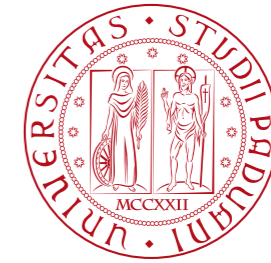




1222·2022  
**800**  
ANNI



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

CSAMA 2022 - BRIXEN/BRESSANONE

---

**SINGLE-CELL RNA-SEQ**

Davide Risso

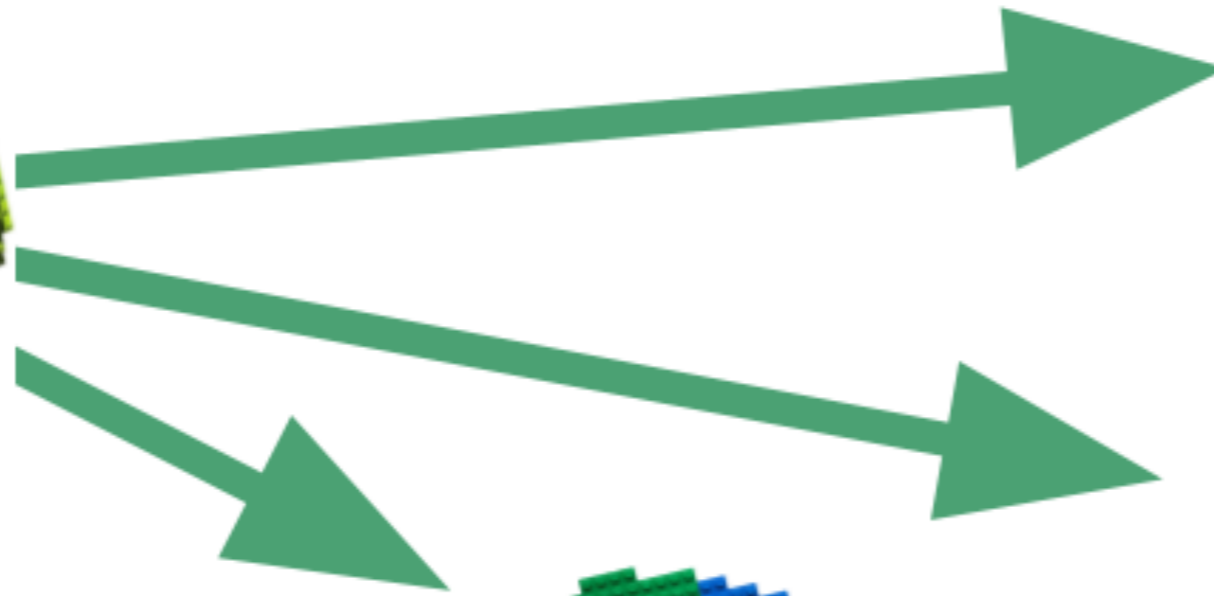
 @drisso1893

 @drisso

# The evolution of gene expression measurements



Original organ



Spatial transcriptomics

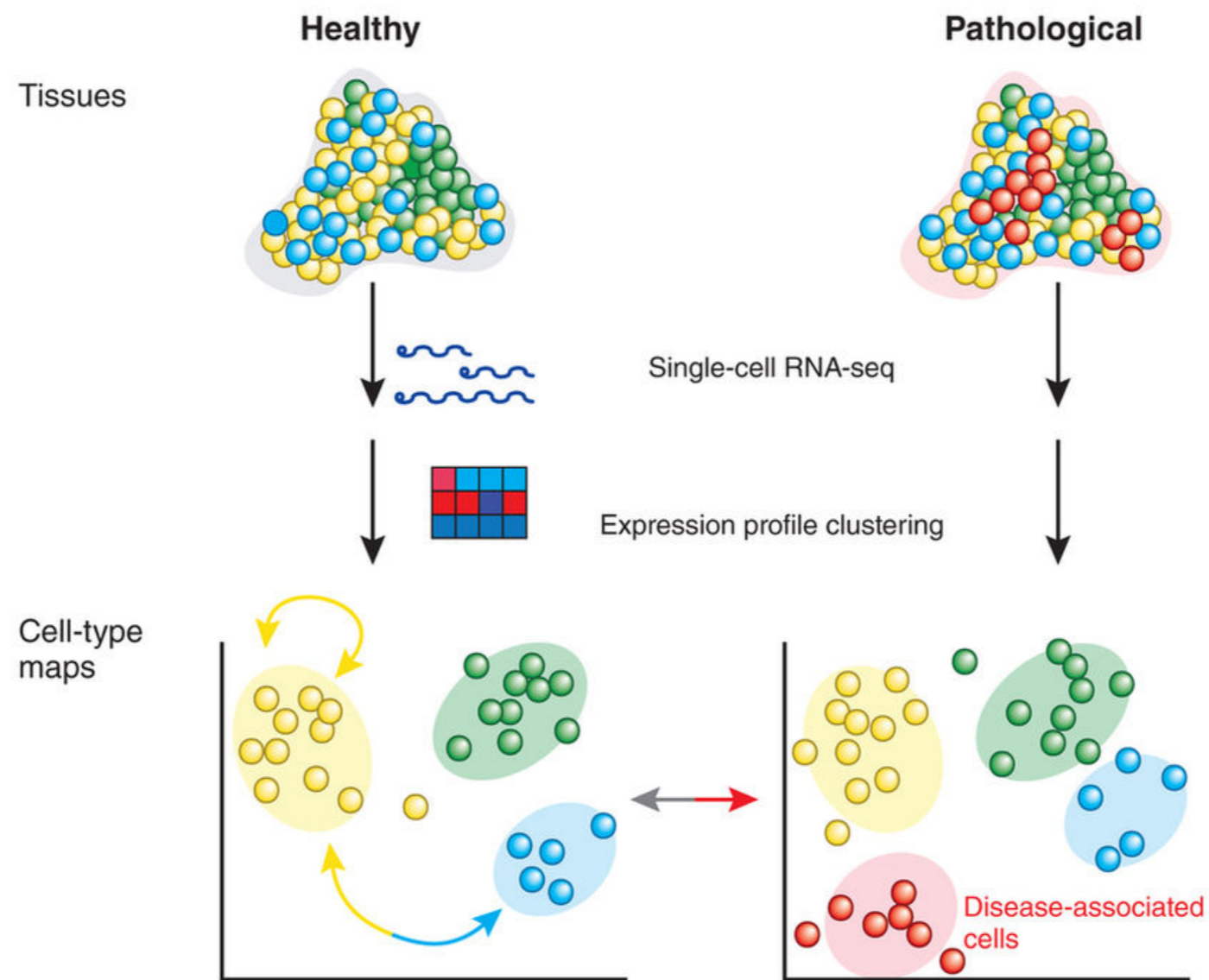


Bulk RNA-seq



Single-cell RNA-seq

# SINGLE-CELL RNA-SEQ



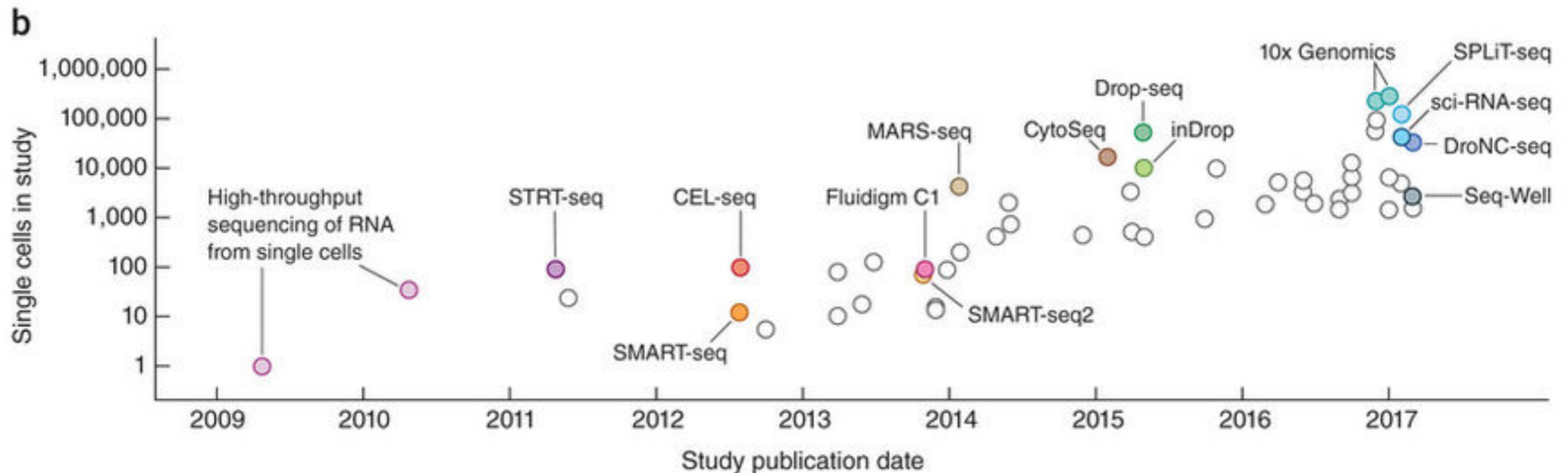
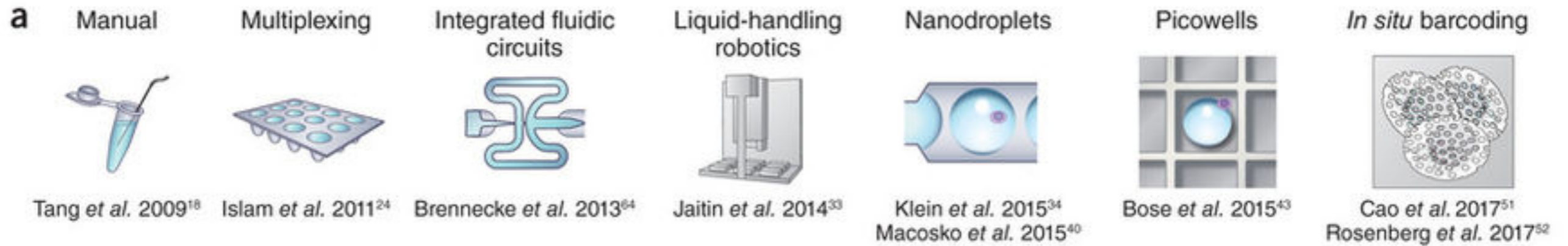
## Types of analyses

- Within cell type**
- Stochasticity, variability of transcription
  - Regulatory network inference
  - Allelic expression patterns
  - Scaling laws of transcription

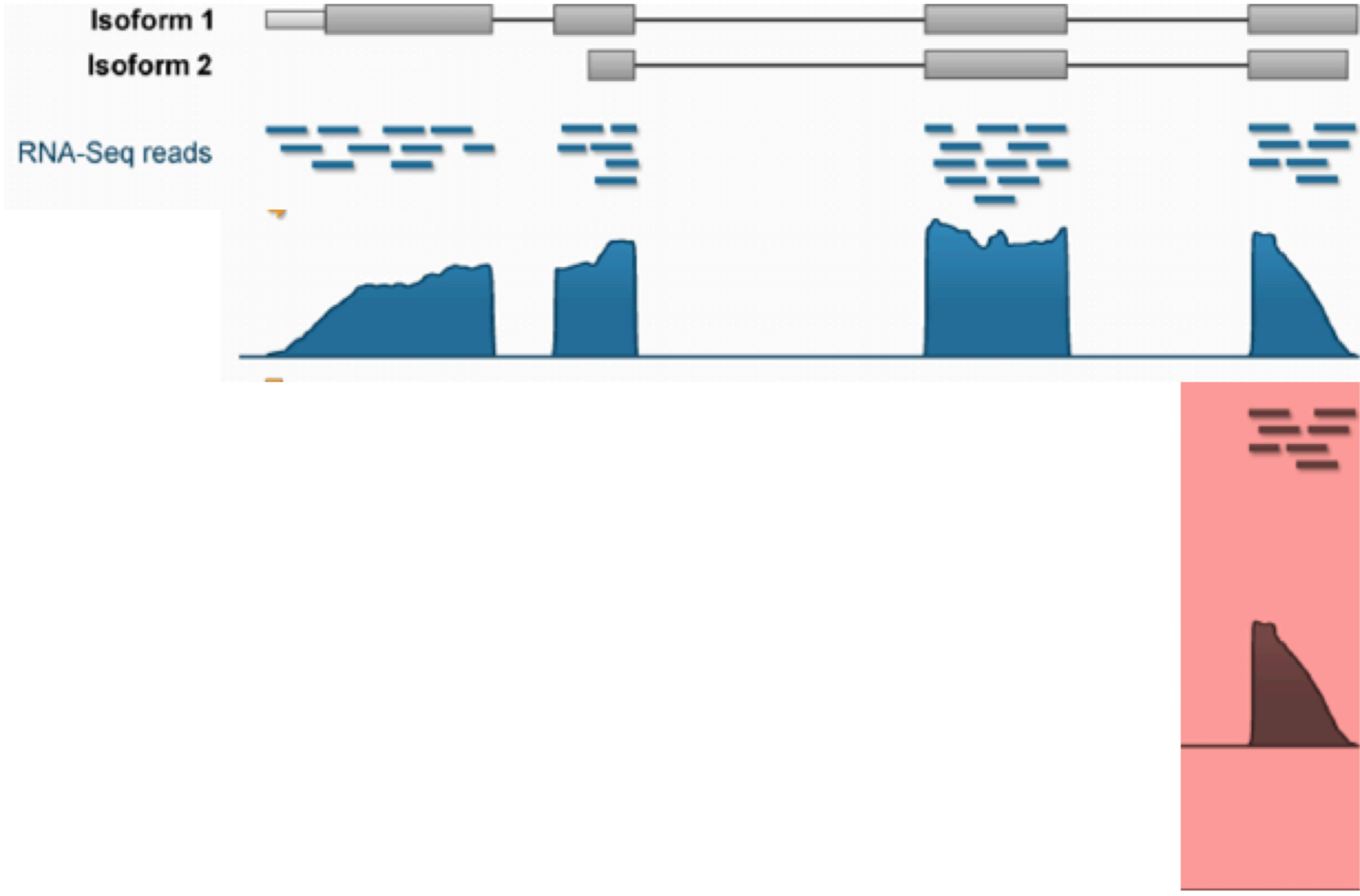
- Between cell types**
- Identify biomarkers
  - (Post)-transcriptional differences

- Between tissues**
- Cell-type compositions
  - Altered transcription in matched cell types

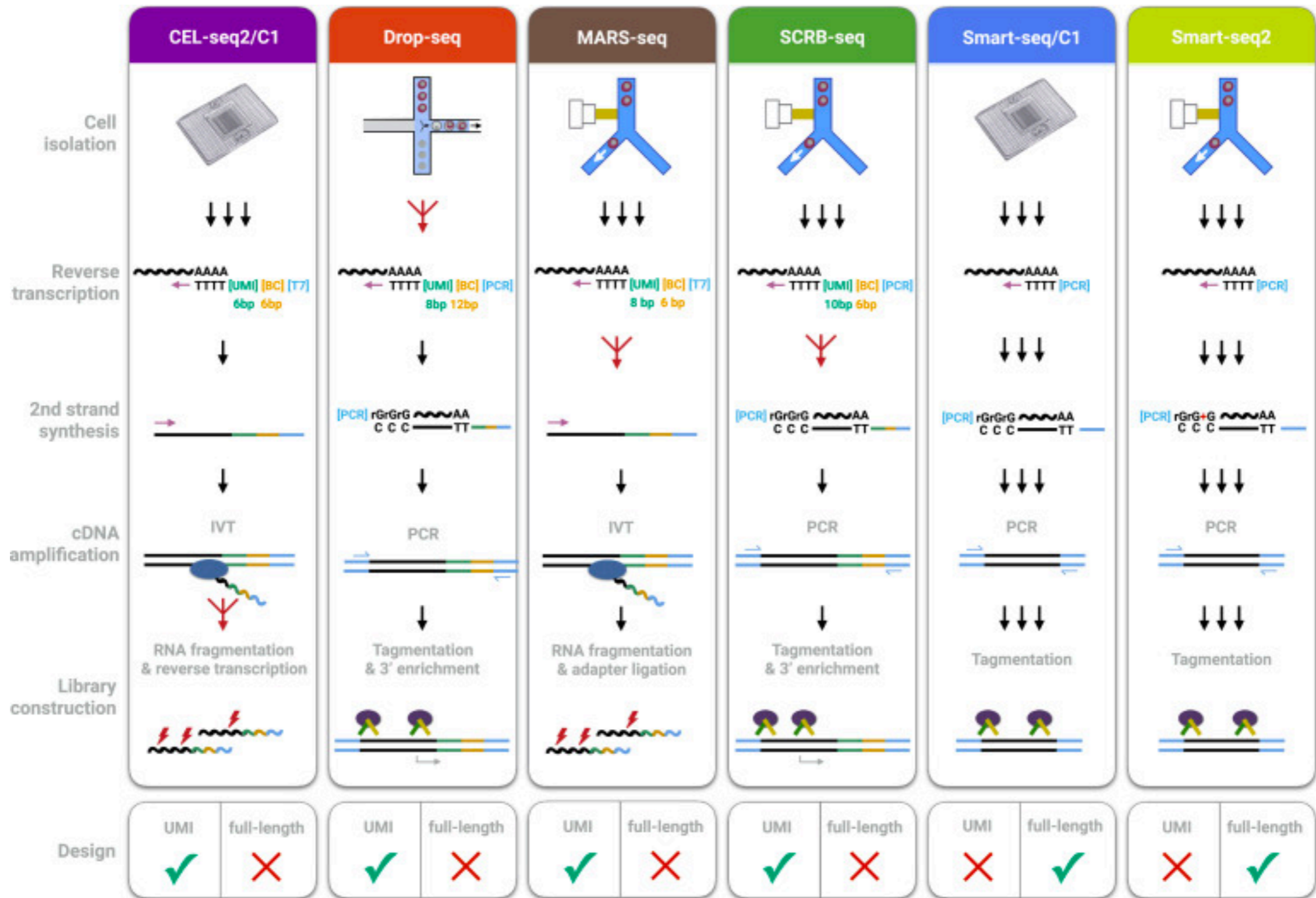
# SEVERAL PROTOCOLS AND PLATFORMS



# DIFFERENT PROTOCOLS HAVE DIFFERENT PROPERTIES

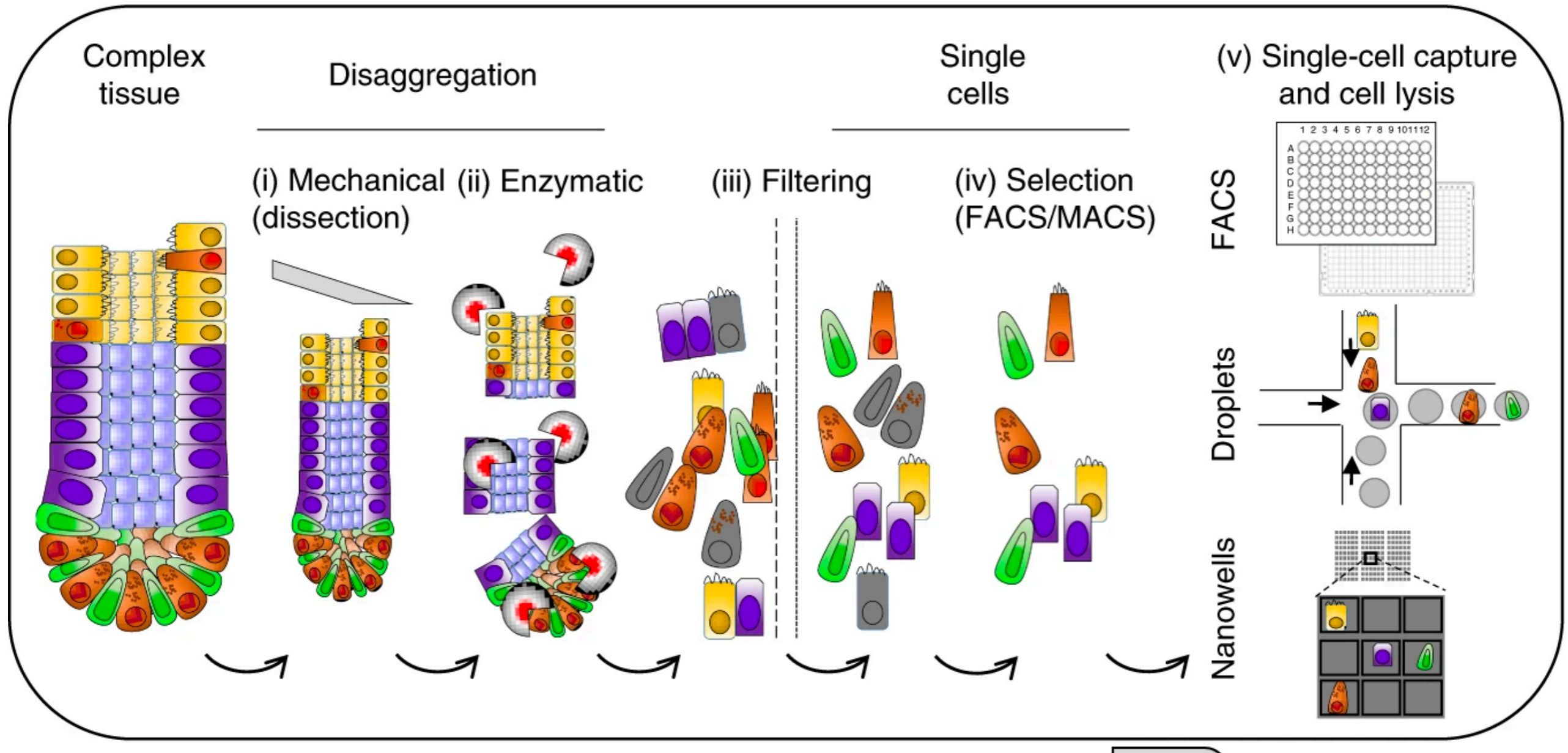


# DIFFERENT PROTOCOLS HAVE DIFFERENT PROPERTIES



# SINGLE-CELL RNA-SEQ IN A NUTSHELL

## (1) Sample preparation

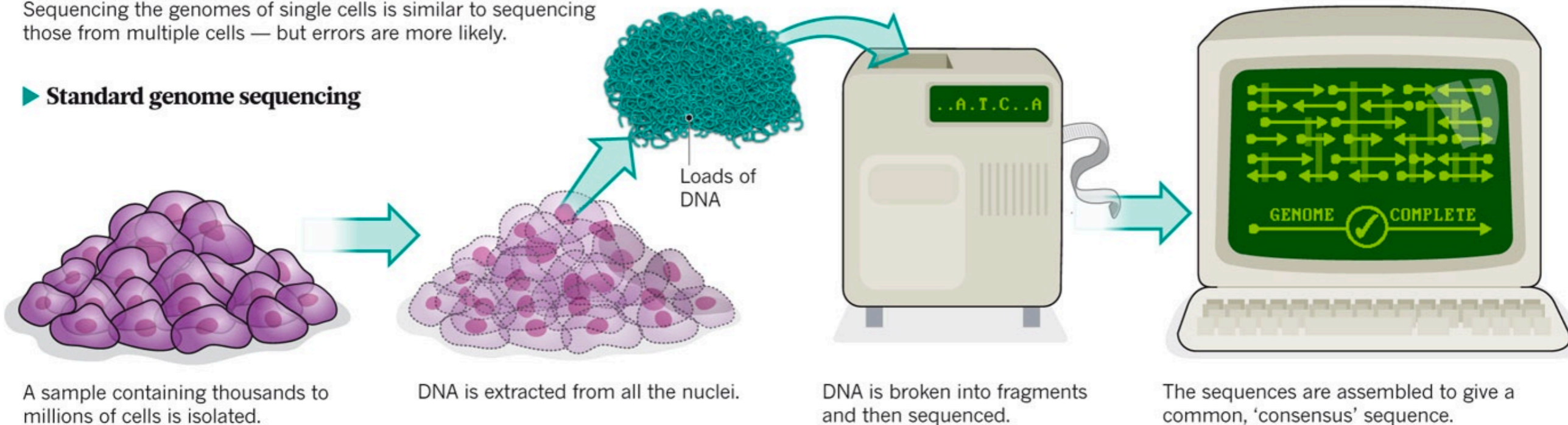


# AMPLIFICATION BIAS LEADS TO...

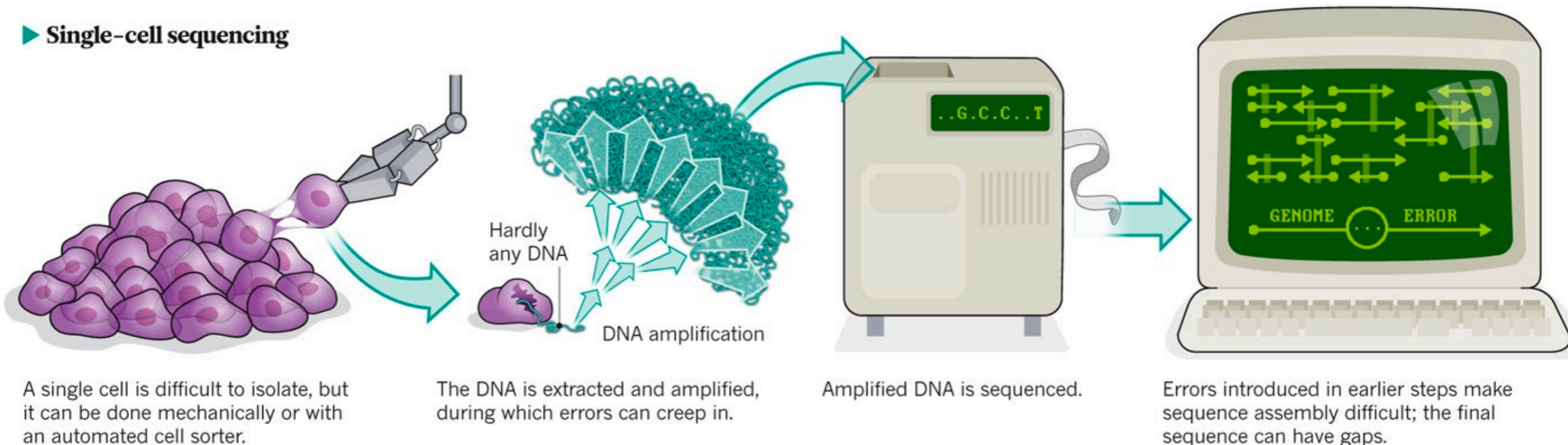
## ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

### ► Standard genome sequencing

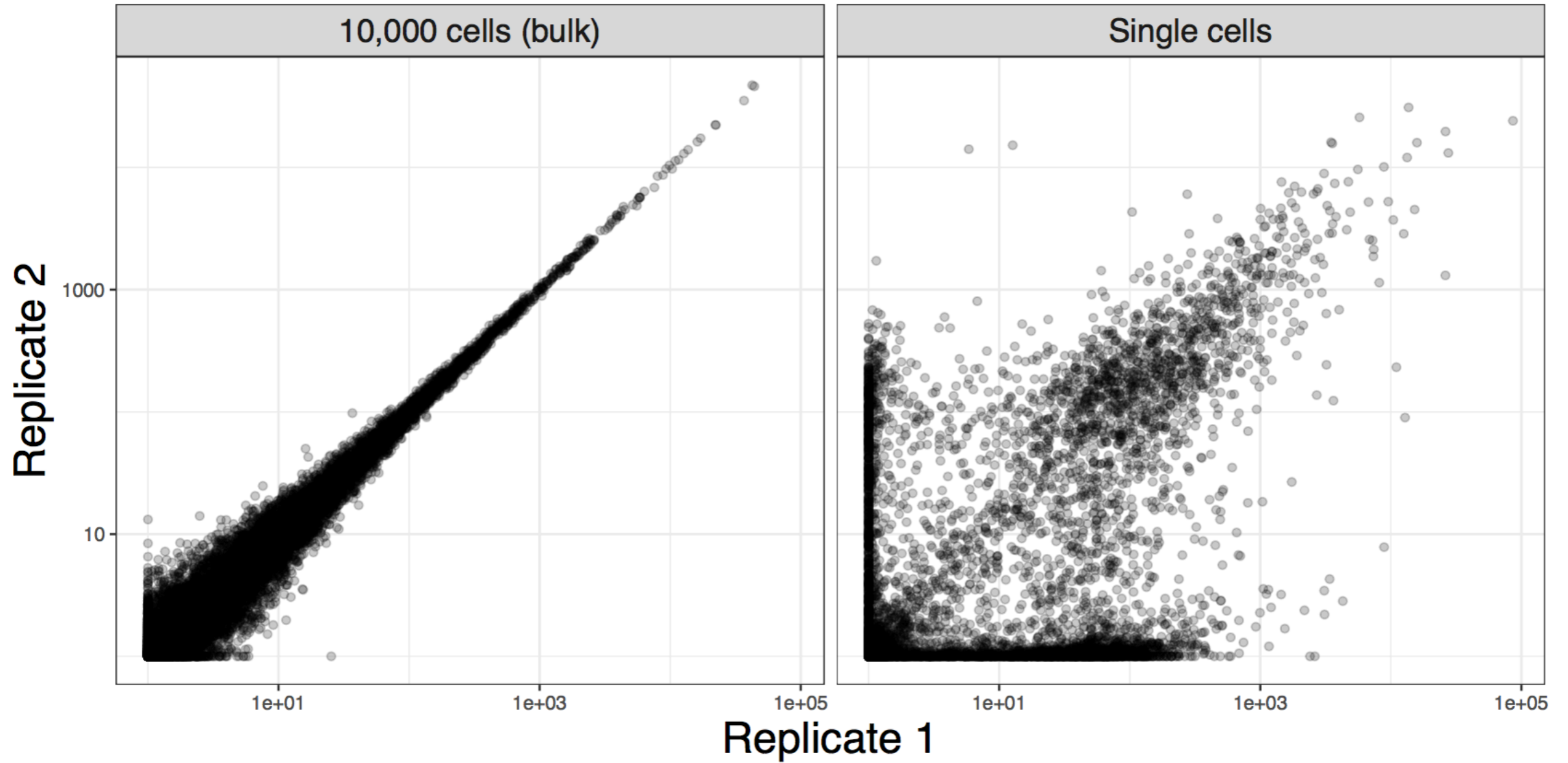


### ► Single-cell sequencing

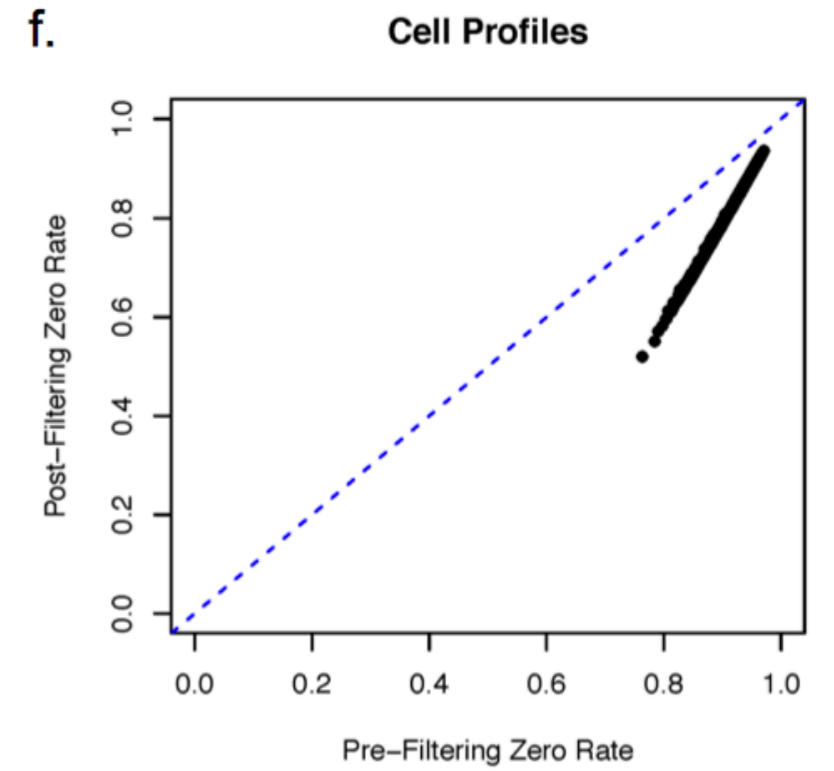
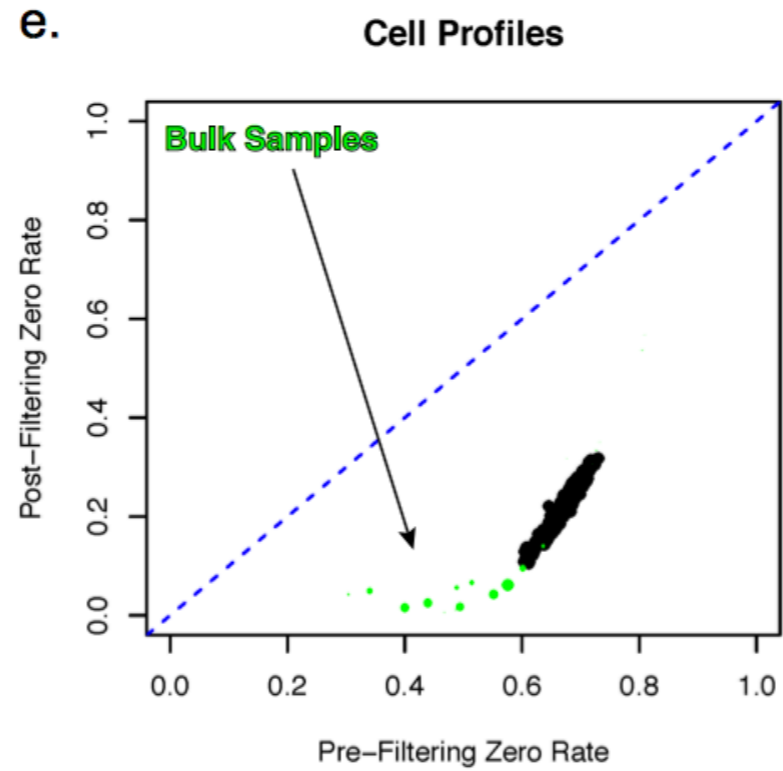
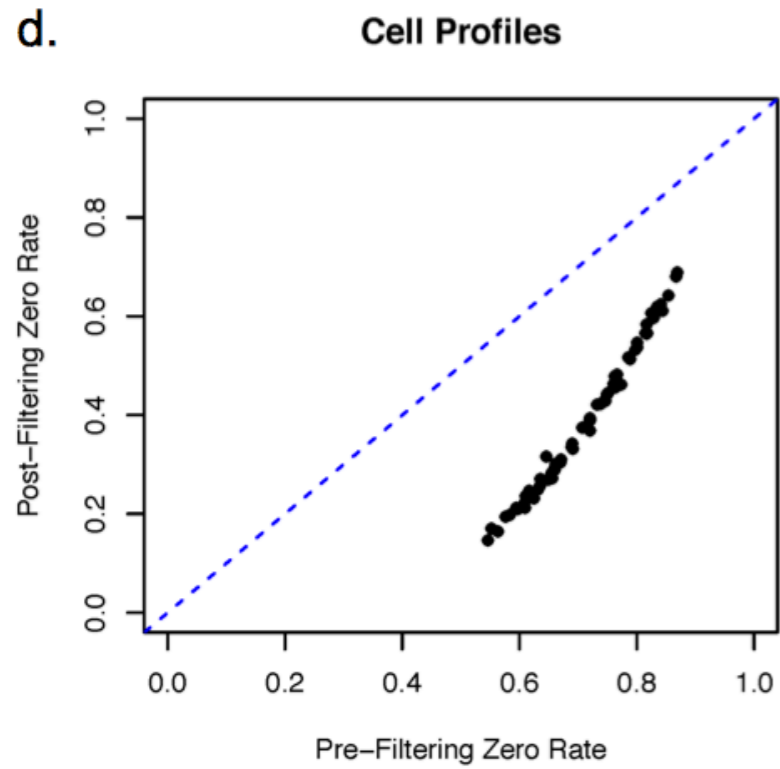




# INCREASED VARIABILITY COMPARED TO “BULK” RNA-SEQ



# EXCESS OF ZERO COUNTS



# UNIQUE MOLECULAR IDENTIFIERS

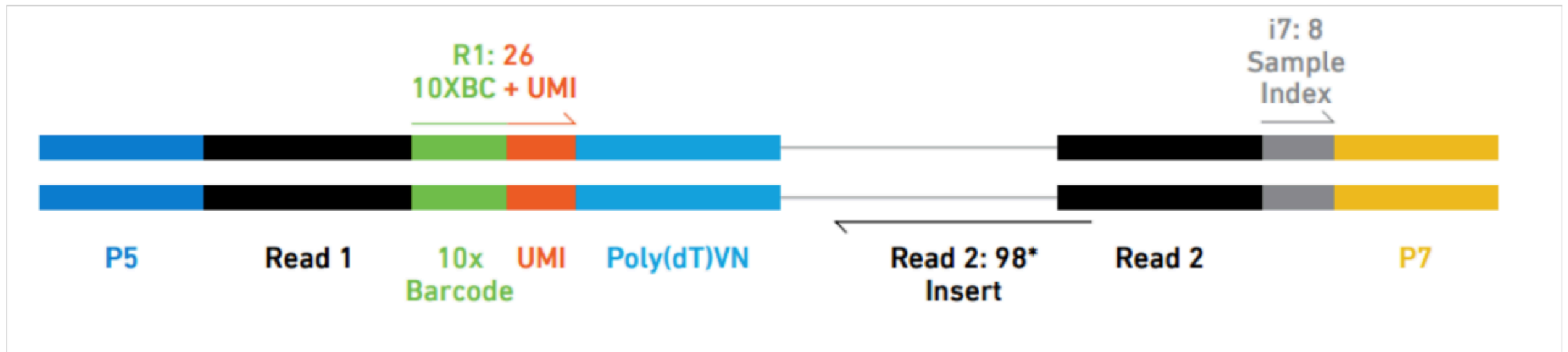
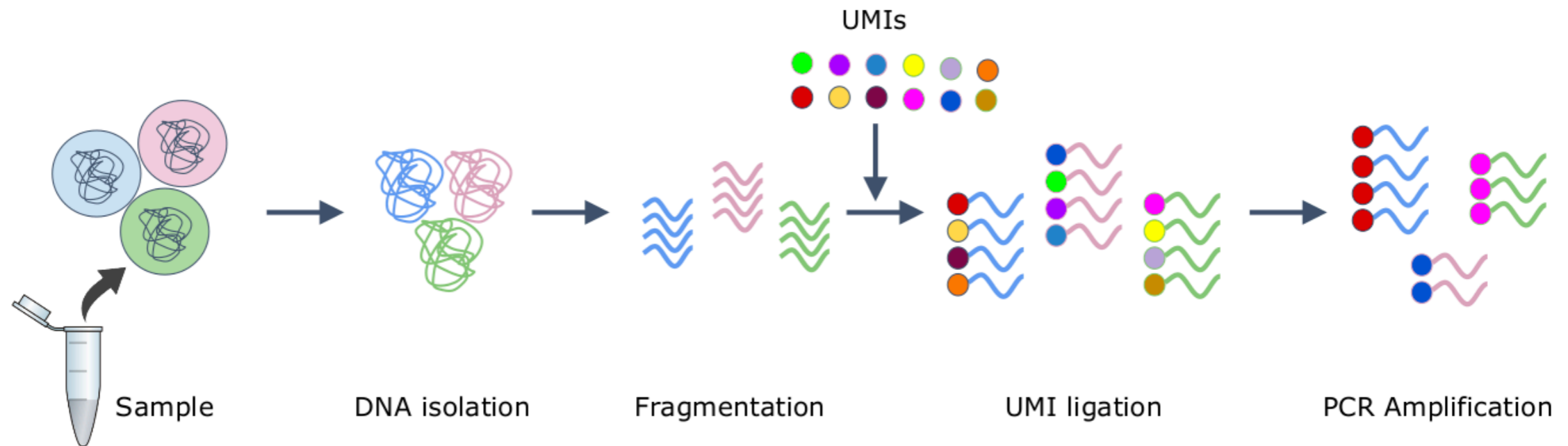
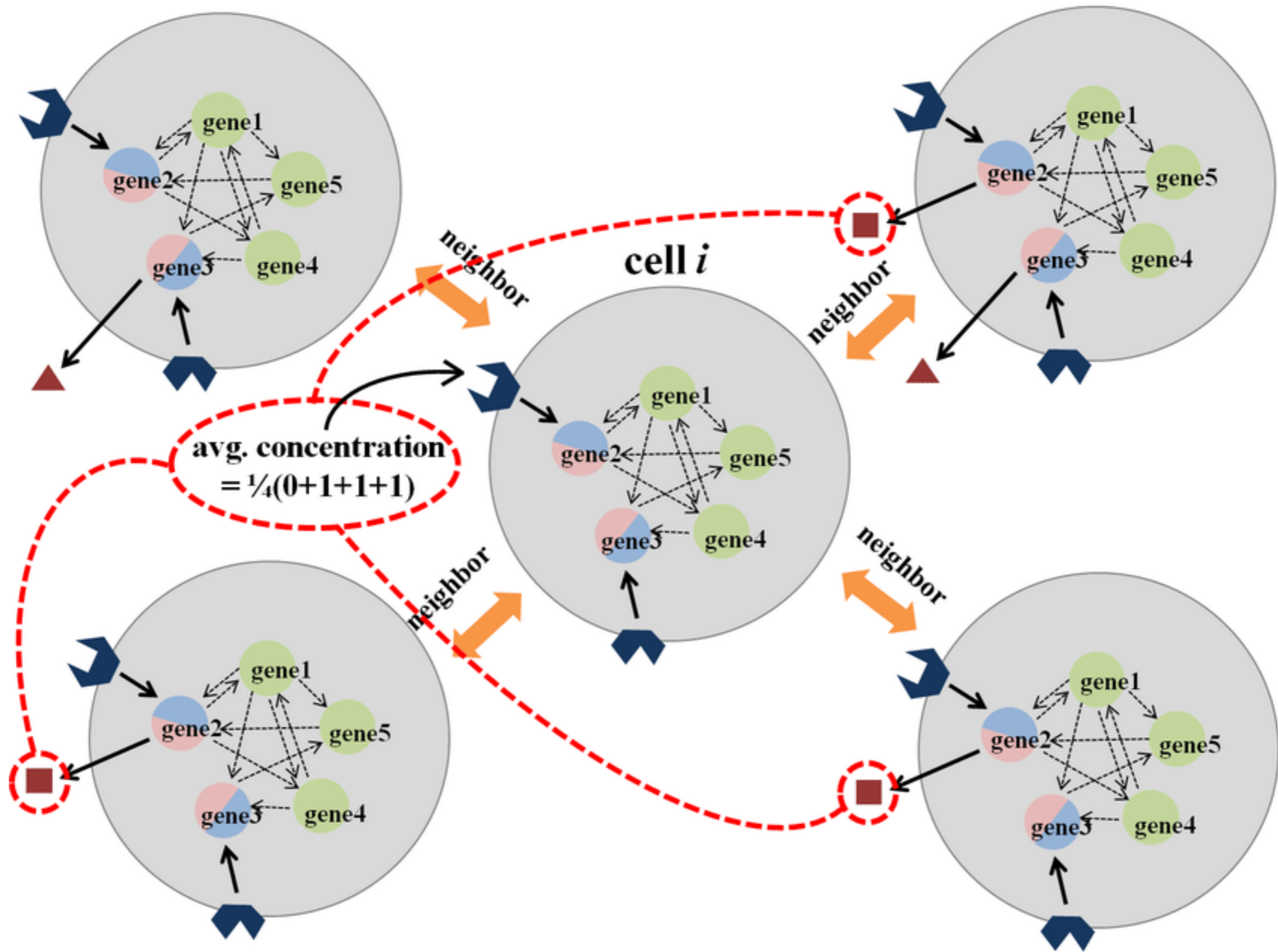


Fig. 2. Schematic of a fragment from a final Chromium™ Single Cell 3' v2 library. \*Can be adjusted.



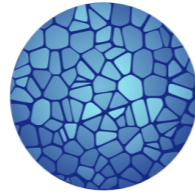
# HIGH COMPLEXITY



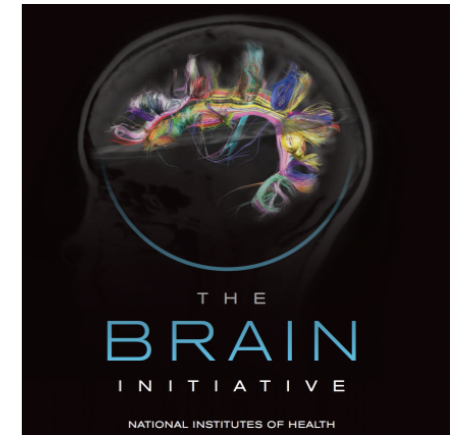
# HIGH DIMENSIONALITY AND SAMPLE SIZE

20M Cells

ALL CELLS



HUMAN  
CELL  
ATLAS



Blood	Kidney
Bone	Liver
Brain	Lung
Pancreas	Heart
Immune System	Skin

## Article

### A transcriptomic and epigenomic cell atlas of the mouse primary motor cortex

<https://doi.org/10.1038/s41586-021-03500-8>

Received: 5 March 2020

Accepted: 26 March 2021

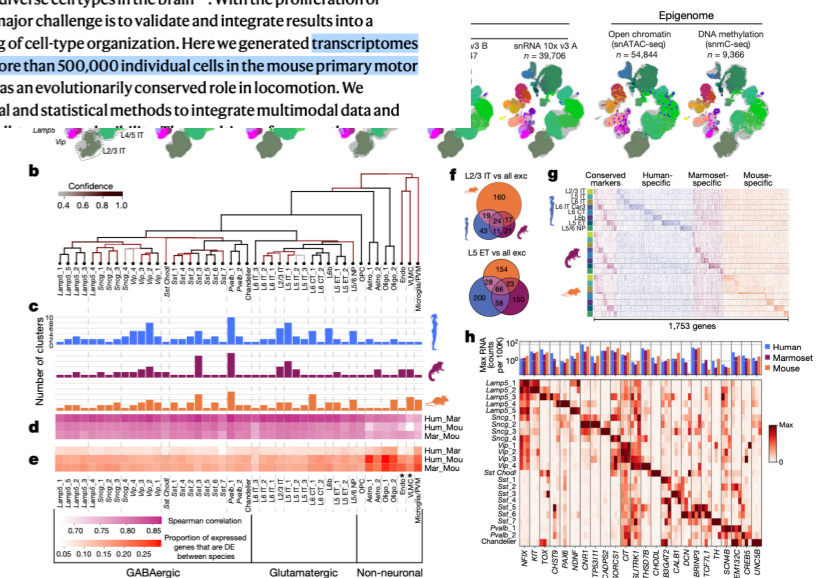
Published online: 6 October 2021

Open access

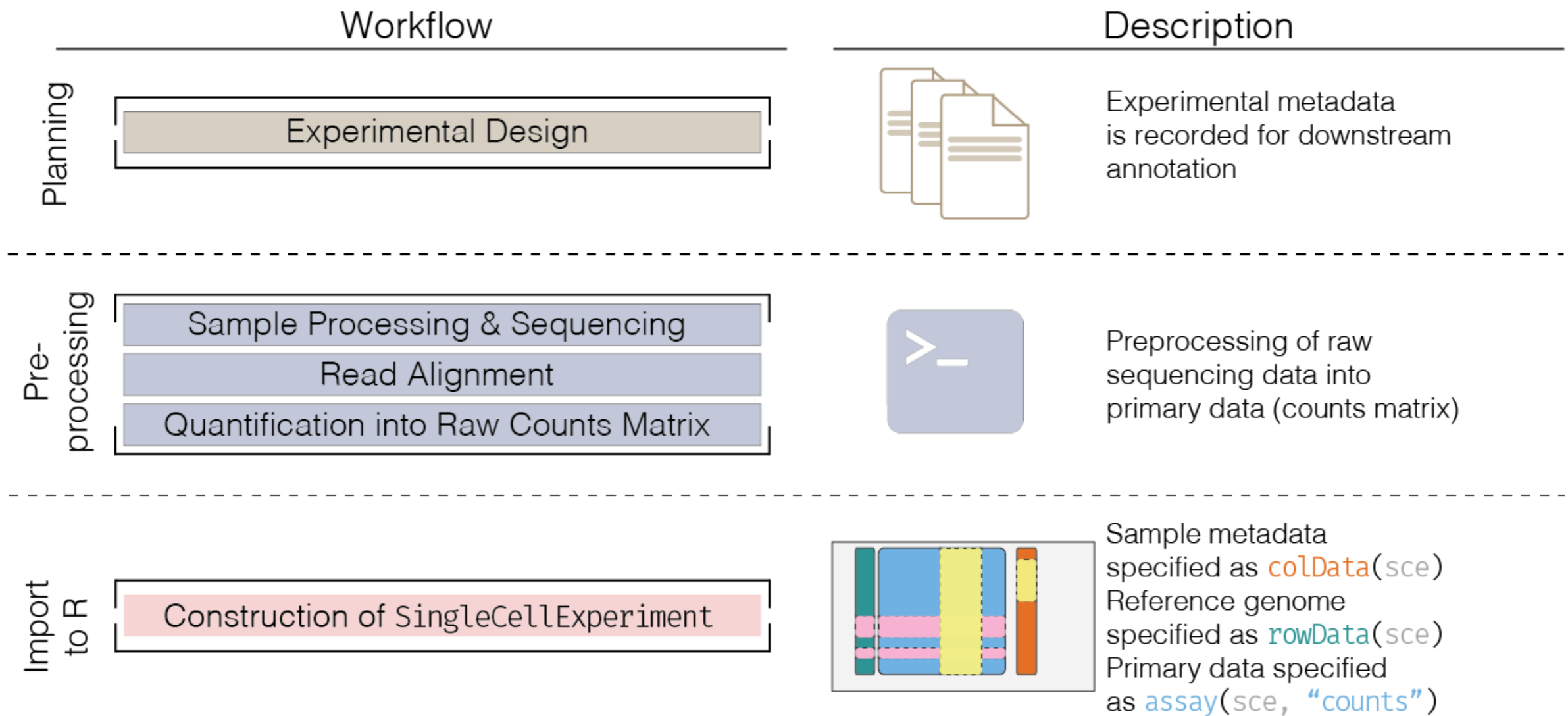
Check for updates

Zizhen Yao<sup>1,32</sup>, Hanqing Liu<sup>2,32</sup>, Fangming Xie<sup>3,32</sup>, Stephan Fischer<sup>4,32</sup>, Ricky S. Adkins<sup>5</sup>, Andrew I. Aldridge<sup>6</sup>, Seth A. Ament<sup>6</sup>, Anna Bartlett<sup>6</sup>, M. Margarita Behrens<sup>6</sup>, Koen Van den Berge<sup>6,8</sup>, Darren Bertagnoli<sup>7</sup>, Hector Roux de Bézieux<sup>8</sup>, Tommaso Biancalani<sup>10</sup>, A. Sina Boeshaghi<sup>11</sup>, Héctor Corrada Bravo<sup>12</sup>, Tamara Casper<sup>1</sup>, Carlo Colantuoni<sup>13,14,15</sup>, Jonathan Crabtree<sup>6</sup>, Heather Creasy<sup>2</sup>, Kirsten Crichton<sup>1</sup>, Megan Crow<sup>6</sup>, Nick Dee<sup>1</sup>, Elizabeth L. Dougherty<sup>10</sup>, Wayne I. Doyle<sup>16</sup>, Sandrine Dudoit<sup>7</sup>, Rongxin Fang<sup>17</sup>, Victor Felix<sup>5</sup>, Olivia Fong<sup>1</sup>, Michelle Giglio<sup>5</sup>, Jeff Goldy<sup>1</sup>, Mike Hawrylycz<sup>1</sup>, Brian R. Herb<sup>5</sup>, Ronna Hertzano<sup>5,18</sup>, Xiaomeng Hou<sup>19</sup>, Qiwen Hu<sup>20</sup>, Jayaram Kancherla<sup>12</sup>, Matthew Kroll<sup>1</sup>, Kanan Lathia<sup>1</sup>, Yang Eric Li<sup>21</sup>, Jacinta D. Lucero<sup>6</sup>, Chongyuan Luo<sup>22,23</sup>, Anup Mahurkar<sup>2</sup>, Delissa McMillen<sup>1</sup>, Naeem M. Nadaf<sup>20</sup>, Joseph R. Nery<sup>2</sup>, Thuc Nghi Nguyen<sup>1</sup>, Sheng-Yong Niu<sup>2</sup>, Vasilis Ntranos<sup>24</sup>, Joshua Orvis<sup>5</sup>, Julia K. Osteen<sup>6</sup>, Thanh Pham<sup>1</sup>, Antonio Pinto-Duarte<sup>6</sup>, Olivier Poirion<sup>19</sup>, Sebastian Preissl<sup>19</sup>, Elizabeth Purdom<sup>7</sup>, Christine Rimorin<sup>1</sup>, Davide Risso<sup>25</sup>, Angeline C. Rivkin<sup>23</sup>, Kimberly Smith<sup>1</sup>, Kelly Street<sup>26</sup>, Josef Sulc<sup>1</sup>, Valentine Svensson<sup>1</sup>, Michael Tieu<sup>1</sup>, Amy Torkelson<sup>1</sup>, Herman Tung<sup>1</sup>, Eeshit Dhaval Vaishnav<sup>10</sup>, Charles R. Vanderburg<sup>10</sup>, Cindy van Velthoven<sup>1</sup>, Xinxin Wang<sup>10,31</sup>, Owen R. White<sup>5</sup>, Z. Josh Huang<sup>27</sup>, Peter V. Kharchenko<sup>20</sup>, Lior Pachter<sup>1</sup>, John Ngai<sup>19</sup>, Aviv Regev<sup>10,29</sup>, Bosiljka Tasic<sup>1</sup>, Joshua D. Welch<sup>20</sup>, Jesse Gillis<sup>4</sup>, Evan Z. Macosko<sup>10</sup>, Bing Ren<sup>10,21</sup>, Joseph R. Ecker<sup>2,23</sup>, Hongkui Zeng<sup>10,32</sup> & Eran A. Mukamel<sup>10,32</sup>

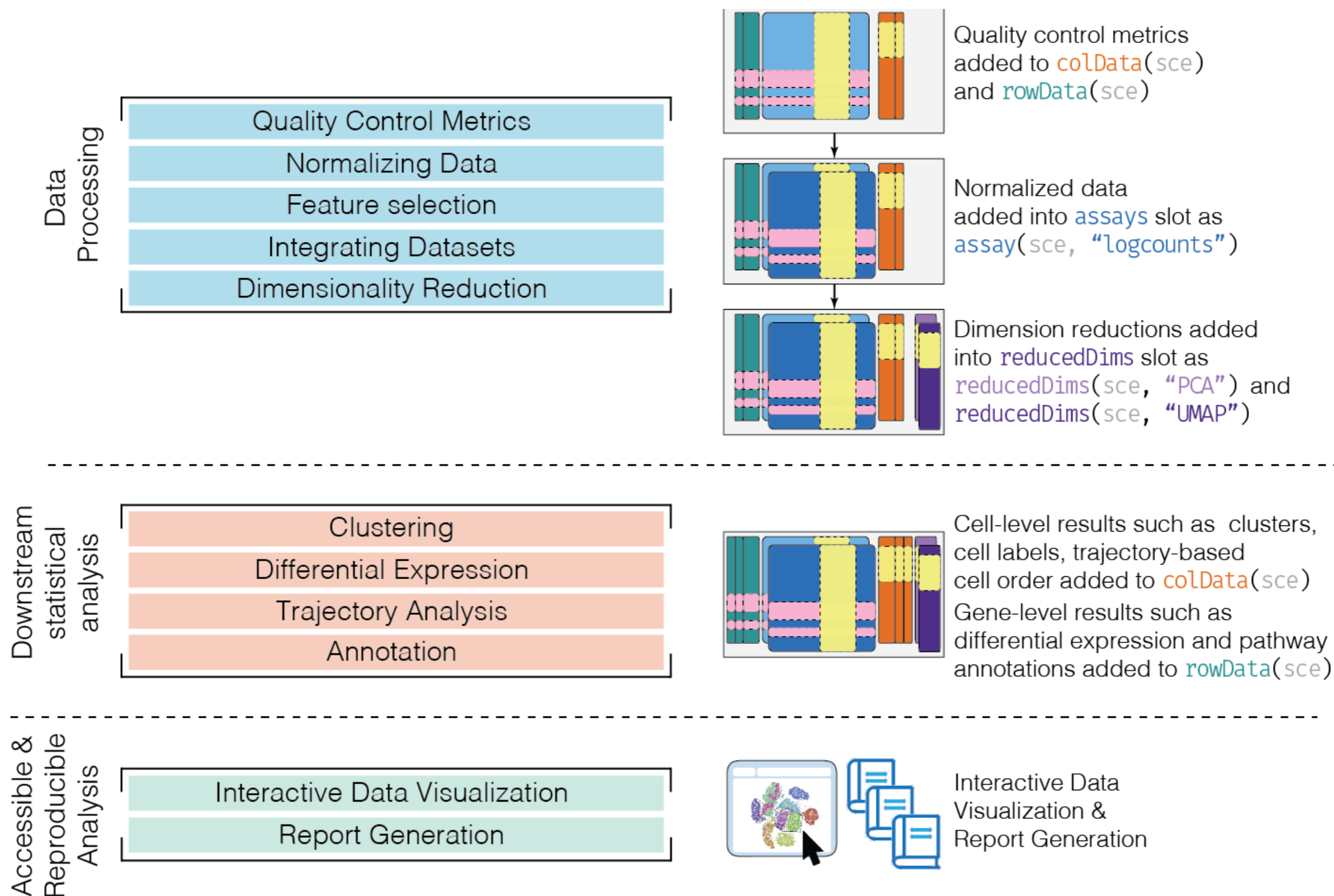
Single-cell transcriptomics can provide quantitative molecular signatures for large, unbiased samples of the diverse cell types in the brain<sup>1–3</sup>. With the proliferation of multi-omics datasets, a major challenge is to validate and integrate results into a biological understanding of cell-type organization. Here we generated **transcriptomes and epigenomes** from more than 500,000 individual cells in the mouse primary motor cortex, a structure that has an evolutionarily conserved role in locomotion. We developed computational and statistical methods to integrate multimodal data and



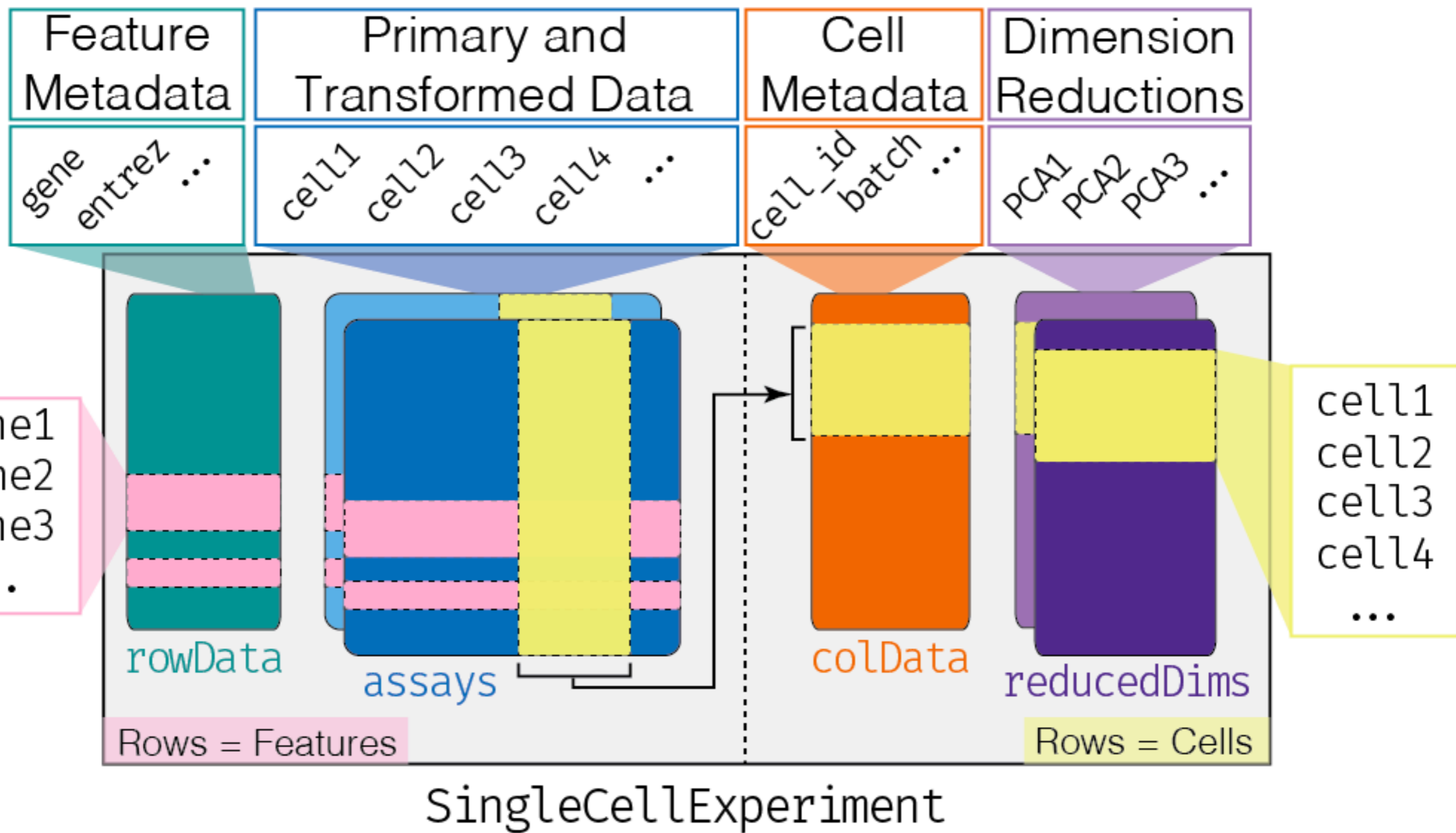
# A TYPICAL WORKFLOW



# A TYPICAL WORKFLOW

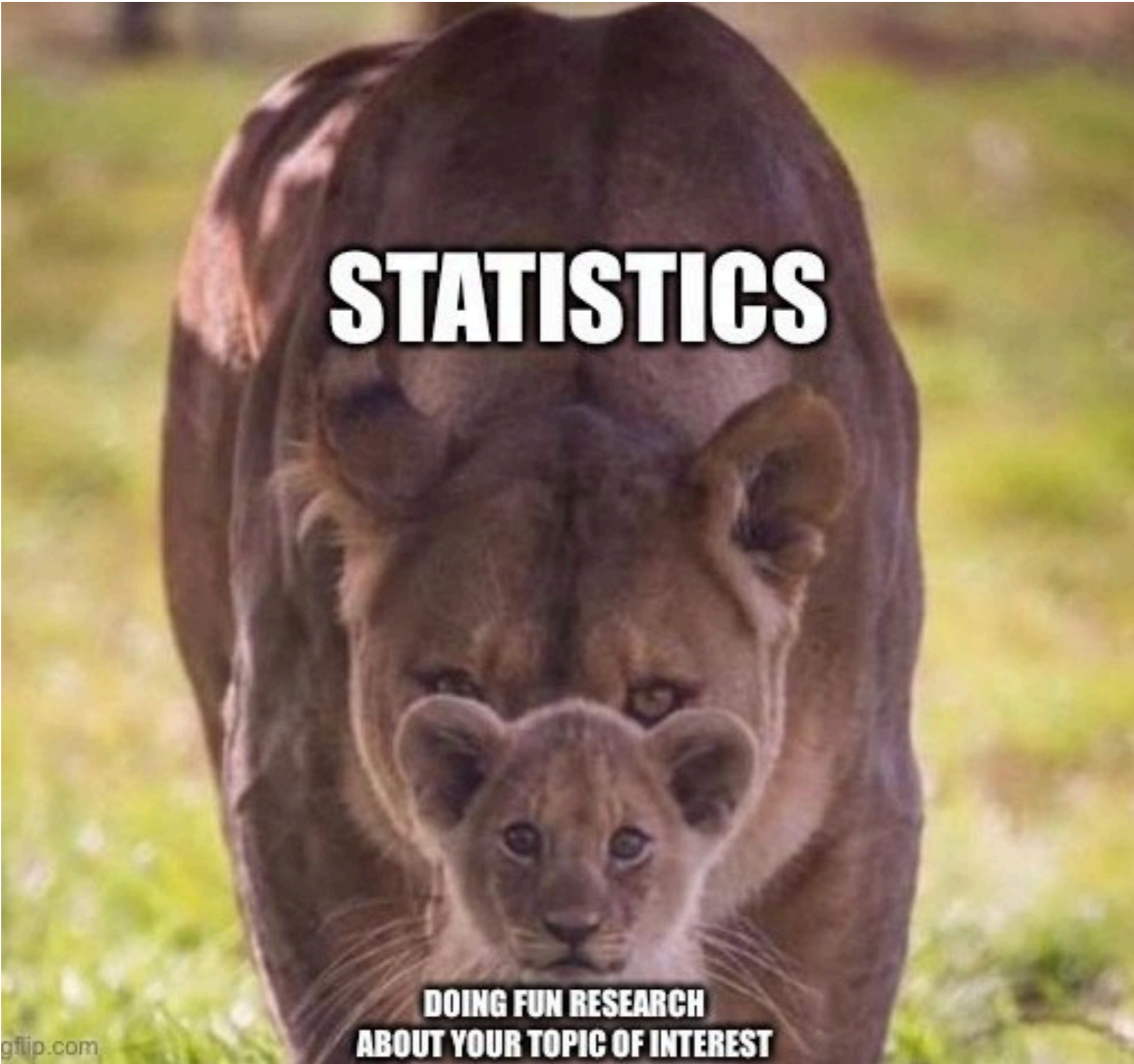


# THE SINGLECELLEXPERIMENT CLASS

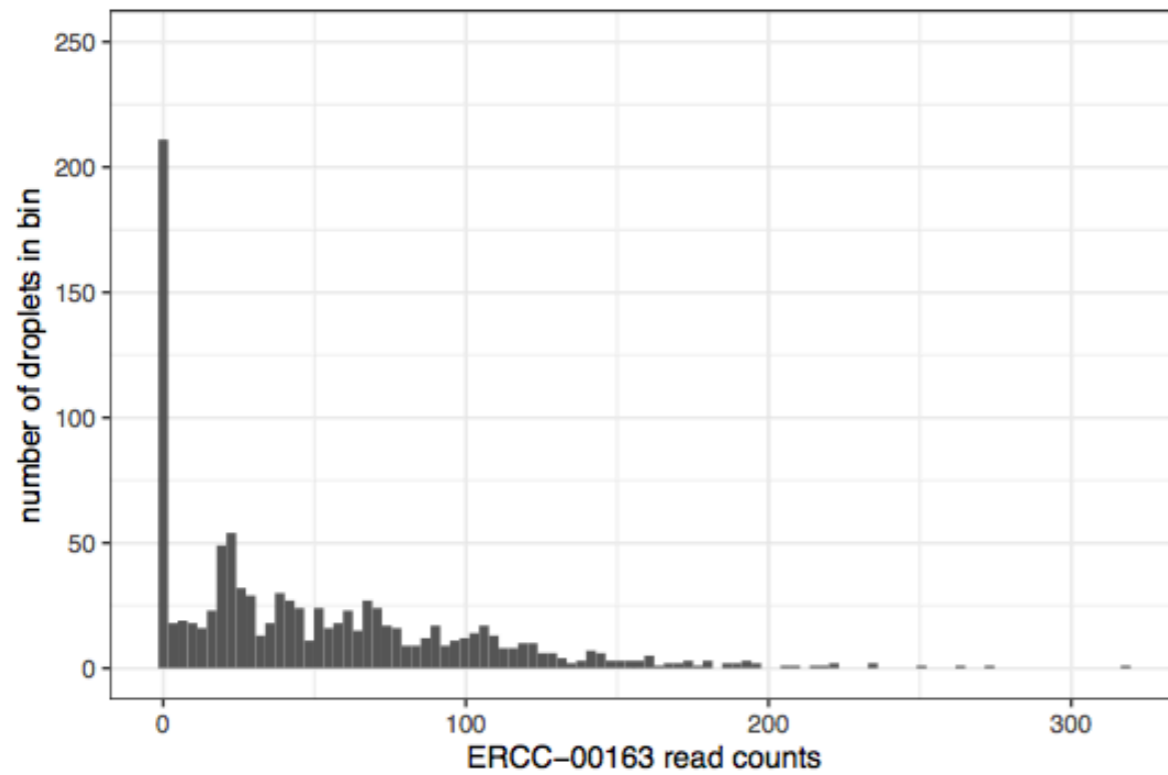




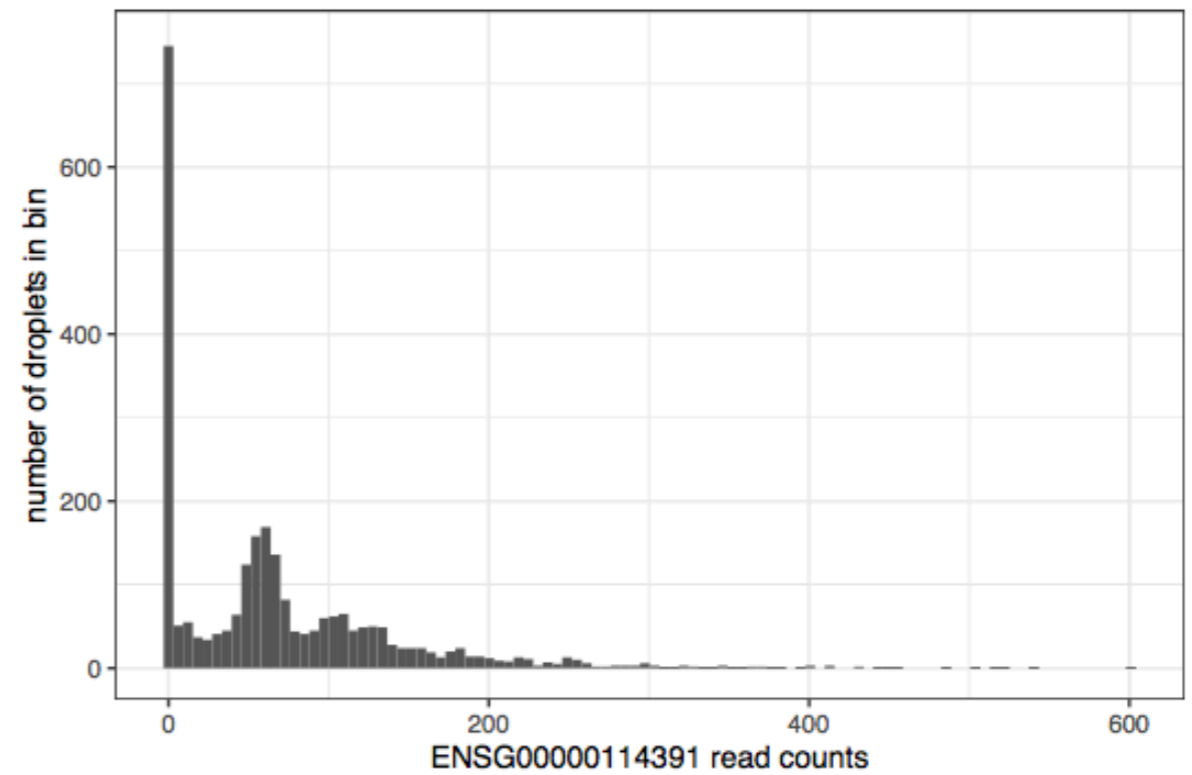
# DATA PROPERTIES



# READ COUNT DISTRIBUTION

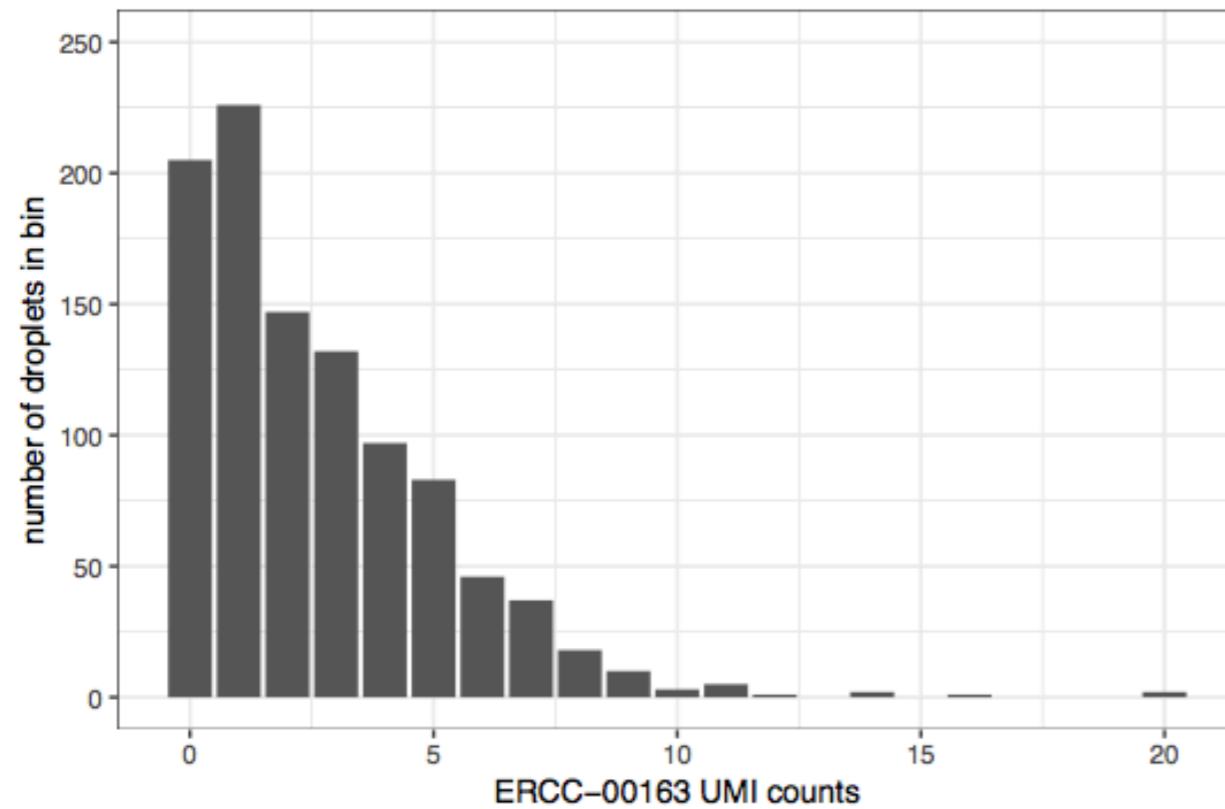


(a) Read counts- technical replicates

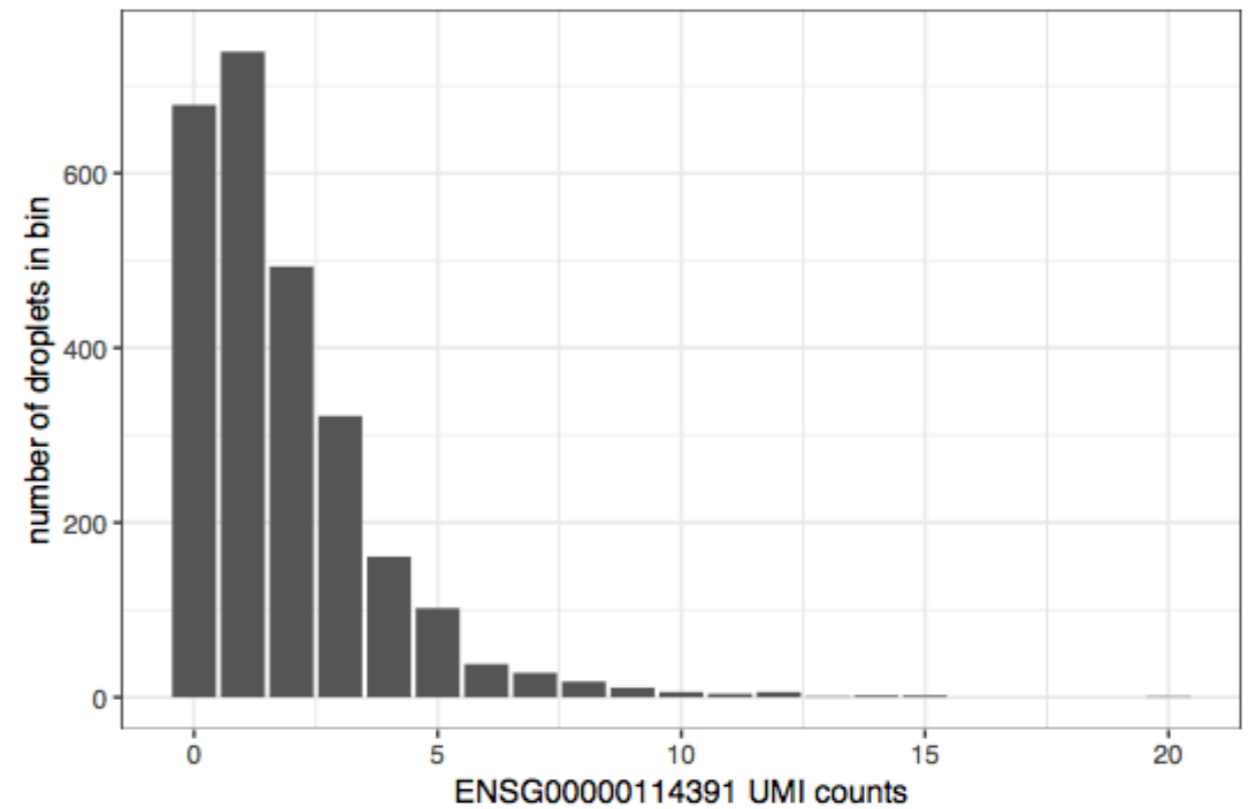


(b) Read counts- biological replicates

# UMI COUNT DISTRIBUTION

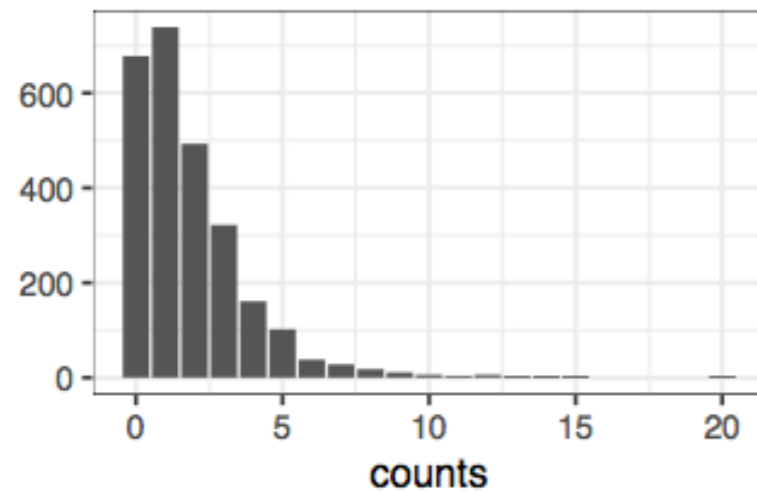


(c) UMI counts- technical replicates

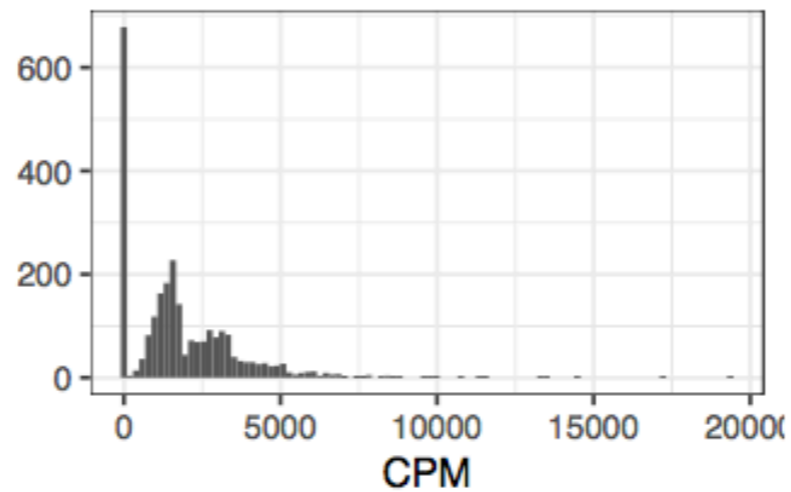


(d) UMI counts- biological replicates

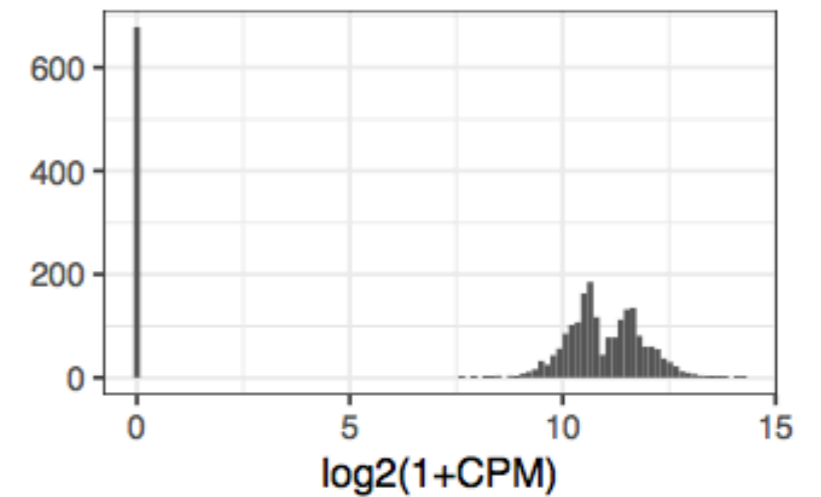
# LOG TRANSFORMATION DOES NOT HELP!



(a) UMI counts

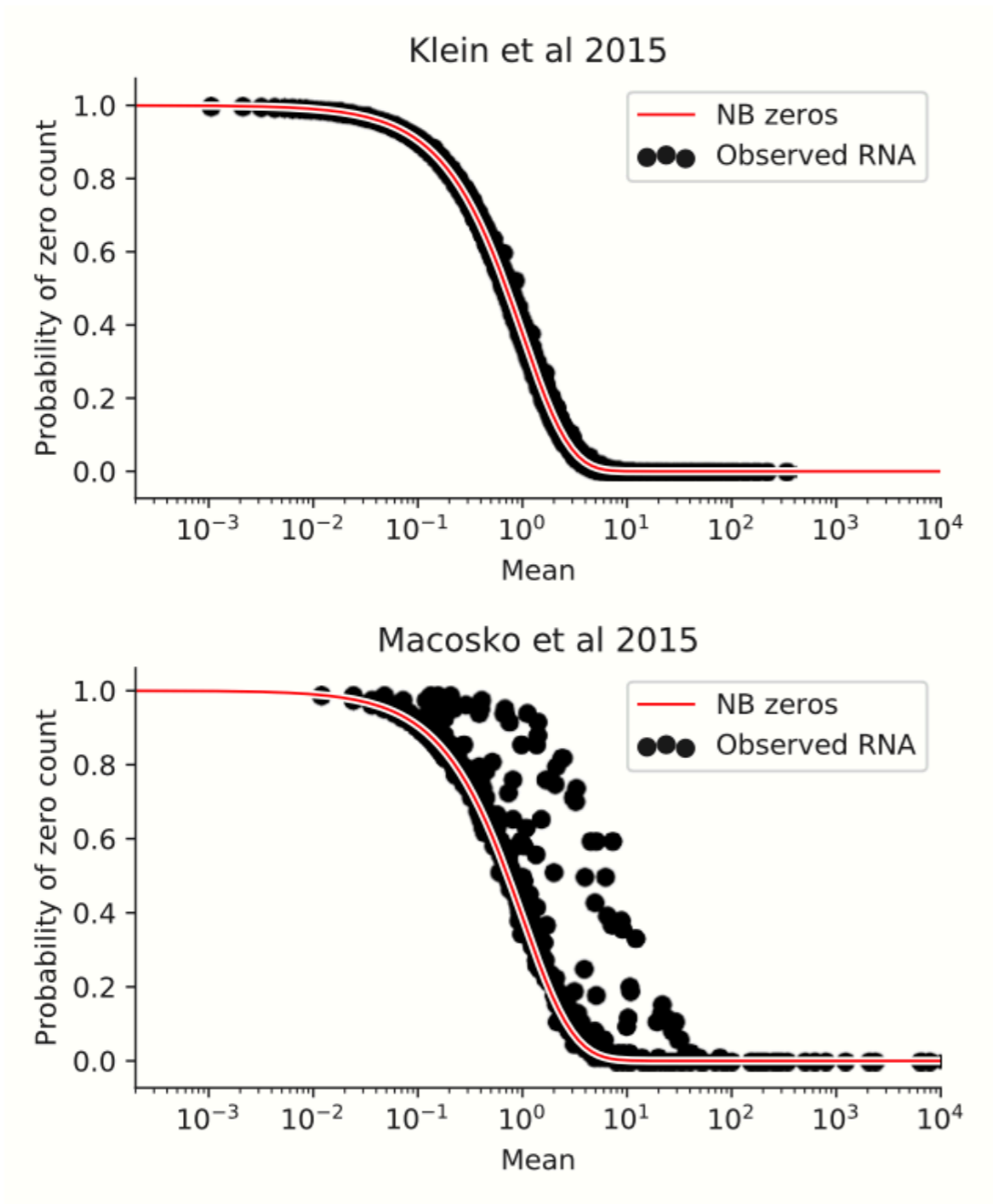


(b) counts per million (CPM)



(c) log of CPM

# SHOULD WE MODEL ZERO INFLATION?



# SHOULD WE MODEL ZERO INFLATION?

- ▶ Non-UMI data: very likely.
- ▶ UMI data: probably not.

## EXPLORATORY DATA ANALYSIS!

- ▶ Measurement vs expression models ([Sarkar & Stephens 2021](#))

**Table 1 | Single-gene models for scRNA-seq data**

Expression model	Observation model	Method <sup>a</sup>
Point mass (no variation)	Poisson	Analytic
Gamma	Negative binomial	MASS <sup>41</sup> , edgeR <sup>42</sup> , DESeq2 (ref. <sup>43</sup> ), SAVER <sup>20</sup> , BASICS <sup>44</sup>
Point-Gamma	ZINB	PSCL <sup>45</sup>
Unimodal (nonparametric)	Unimodal	ashr <sup>24,46</sup>
Point-exponential family	Flexible	DESCEND <sup>4</sup>
Fully nonparametric <sup>47</sup>	Flexible	ashr

Different expression models, when combined with the Poisson measurement model, yield different observation models. <sup>a</sup>Previously published methods and software packages that use the corresponding observation model to analyze data.

**Table 2 | Multigene models for scRNA-seq data**

Link function	Noise distribution	Method <sup>a</sup>
Identity	None	NMF <sup>48</sup> , scHPF <sup>49</sup>
Identity	Gamma	NBMF <sup>50</sup>
log	None	GLM-PCA <sup>19</sup>
log	Gamma	scNBMF <sup>51</sup> , GLM-PCA <sup>19</sup>
log	Point-Gamma	ZINB-WaVE <sup>52</sup>
Neural network	Point-Gamma	scVI <sup>29</sup> , DCA <sup>21</sup>

Multigene models partition variation in true expression into structured and stochastic components. The link function describes a transformation and the noise distribution indicates an assumption about the stochastic component. <sup>a</sup>Previously published methods and software packages that use the corresponding observation model to analyze data.

# SHOULD WE MODEL ZERO INFLATION?

Table 1: Hellinger distance between zinb and NB distribution

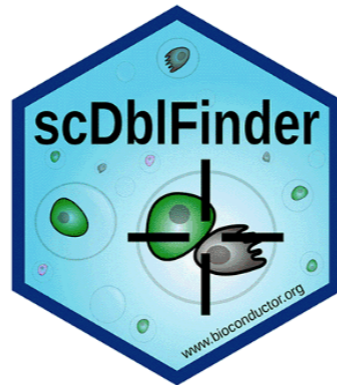
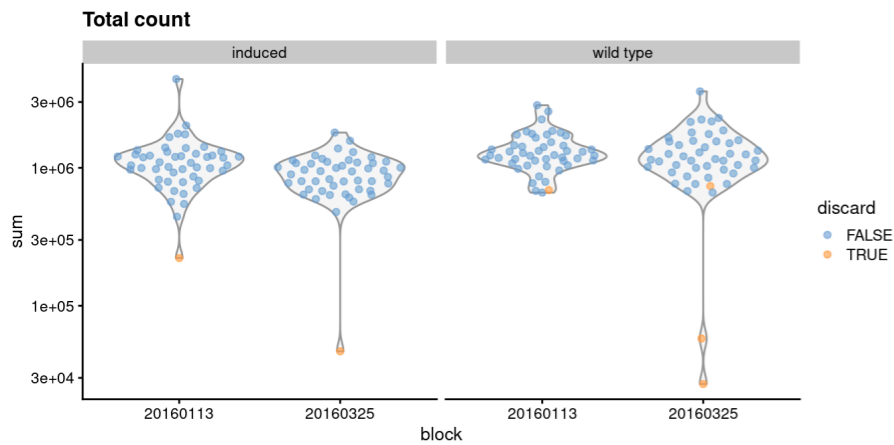
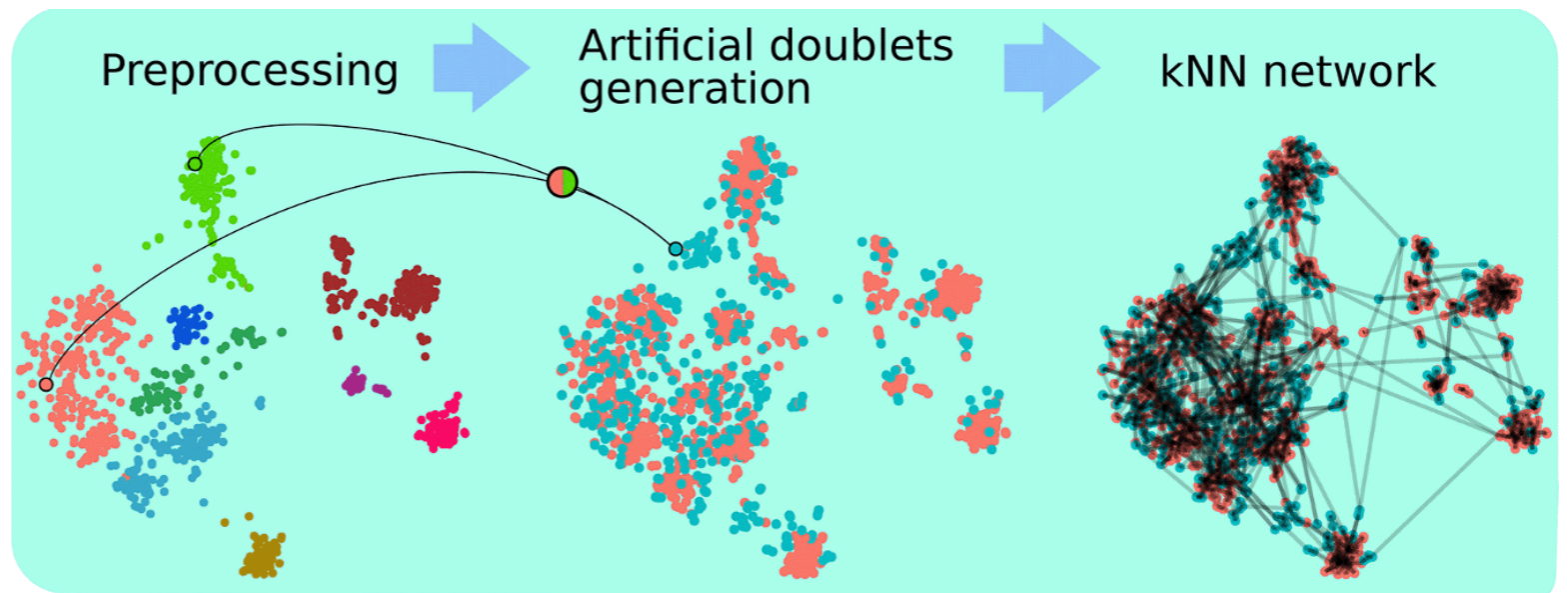
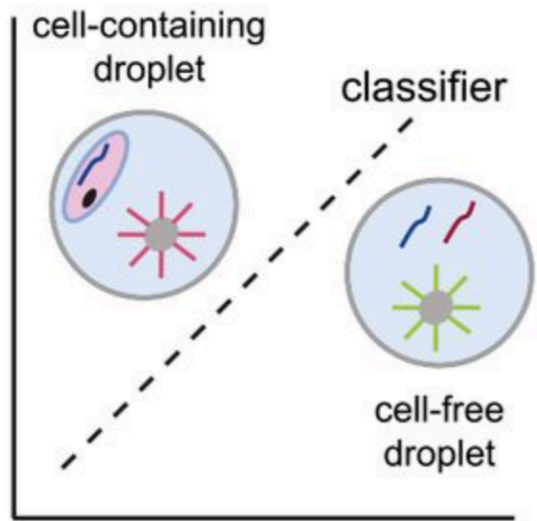
$\theta_0$	$\mu_0$	$\pi$									
		0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
0.5	0.5	0.00	0.02	0.05	0.07	0.10	0.13	0.16	0.19	0.23	0.28
	5	0.00	0.05	0.10	0.15	0.20	0.25	0.30	0.36	0.42	0.50
	10	0.00	0.06	0.12	0.18	0.23	0.29	0.34	0.40	0.47	0.55
	15	0.00	0.07	0.13	0.19	0.25	0.31	0.37	0.43	0.50	0.58
	20	0.00	0.07	0.14	0.20	0.26	0.32	0.38	0.44	0.51	0.60
	25	0.00	0.08	0.15	0.21	0.27	0.33	0.39	0.46	0.52	0.61
5	0.5	0.00	0.03	0.06	0.09	0.12	0.15	0.19	0.23	0.27	0.33
	5	0.00	0.13	0.22	0.30	0.37	0.43	0.50	0.57	0.64	0.72
	10	0.00	0.19	0.28	0.36	0.43	0.50	0.57	0.63	0.70	0.79
	15	0.00	0.21	0.30	0.38	0.45	0.52	0.59	0.65	0.72	0.81
	20	0.00	0.21	0.31	0.39	0.46	0.53	0.59	0.66	0.73	0.82
	25	0.00	0.22	0.32	0.40	0.47	0.53	0.60	0.67	0.74	0.82

# QUALITY CONTROL AND FILTERING



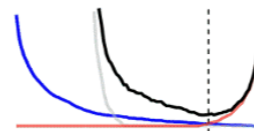
**DO YOUR DATA SPARK JOY?**





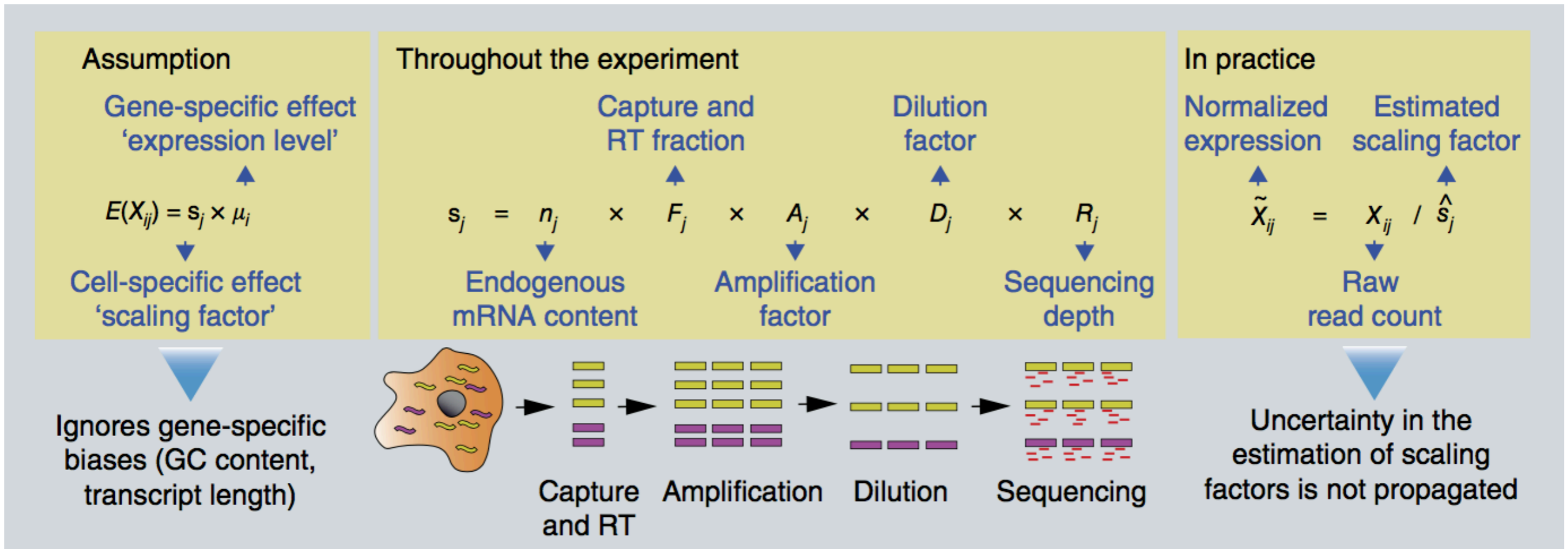
Gradient boosted trees

Thresholding

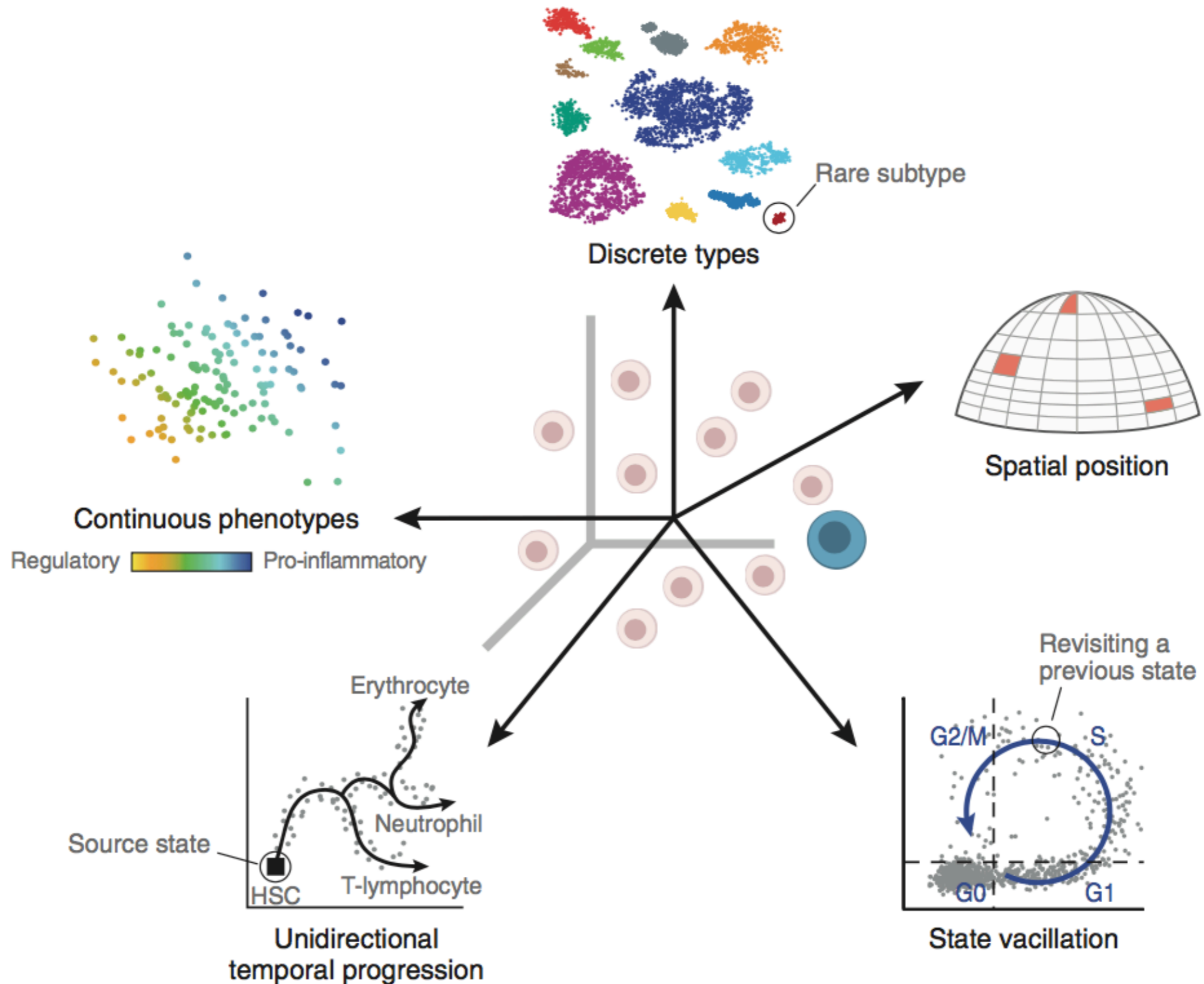


Cell-wise predictors

- library size
- cxds score
- marker coexpression
- ratio doublets in neighborhood
- principal components
- etc...



# DIMENSIONALITY REDUCTION



# DIMENSIONALITY REDUCTION

We talk about “dimensionality reduction” when referring to two different goals:

## 1. **Visualize** high-dimensional data

- ▶ Usually 2-3 dimensions
- ▶ Non-linear, local techniques

## 2. **Infer** low-rank signal from high-dimensional data

- ▶ Usually 10-50 dimensions
- ▶ Factor analysis models

# PRINCIPAL COMPONENT ANALYSIS (PCA)

- ▶ PCA is the starting point and baseline approach for both types of analysis.
- ▶ PCA can be used to visualize high-dimensional data in 2-3 dimensions.
- ▶ PCA can be seen as a solution of a factor analysis model for Gaussian data.

# DESIRED PROPERTIES OF DIMENSIONALITY REDUCTION MODELS

- ▶ Accounting for the count nature of the data, overdispersion, and possibly zero inflation.
- ▶ General and flexible.
- ▶ Extract low-dimensional signal from the data.
- ▶ Adjust for complex effects (batch effects, sample quality).
- ▶ Scalable.

# EXAMPLE: TABULA MURIS DATA

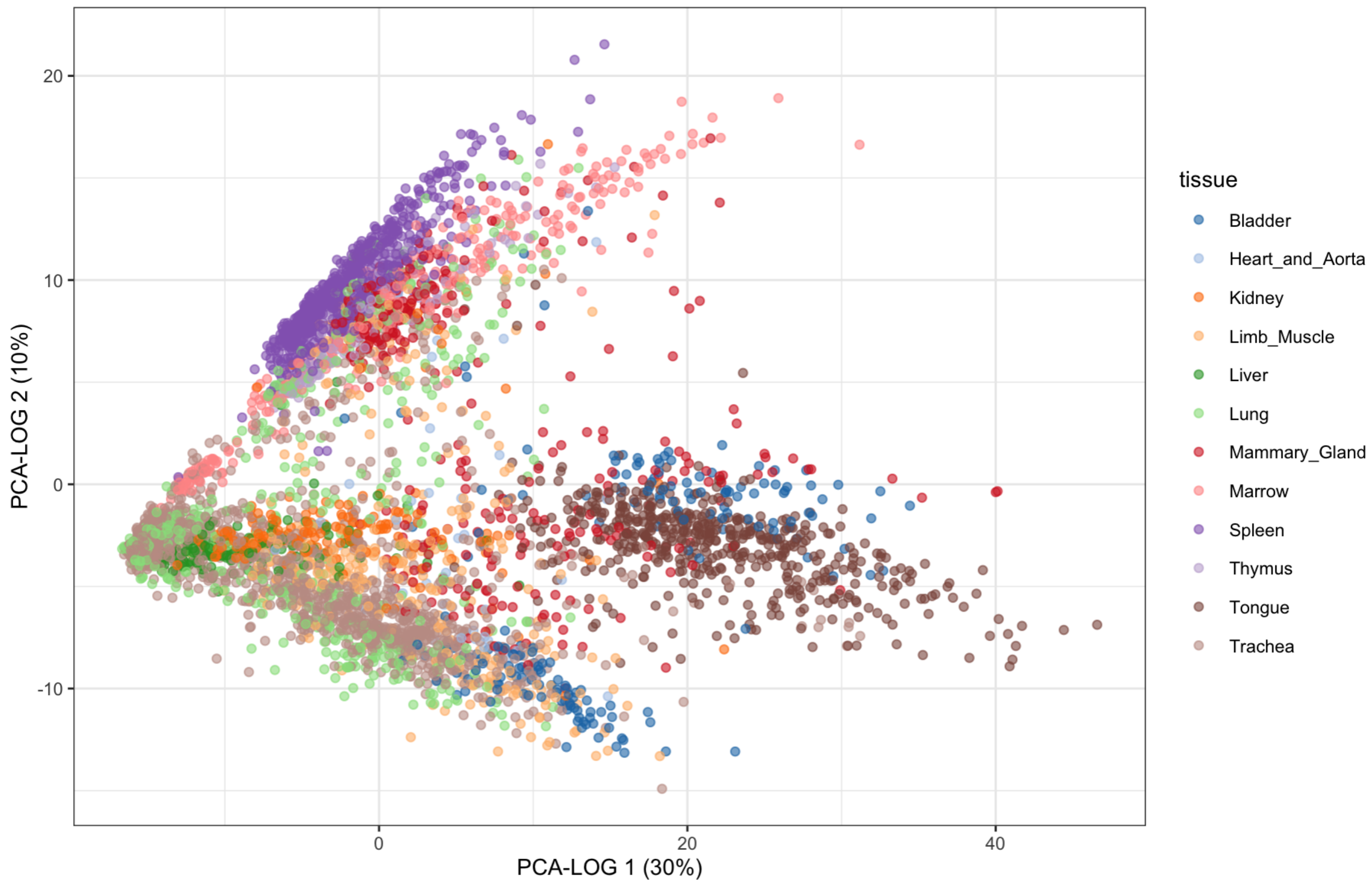
Tabula Muris is a project aimed at characterizing all cell types in the mouse.

The droplet dataset comprised 70,000 cells from 12 tissues.

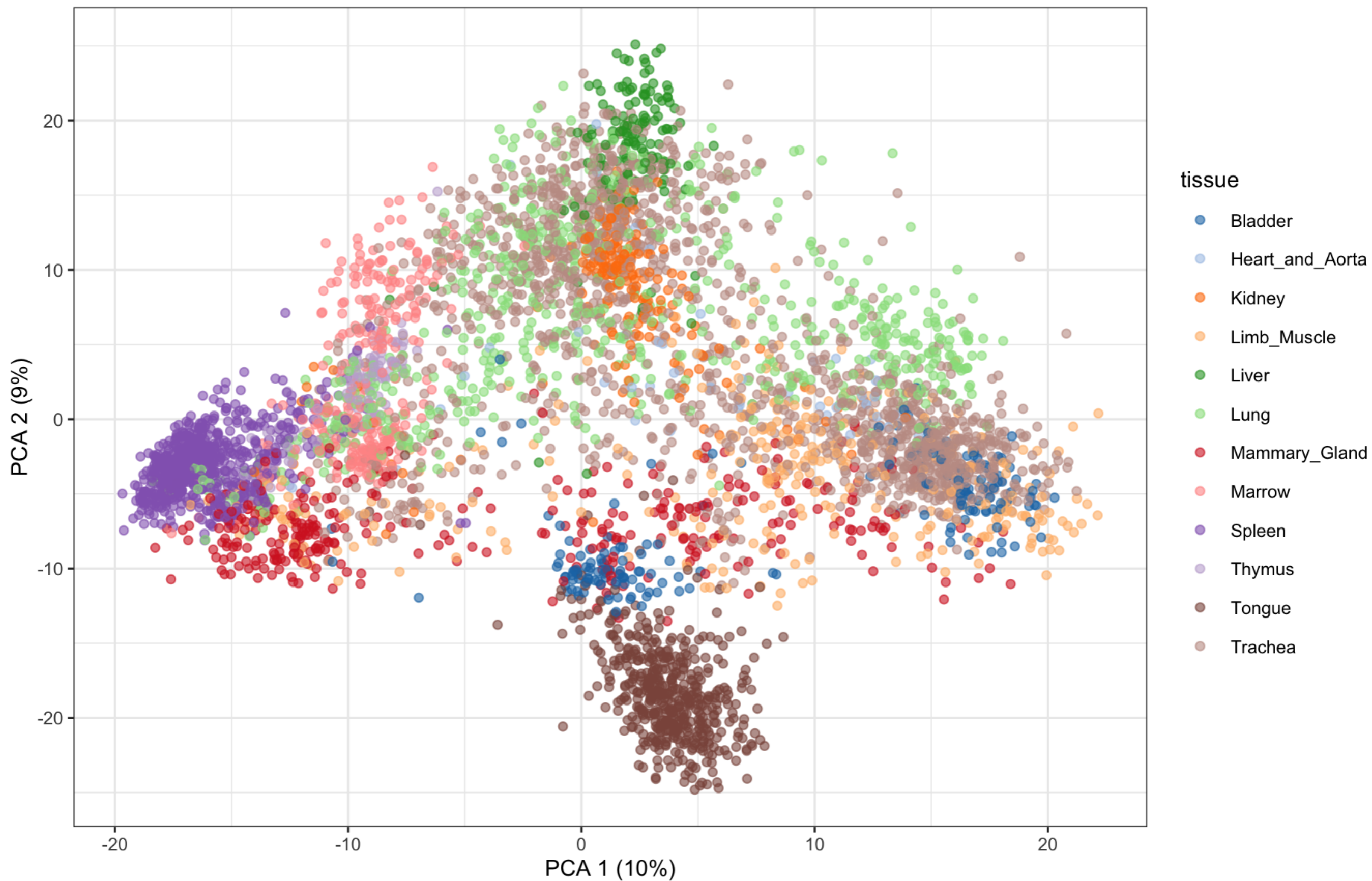
We see here a random subset of 5,000 cells, limiting the dataset to the 1,000 most variable genes.

[TabulaMurisData Bioconductor package.](#)

# TABULA MURIS: PCA (LOG SCALE)

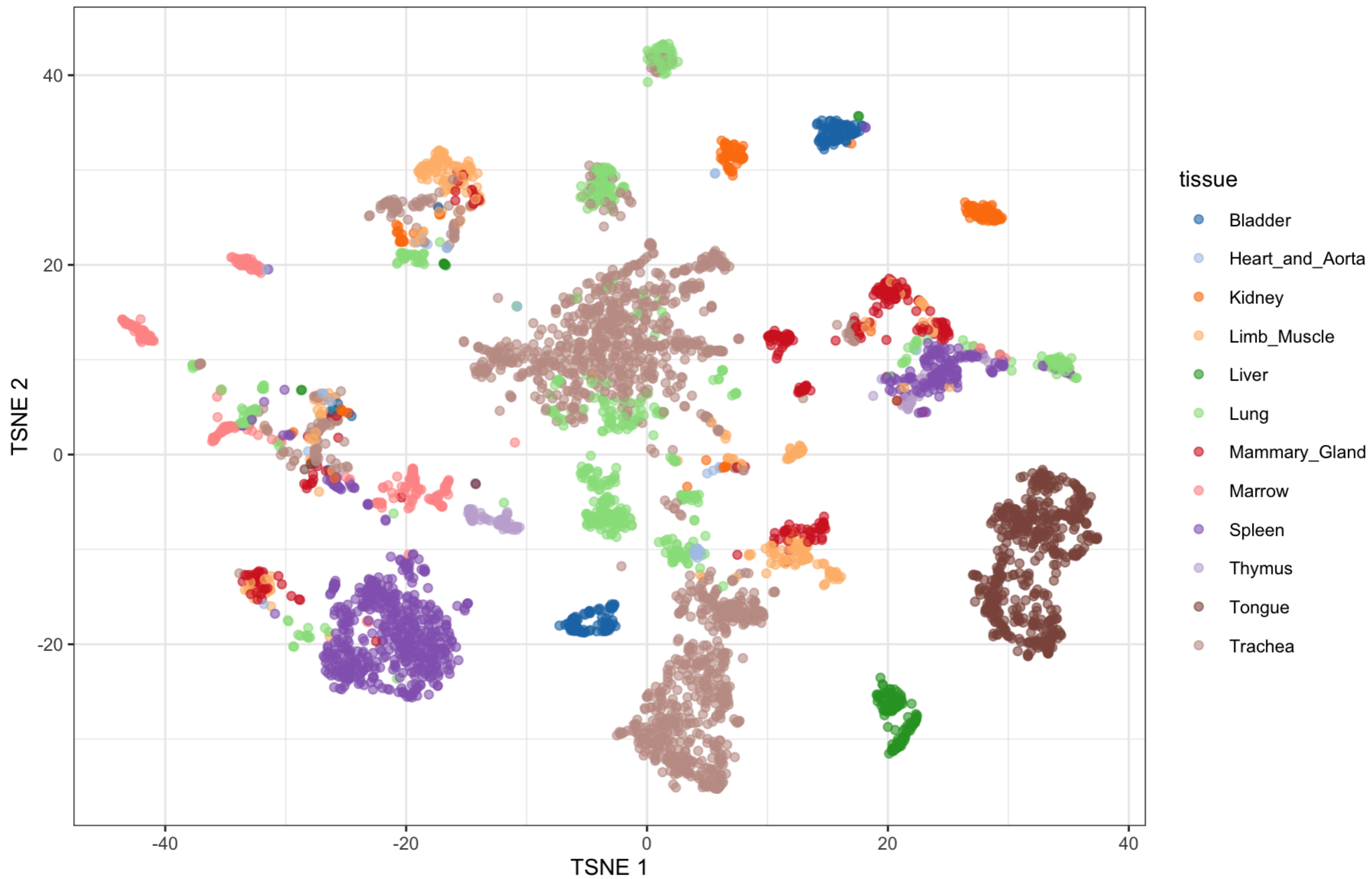


# TABULA MURIS: PCA AFTER SCRAN NORMALIZATION (LOG SCALE)

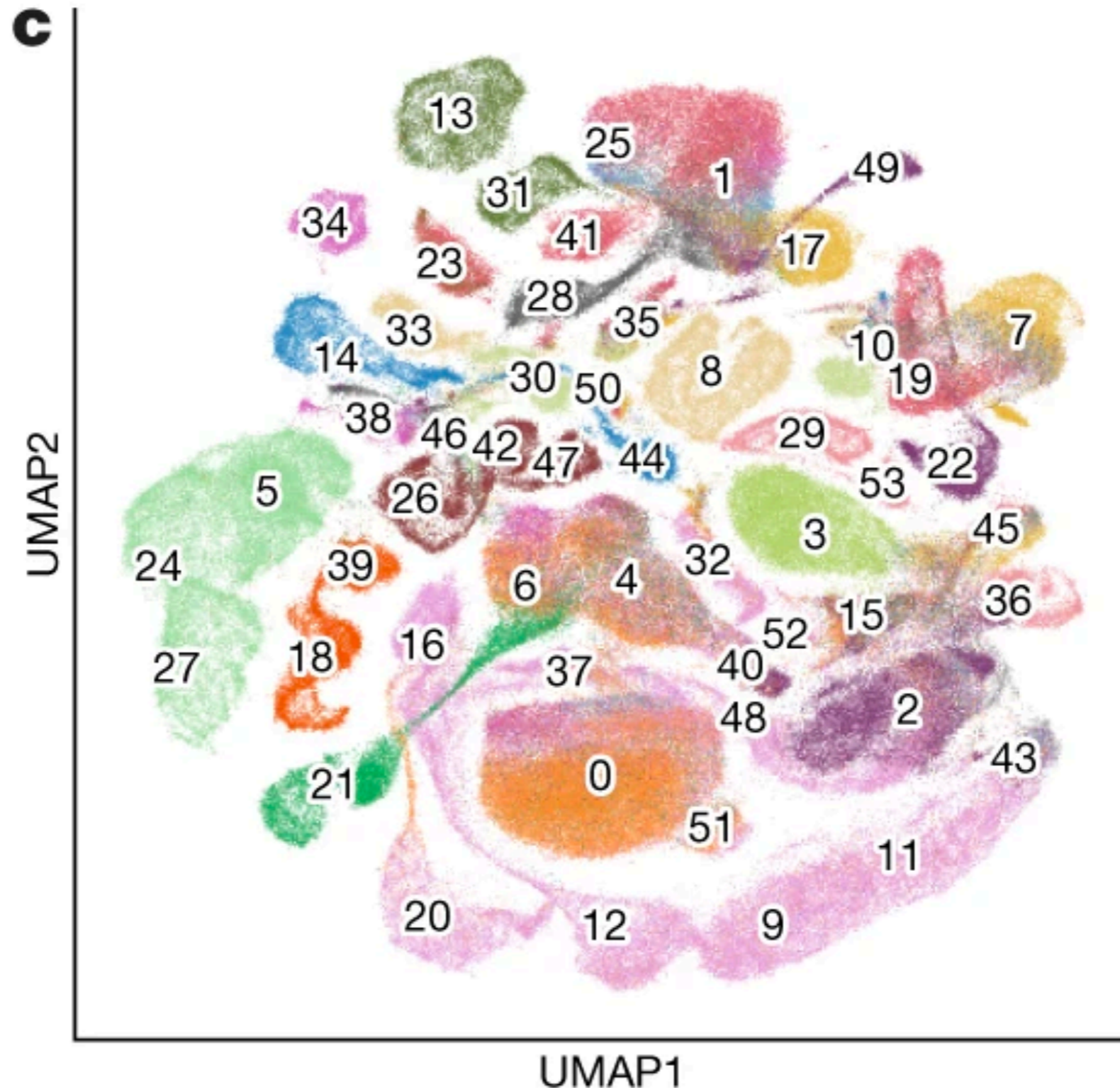




# TABULA MURIS: T-SNE



# TABULA MURIS: UMAP (ALL CELLS)



# PCA IS A LINEAR METHOD

- ▶ One way to define the first principal component is: the **linear combination** of the original variables that explain the most variability in the data.
- ▶ Similarly, subsequent PCs are linear combinations of the original variables that are orthogonal to the first and explain the most variance among the orthogonal combinations.
- ▶ Are we limiting ourselves by only looking at **linear** combinations?
- ▶ Would a non-linear method have more flexibility in explaining our data?

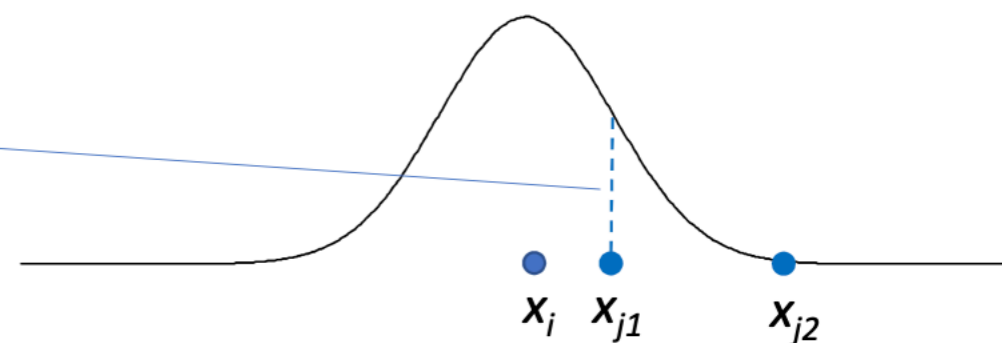
# T-DISTRIBUTED STOCHASTIC NEIGHBOR EMBEDDING (T-SNE)

- ▶ One option, very popular in single-cell genomics, is t-distributed Stochastic Neighbor Embedding (t-SNE).
- ▶ Briefly, the problem that we want to solve is to represent in a 2-3 dimensional map (*embedding*) the observations from a high-dimensional space **preserving as much as possible the distance between points.**

# STOCHASTIC EMBEDDING: PROBABILISTIC REPRESENTATION OF DISTANCES

- ▶ Similarity between two points,  $x_i$  and  $x_j$  in the **original high-dimensional space** is given by

$$p_{j|i} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|x_i - x_k\|^2 / 2\sigma_i^2)},$$



The denominator scales the sum of all the scores to 1

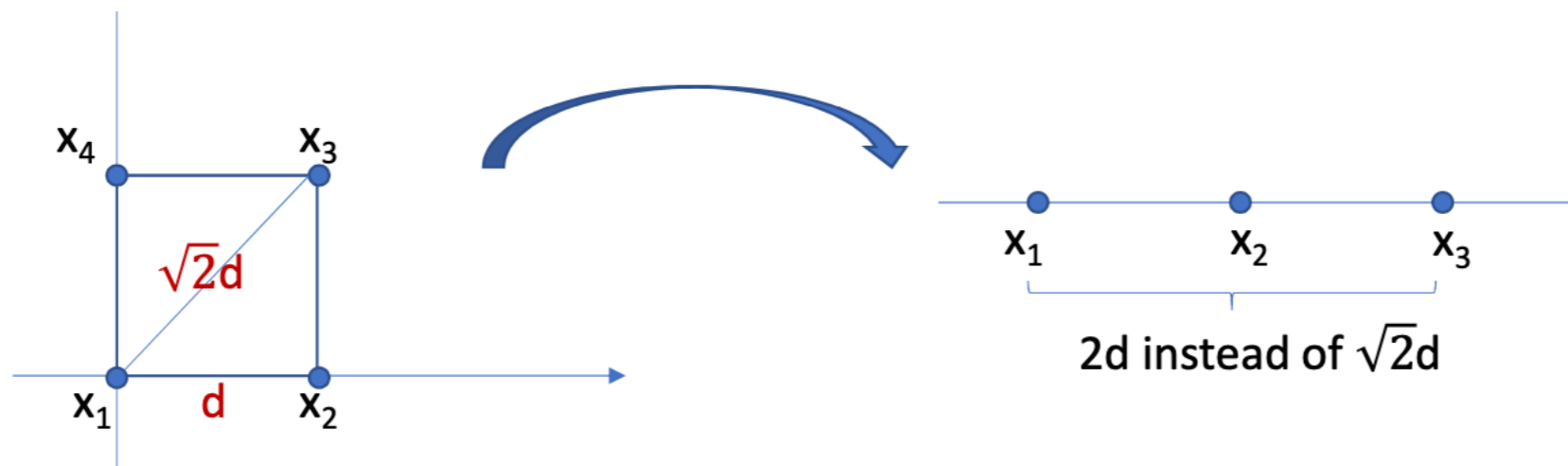
- ▶ Essentially, the probability that  $x_i$  would pick  $x_j$  as its neighbor.
- ▶ We set  $p_{i|i} = 0$  and actually use a symmetrized version that ensures  $p_{ij} = p_{ji}$ .

# STOCHASTIC EMBEDDING: PROBABILISTIC REPRESENTATION OF DISTANCES

- ▶ We could define a similar density in the *low-dimensional space*, but we use a t-distribution instead of a Gaussian kernel

$$\cancel{q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)},} \quad q_{ij} = \frac{(1 + \|y_i - y_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|y_k - y_l\|^2)^{-1}}.$$

- ▶ The t distribution has heavier tails and partially account for the *crowding problem*.



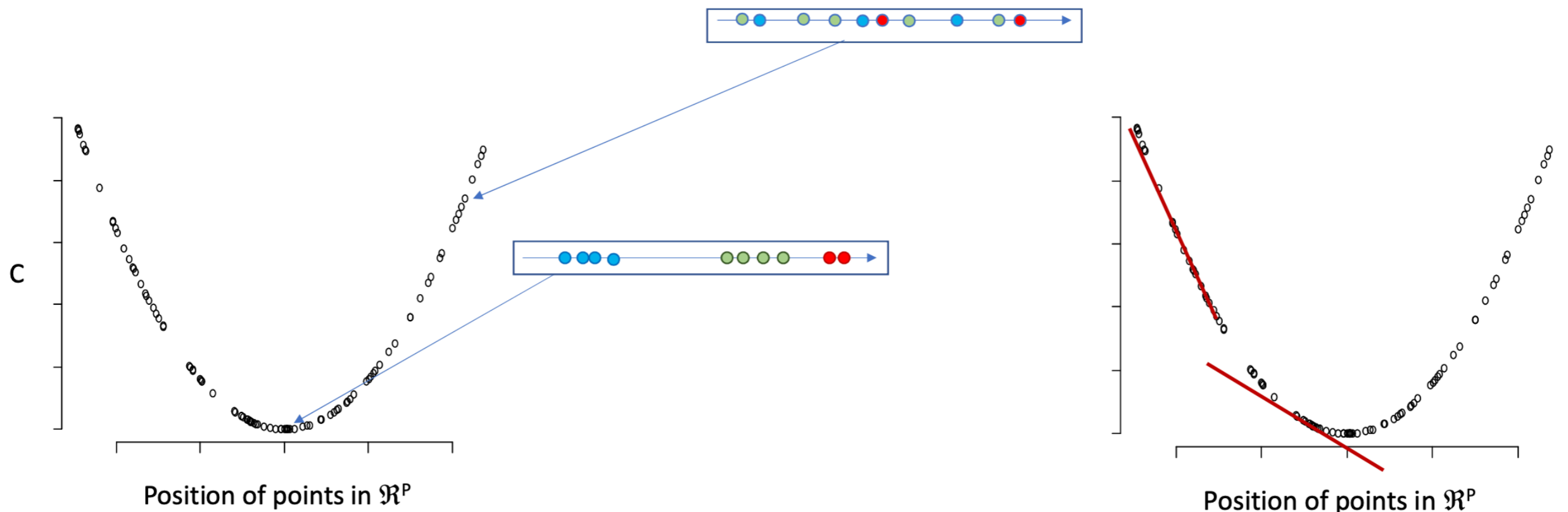
# T-SNE ALGORITHM

$$C = KL(P||Q) = \sum_i \sum_j p_{ij} \log \frac{p_{ij}}{q_{ij}}.$$

$$q_{ij} = \frac{\exp(-\|y_i - y_j\|^2)}{\sum_{k \neq l} \exp(-\|y_k - y_l\|^2)},$$

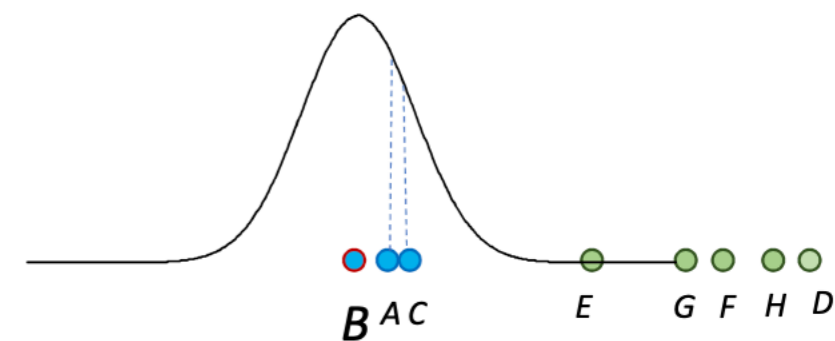
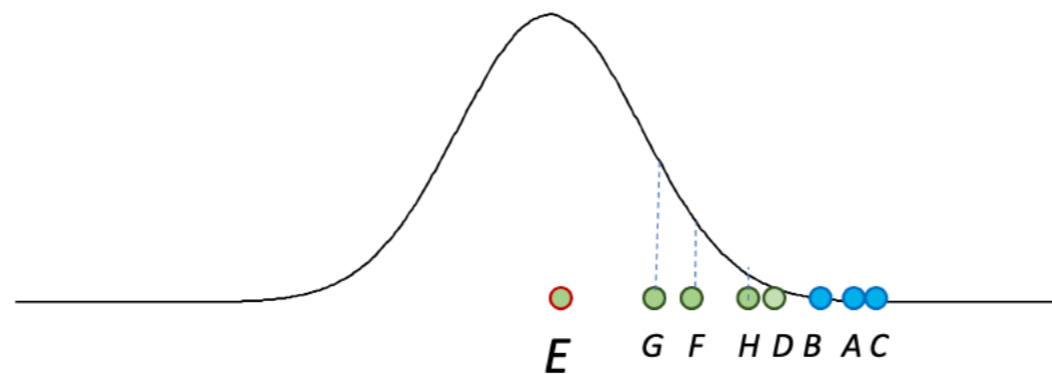
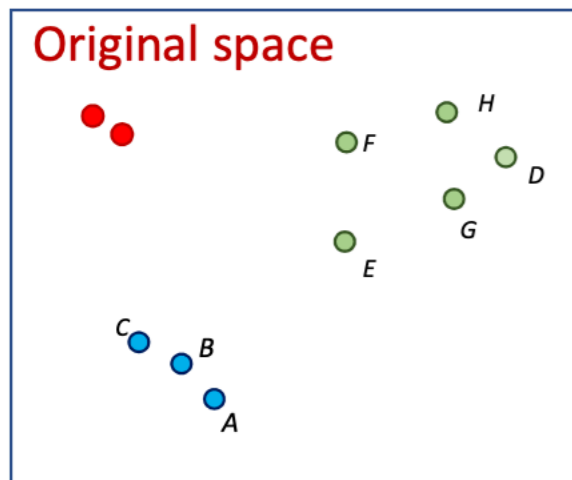
$$p_{ij} = \frac{\exp(-\|x_i - x_j\|^2 / 2\sigma^2)}{\sum_{k \neq l} \exp(-\|x_k - x_l\|^2 / 2\sigma^2)},$$

- ▶ We minimize the Kullback-Leibler (KL) divergence between the two distributions with *gradient descent*.



# CHOICE OF $\sigma^2$

- ▶ It's not appropriate to have a single value of  $\sigma^2$  as you need a smaller value in more dense regions.
- ▶ The user controls it through a parameter called **perplexity**
- ▶ **Perplexity can have a big impact on the result!**



$$\text{Perp}(P_i) = 2^{H(P_i)},$$

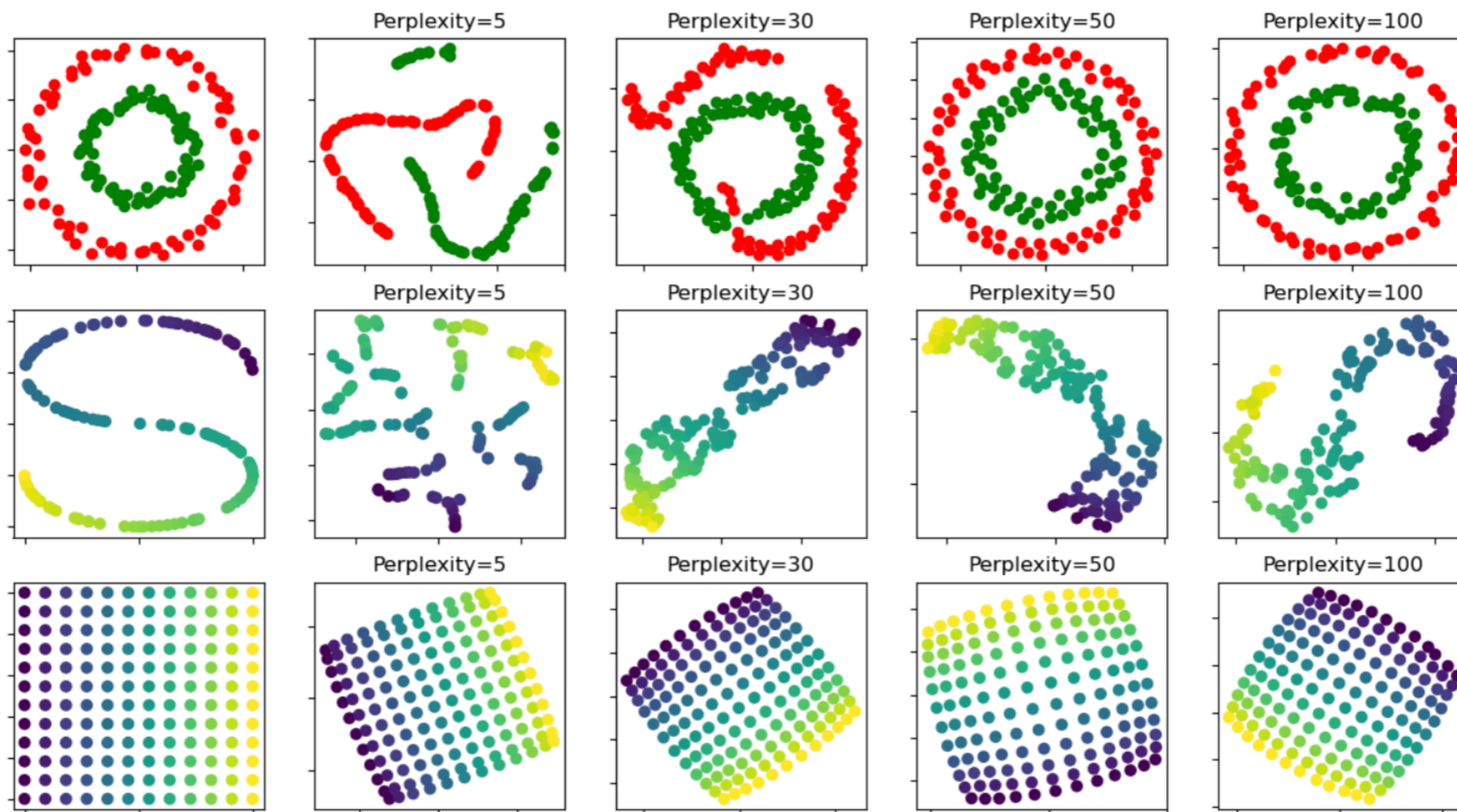
$$H(P_i) = -\sum_j p_{j|i} \log_2 p_{j|i}.$$

$H(P_i)$  is the Shannon entropy of  $P_i$  measured in bits



# T-SNE ART, OR THE CHOICE OF THE PERPLEXITY PARAMETER

<https://distill.pub/2016/misread-tsne/>

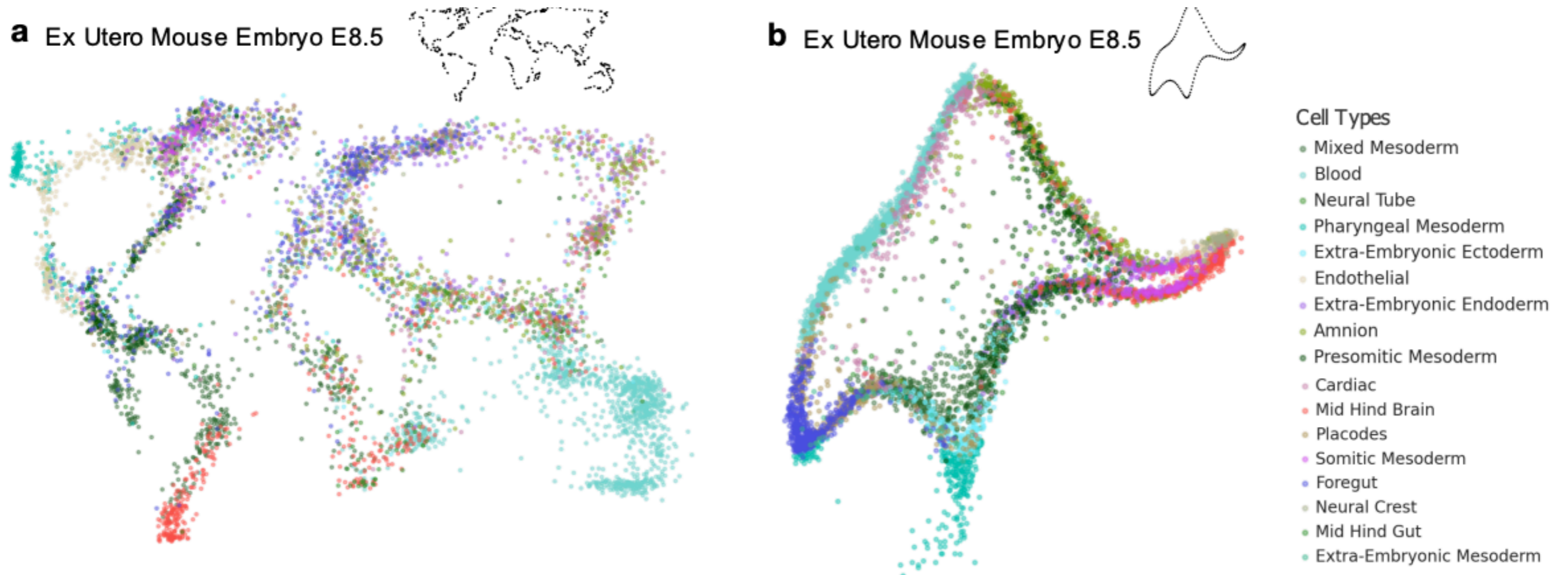


# LIMITATIONS OF T-SNE (AND UMAP)

- ▶ Unlike PCA, we do not have a simple interpretation for our low-dimensional embedding (the axes have “no meaning”).
- ▶ t-SNE preserves only the local structure (who is neighbor of whom) but not the global structure
- ▶ There is no guarantee of convergence to the global minimum (non-convex problem), hence different runs will lead to different embeddings.
- ▶ Some argue that t-SNE and UMAP do not even preserve the local structure or the neighbors (Chari et al. 2021)

# LIMITATIONS OF T-SNE (AND UMAP)

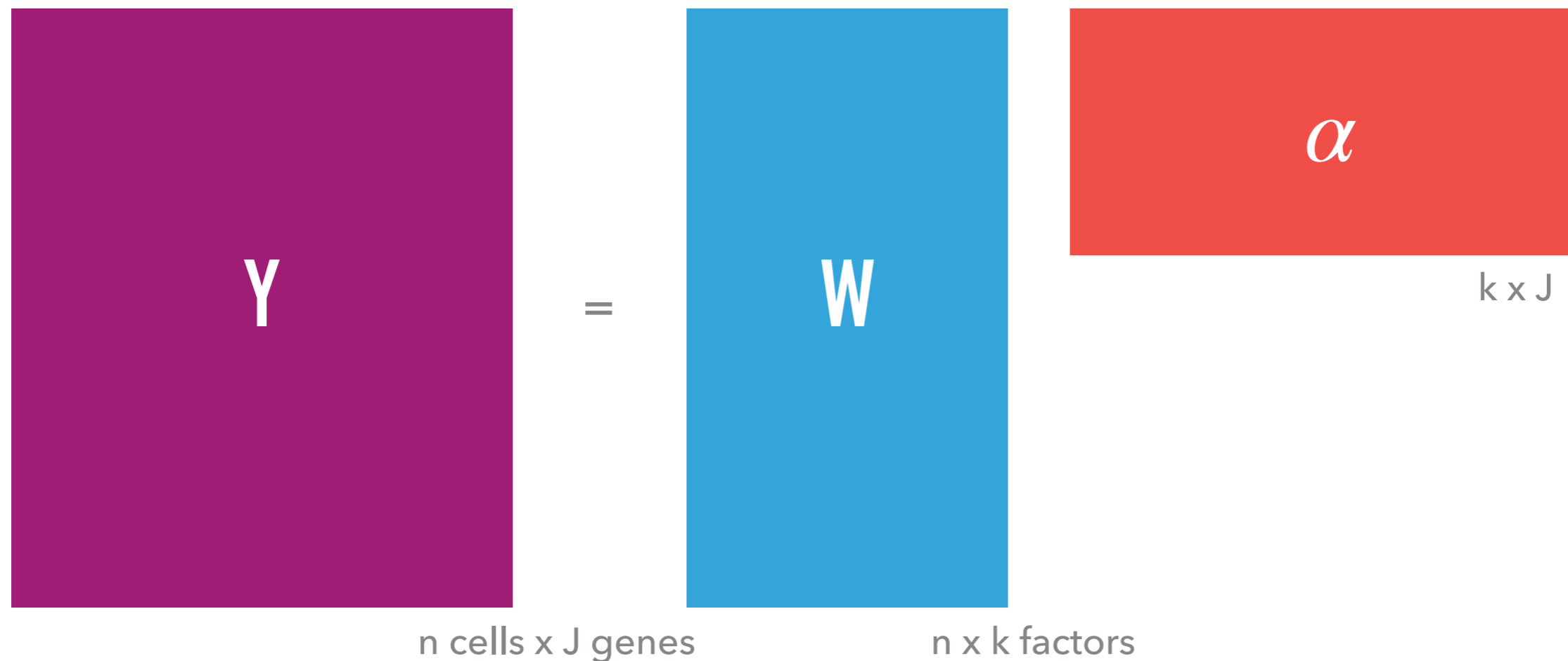
- ▶ The “shape” of the data in the embedding is arbitrary.



# FACTOR ANALYSIS

From a statistical model's perspective, we can state the problem using the following model

$$Y = W\alpha + \varepsilon$$



# FACTOR ANALYSIS

The goal is to find  $k \ll J$  factors that describe, with the minimum possible loss of information, the  $J$  original variables (genes).

We can show that if  $\varepsilon$  (or equivalently  $Y$ ) is Gaussian, a solution of the model is PCA.

# ADVANTAGES OF PCA

In one word: **interpretability!**

- ▶ The first principal component is the direction of greater variability in the data.
- ▶ It is easy to compute the variance explained by the first  $m$  principal components.
- ▶ Very computationally efficient.

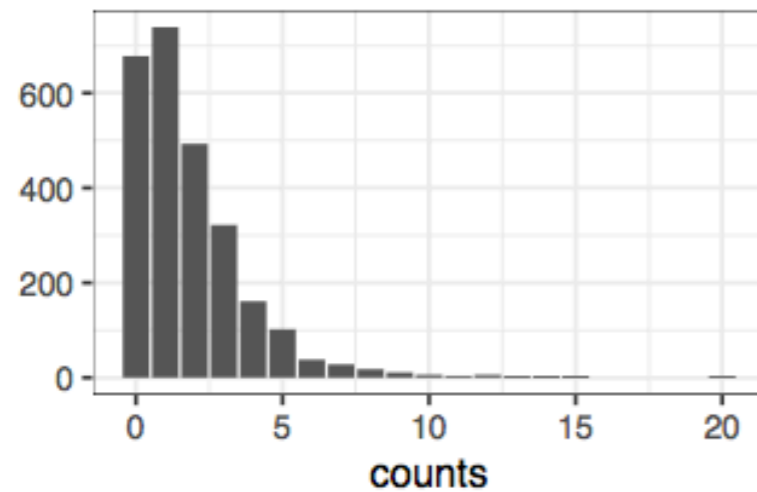
# SO... WHY NOT PCA?

The main issue of PCA for single-cell data is that the data are non-negative integer counts, which exhibit skewed distributions and are not well fit by a Gaussian model.

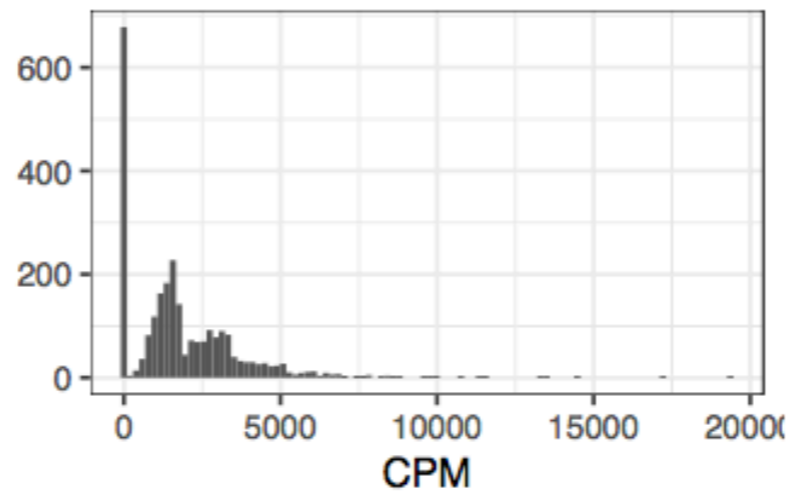
A simple, and somewhat effective, solution is to transform the data, e.g., by  $\log(x + 1)$ , but this is not always straightforward:

- ▶ Which transformation to use?
- ▶ Do we need to normalize the data for sequencing depth and other cell-specific effects?
- ▶ Zero counts complicate the analysis.

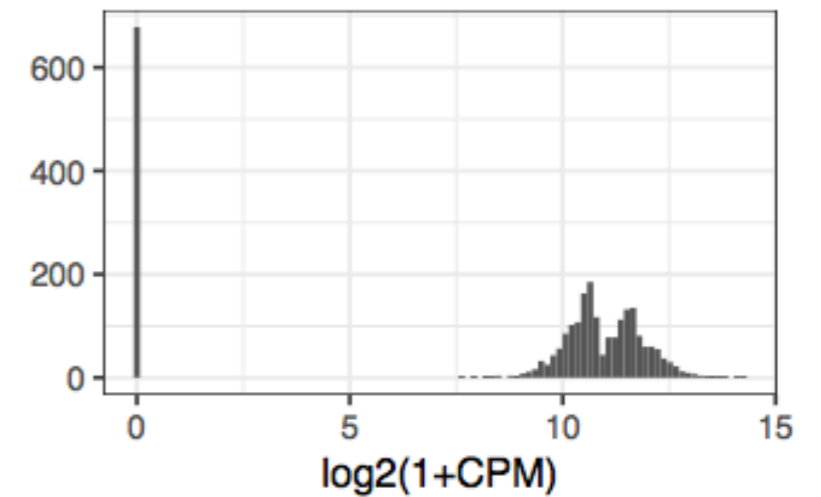
# REMEMBER THE EFFECT OF LOG TRANSFORMATION...



(a) UMI counts



(b) counts per million (CPM)



(c) log of CPM



# GLM-PCA

One alternative to transforming the data, is to generalize our model to non-Gaussian data.

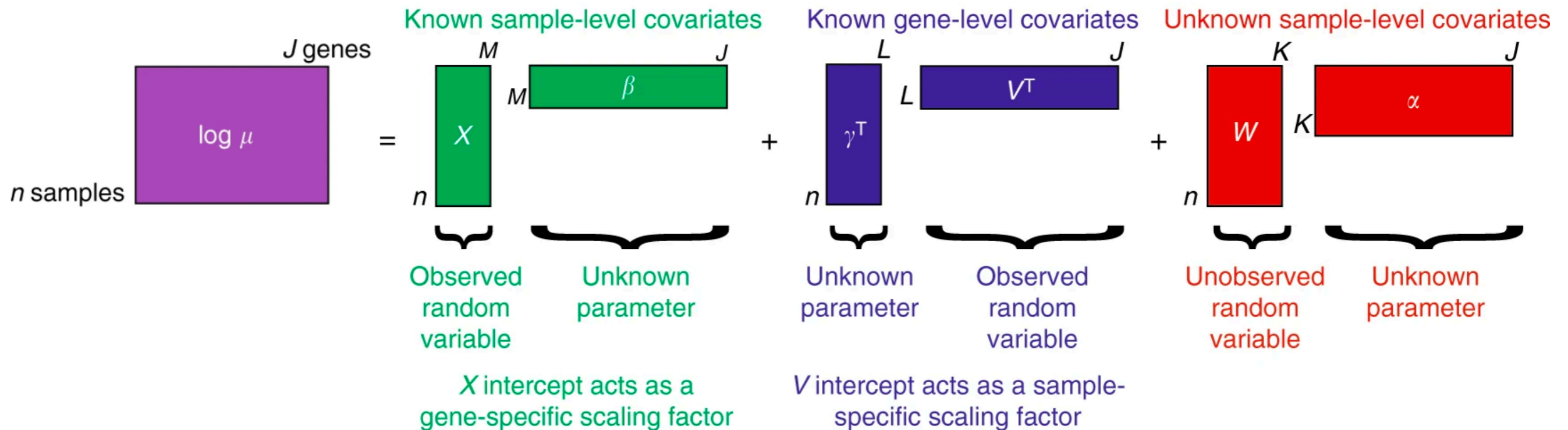
This can be done by defining a set of models, known as GLM-PCA () that extend the framework to a set of well behaving distributions (exponential family) similar to how GLM extends the linear model.

In particular, since we have count data, we can use the Poisson or negative binomial model, which has a log link function.

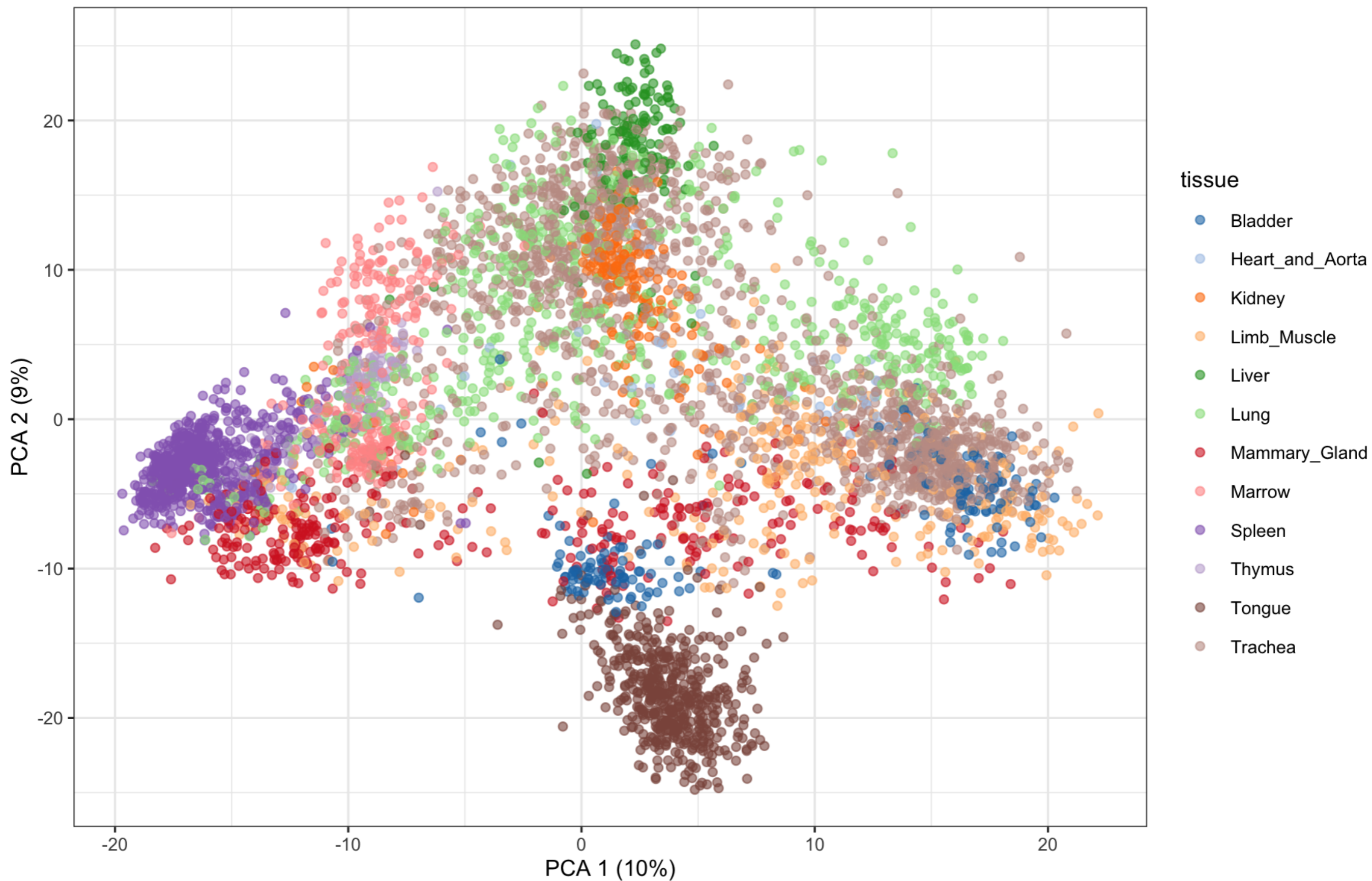
$$E[Y | W] = \mu, \quad \log \mu = W\alpha$$

# GLM-PCA / WANTED VARIATION EXTRACTION

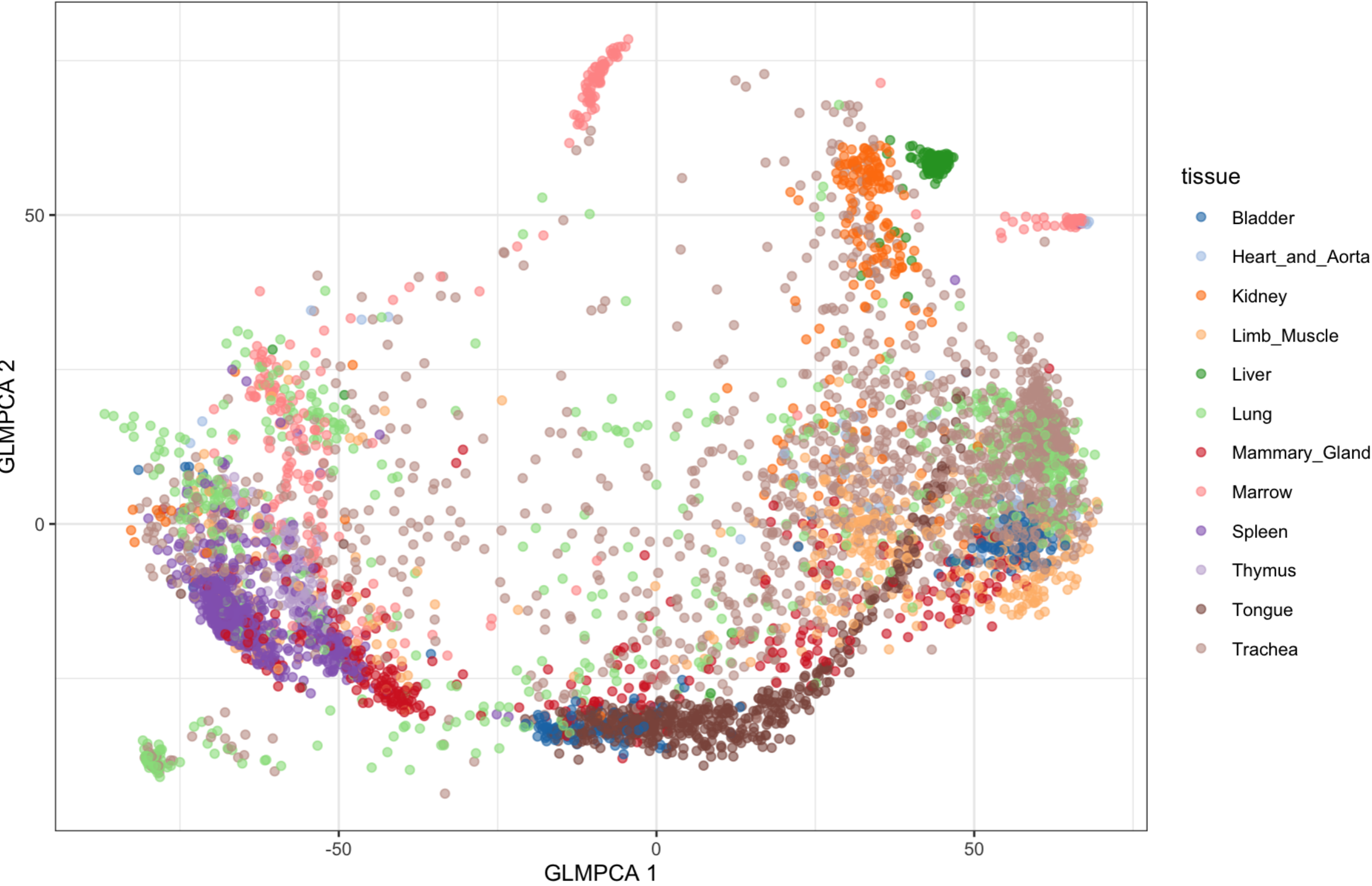
A further generalization allows us to include *observed covariates* in the model. These can be covariates at the cell and gene level and it is useful for normalization and batch effect correction.



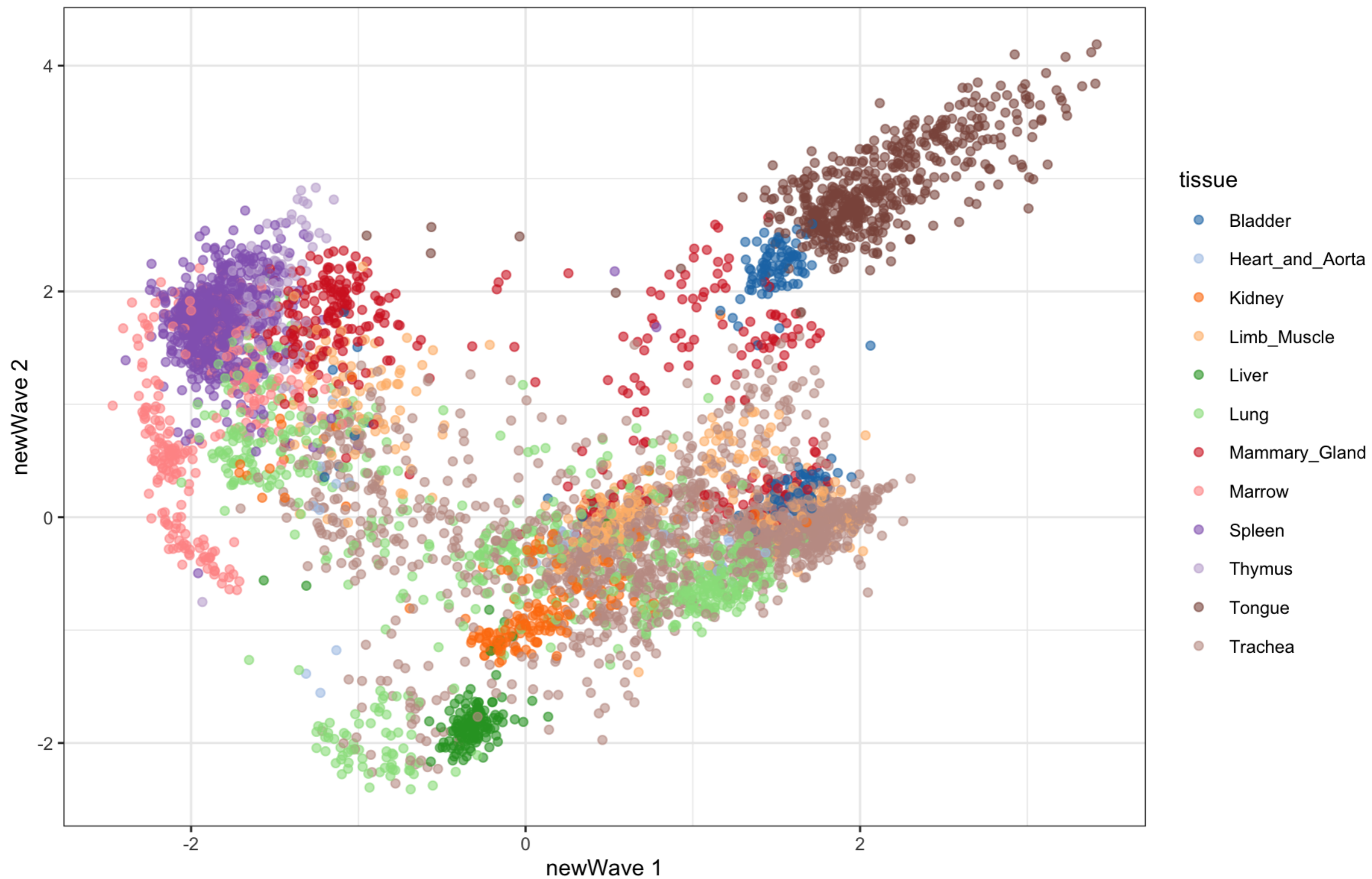
# TABULA MURIS: PCA AFTER SCRAN NORMALIZATION (LOG SCALE)



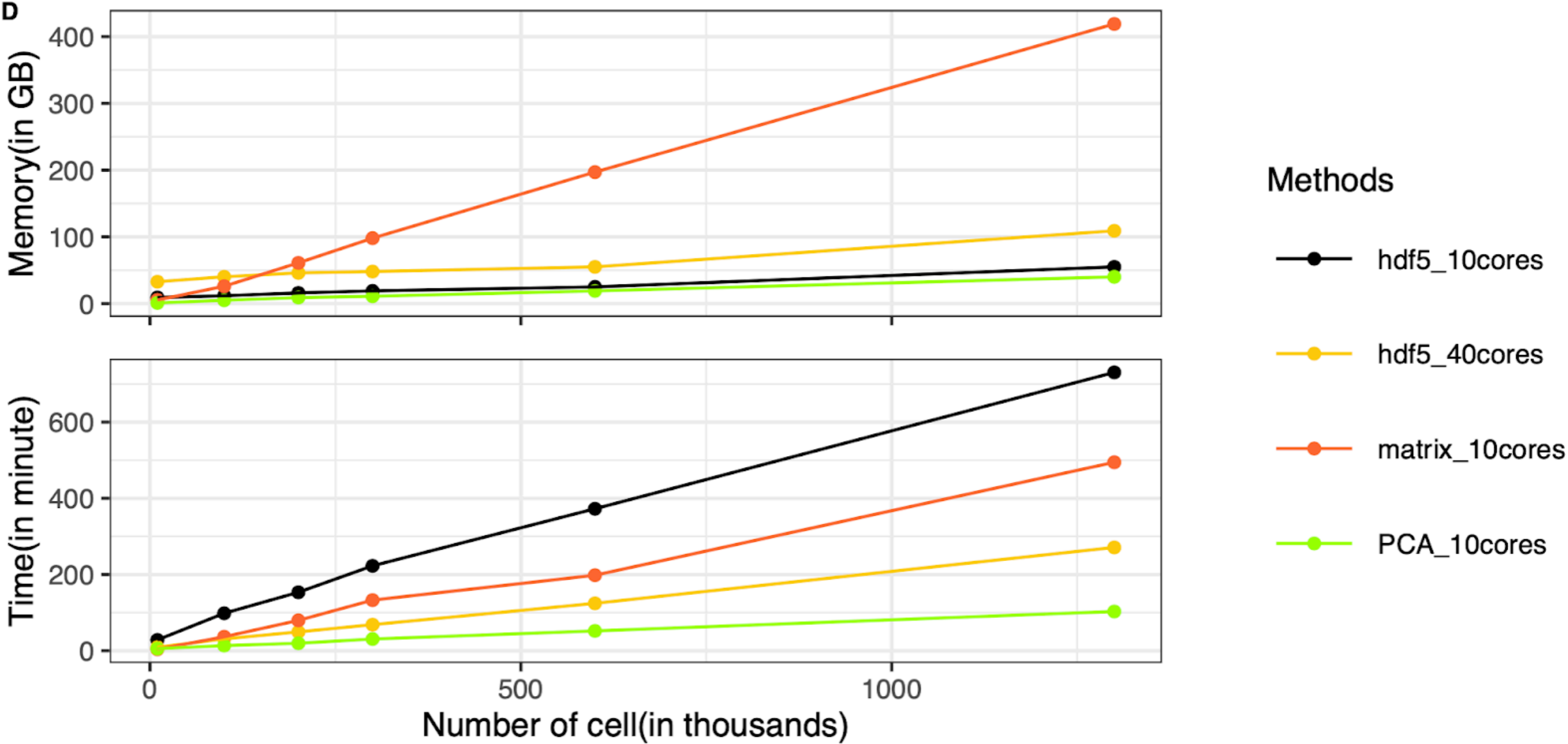
# TABULA MURIS: GLM-PCA (POISSON)



# TABULA MURIS: NEWAVE (NEGATIVE BINOMIAL)



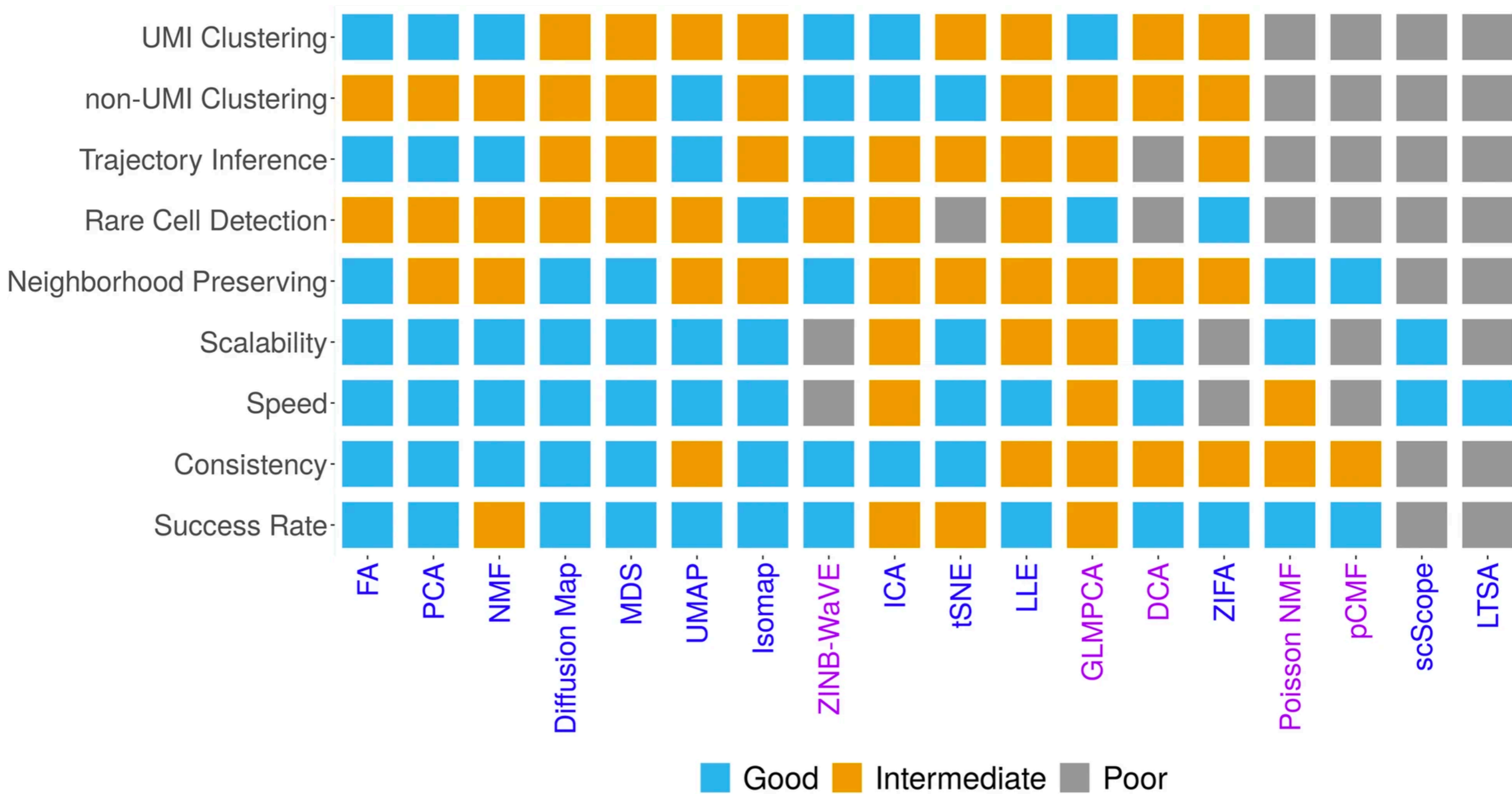
# SCALABILITY



# SCALABILITY

- ▶ Townes et al. (2019) propose an *approximate approach* to speed up computations.
- ▶ Essentially, they compute Pearson or deviance residuals of a *GLM* fit on each gene independently, and then compute PCA of the residuals.
- ▶ A similar approach, *correspondence analysis*, uses chi-squared Pearson residuals + PCA/SVD.
- ▶ These methods are implemented in the [scry](#) and [corral](#) Bioconductor packages, respectively.

# WHICH SHOULD I USE?





# MORE QUESTIONS THAN ANSWERS

- ▶ How many factors should I estimate?
- ▶ Should I include covariates? Which ones?
- ▶ If PCA, should I center/scale?
- ▶ Which data transformations should I use?
- ▶ Which normalization should I use?
- ▶ Why not deep neural networks? (That should take care of it!)
- ▶ Importance of simple models and interpretability of the solutions.

# TAKE-HOME MESSAGE

- ▶ t-SNE / UMAP are fine for visualization
- ▶ Do not use them for inference (e.g., clustering)
- ▶ Linear/more interpretable techniques should be preferred



A photograph of a building belonging to the University of Padua. The building has a brick base and a light-colored upper section. A sign on the wall features the university's crest and the text 'UNIVERSITA' DI PADOVA'. A large white text overlay is at the bottom.

UNIVERSITA'  
DI PADOVA

THANKS FOR YOUR ATTENTION!