

# Single Cell RNA-seq Data Integration

---

Ahmed Mahfouz

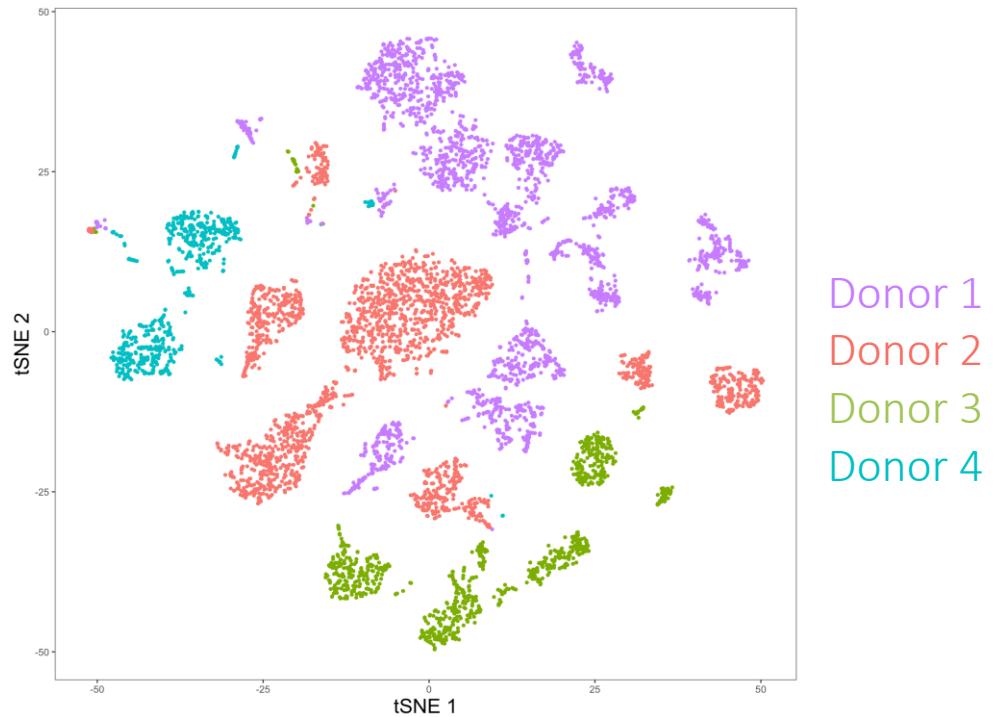
Human Genetics Department, LUMC

Leiden Computational Biology Center, LUMC

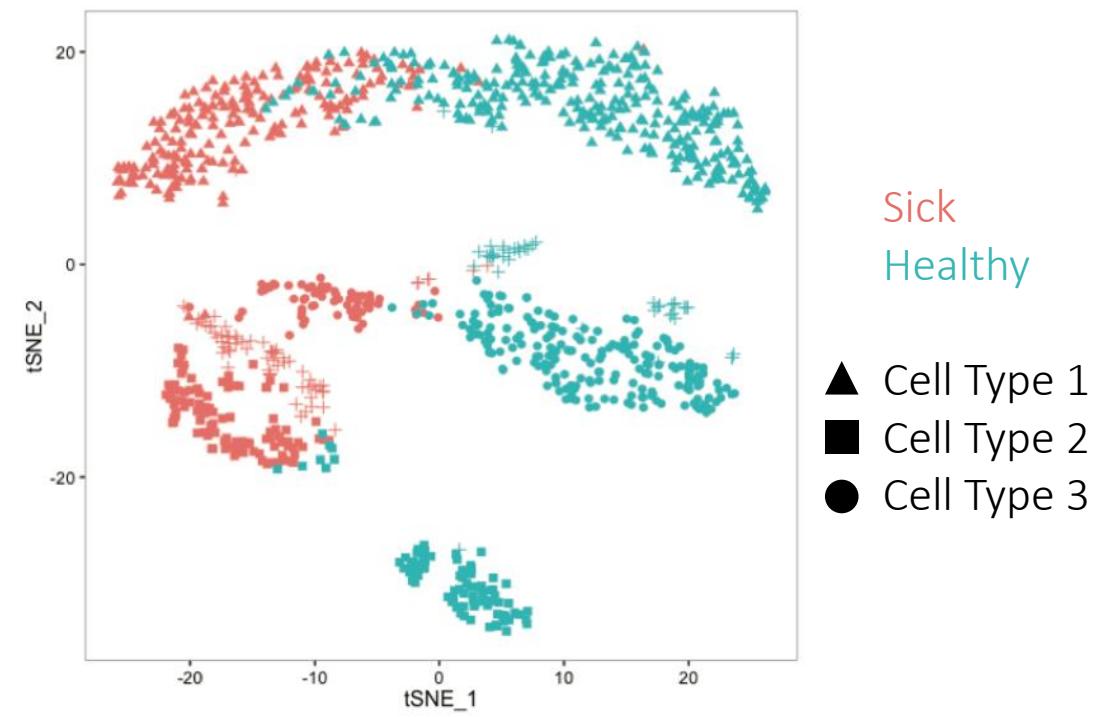
Delft Bioinformatics Lab, TU Delft



# Why integrate?



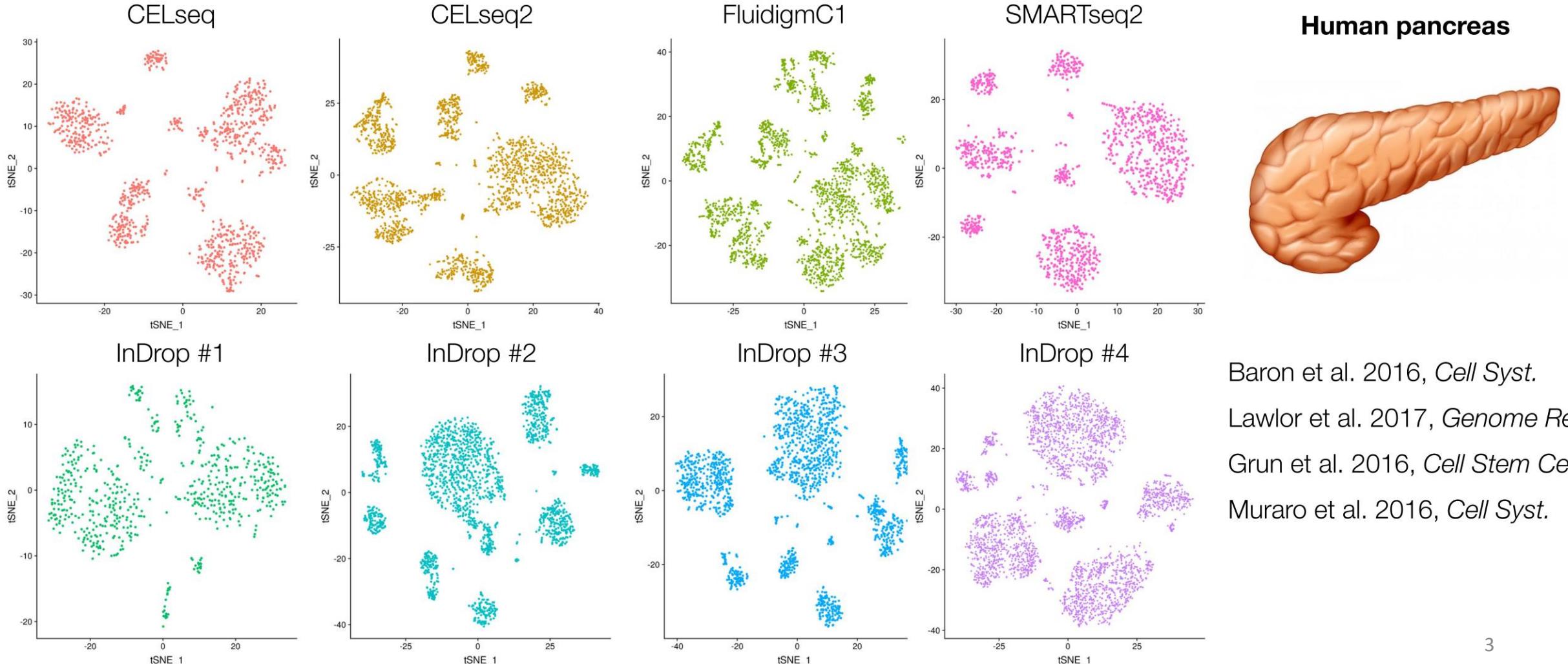
Same tissue from different donors



Cross condition comparisons

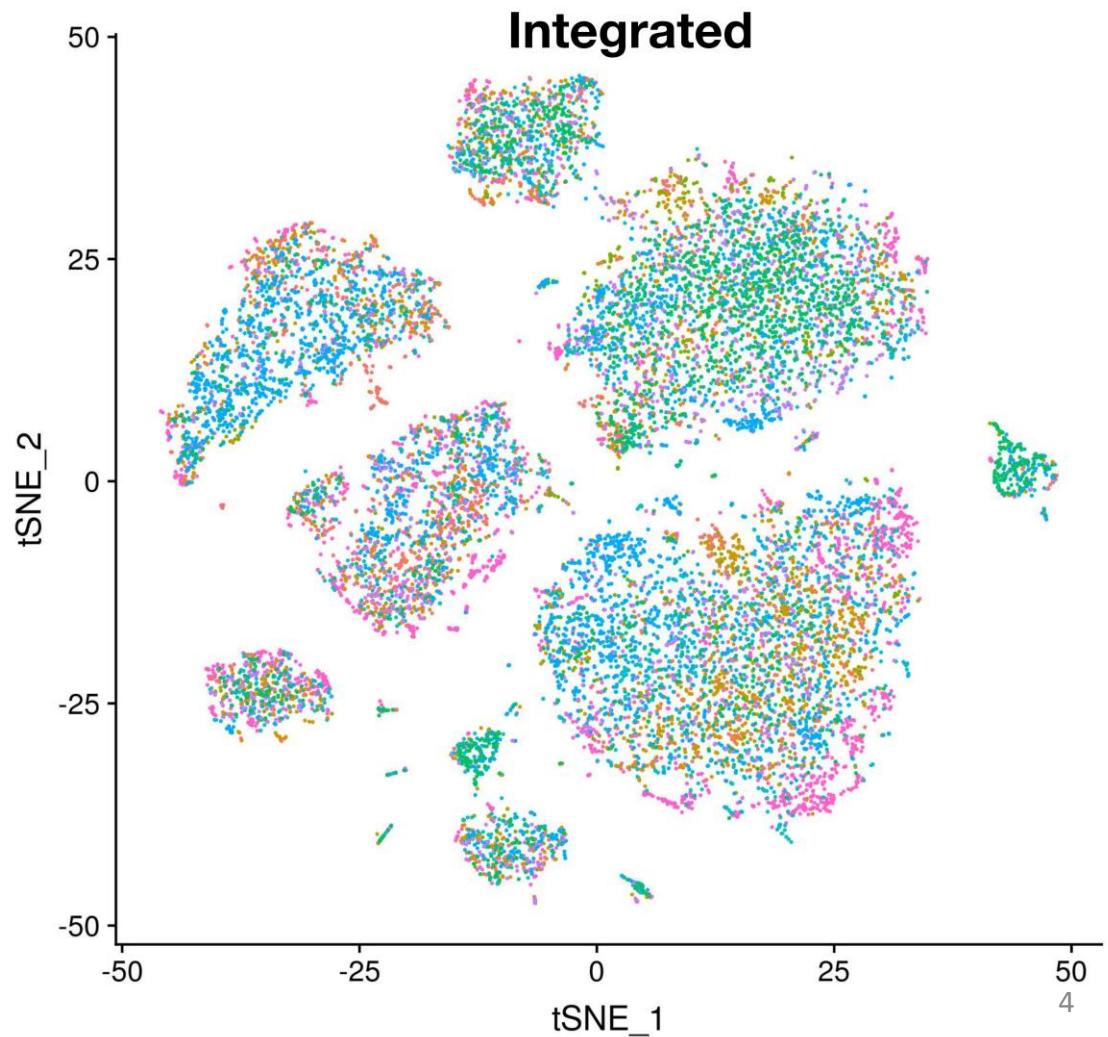
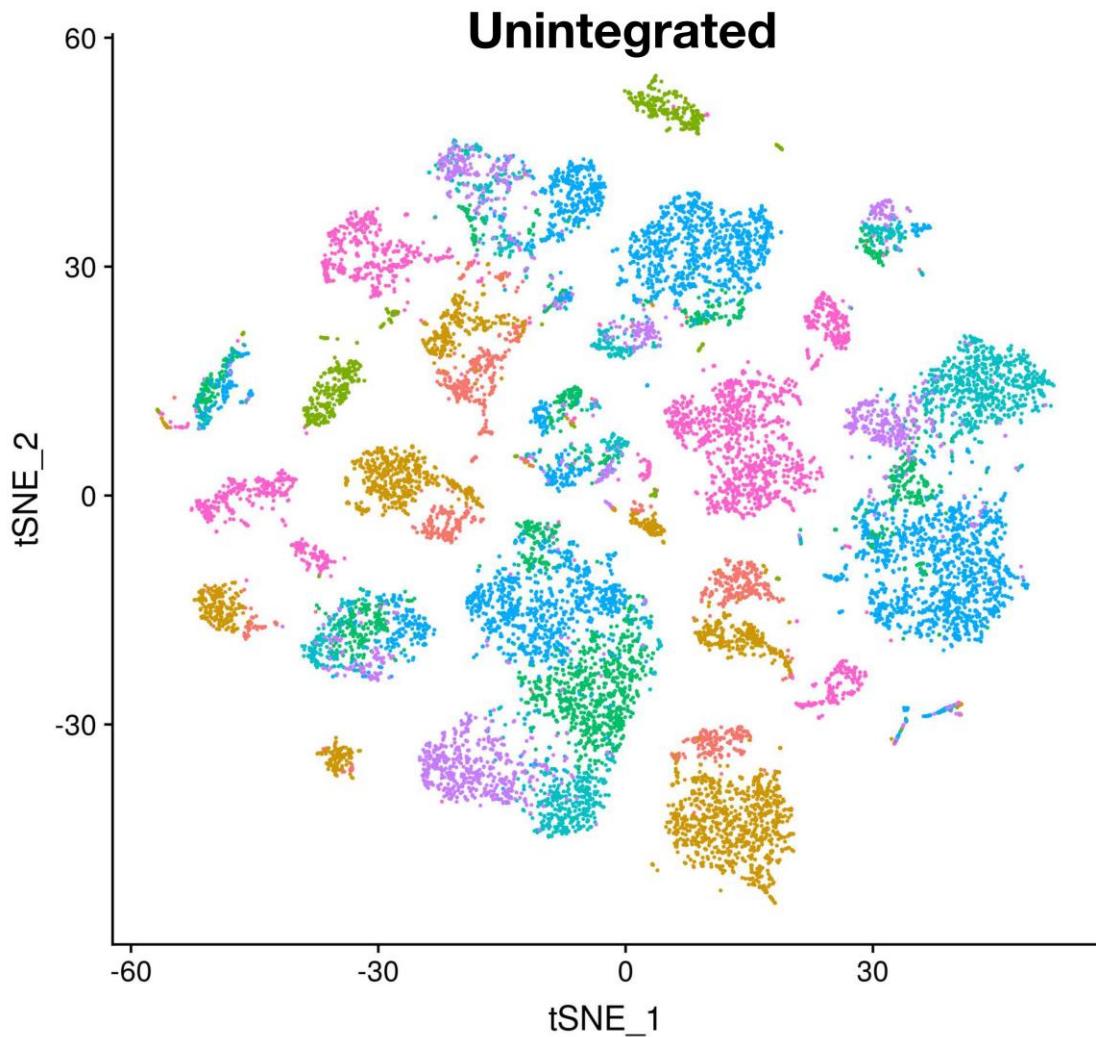
# Building a cell atlas

## 8 maps of the human pancreas



# Building a cell atlas

## 8 maps of the human pancreas

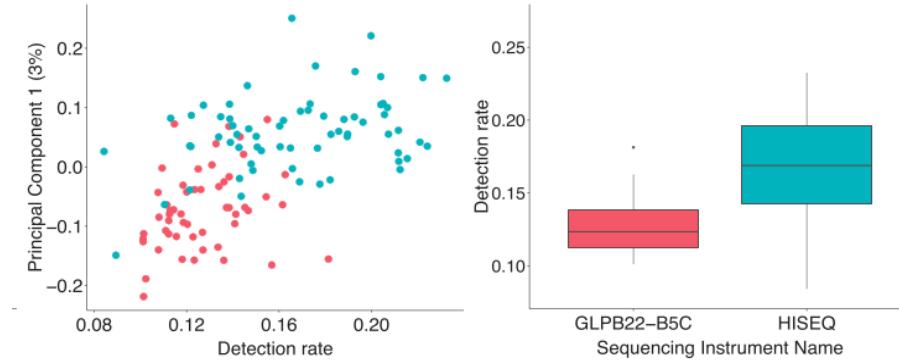


# Confounders and batch effects

## 1. Technical variability

- Changes in sample quality/processing
- Library prep or sequencing technology
- ‘Experimental reality’

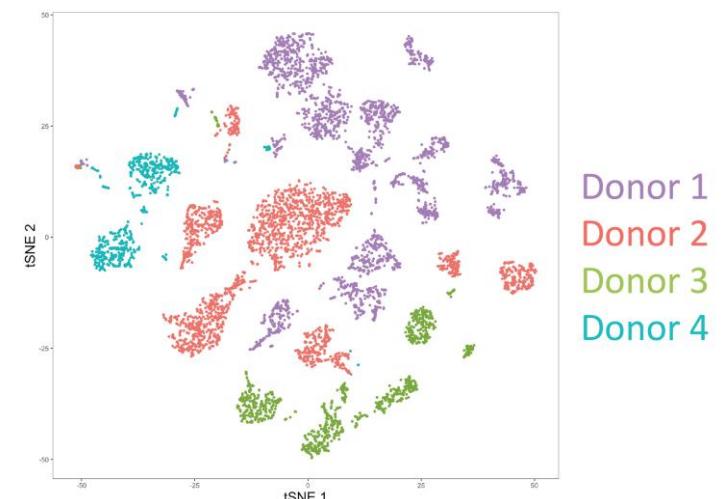
Technical ‘batch effects’ confound downstream analysis



## 2. Biological variability

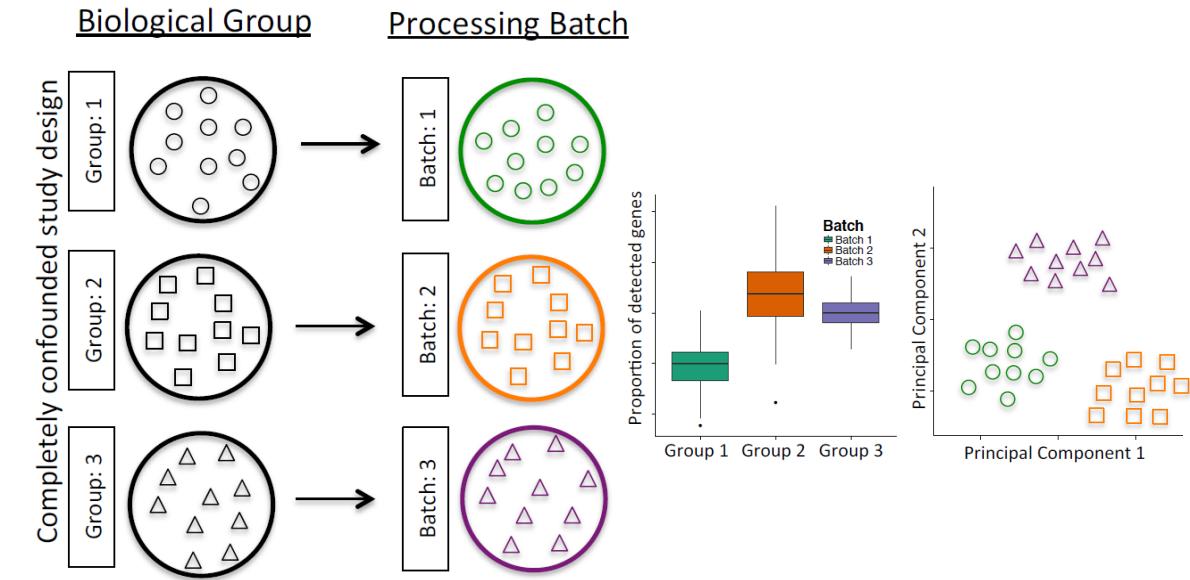
- Patient differences
- Environmental/genetic perturbation
- Evolution! (cross-species analysis)

Biological ‘batch effects’ confound comparisons of scRNA-seq data



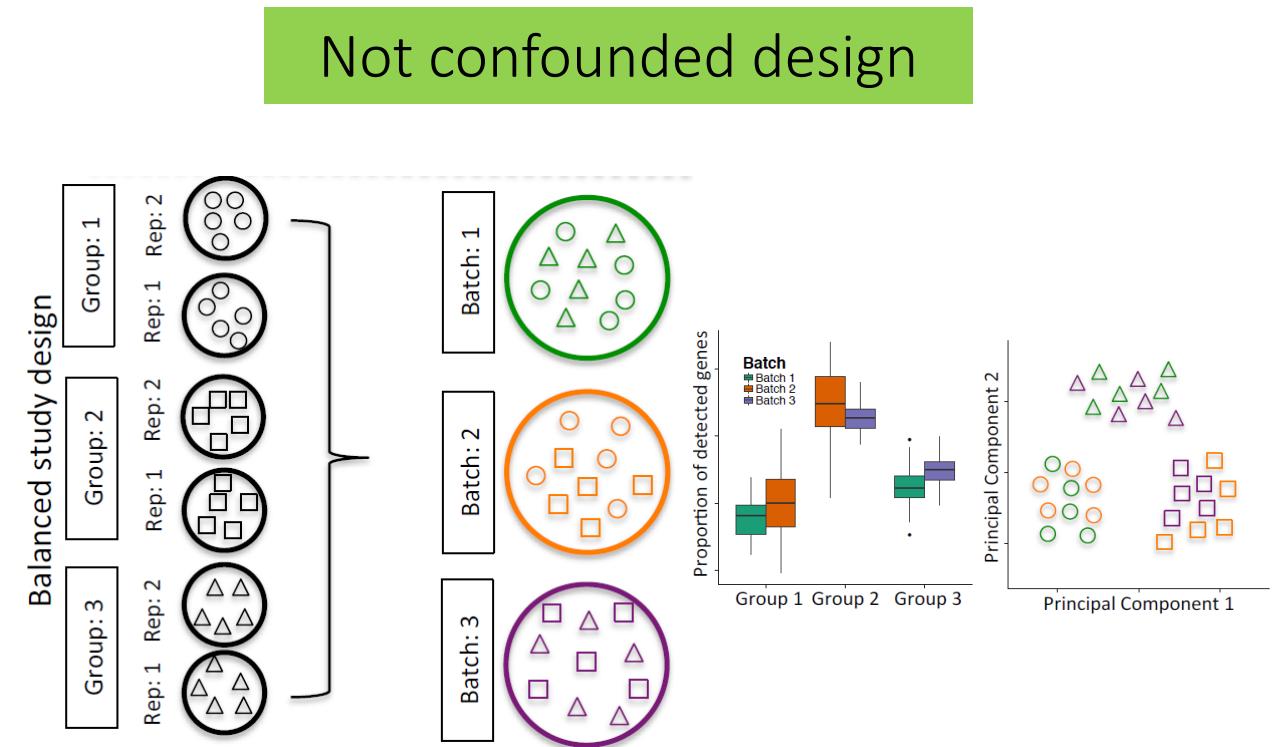
# Confounders and batch effects

## Confounded design



Don't design your experiment like this!!!

## Not confounded design



Good experimental design *does not remove batch effects*, it prevents them from biasing your results.

# Our agenda

- Single cell batch Correction methods
- Performance assessment

# Batch correction methods

- Many good options have been developed for bulk RNA-seq data:
  - RUVseq() or svaseq()
  - Linear models with e.g. removeBatchEffect() in limma or scater
  - ComBat() in sva
  - ...
- But bulk RNA-seq methods make modelling assumptions that are likely to be violated in scRNAseq data
  - The composition of cell populations are either known or the same across batches
  - Batch effect is additive: batch-induced fold-change in expression is the same across different cell subpopulations for any given gene

# Batch correction methods

- MNNcorrect (<https://doi.org/10.1038/nbt.4091>)
- CCA + anchors (Seurat v3) (<https://doi.org/10.1101/460147>)
- CCA + dynamic time warping (Seurat v2) (<https://doi.org/10.1038/nbt.4096>)
- LIGER (<https://doi.org/10.1101/459891>)
- Harmony (<https://doi.org/10.1101/461954>)
- Conos (<https://doi.org/10.1101/460246>)
- Scanorama (<https://doi.org/10.1101/371179>)
- scMerge (<https://doi.org/10.1073/pnas.1820006116>)
- ...

**Two broad strategies:**

- Joint dimension reduction
- Graph-based joint clustering

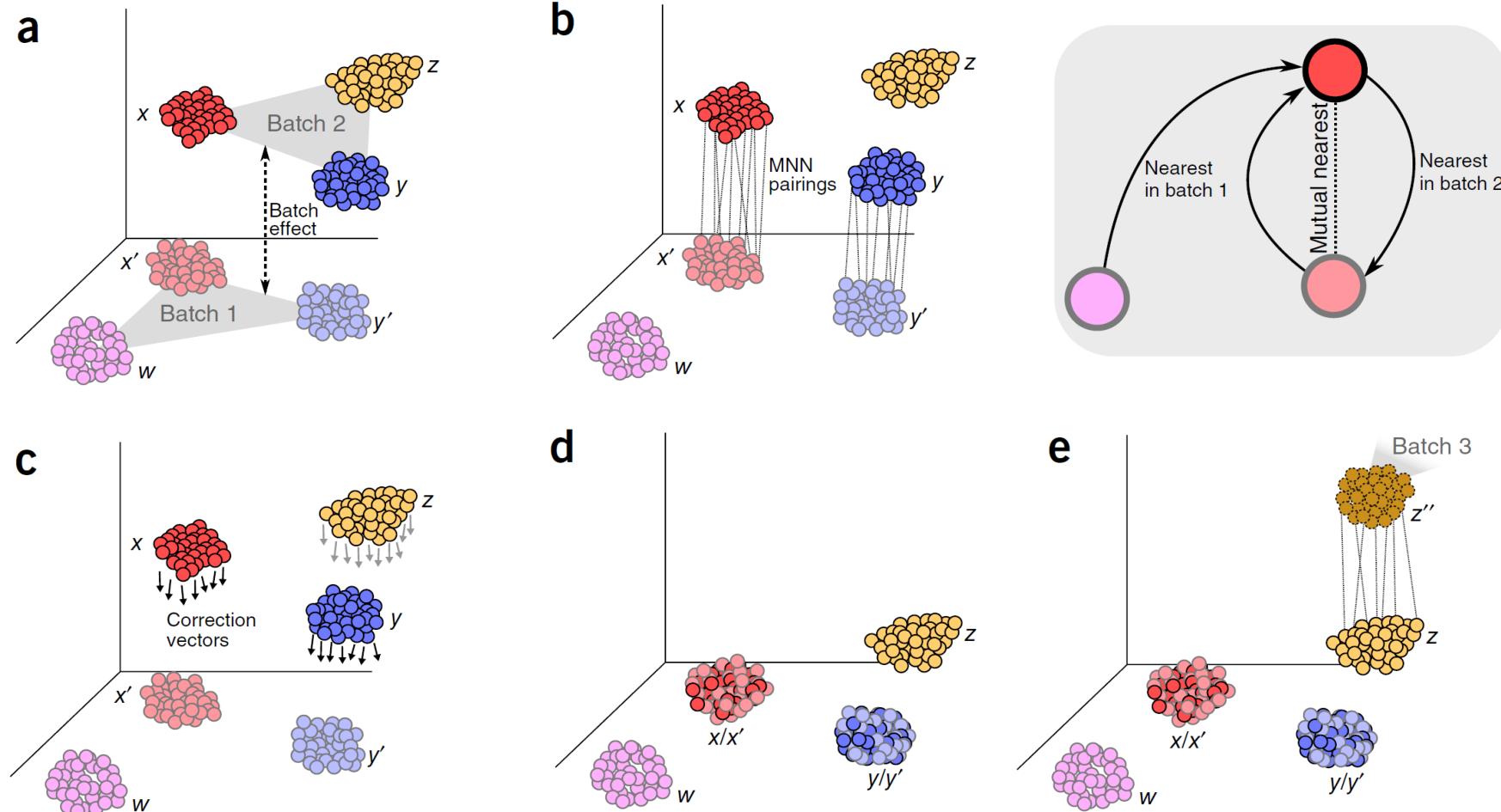
# Batch correction methods

- MNcorrect (<https://doi.org/10.1038/nbt.4091>)
- CCA + anchors (Seurat v3) (<https://doi.org/10.1101/460147>)
- CCA + dynamic time warping (Seurat v2) (<https://doi.org/10.1038/nbt.4096>)
- LIGER (<https://doi.org/10.1101/459891>)
- Harmony (<https://doi.org/10.1101/461954>)
- Conos (<https://doi.org/10.1101/460246>)
- Scanorama (<https://doi.org/10.1101/371179>)
- scMerge (<https://doi.org/10.1073/pnas.1820006116>)
- ...

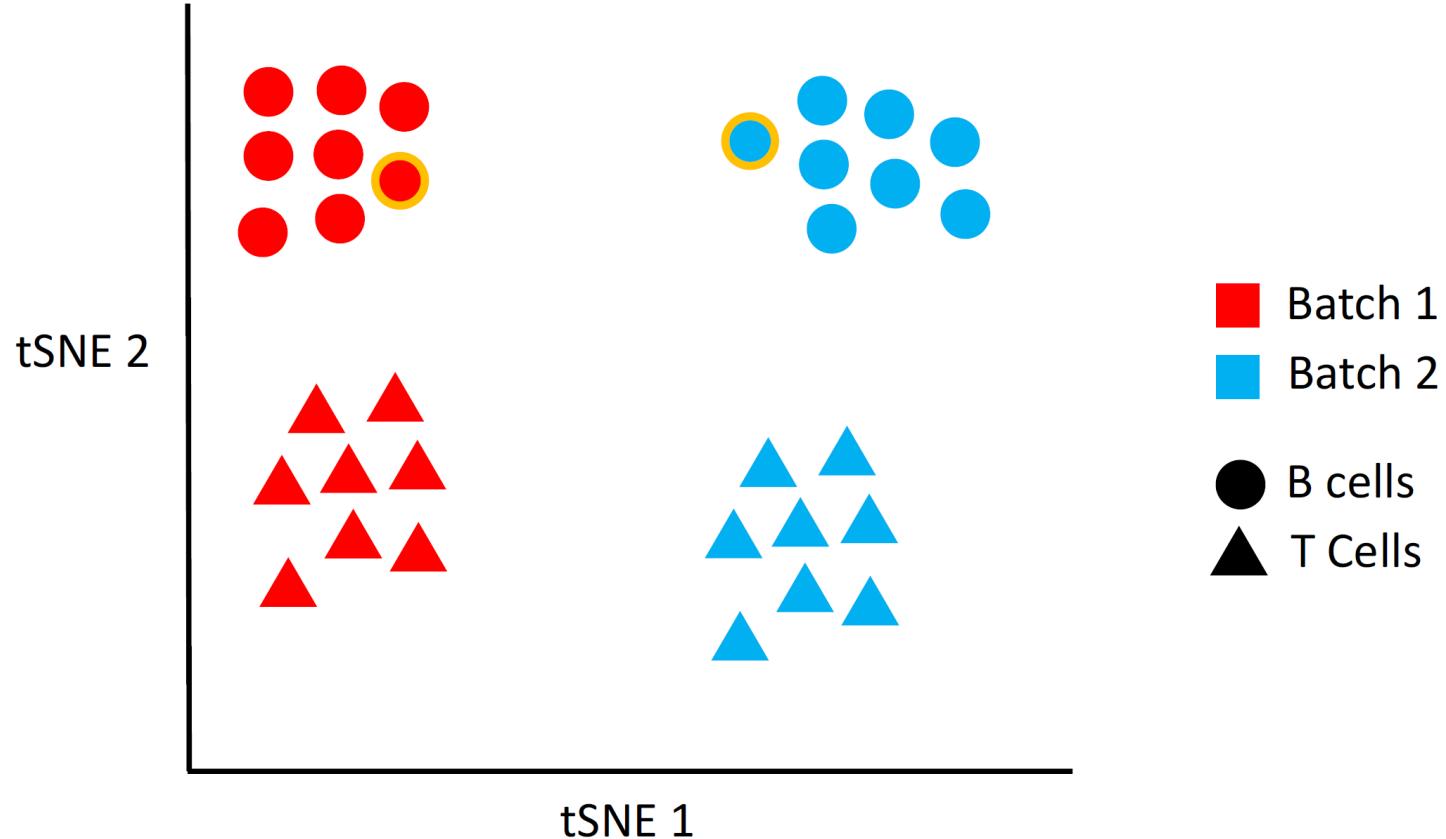
Two broad strategies:

- Joint dimension reduction
- Graph-based joint clustering

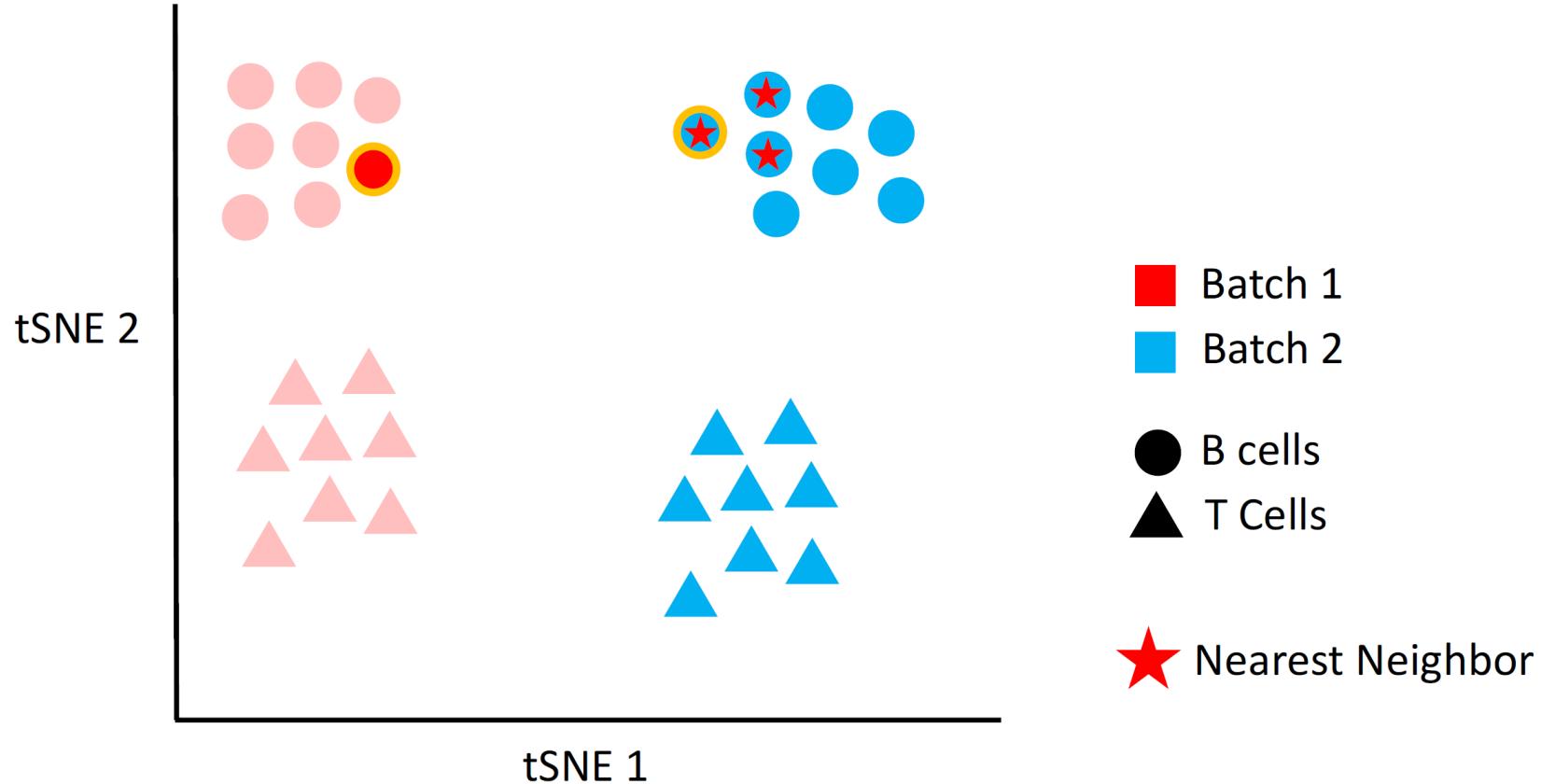
# Mutual Nearest Neighbors (MNN)



# Mutual Nearest Neighbors (MNN)



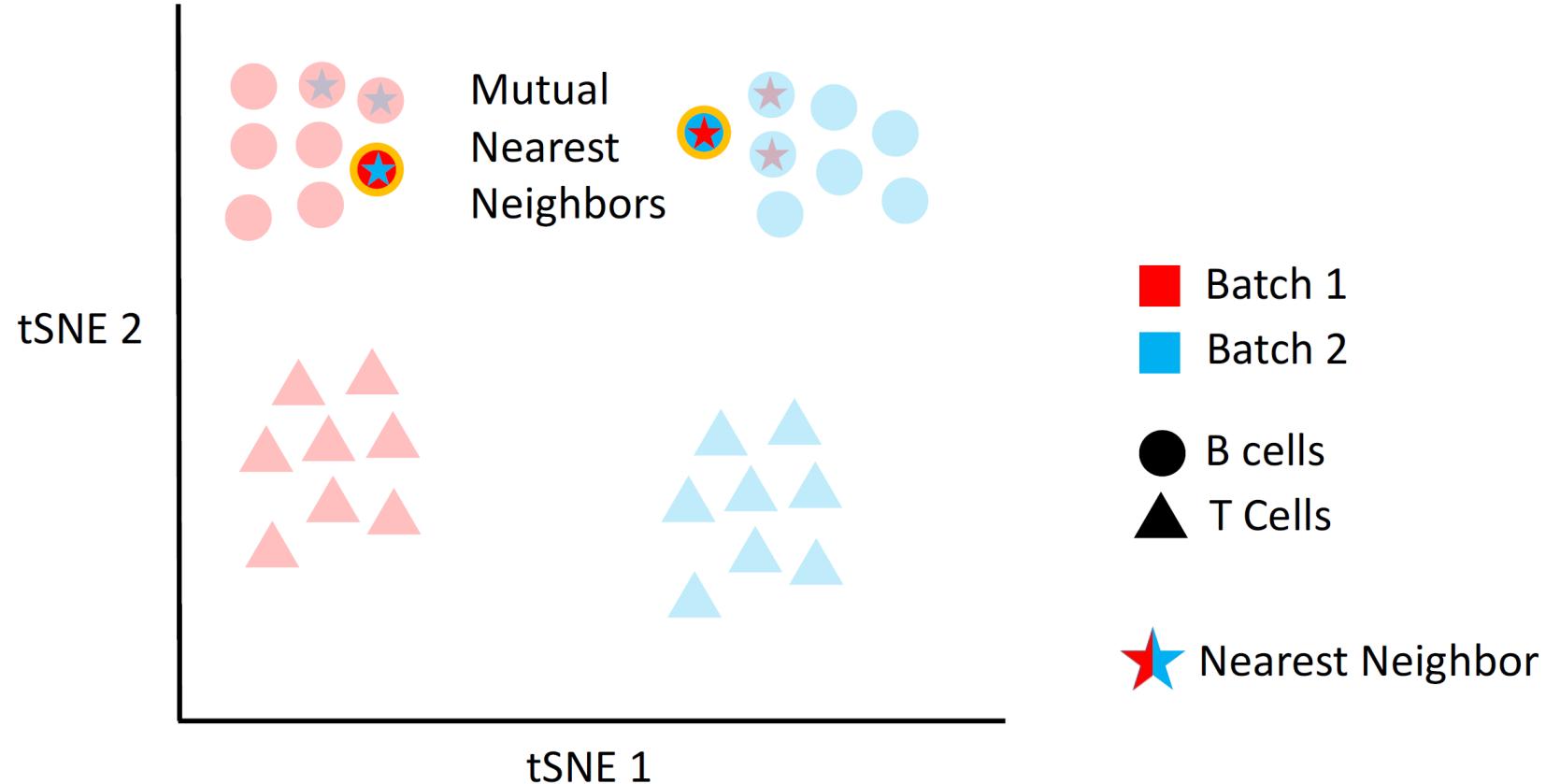
# Mutual Nearest Neighbors (MNN)



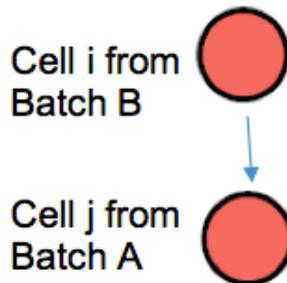
# Mutual Nearest Neighbors (MNN)



# Mutual Nearest Neighbors (MNN)

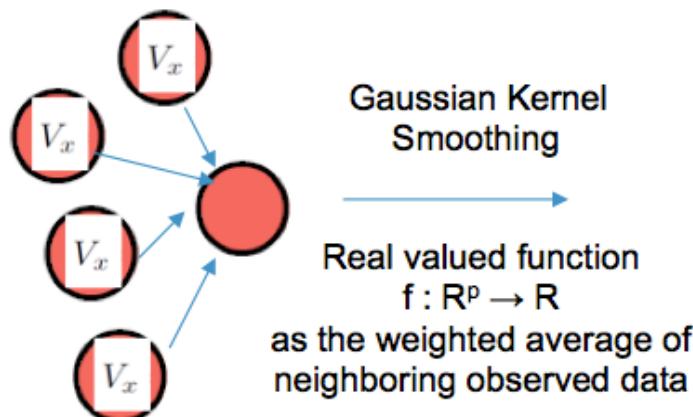


# Mutual Nearest Neighbors (MNN)



1) For each MNN pair, a pair-specific batch-correction vector is computed as the vector difference between the expression profiles of the paired cells.

2) A cell-specific batch-correction vector is then calculated as a weighted average of these pair-specific vectors, as computed with a Gaussian kernel.

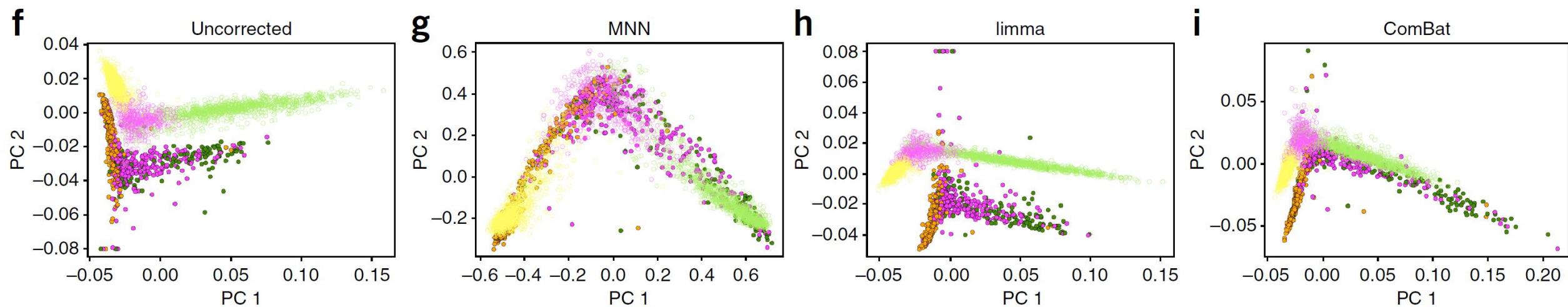


$$V_x = \begin{pmatrix} gene1_a - gene1_b \\ gene2_a - gene2_b \\ gene3_a - gene3_b \\ \dots \\ geneN_a - geneN_b \end{pmatrix}$$

Batch Correction vector for each cell



# Mutual Nearest Neighbors (MNN)



SMART-seq2

- MEP
- GMP
- CMP

MARS-seq

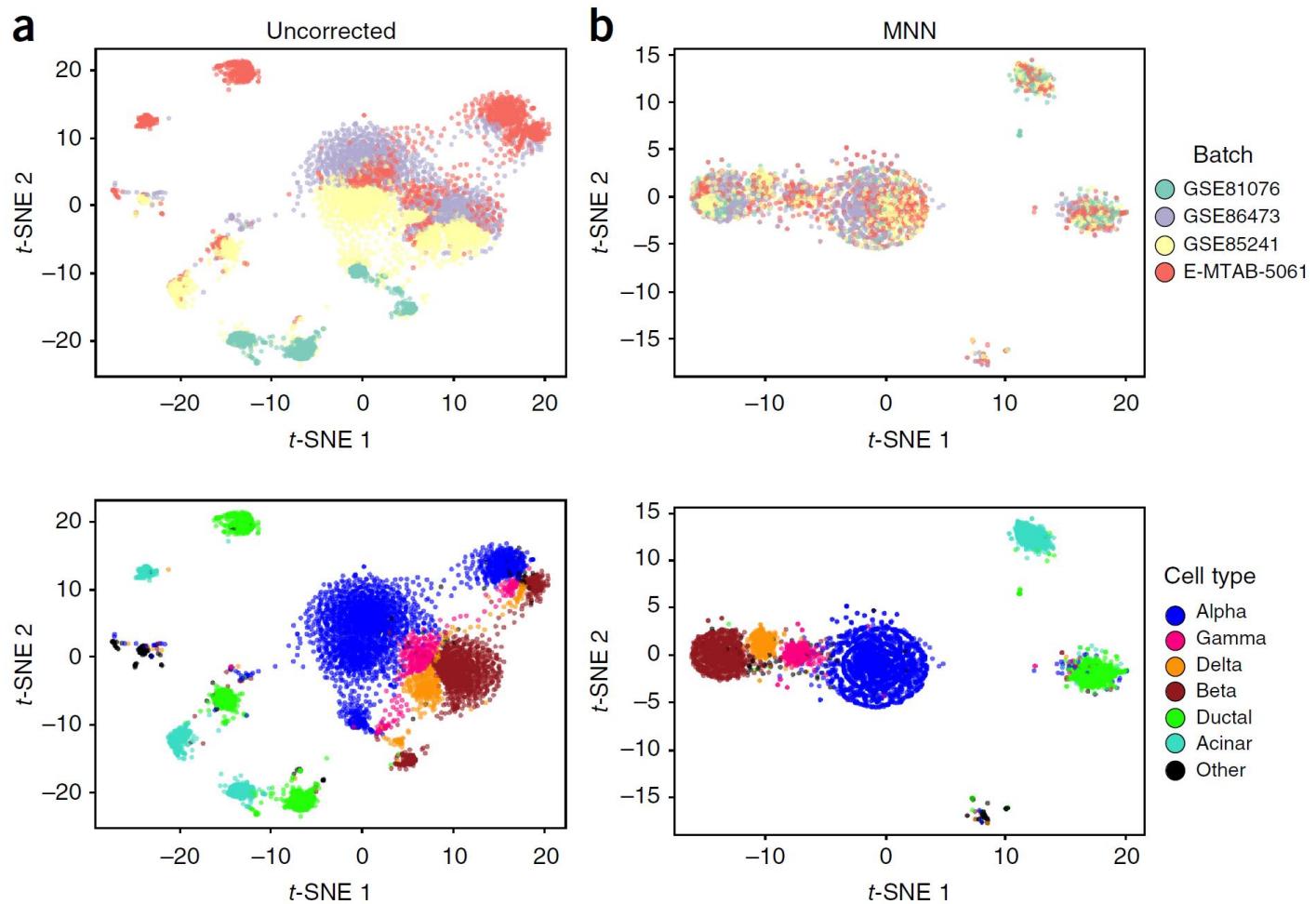
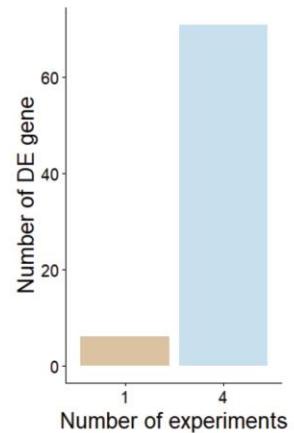
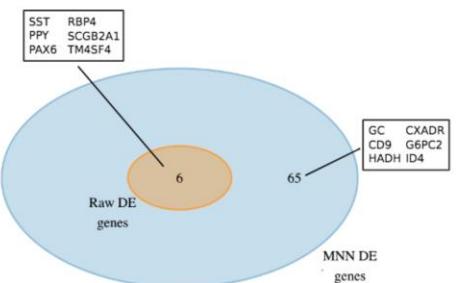
- MEP
- GMP
- CMP

MEPs: megakaryocyte–erythrocyte progenitors  
GMPs: granulocyte–monocyte progenitors  
CMPs: common myeloid progenitors

# Mutual Nearest Neighbors (MNN)

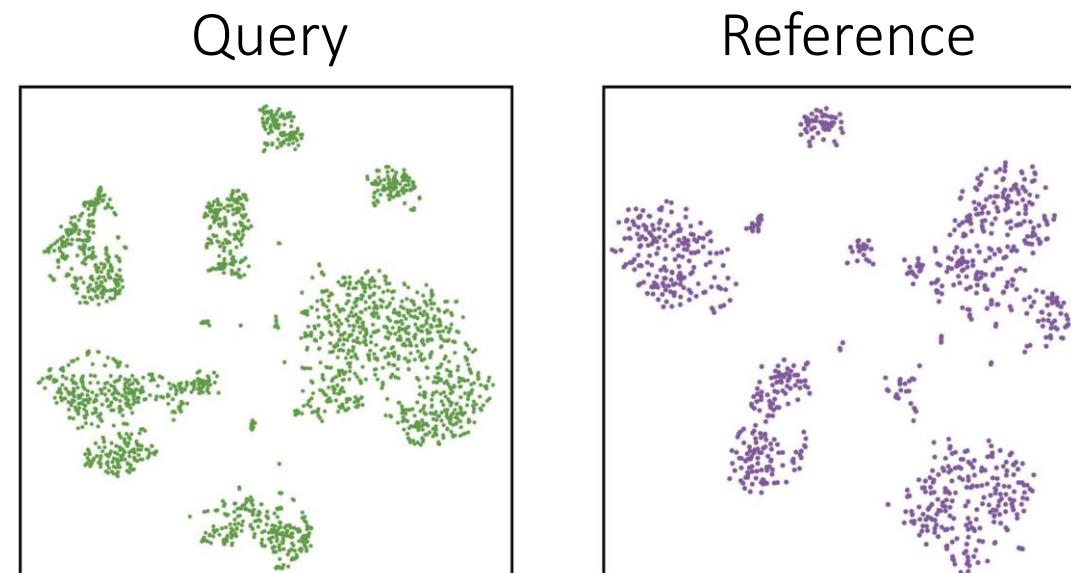
- Pooling experiments -> increased statistical power

Delta vs Gamma Islet Cells

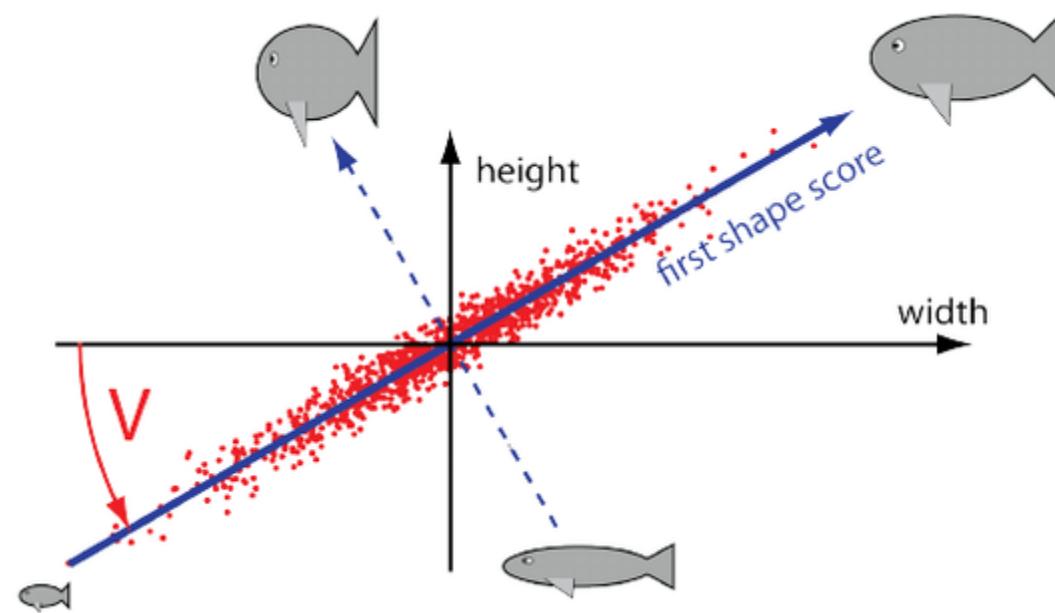
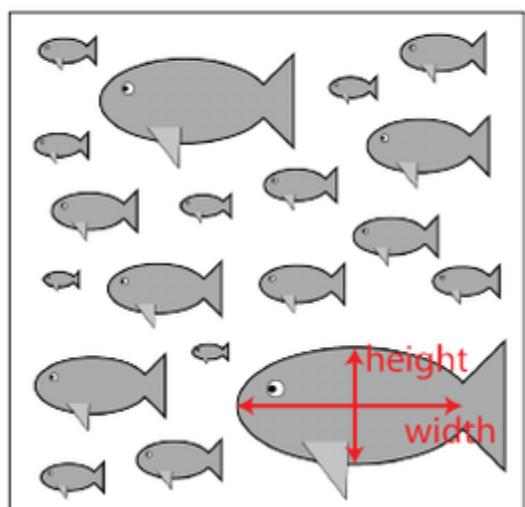


# CCA + anchors (Seurat v3)

1. Find corresponding cells across datasets
2. Compute a data adjustment based on correspondences between cells
3. Apply the adjustment

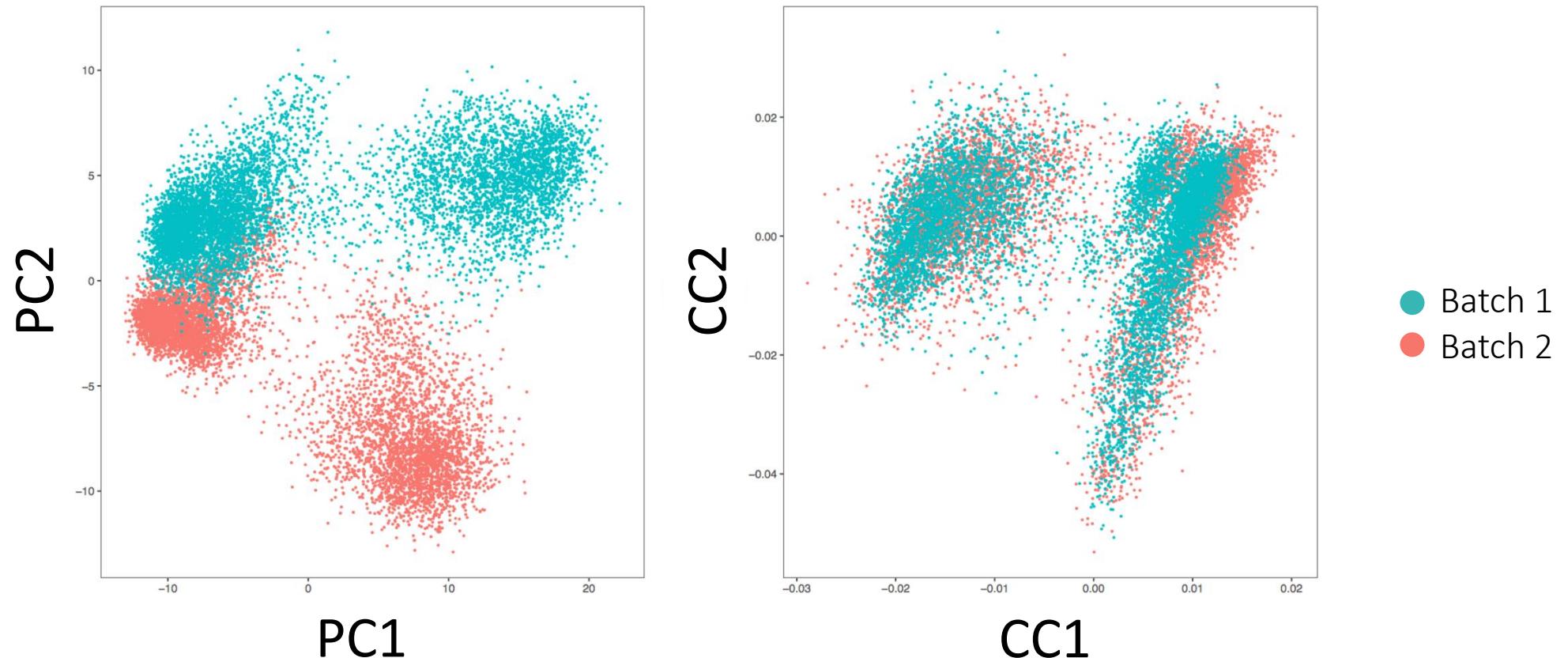


# Principle component analysis



# Finding corresponding cells

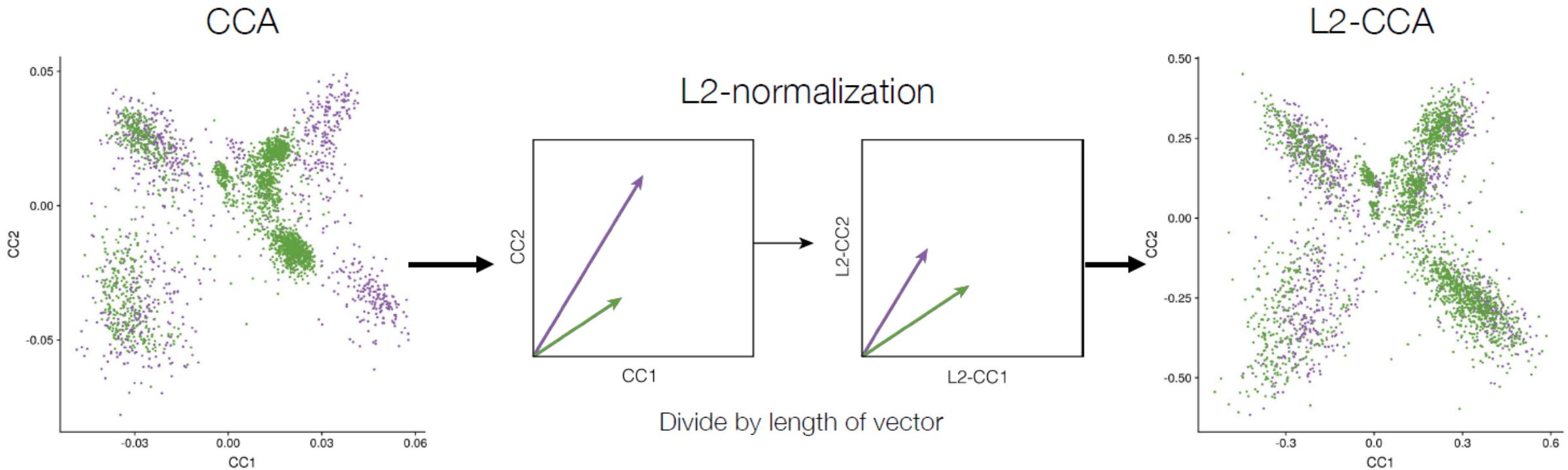
## Canonical correlation analysis and normalization



CCA captures correlated sources of variation between two datasets

# Finding corresponding cells

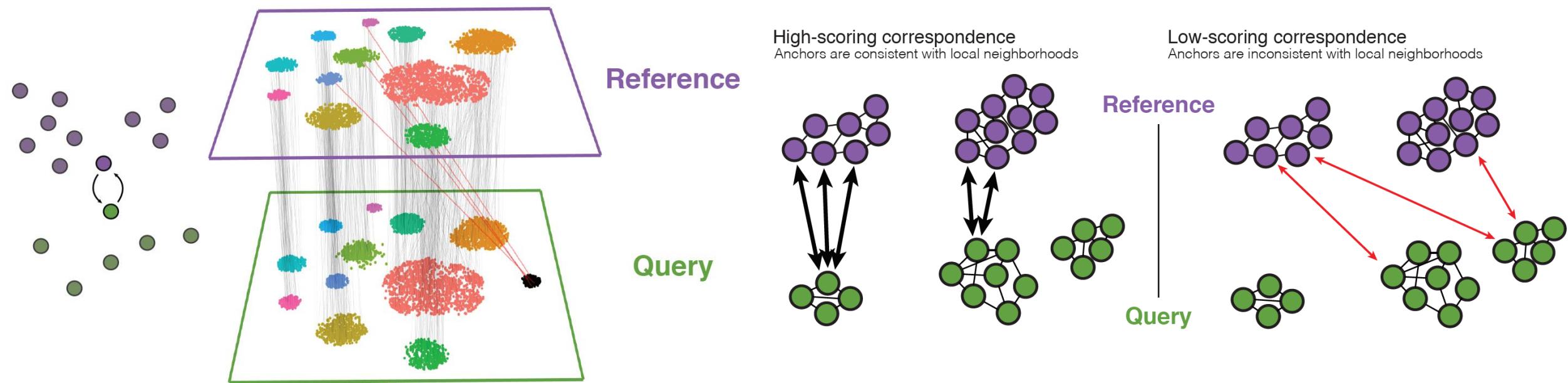
Canonical correlation analysis and normalization



L2-normalization corrects for differences in scale

# Finding corresponding cells

Anchors: mutual nearest neighbors



# Finding corresponding cells

## Data integration

1. Calculate the matrix  $B$ , where each column represents the difference between the two expression vectors for every pair of anchor cells  $a$
2. Construct a weight matrix  $W$  that defines the strength of association between each query cell  $c$ , and each anchor  $i$
3. Calculate a transformation matrix  $C$  using the previously computed weights matrix and the integration matrix as
4. Subtract the transformation matrix  $C$  from the original expression matrix  $Y$  to produce the integrated expression matrix  $\hat{Y}$

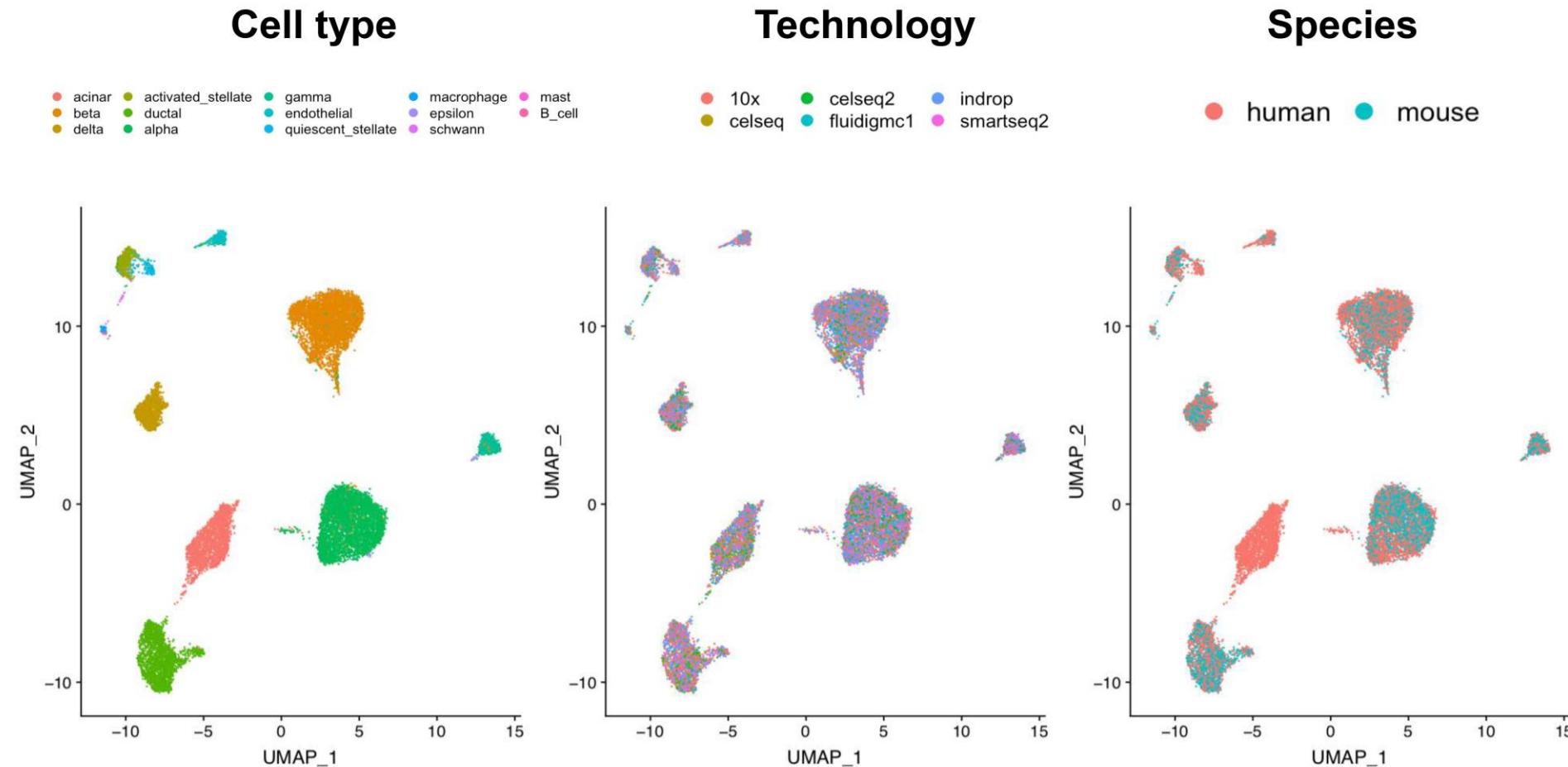
$$B = X[, a] - Y[, a]$$

$$W_{c,i} = \frac{\tilde{D}_{c,i}}{\sum_1^{j=k.weight} \tilde{D}_{c,j}}$$

$$C = BW^T$$

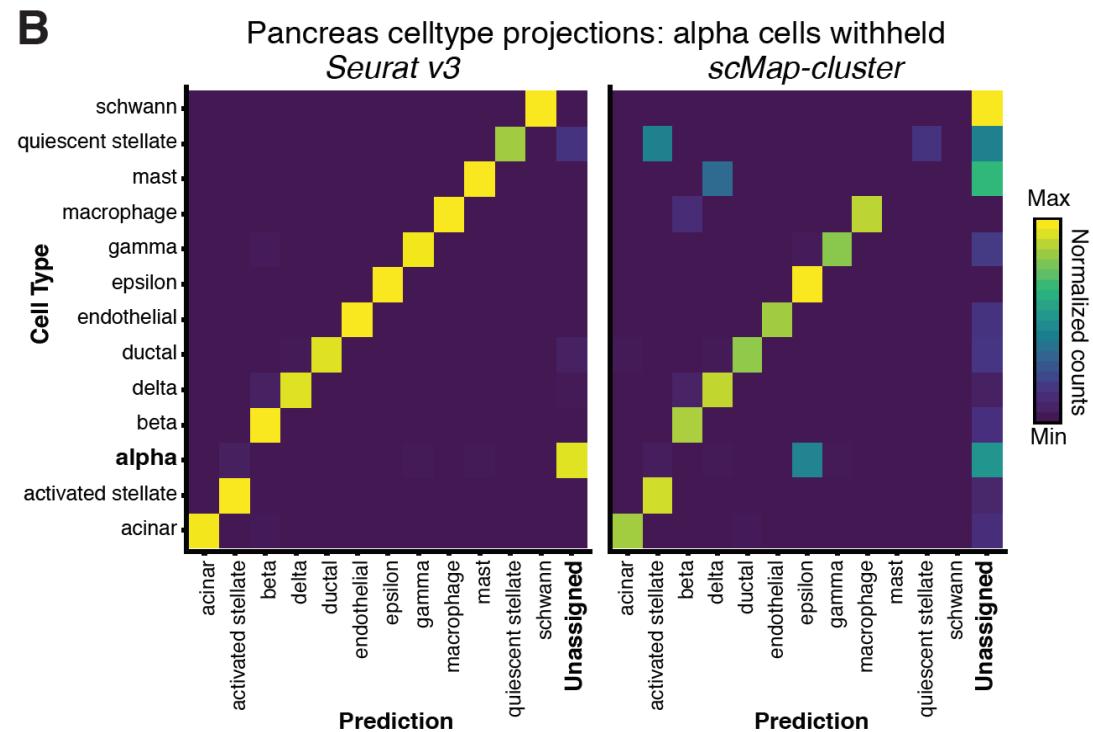
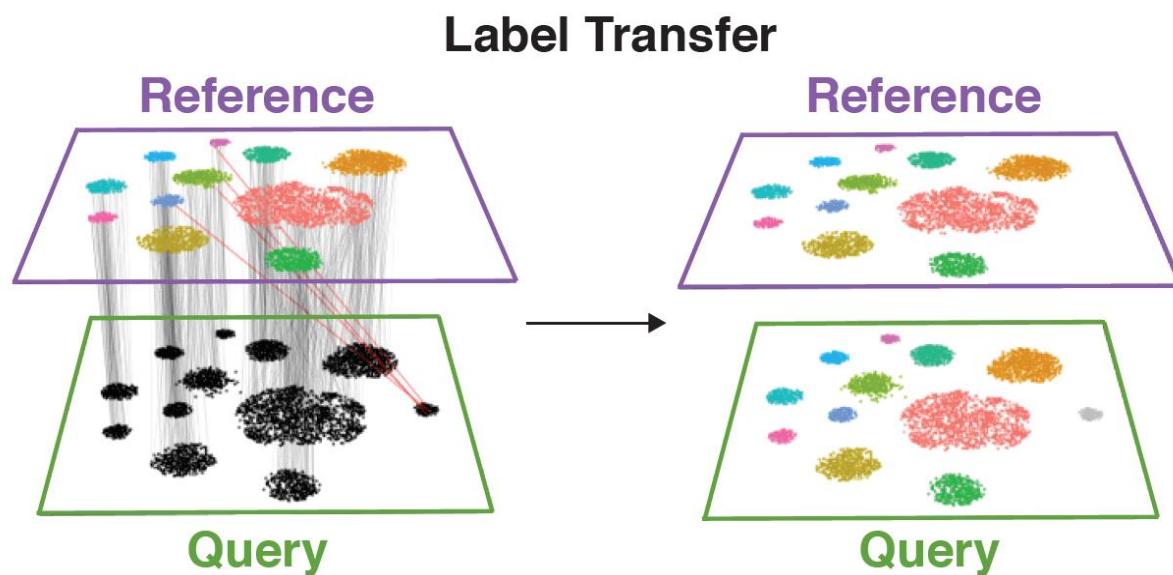
$$\hat{Y} = Y - C$$

# CCA + anchors (Seurat v3)



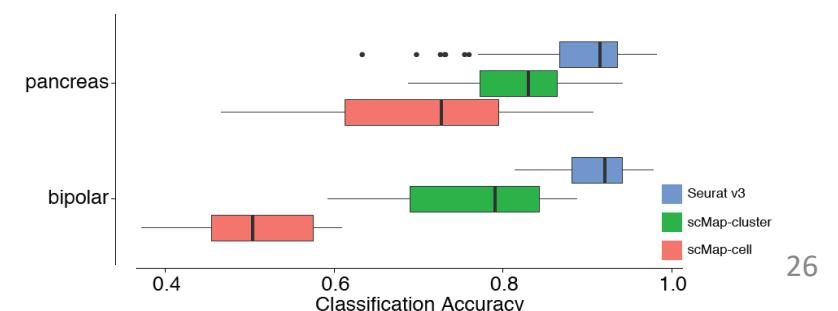
Retinal bipolar datasets: 51K cells, 6 technologies, 2 Species

# Label transfer (classification)

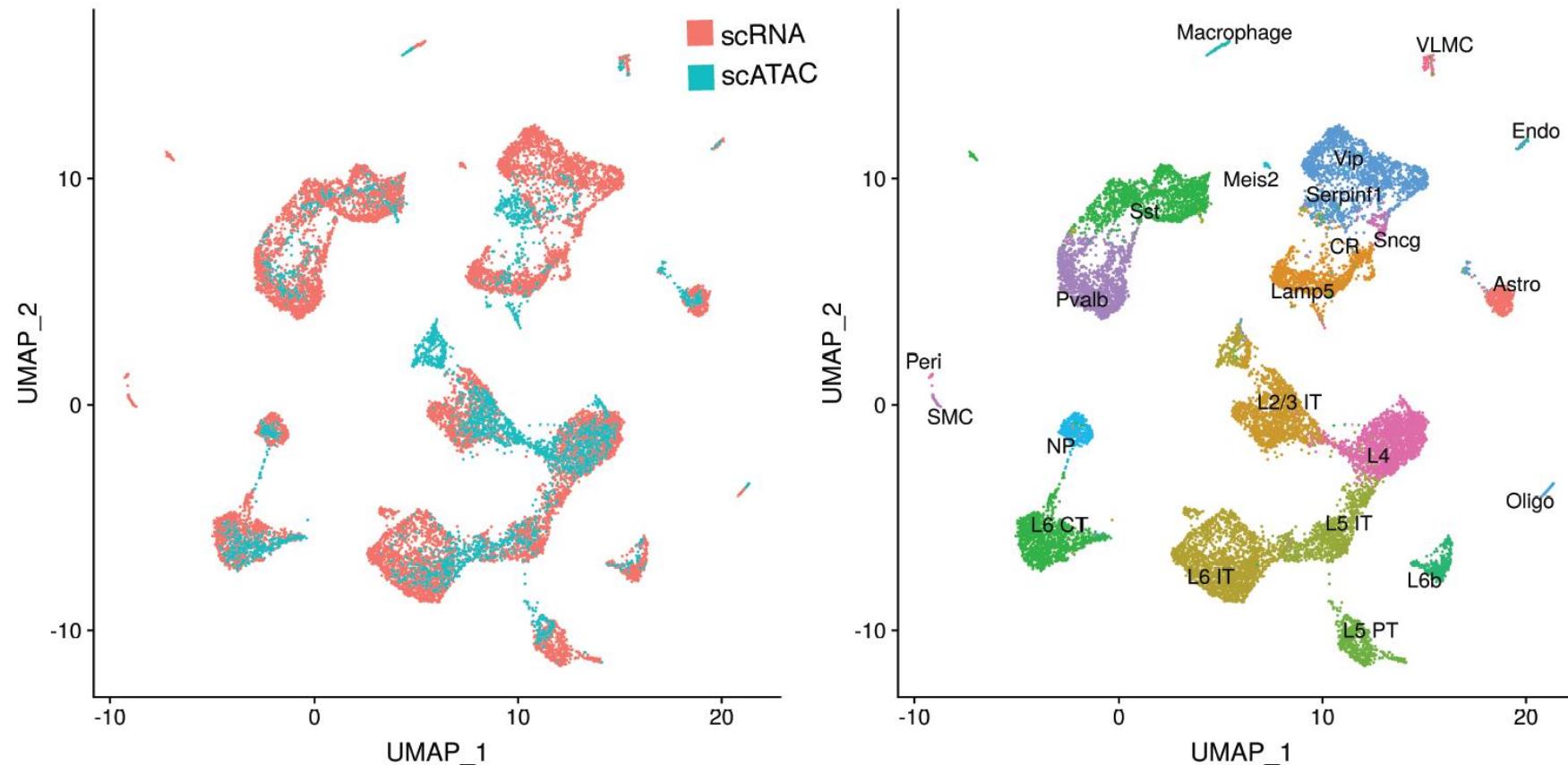


## Weighted vote classifier

What is the classification of each cells nearest anchors?



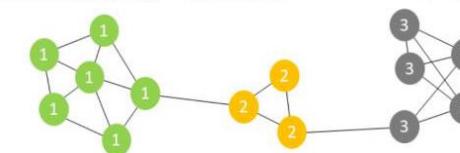
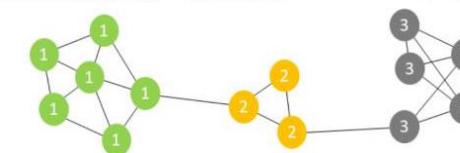
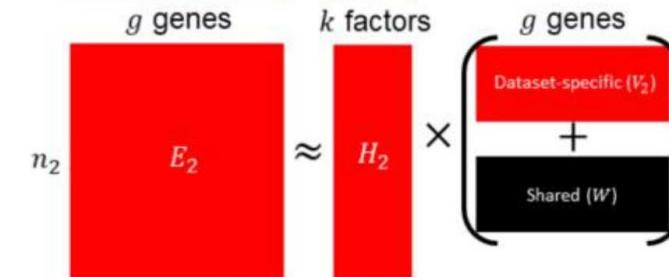
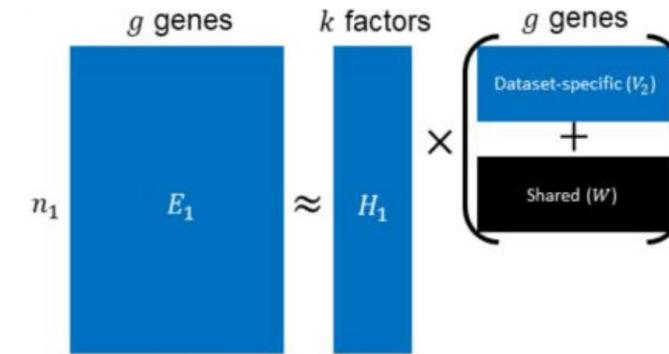
# Integration across modalities



# LIGER

## Linked Inference of Genomic Experimental Relationships

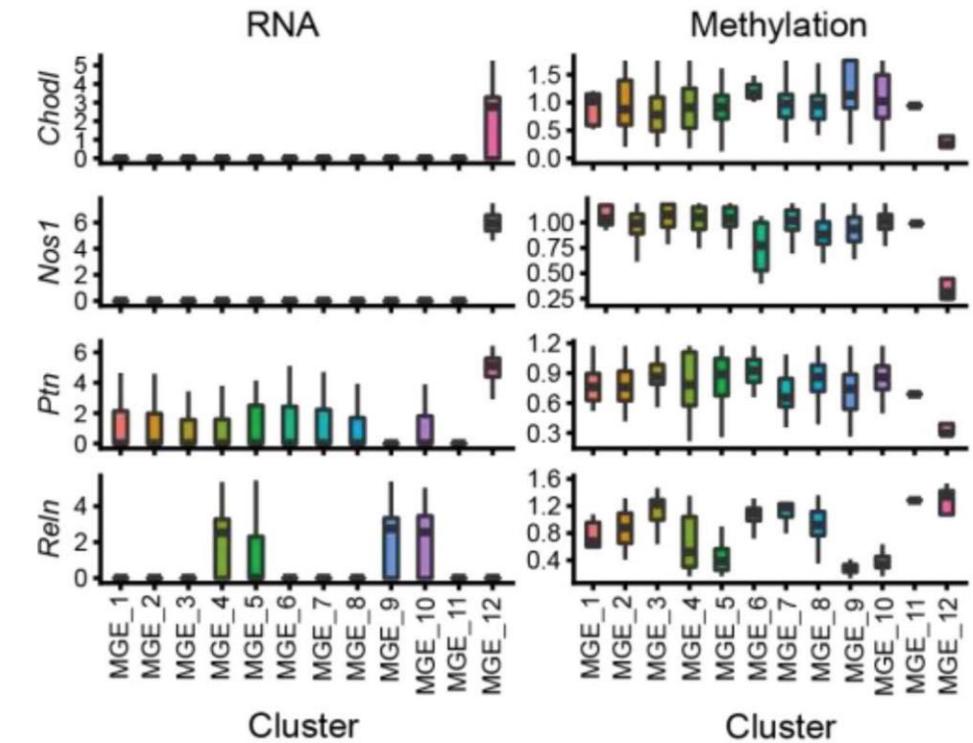
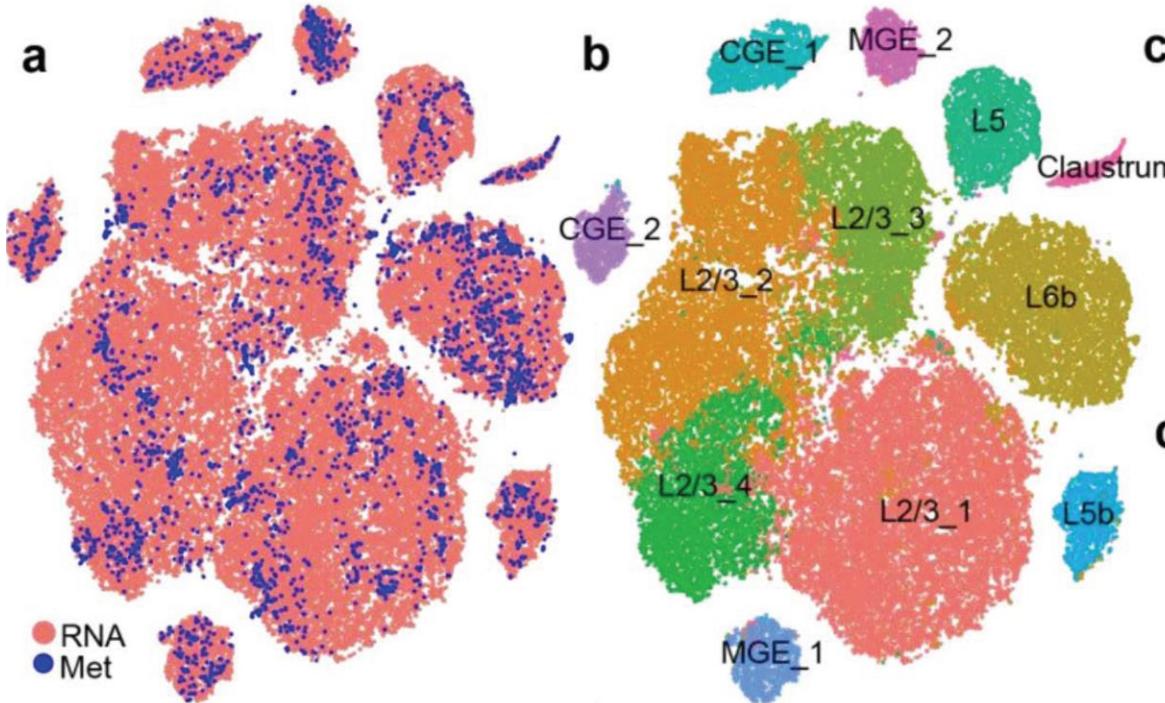
- 1) Integrative non-negative matrix factorization (iNMF) to learn a shared low-dimensional space
- 2) Perform joint clustering on the shared factor neighborhood graph
  - Factors are interpretable due to non-negative constraint
  - Finds set of dataset-specific factors and a set of shared factors



# LIGER

## Linked Inference of Genomic Experimental Relationships

- Joint clustering of gene expression and DNA methylation data



# Using the corrected values

- Batch correction facilitates cell-based analysis of population heterogeneity in a consistent manner across batches.
  - No need to identify mappings between separate clusterings
  - Increased number of cells allows for greater resolution of population structure
- BUT...
- It is not recommended to use the corrected expression values for gene-based analyses (e.g. differential expression)
- Arbitrary correction algorithms are not obliged to preserve the magnitude (or even direction) of differences in per-gene expression when attempting to align multiple batches

# Performance assessment

- Qualitative (visualization)
- Quantitative:
  - Silhouette score
  - kBET: k-nearest-neighbor batch-effect test
  - ...

# Silhouette score

A score for each cell that assesses the separation of cell types, with a high score suggesting that cells of the same cell type are close together and far from other cells of a different type.

$a(i)$  is the average distance of cell  $i$  to all other cells within  $i$ 's cluster.

$b(i)$  is the average distance of  $i$  to all cells in the nearest cluster to which  $i$  does not belong.

Silhouette score:

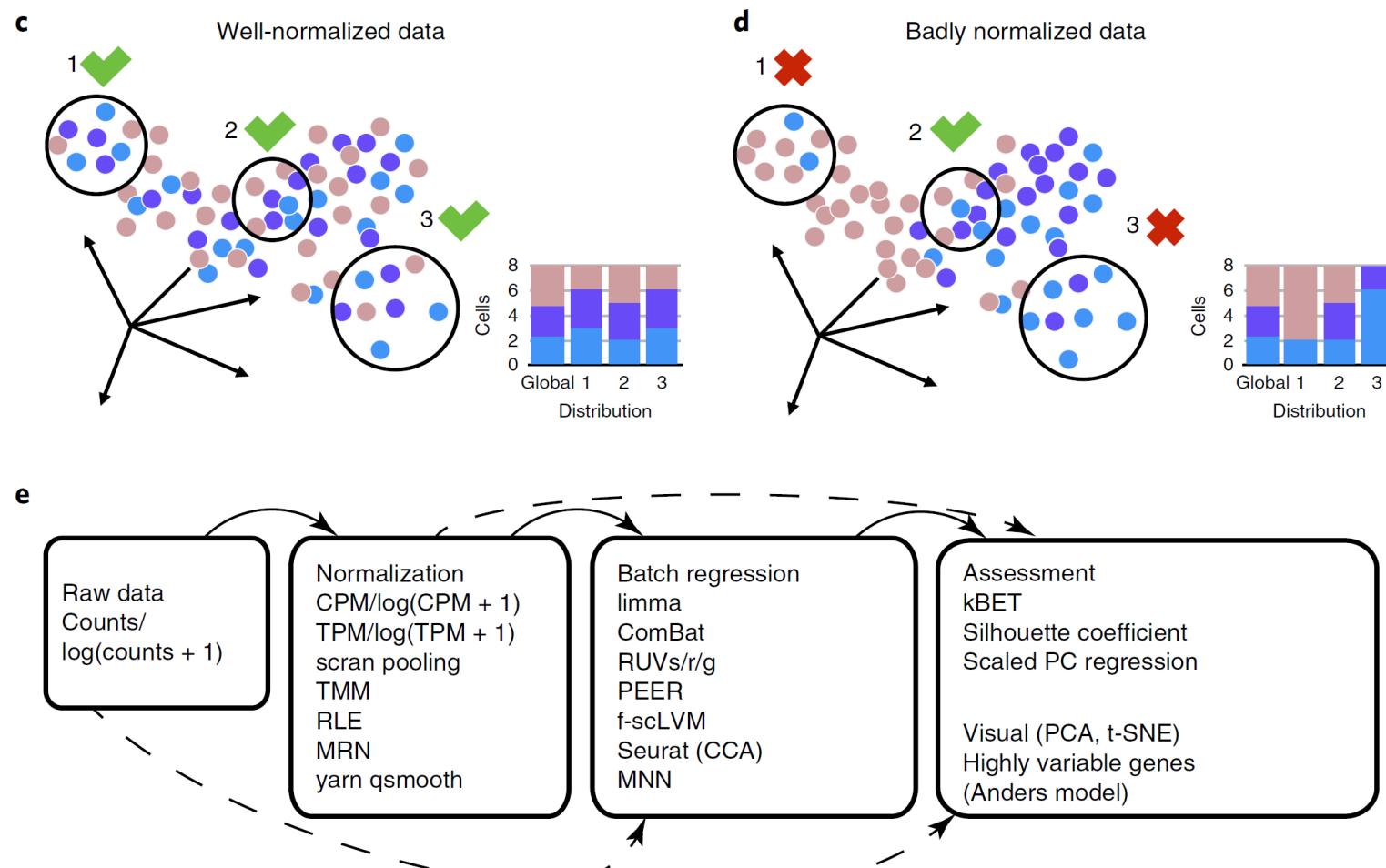
$$S = \frac{1}{N} \sum s(i)$$

$$s(i) = \frac{b(i) - a(i)}{\max(a(i), b(i))}$$

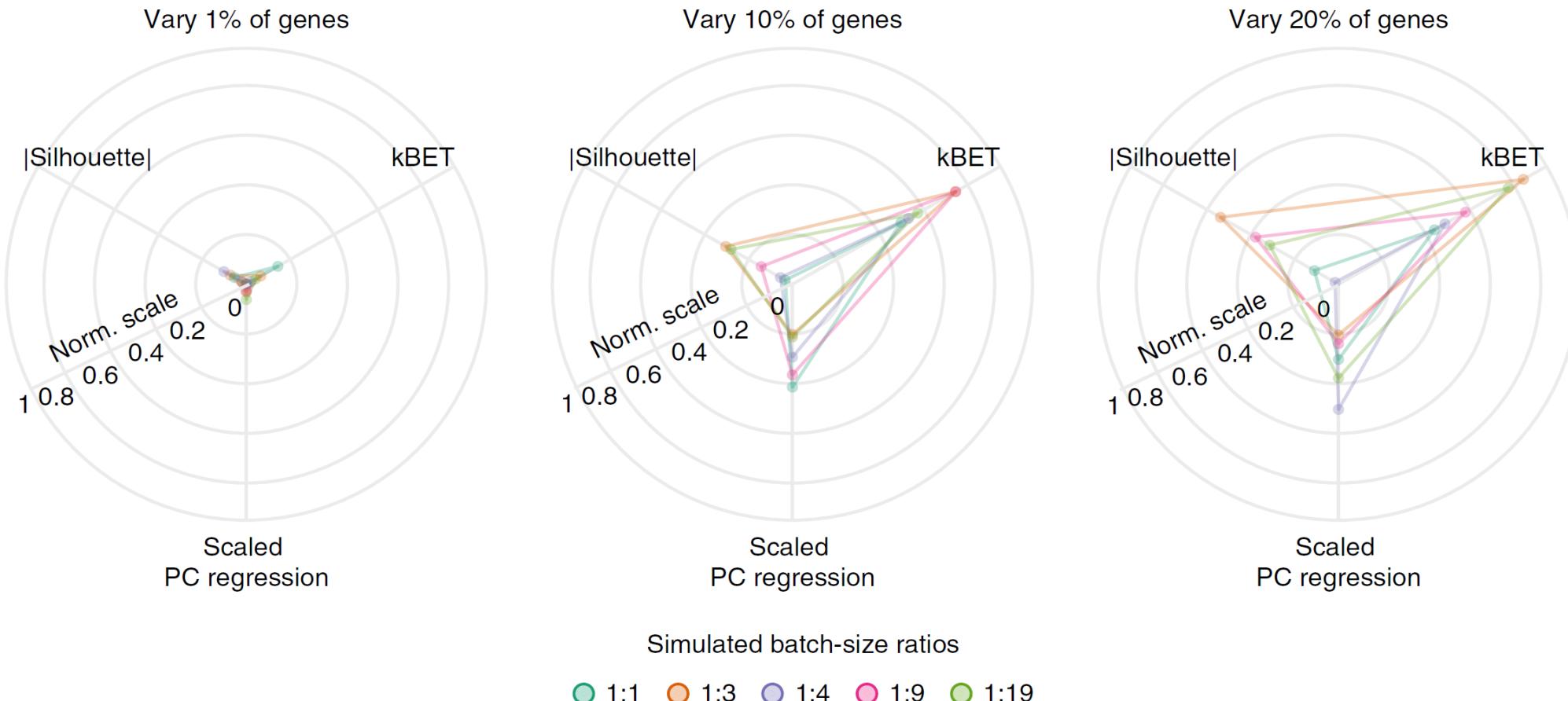
$$a(i) = \frac{1}{|C_i|} \sum_{\forall j} d(x_i, x_j)$$

$$b(i) = \min_{\forall j, j \notin C_i} d(x_i, x_j)$$

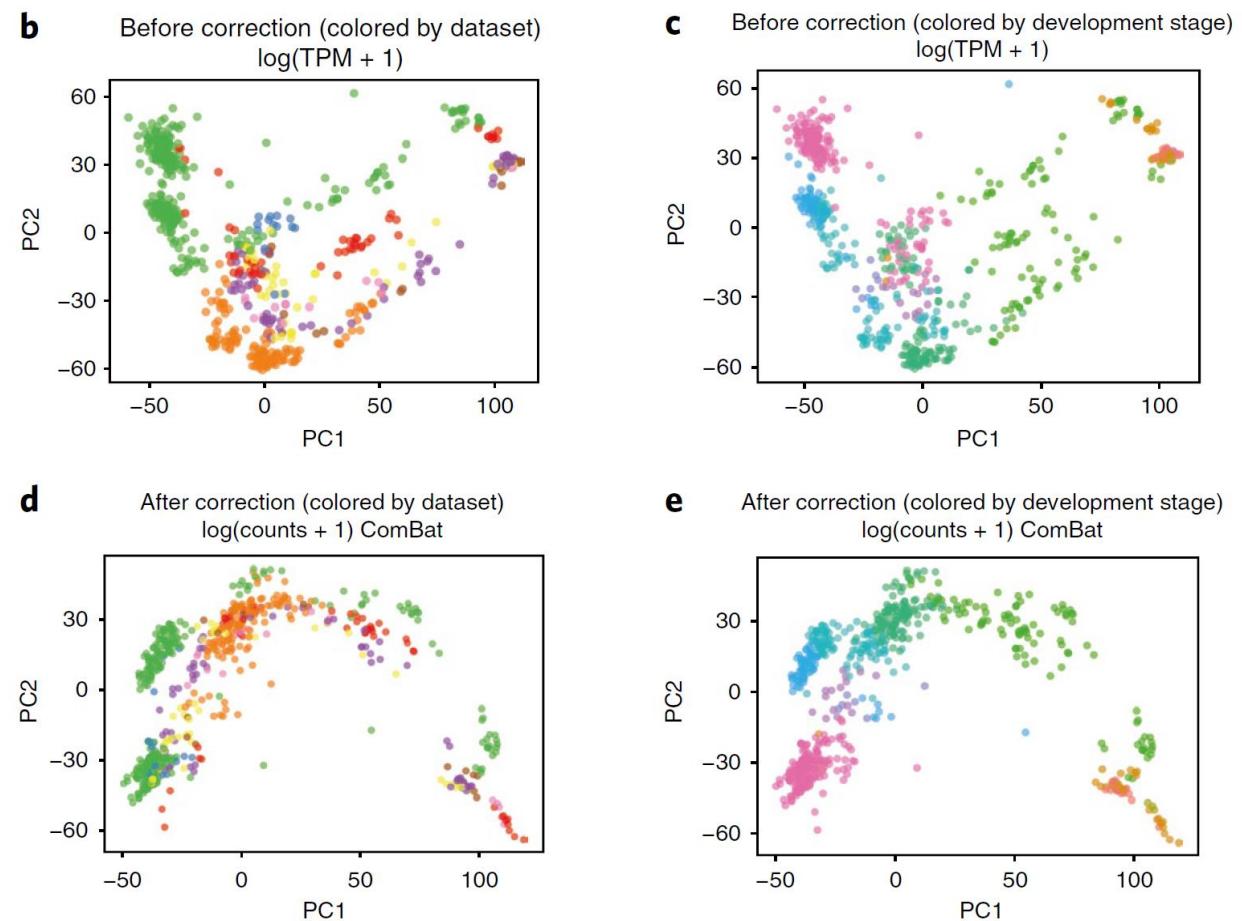
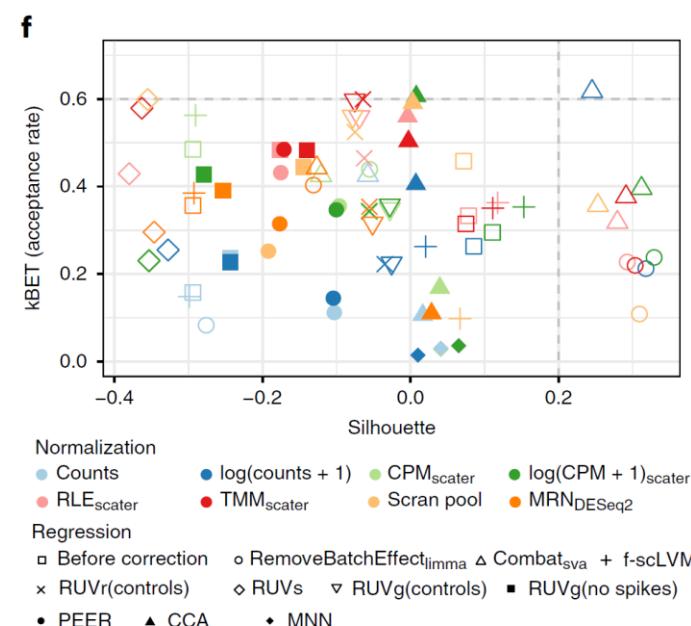
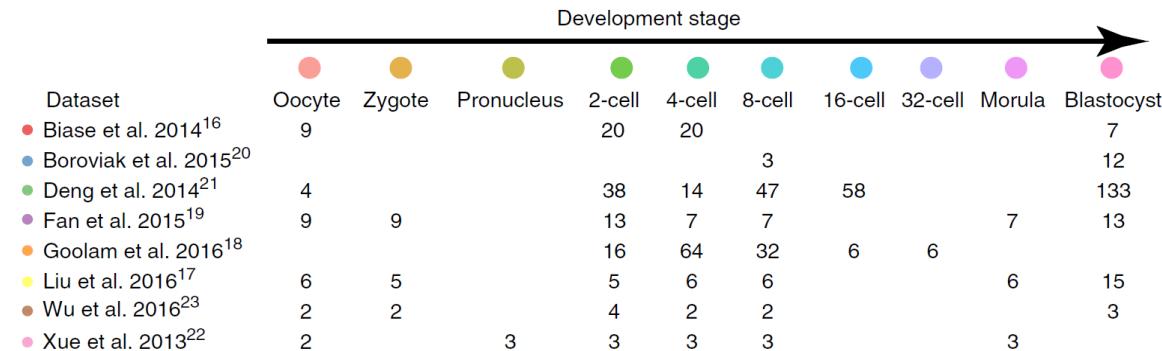
# kBET: $k$ -nearest-neighbor batch-effect test



# kBET is more responsive than other batch tests on simulated data



# kBET assesses data-integration quality

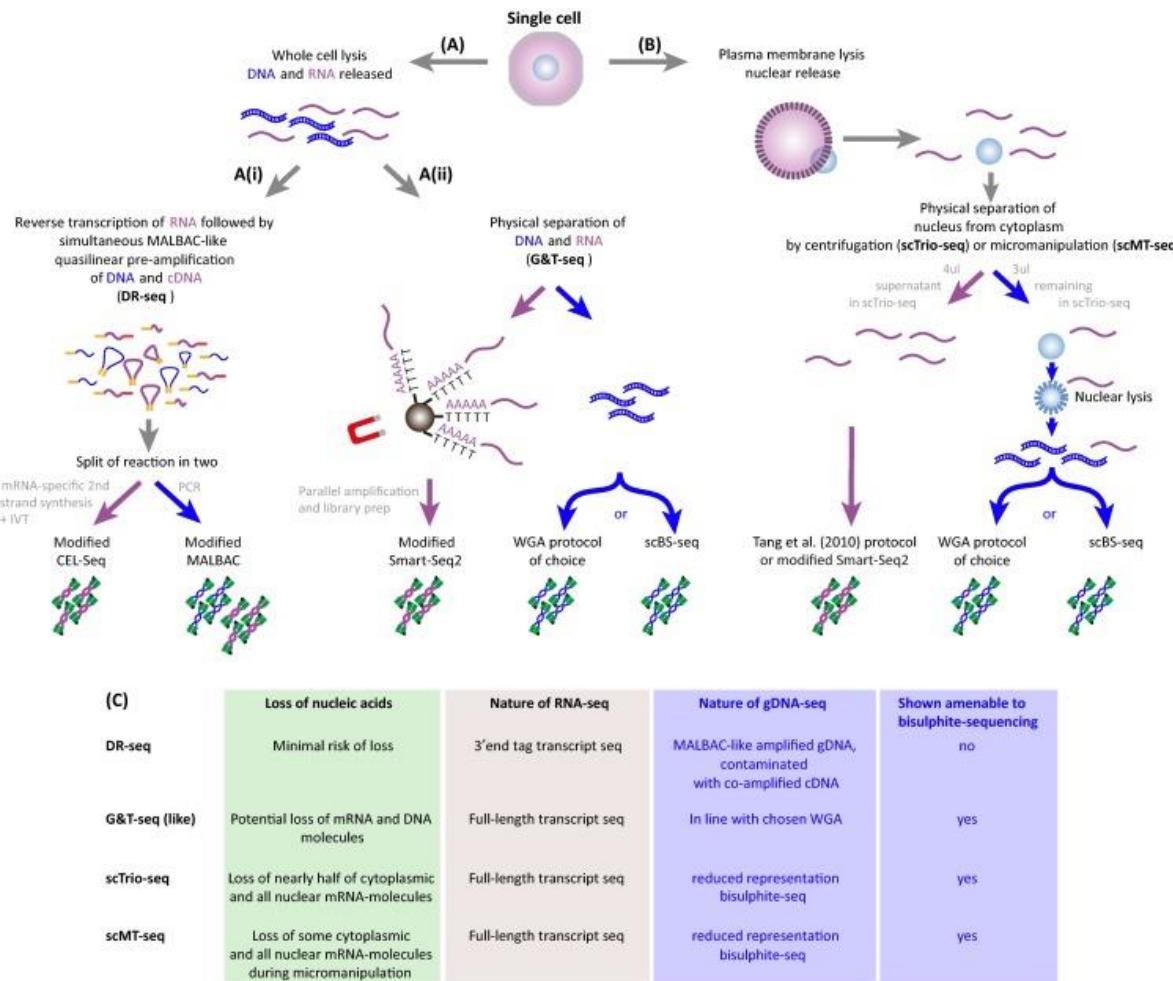


# Summary (so far)

- Integration can allow us to **improve the interpretation** of single-cell data, and build a **multi-modal view** of the tissue
- Numerous methods now available for integration, mainly using **joint dimension reduction**, or **joint clustering**, or a combination of both
- Joint dimension reduction can yield **interpretable factors** and aid in the identification of equivalent states, but is computationally expensive
- Graph-based methods alone can be **extremely fast**, but may struggle when technical differences are on a similar scale to biological differences

# Single cell Multi-omics

Same cell



# High-throughput sequencing of the transcriptome and chromatin accessibility in the same cell

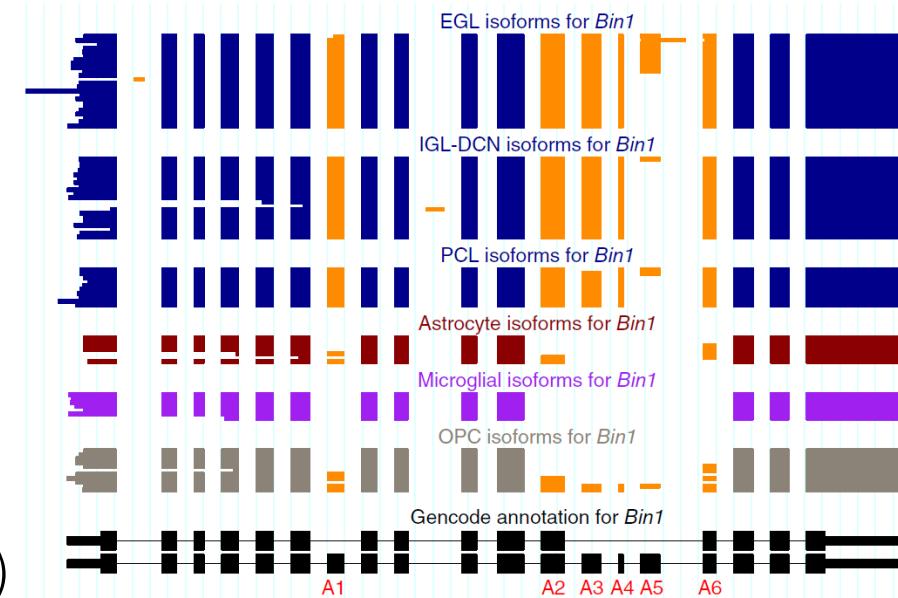
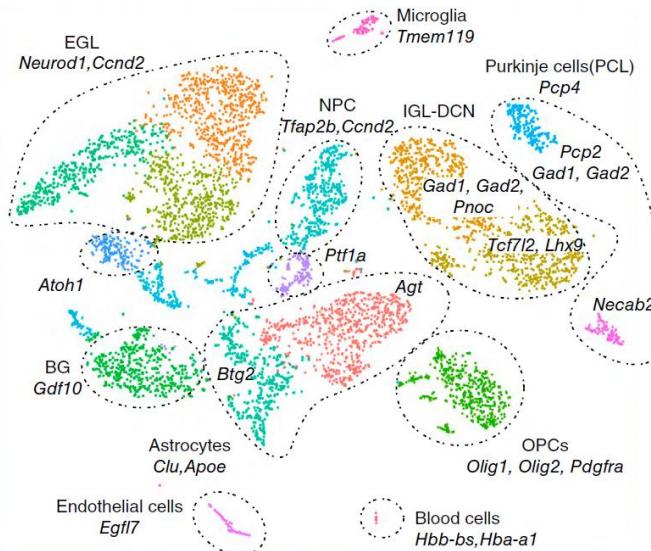
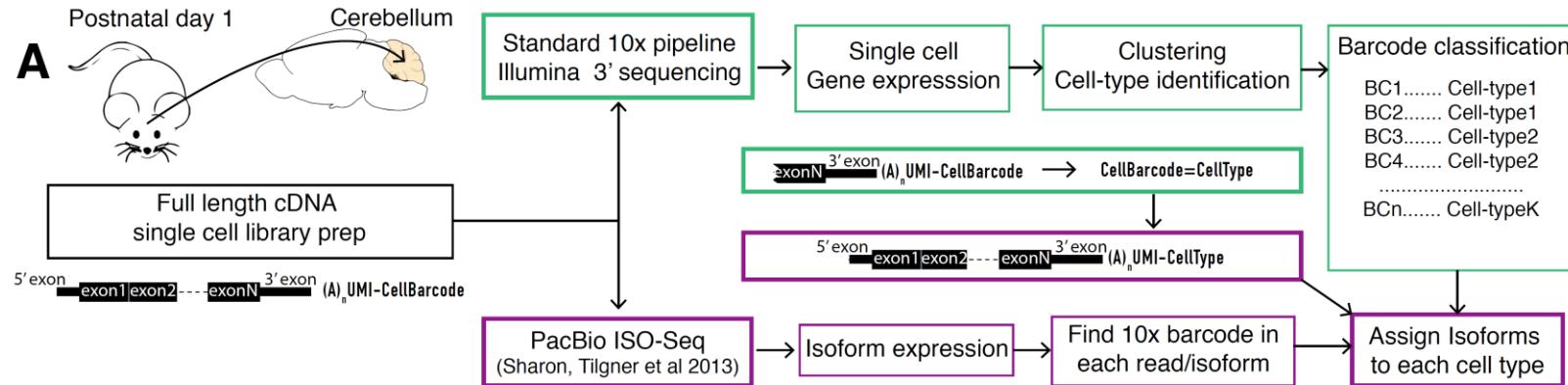
Song Chen , Blue B. Lake  and Kun Zhang \*

# Simultaneous profiling of 3D genome structure and DNA methylation in single human cells

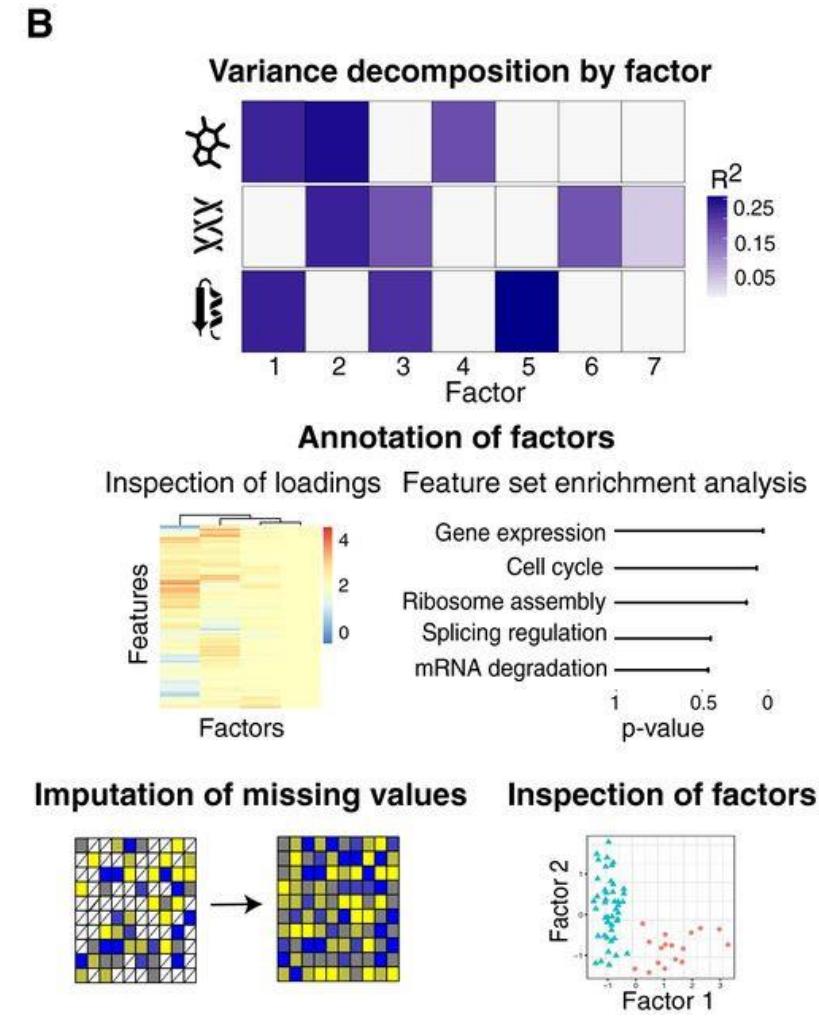
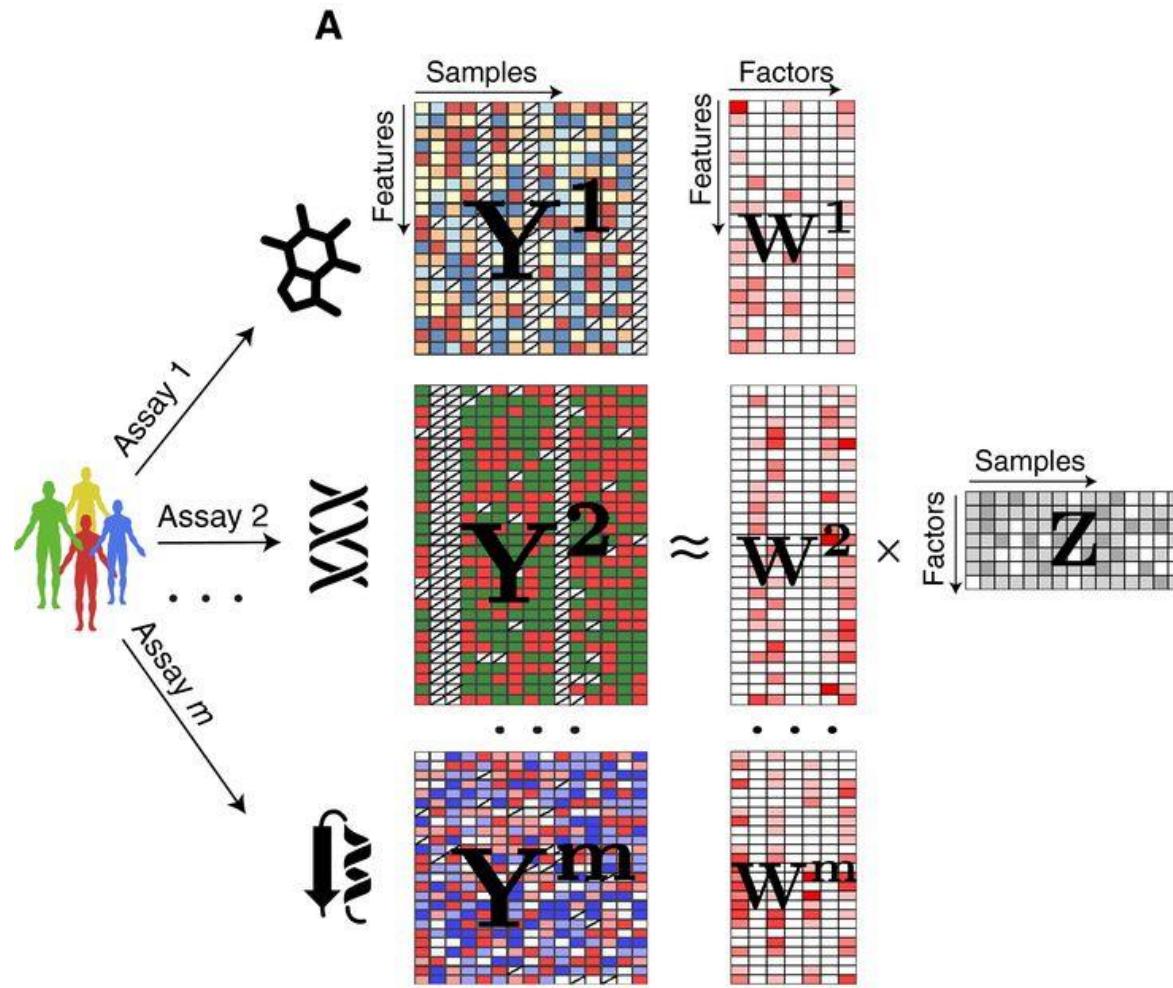
Dong-Sung Lee<sup>1,5</sup>, Chongyuan Luo<sup>2,3,5</sup>, Jingtian Zhou<sup>2,5</sup>, Sahaana Chandran<sup>1</sup>, Angeline Rivkin<sup>2</sup>, Anna Bartlett<sup>2</sup>, Joseph R. Nery  <sup>2</sup>, Conor Fitzpatrick<sup>4</sup>, Carolyn O'Connor<sup>4</sup>, Jesse R. Dixon  <sup>1\*</sup> and Joseph R. Ecker  <sup>2,3\*</sup>

# Single cell isoform RNA sequencing

## ScISOr-seq



# Multi-omics factor analysis



# Data integration practical

- MNN correction
- Seurat v3
- Four pancreatic datasets

# Resources

- Stuart et al. “Comprehensive integration of single cell data”  
<https://www.biorxiv.org/content/10.1101/460147v1>
- Haghverdi et al. “Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors”  
<https://doi.org/10.1038/nbt.4091>
- Tim Stuart “Integration and harmonization of single-cell data” (Satija Lab single cell genomics day 2019)  
<https://satijalab.org/scgd/>
- Andrew Butler “Batch Correction and Data Integration for Single Cell Transcriptomics” (Satija Lab single cell genomics day 2018)  
<https://satijalab.org/scgd18/>
- Orchestrating Single-Cell Analysis with Bioconductor  
<https://osca.bioconductor.org/>
- Hemberg’s group course: Analysis of single cell RNA-seq data  
<https://scrnaseq-course.cog.sanger.ac.uk/website/index.html>
- Seurat Integration and Label Transfer tutorial  
[https://satijalab.org/seurat/v3.0/pancreas\\_integration\\_label\\_transfer.html](https://satijalab.org/seurat/v3.0/pancreas_integration_label_transfer.html)

# Thank You!

✉ a.mahfouz@lumc.nl

🔗 <https://www.lcbc.nl/>

🐦 @ahmedElkoussy