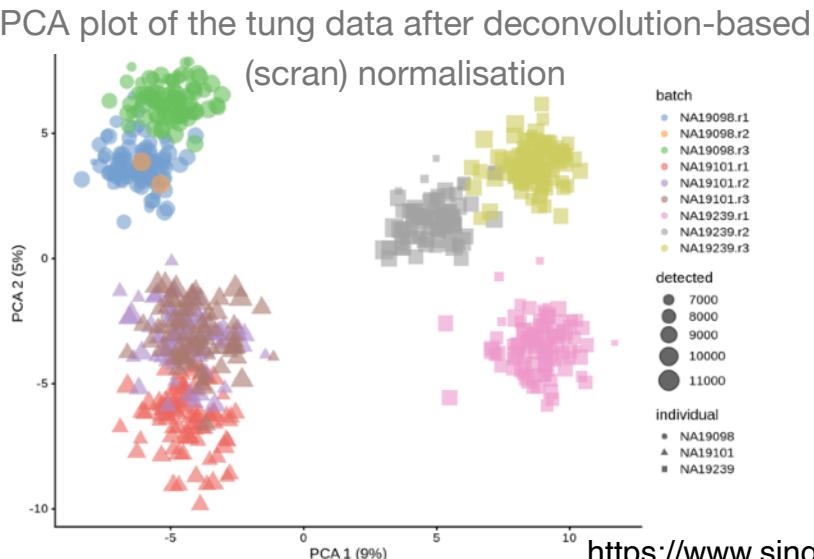
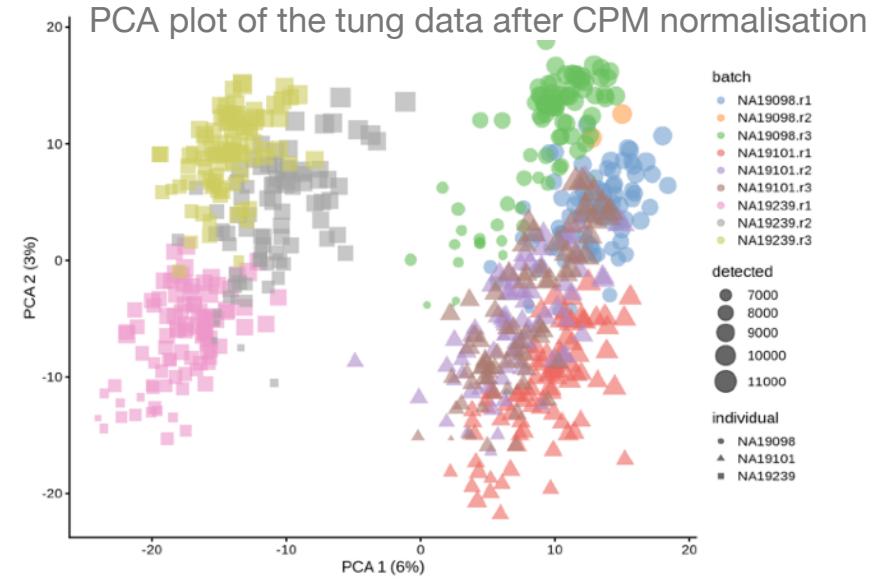
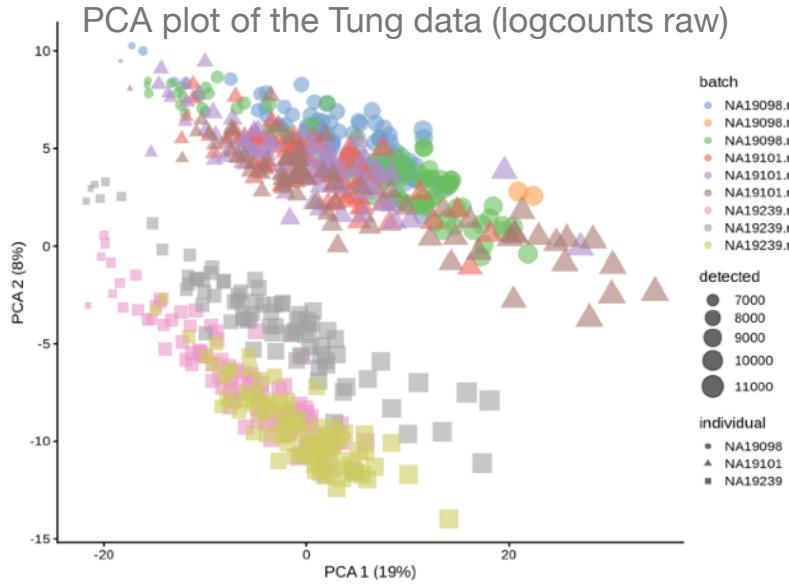


<https://www.singlecellcourse.org/>



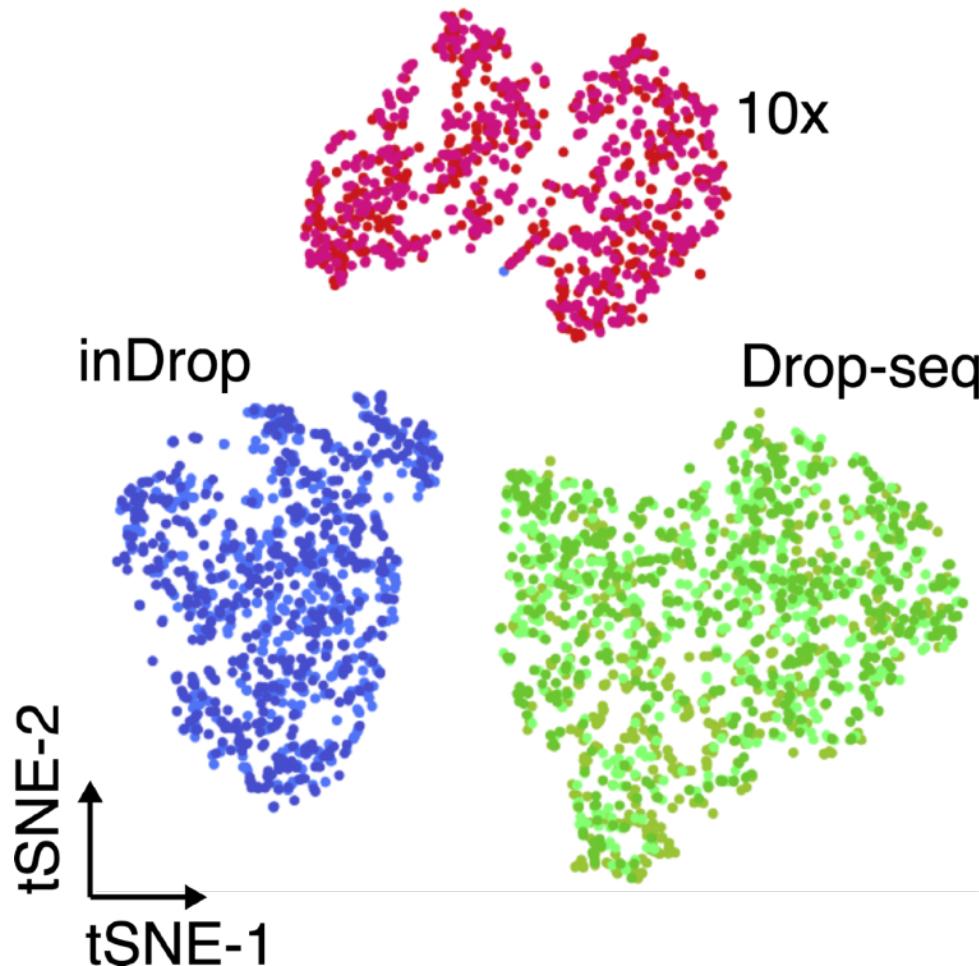
4. Normalisation

The simplest approach would be scaling normalisation by cell-specific size factor such as total transcripts per cell.



5. Batch Correction

Related to normalisation and dealing with sparsity and dropouts

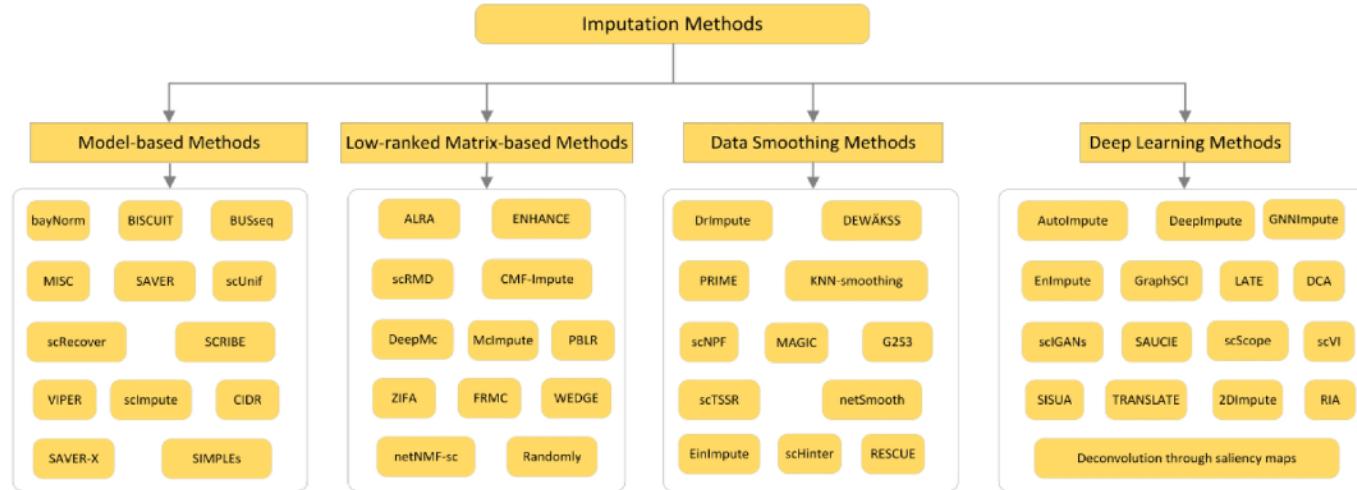


The same cell population was sequenced with three different single-cell protocols (colours).

Adapted from [Zhang et al.. Molecular Cell 73, 130 \(2019\)](#)

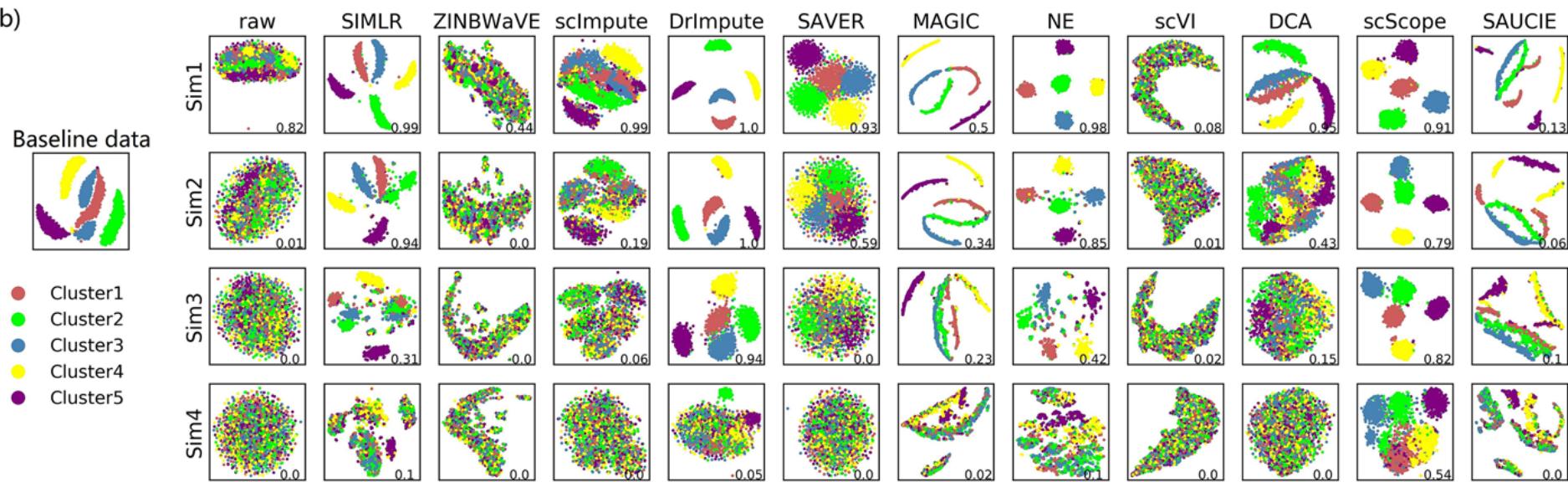
6. Imputation

Related to normalisation and dealing with sparsity and dropout



<https://doi.org/10.3390/app122010684>

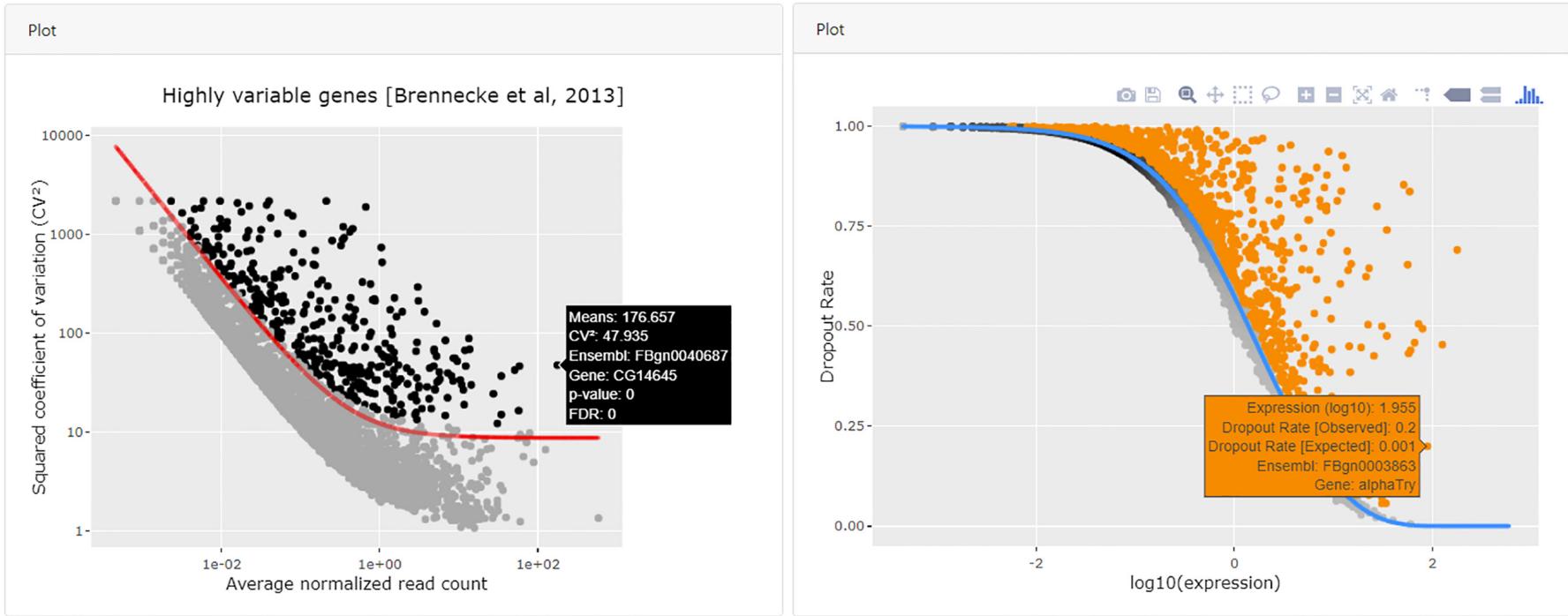
(b)



<https://doi.org/10.1186/s12859-023-05417-7>

7. Future selection

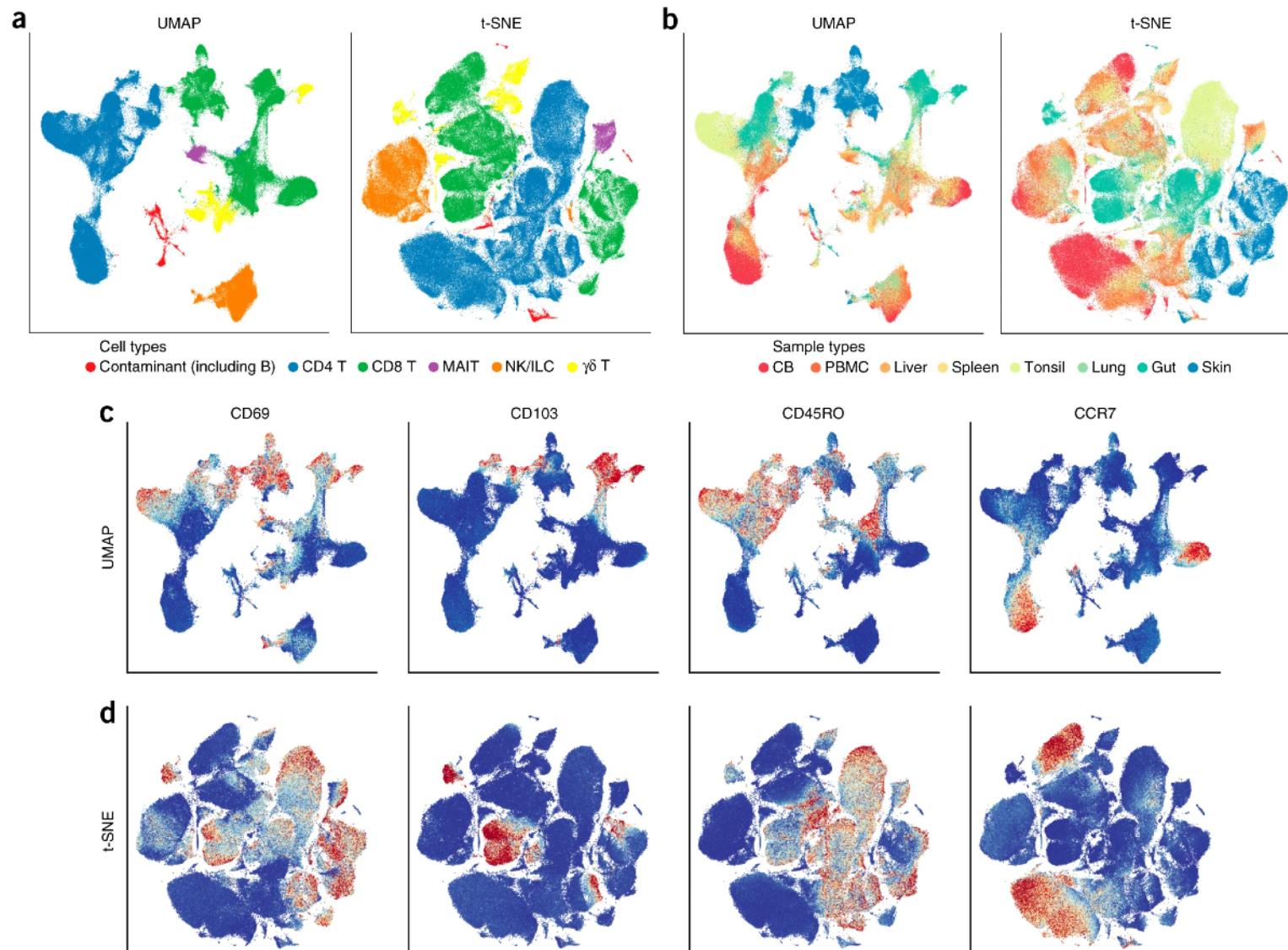
e.g. using Highly variable genes



doi: 10.1093/nar/gkaa412

8. Dimensionality reduction

Promise of UMAP and t-SNE

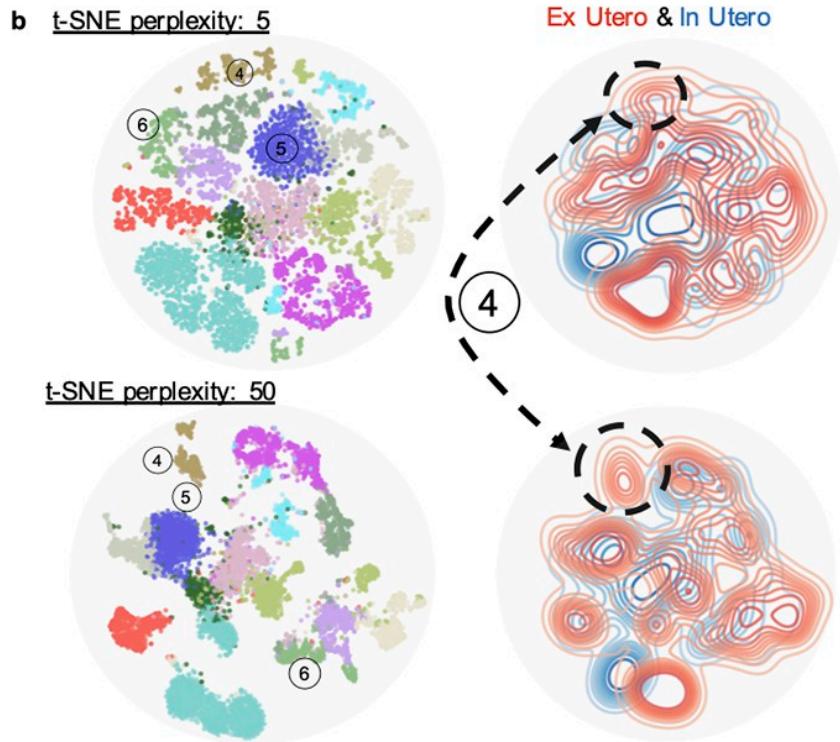
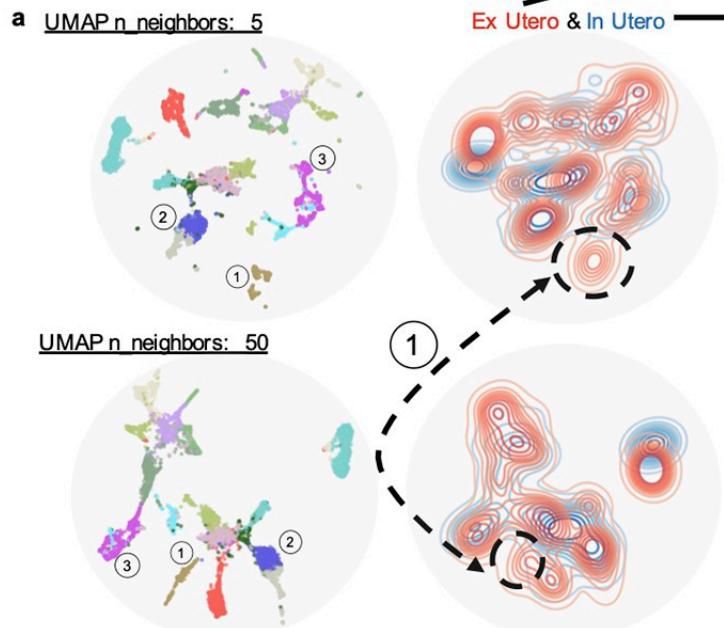


8. Dimensionality reduction

Problems with UMAP and t-SNE

	Necessary Properties		
	Local	Global	Distance
Modality-Mixing, Integration & Reference Mapping	◆	◆	
Cluster Validation & Relationships		◆	◆
Density-Based Visuals & Marker Analysis	◆	◆	◆
Trajectory Inference & Continuous Relationships	◆		◆

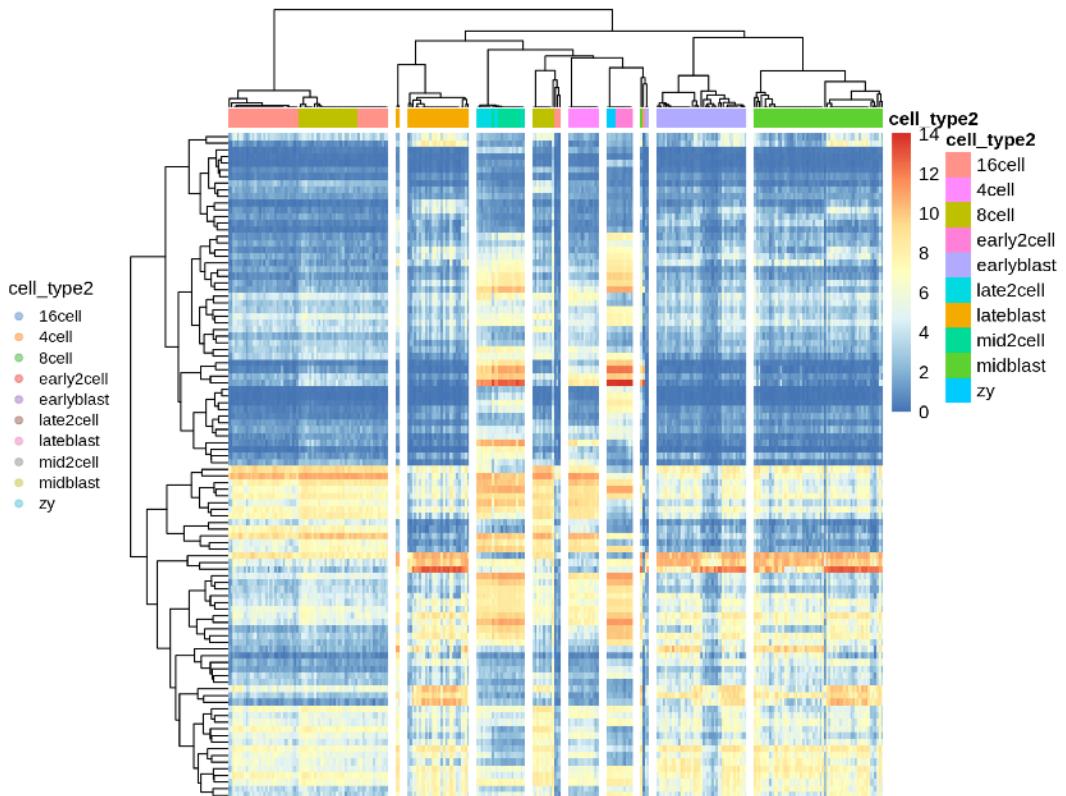
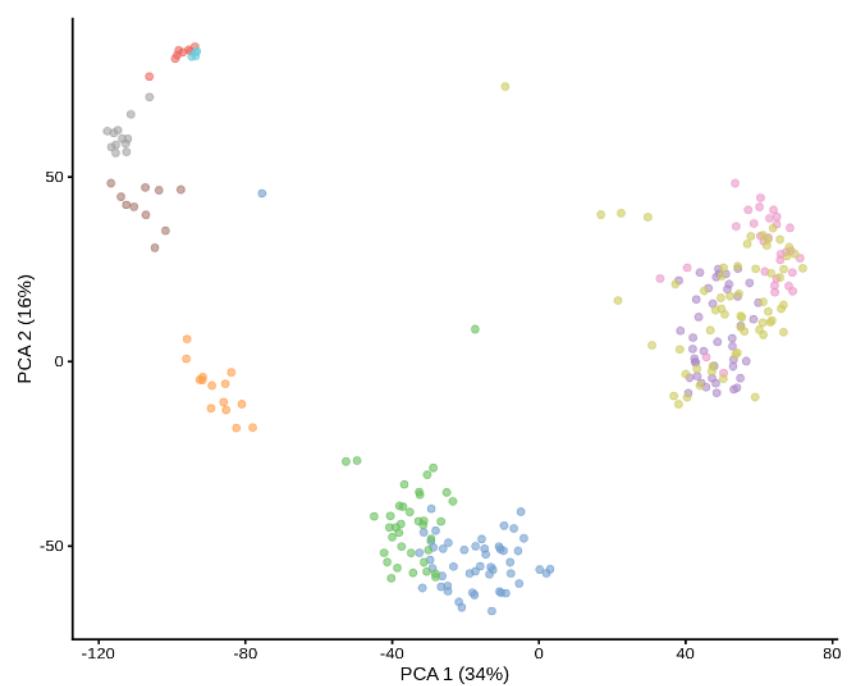
◆ - Yes ◇ - Optional



<https://doi.org/10.1371/journal.pcbi.1011288>

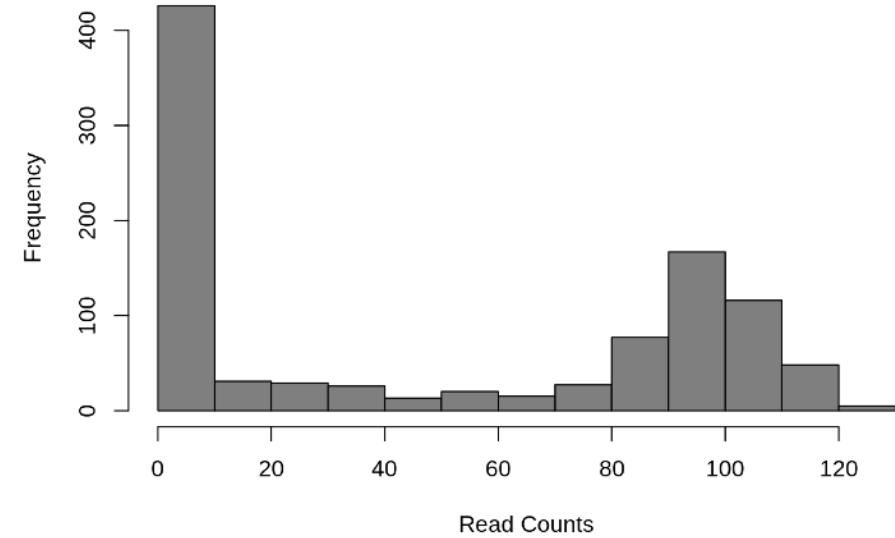
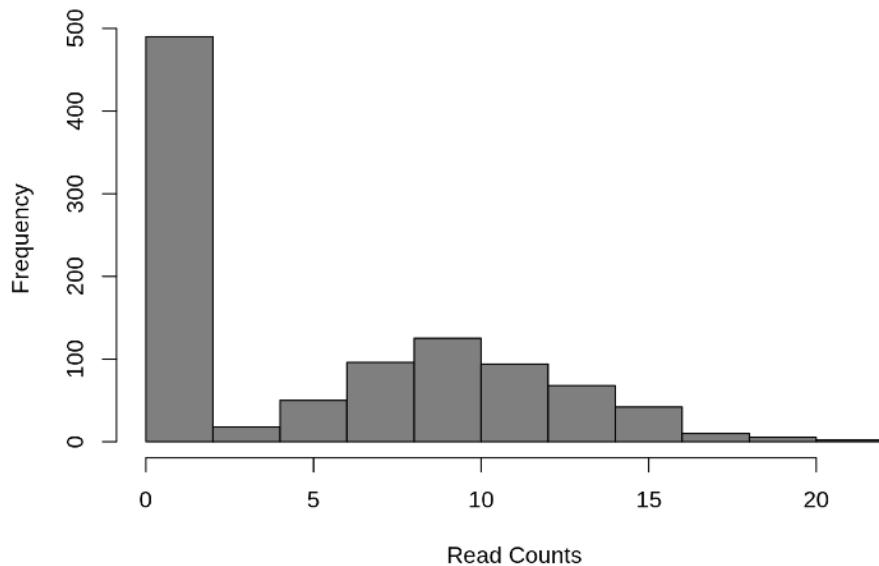
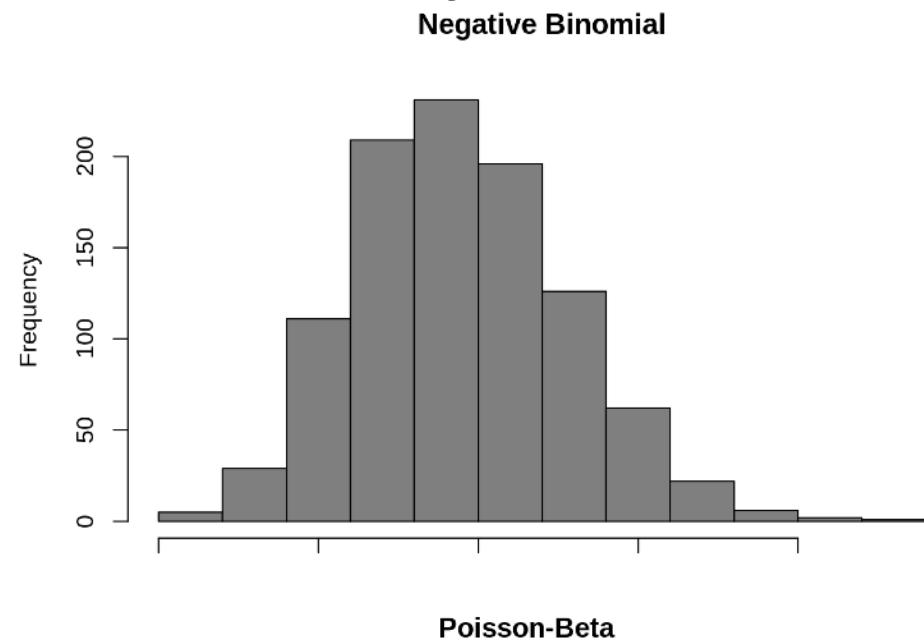
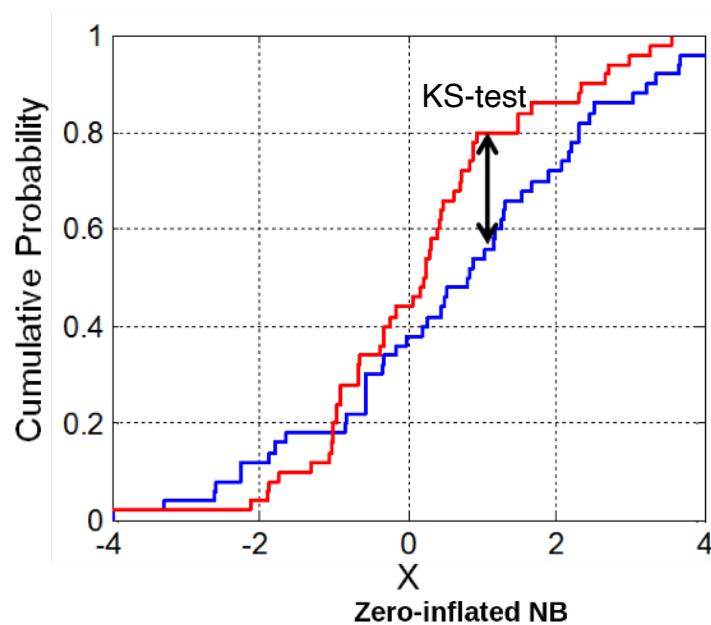
9. Clustering

Classification: cell type allocation-based on marker gene expressions

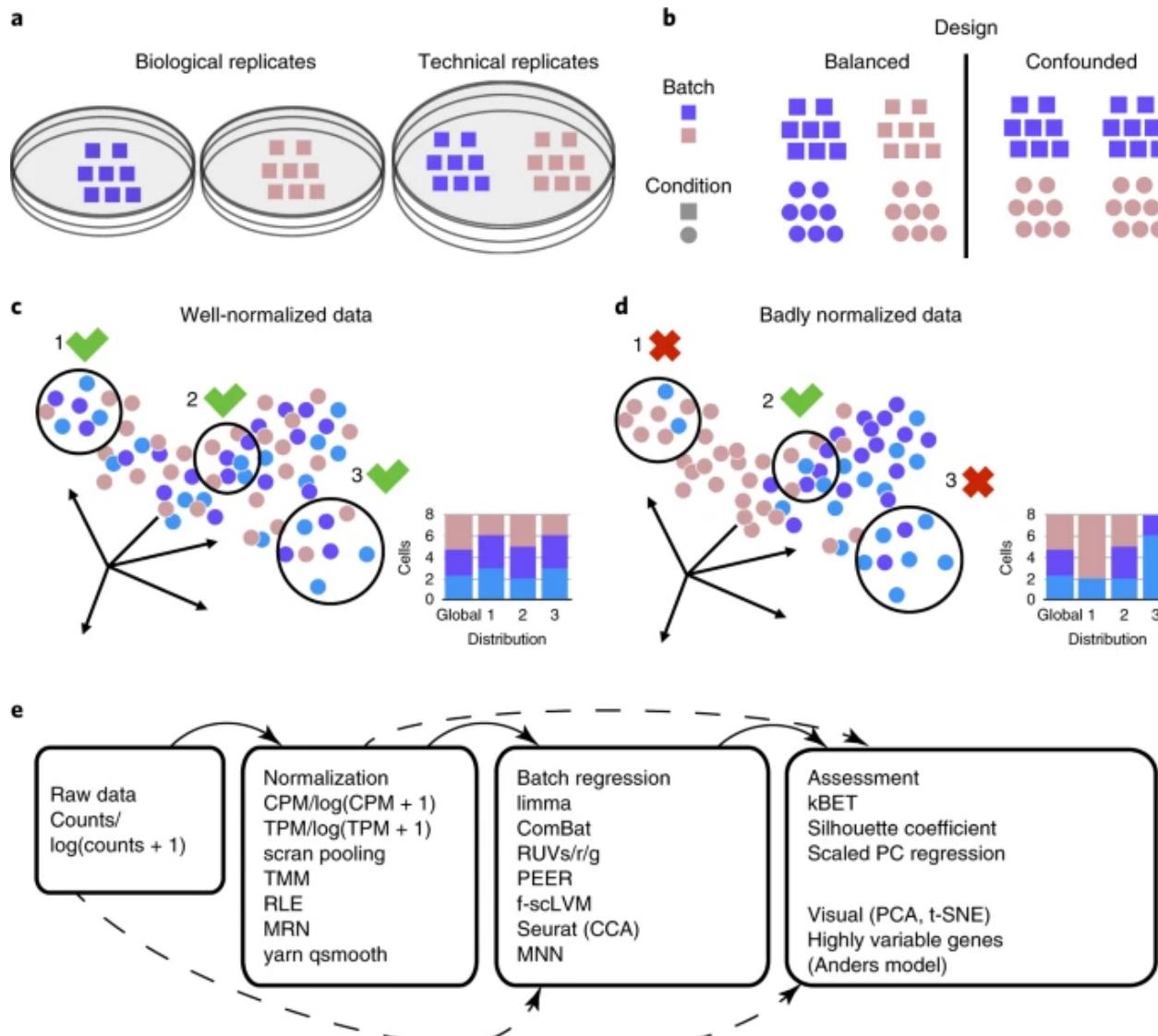


10. Differential Expression Analysis

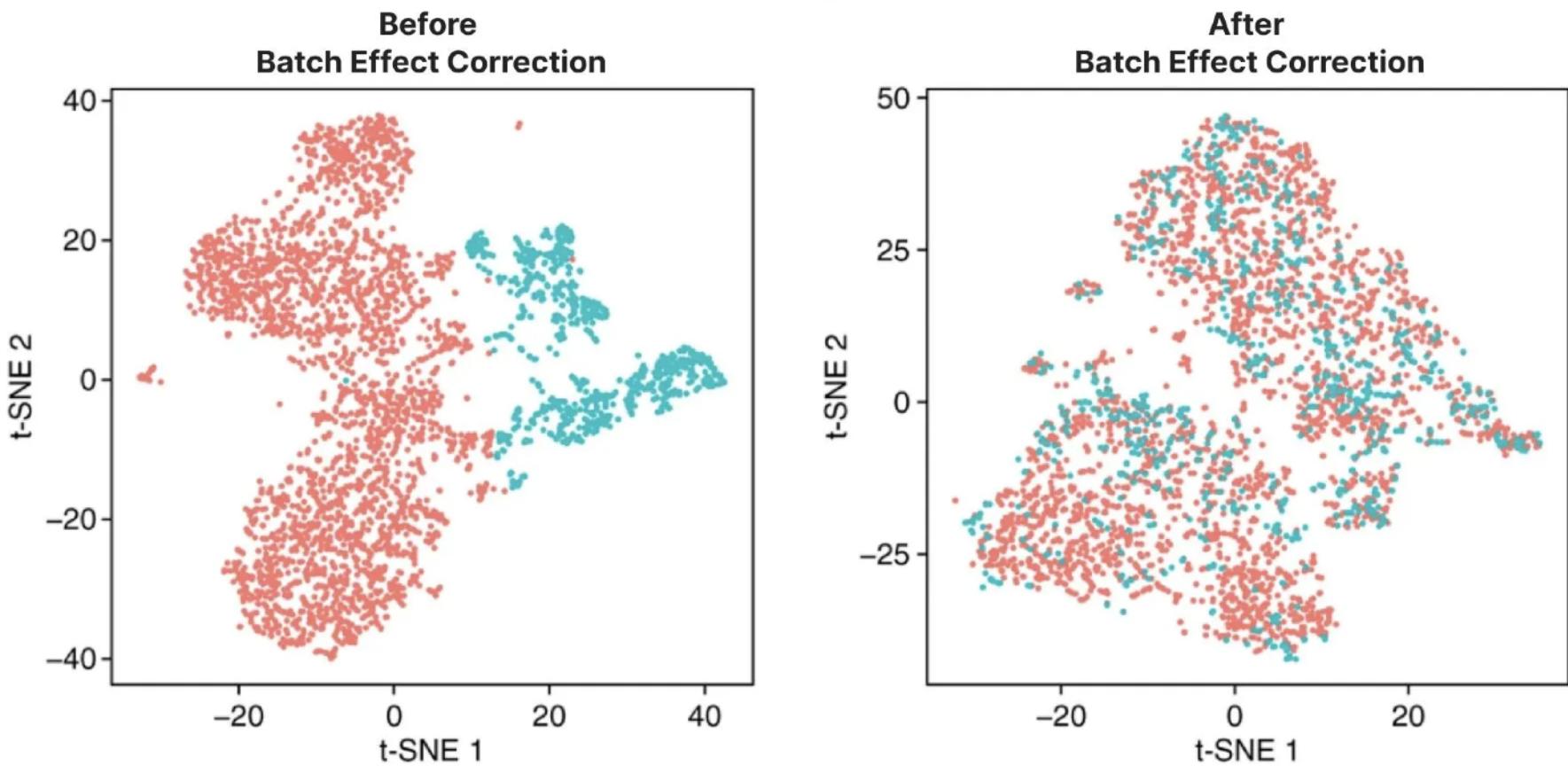
Parametric and non-parametric (+ counting for FDR)



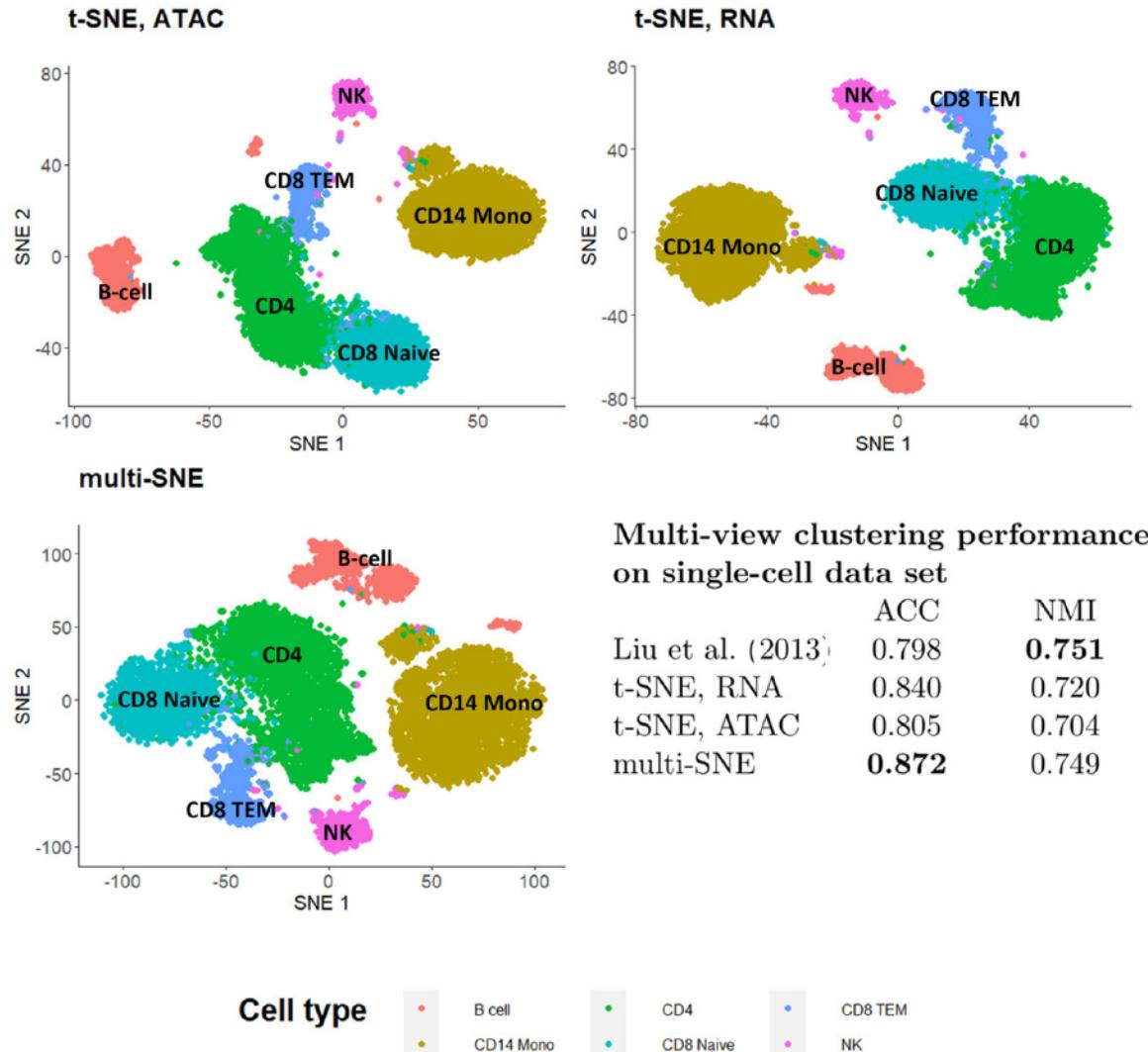
11. Data integration (and batch correction)



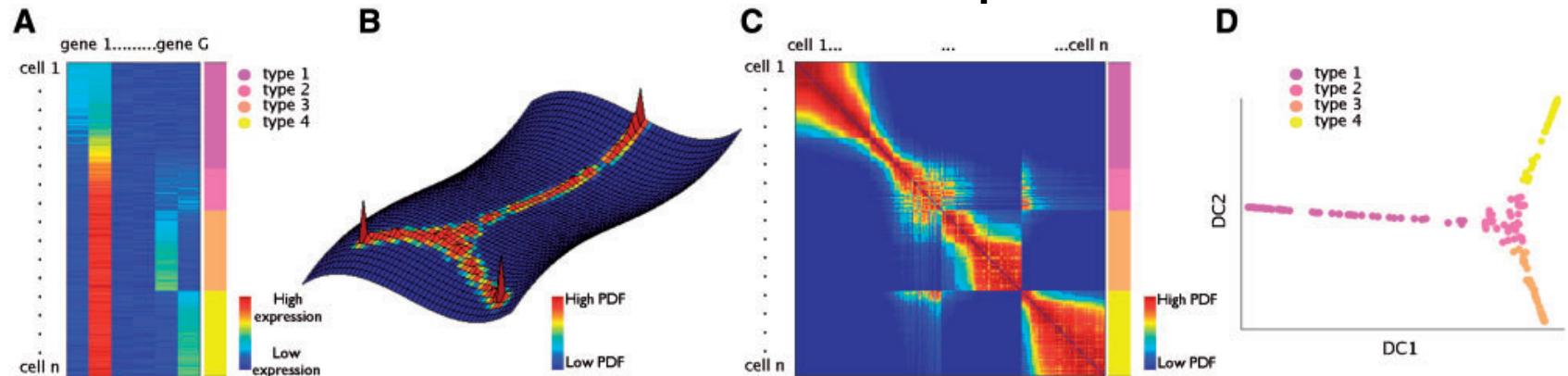
11. Data integration (and batch correction)



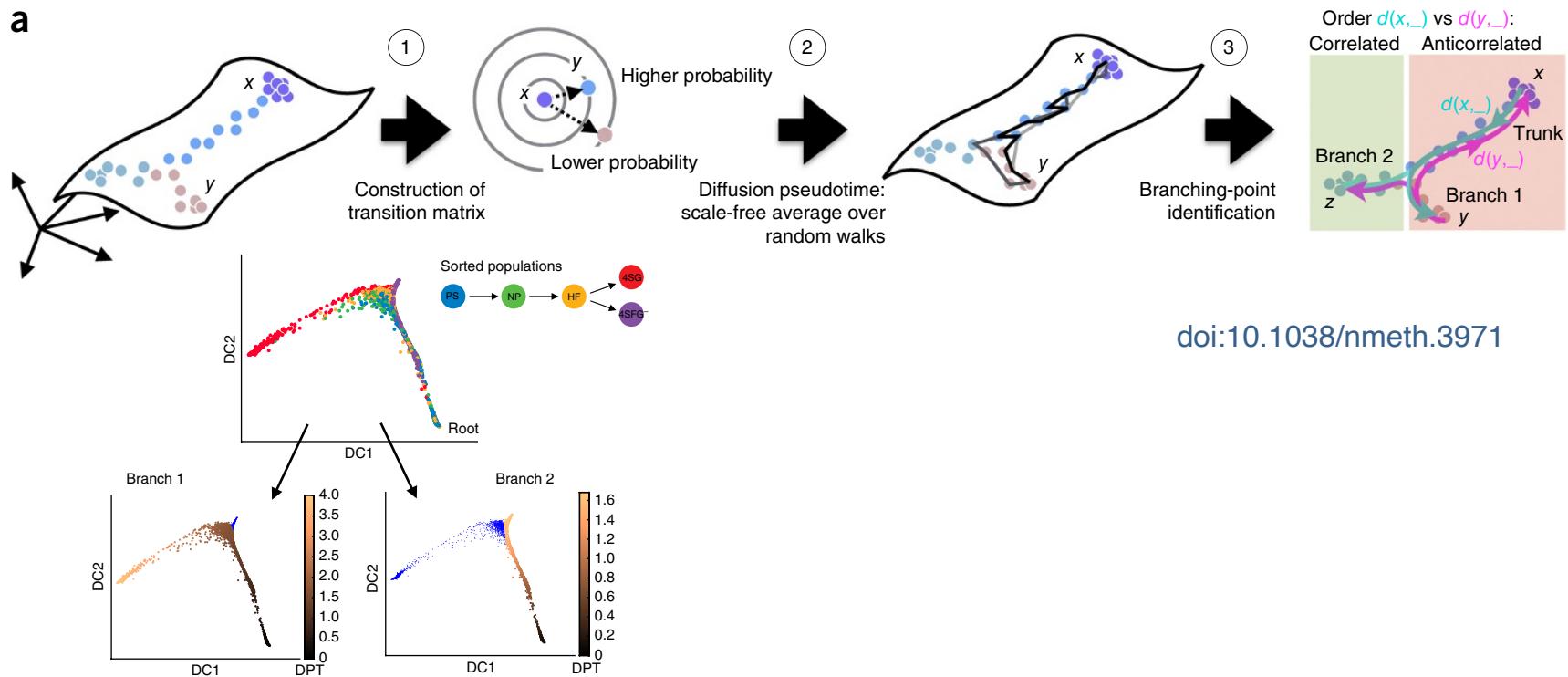
12. Multi-modal data integration



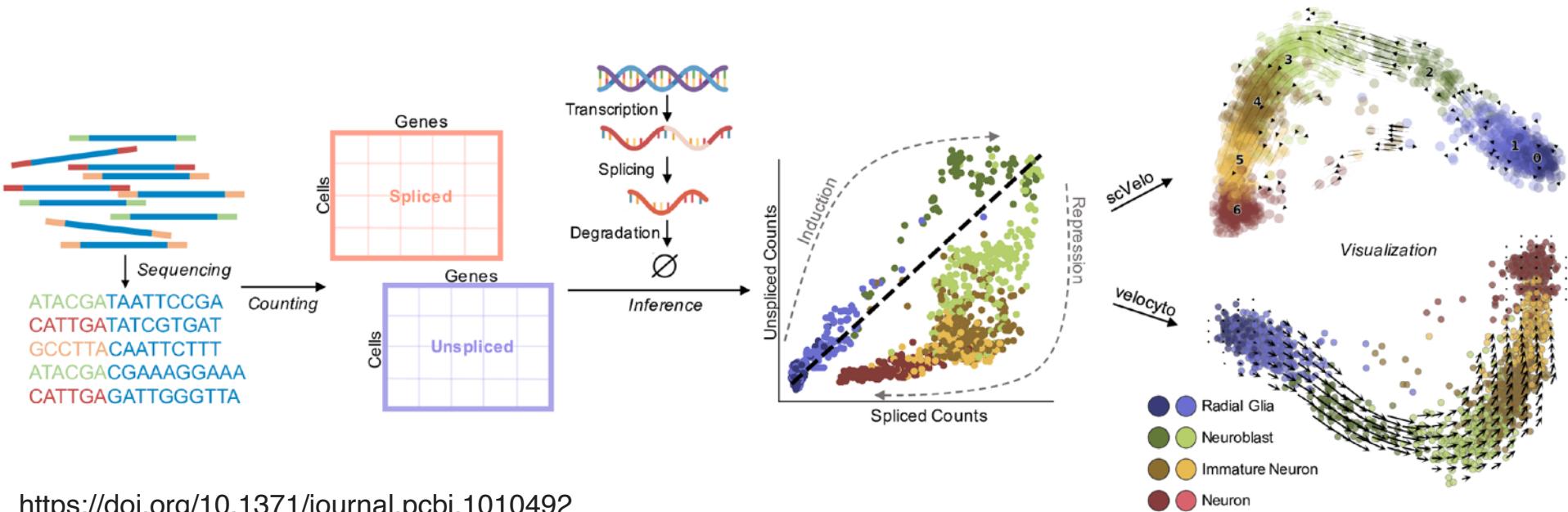
13. Trajectory-inference and pseudo-time ordering Diffusion maps



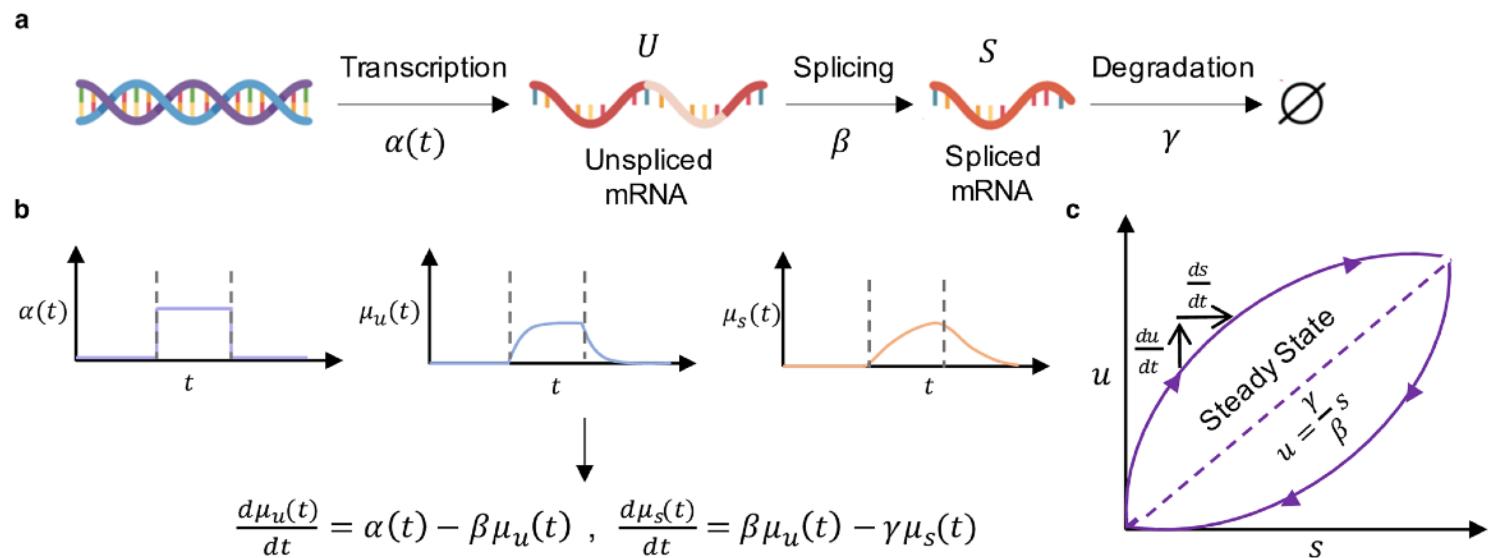
doi: 10.1093/bioinformatics/btv325



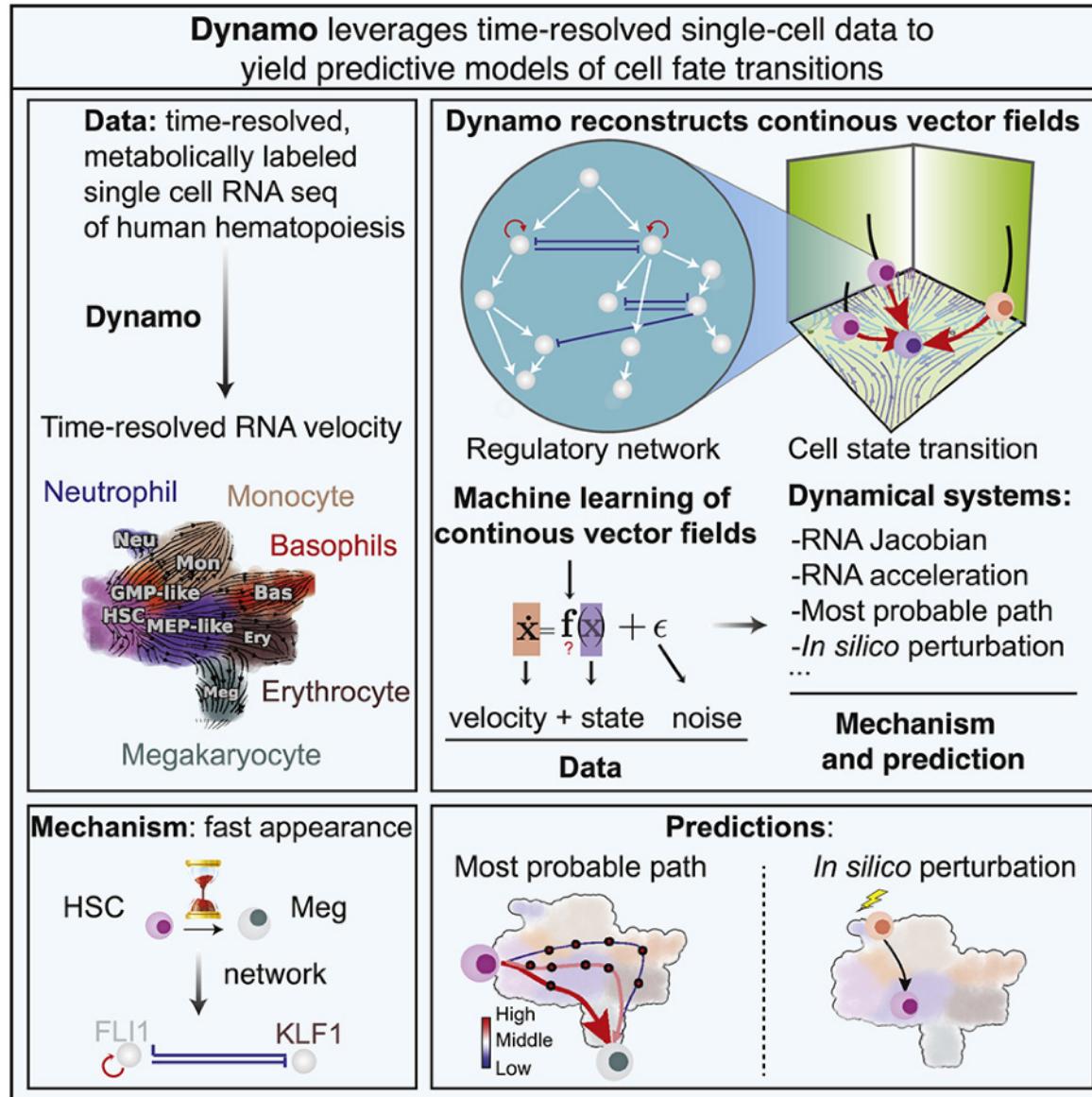
14. RNA-velocity analysis



<https://doi.org/10.1371/journal.pcbi.1010492>

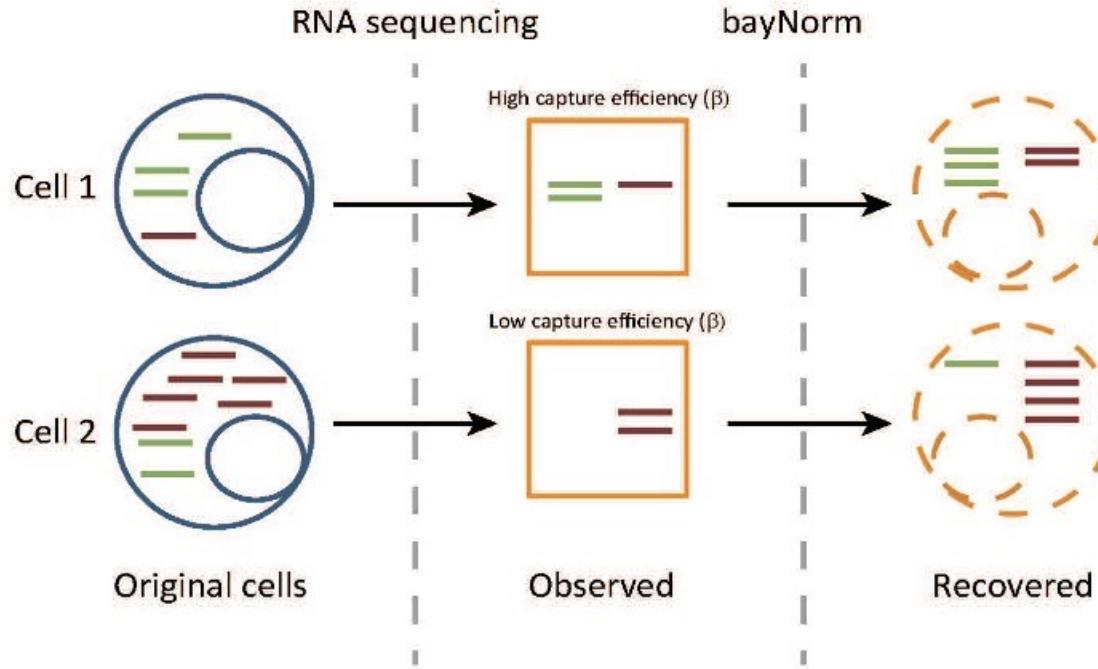


15. Gene-regulatory network inference



Some research examples from the
work in my group

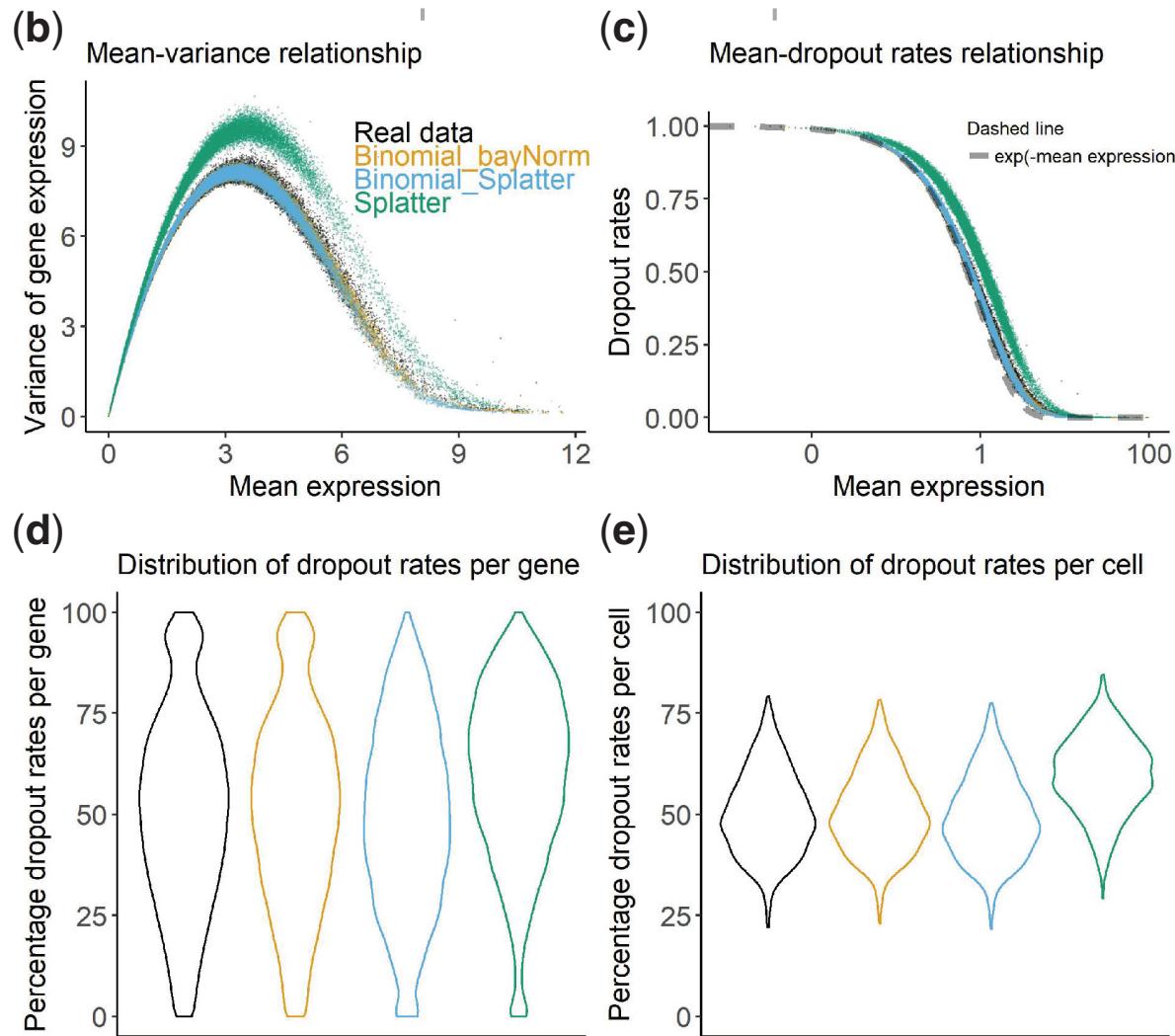
bayNorm: A method for normalisation, imputation and true gene expression recovery



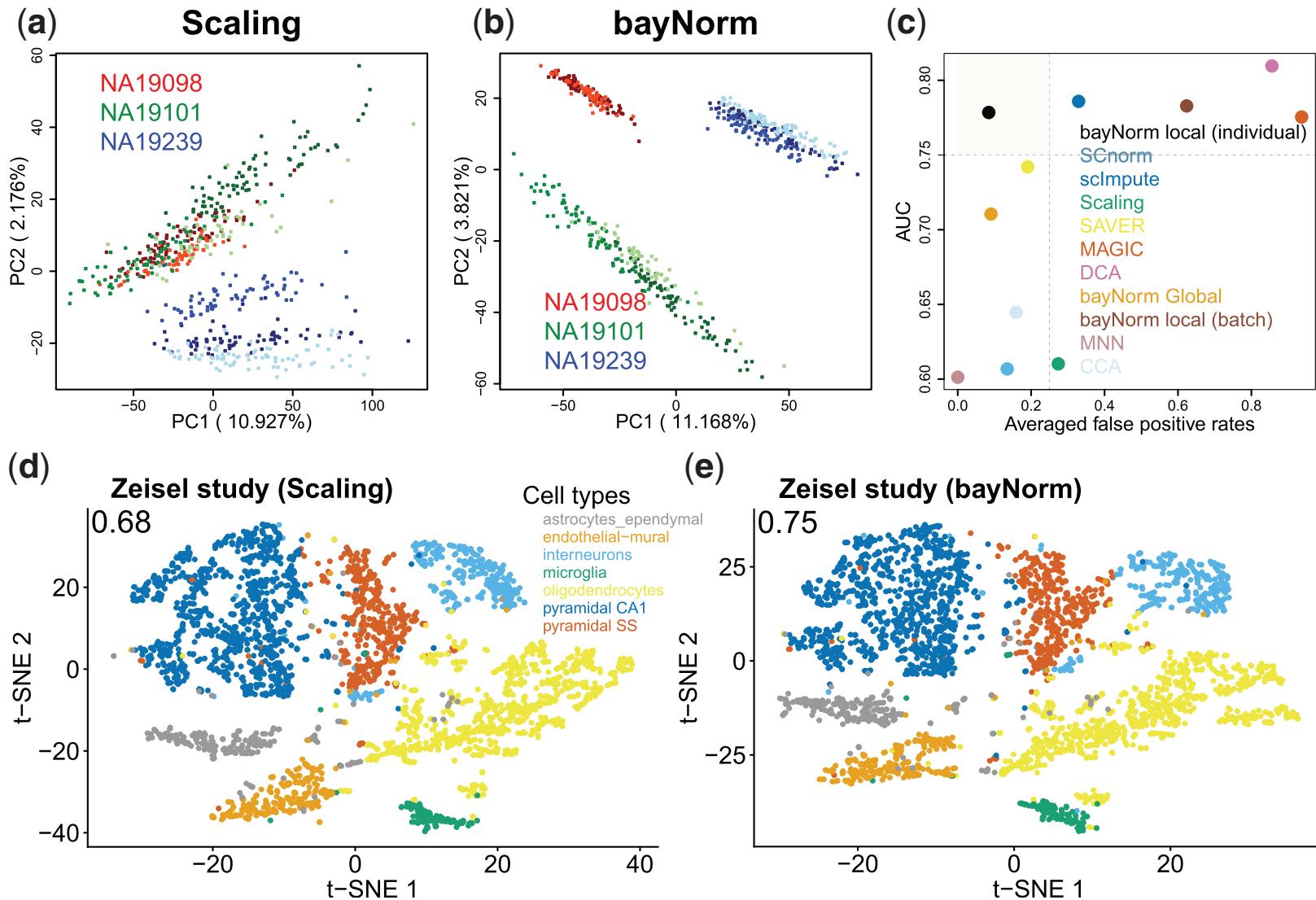
$$\underbrace{\Pr(x_{ij}^0 | x_{ij}, \beta_j, \mu_i, \phi_i)}_{\text{Posterior}} = \frac{\overbrace{\Pr(x_{ij} | x_{ij}^0, \beta_j) \times \Pr(x_{ij}^0 | \mu_i, \phi_i)}^{\text{Likelihood}}}{\underbrace{\Pr(x_{ij} | \mu_i, \phi_i, \beta_j)}_{\text{Marginal likelihood}}}$$

$$x_{ij} | x_{ij}^0 \sim \text{Binom}(x_{ij}^0, \text{prob} = \beta_j),$$
$$x_{ij}^0 \sim \text{NB}(\text{mean} = \mu_i, \text{size} = \phi_i).$$

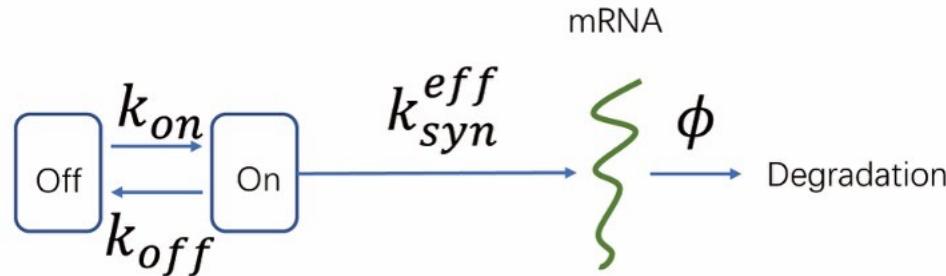
bayNorm posterior predictive explains drop-out (No need for zero-inflated distributions)



bayNorm can remove batch affect and cell type identification



Inference of burst kinetics from scRNA-seq data using bayNorm model

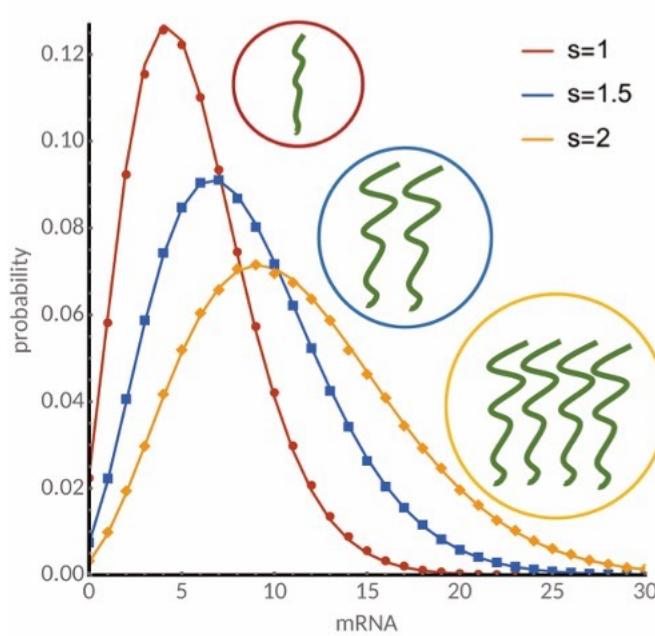


Cells

$$k_{syn}^{eff} = s * k_{syn}$$

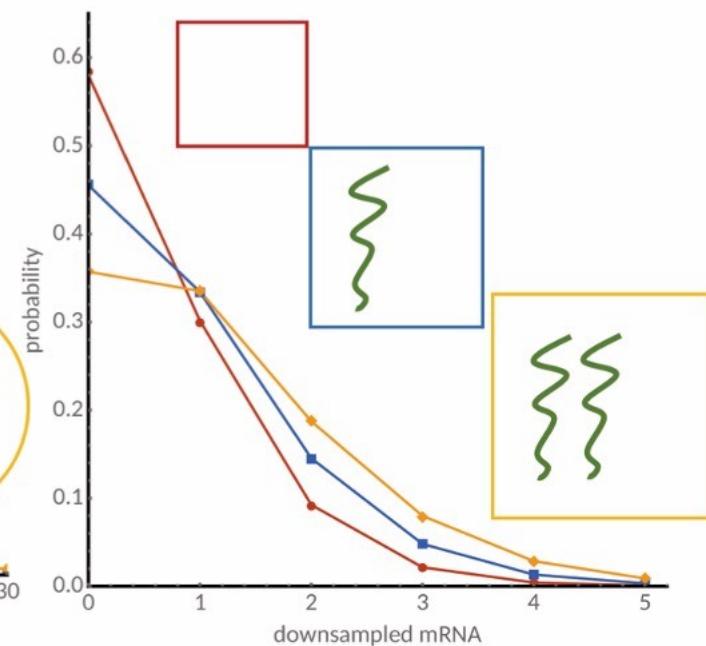
$$X \sim \text{Pois}(k_{syn}^{eff} p)$$

$$p \sim \text{Beta}(k_{on}, k_{off})$$

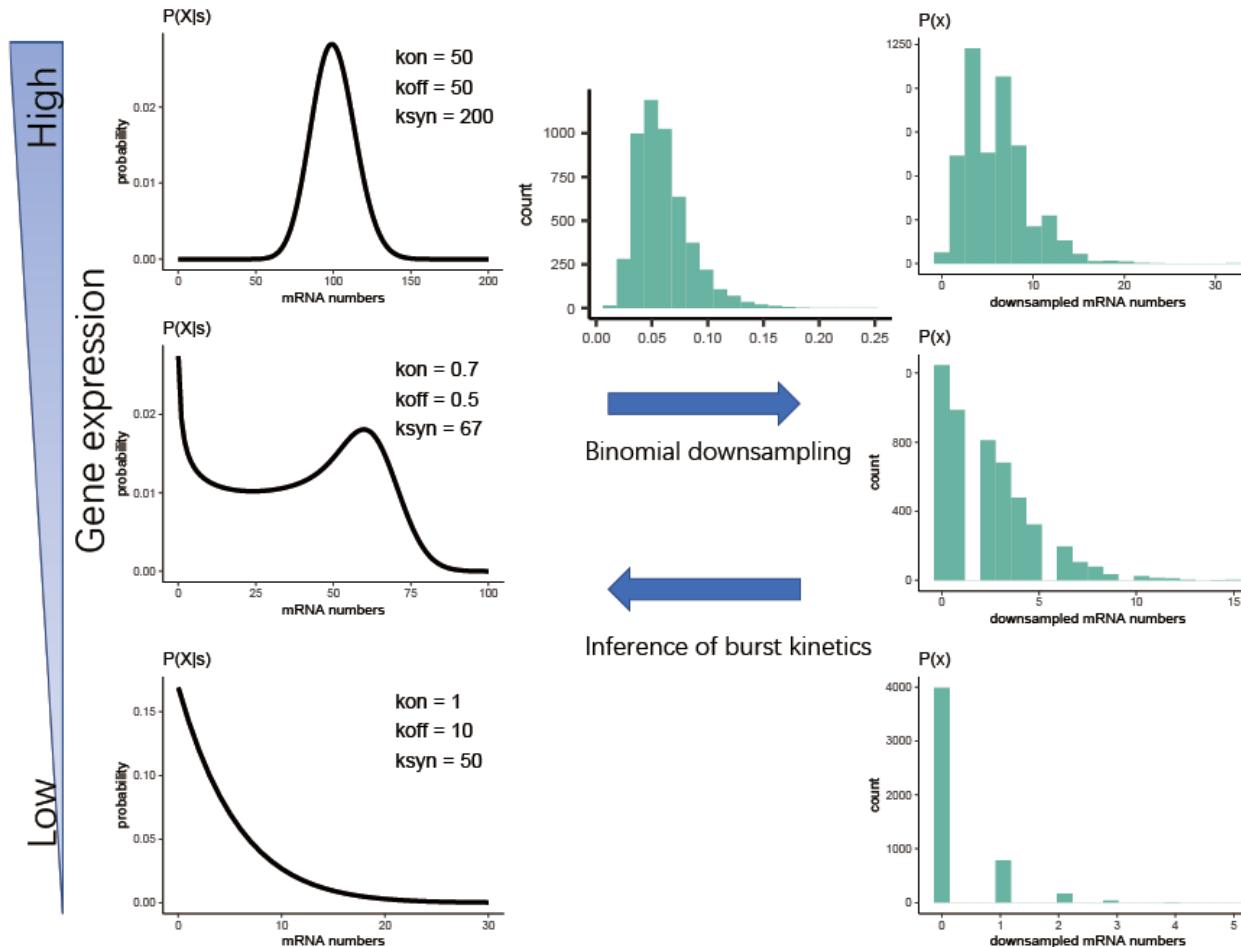


Observations

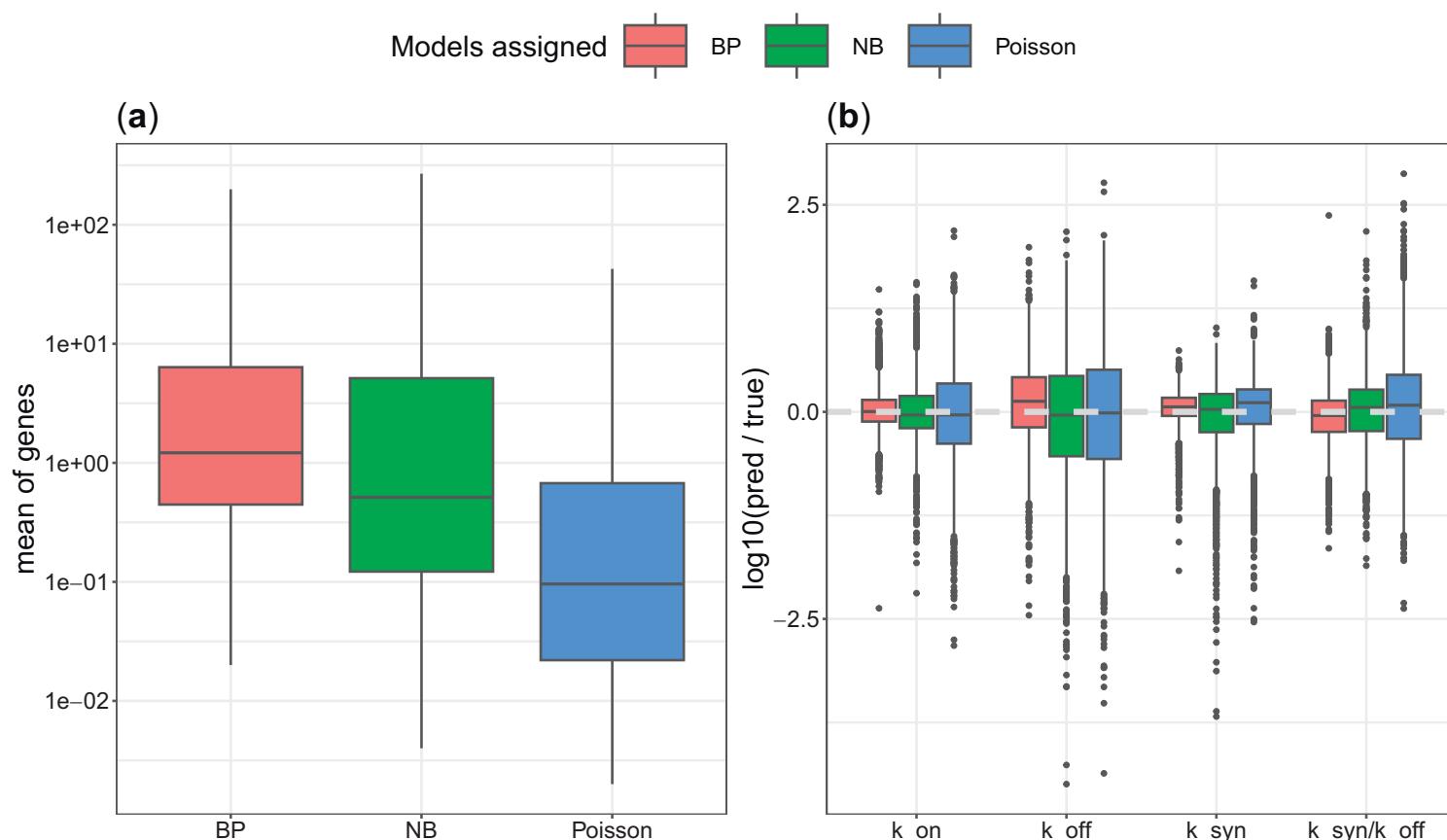
$$k_{syn}^{eff} = \beta * s * k_{syn}$$



Burst kinetics inference using simulation-based inference based on neural networks

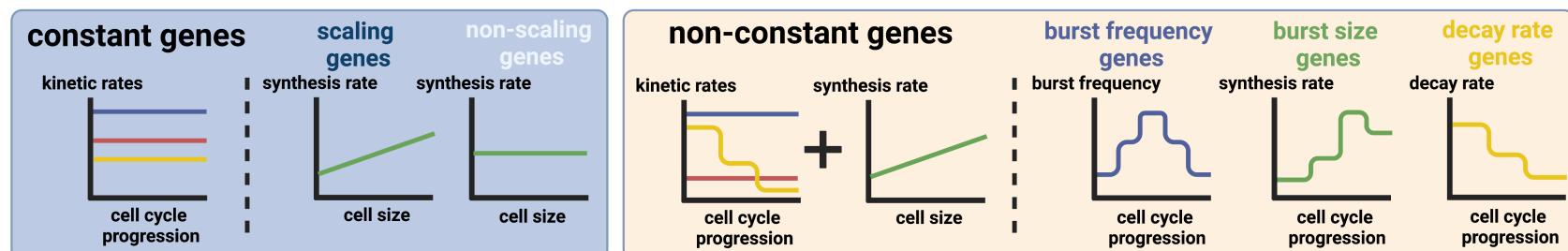


Model selection results chooses simpler models when there is less data

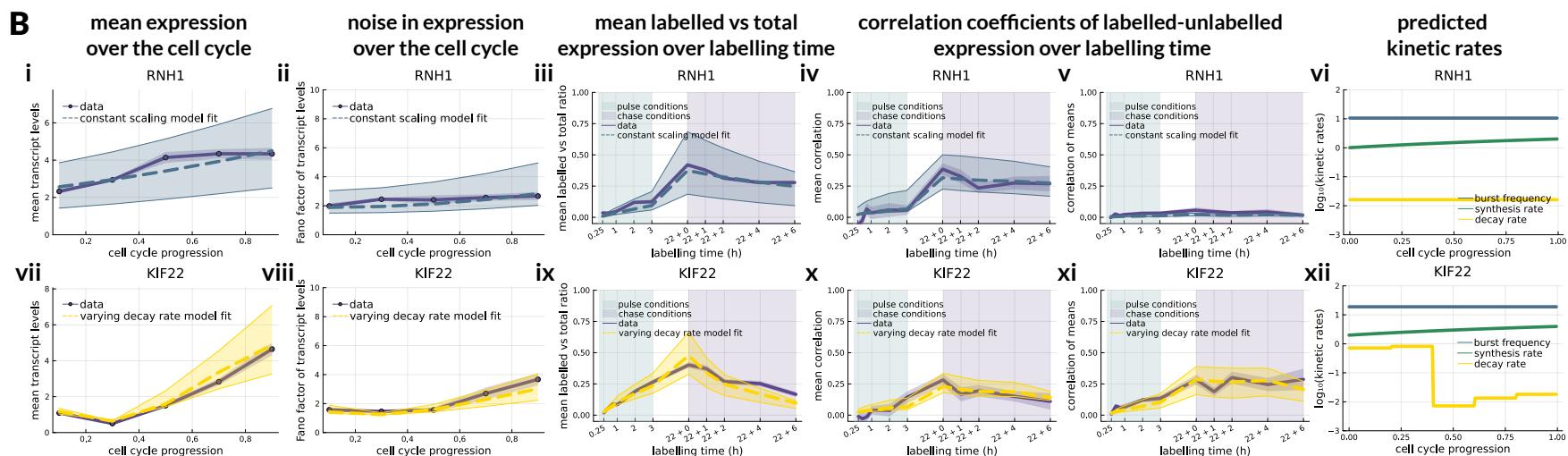


Genome-wide inference of cell size scaling and cell cycle dependence of gene expression Using metabolic labelling time-resolved transcription data (scEU-RNA-seq)

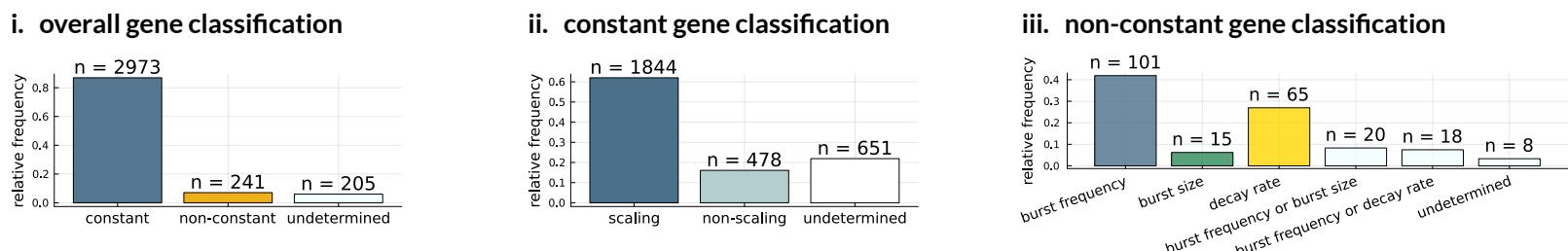
A



B



C

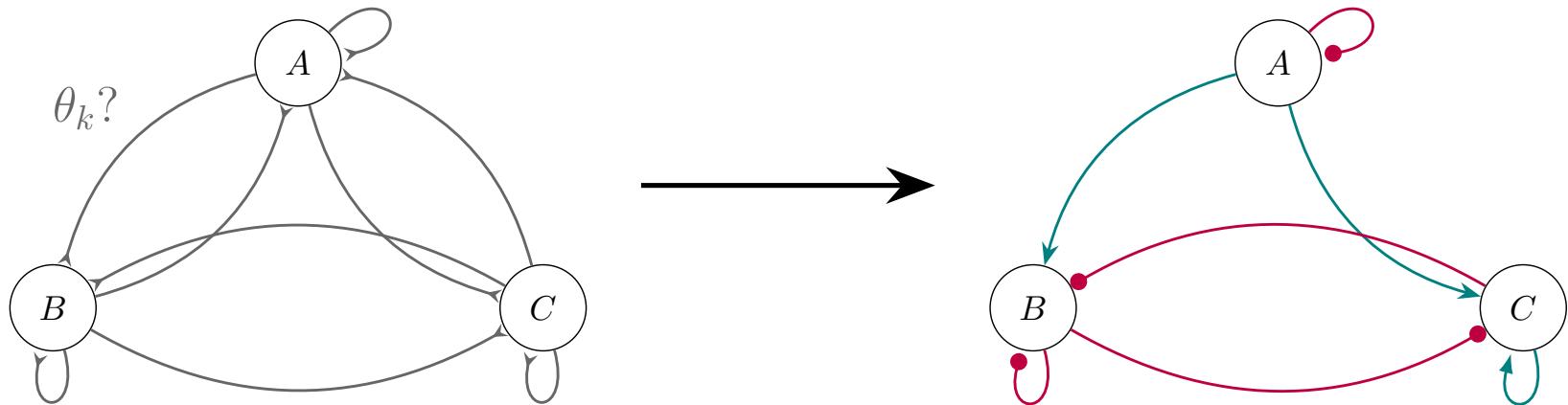


Volteras, Shahrezaei & Thomas Cell Systems 2024

seeDimitris Volteras talk on Friday

Inferring model structure

Model discovery (or equation discovery) for biochemical networks

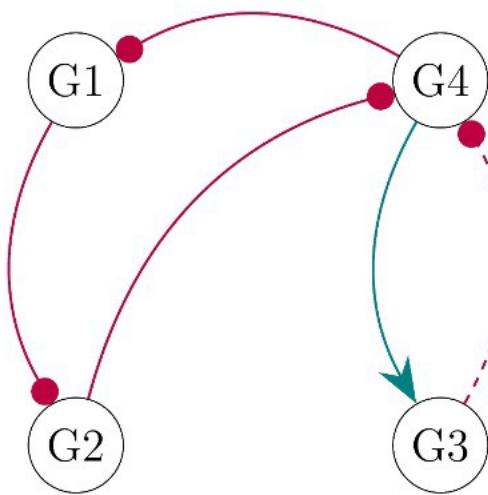


Large scale model selection can be reduced to variable selection

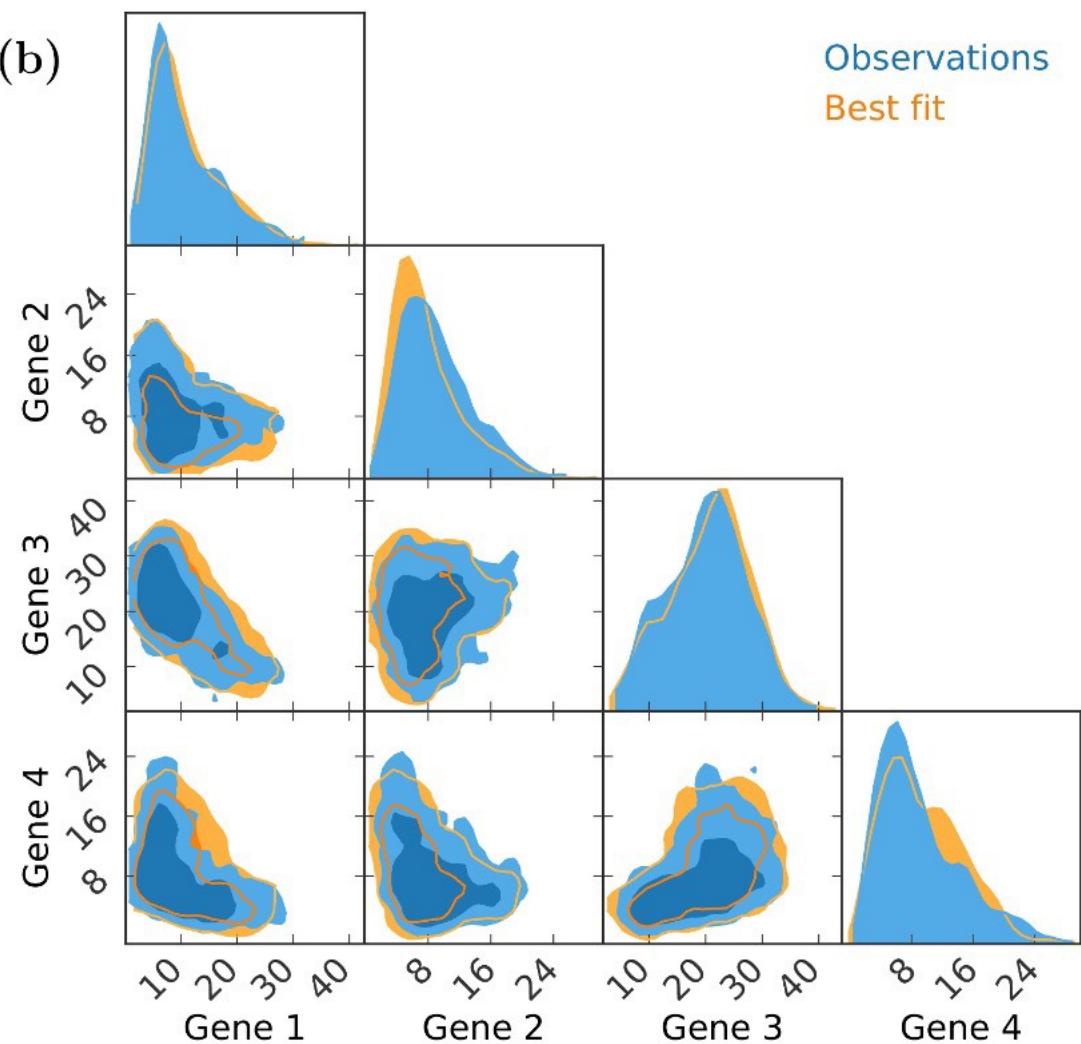
Bayesian Lasso, ABC with sparsity inducing parameters, SINDy

Inference of stochastic genetic networks from single cell data

(a)

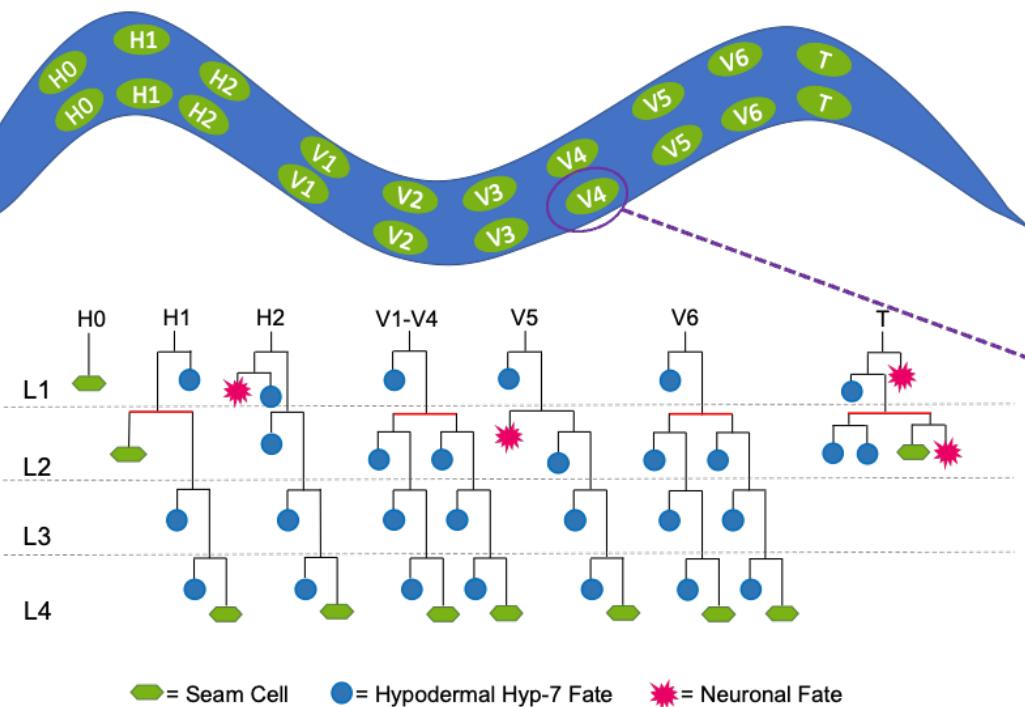


(b)

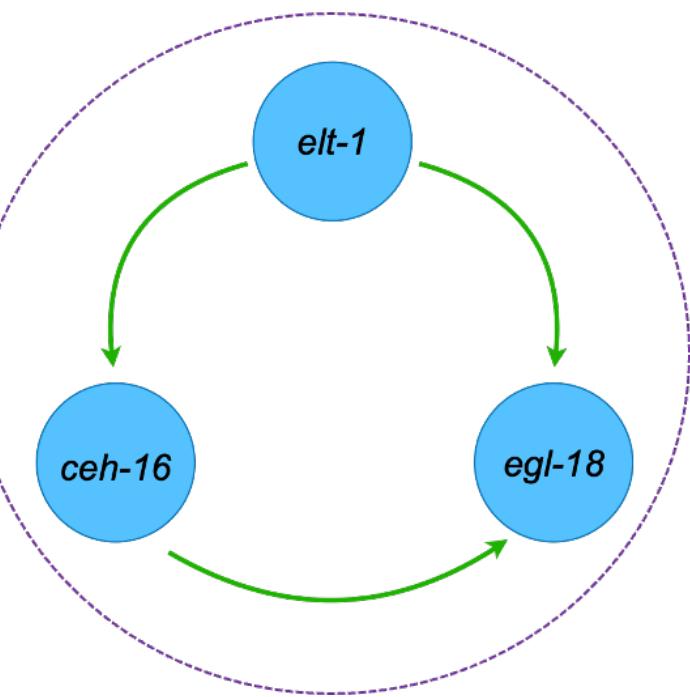


Application to the inference of the genetic network controlling stem cell fate decision

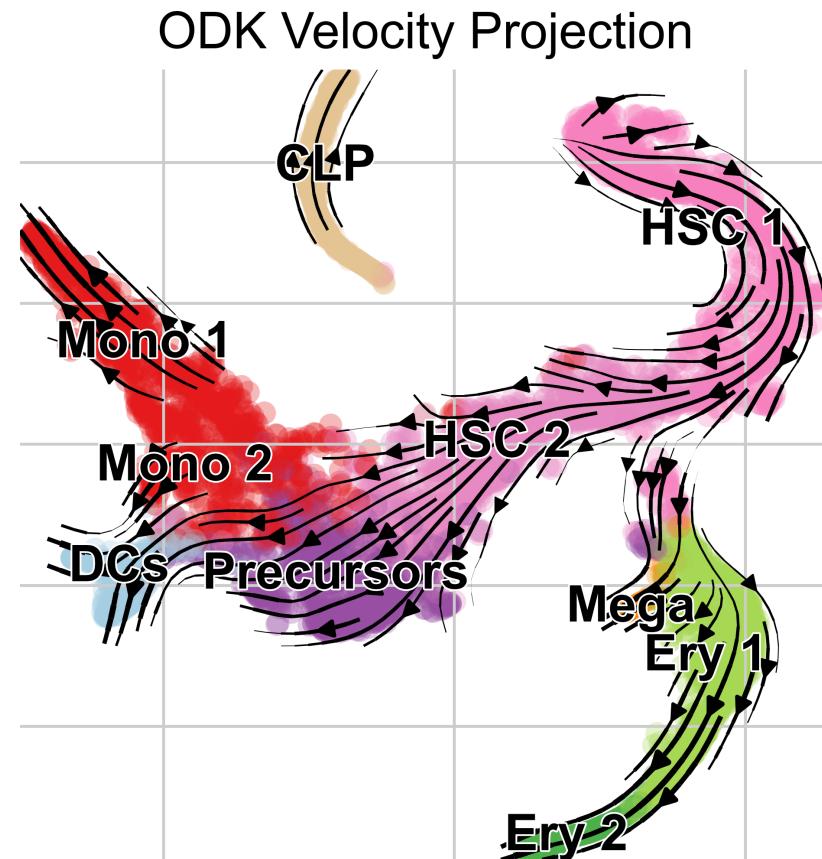
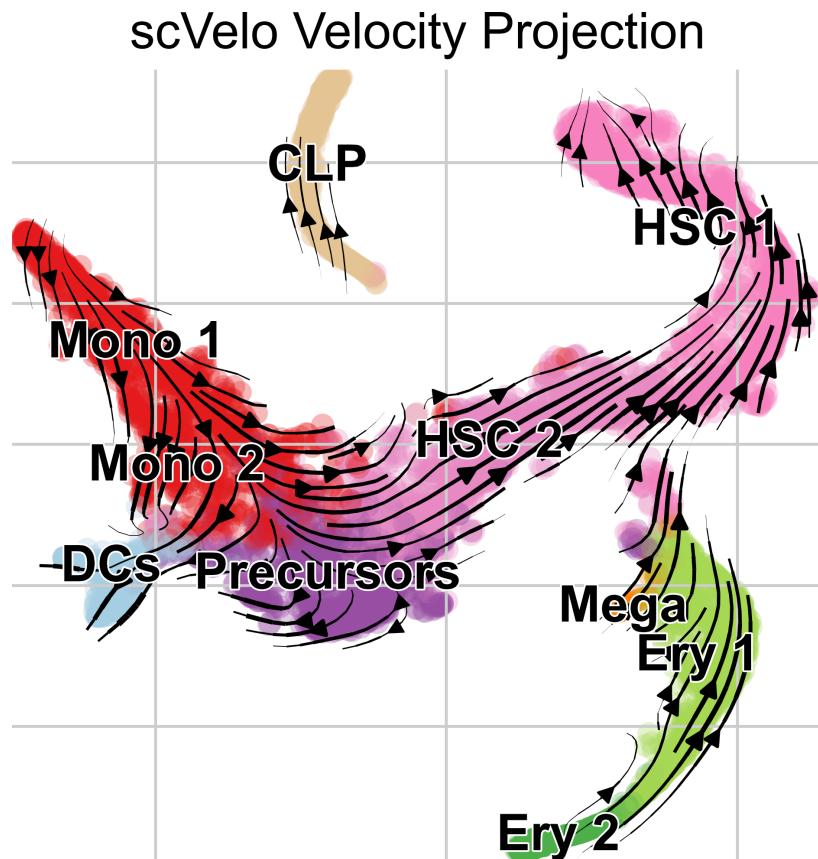
A



B



Inferring the low-dimensional dynamics of underlying stochastic process from the high-dimensional scRNA-seq data using a novel ordered diffusion kernel approach



Jack Soulsby, unpublished work of current PhD student

See talk on Friday

There are long-standing grand challenges for single-cell data science:

Technical errors, sparsity, high-dimensionality and data integration across batches and modalities
Inferring patterns, trajectories and dynamics

