

# Data integration and batch correction in single-cell data

**Paulo Czarnewski**

*ELIXIR Single-Cell Omics community co-lead*  
*National Bioinformatics Infrastructure Sweden (ELIXIR-SE)*

# Sources of variation in single-cell



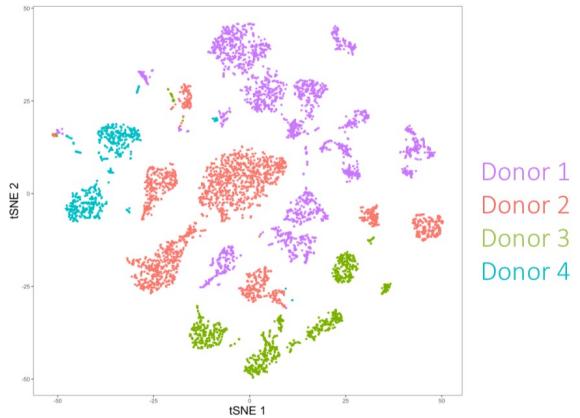
## Biological:

- Cell Type Heterogeneity
- Genetics
- Cell State/Microenvir.
- GExpr Stochasticity
- Cell Cycle Dynamics
- Transcriptional Bursts
- Oscillations
- ...

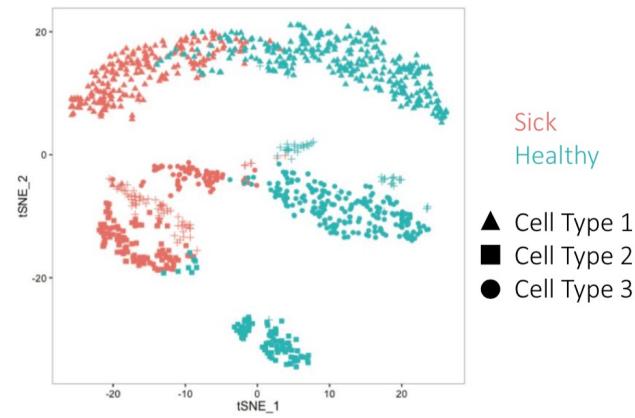
## Technical:

- Capture Efficiency
- Amplification Bias
- PCR artifacts
- Contamination
- Cell Doublets
- Cell Damage
- Sampling (Jackpot Effects)
- ...

# Sources of variation in single-cell: Donors

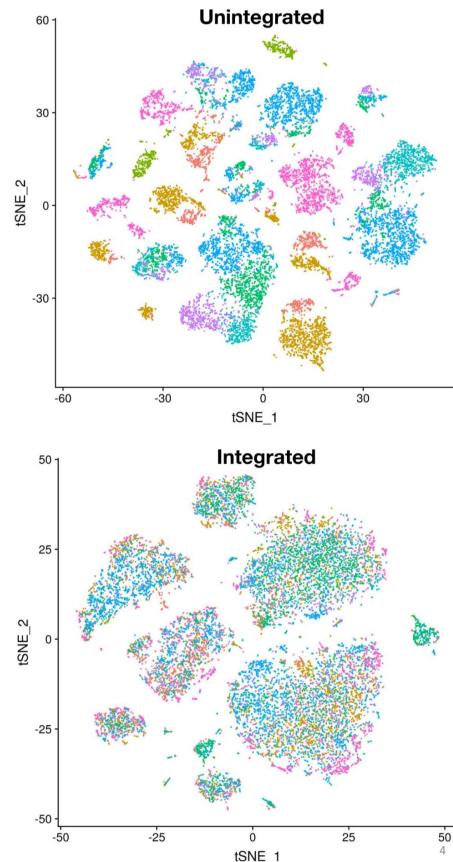
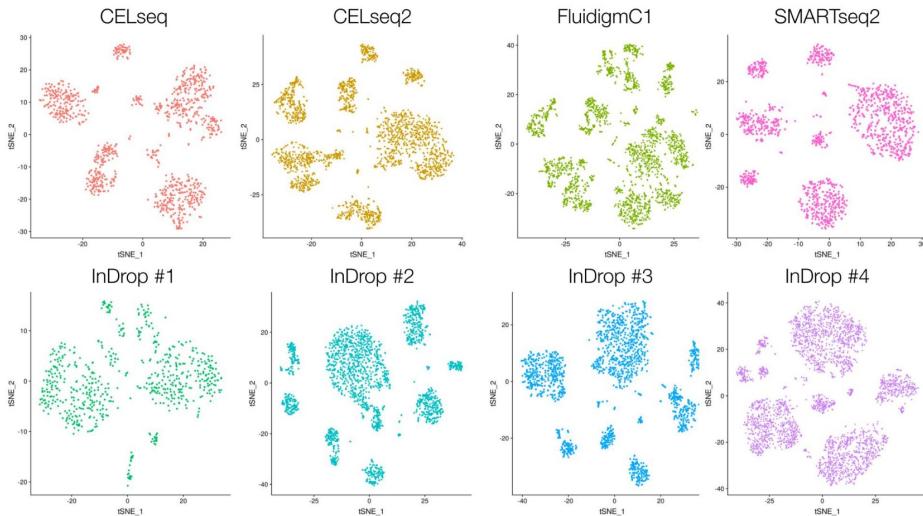


Same tissue from different donors



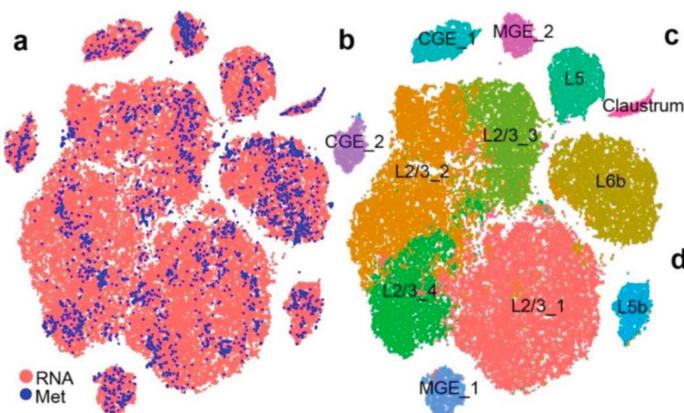
Cross condition comparisons

# Sources of variation in single-cell: Technology

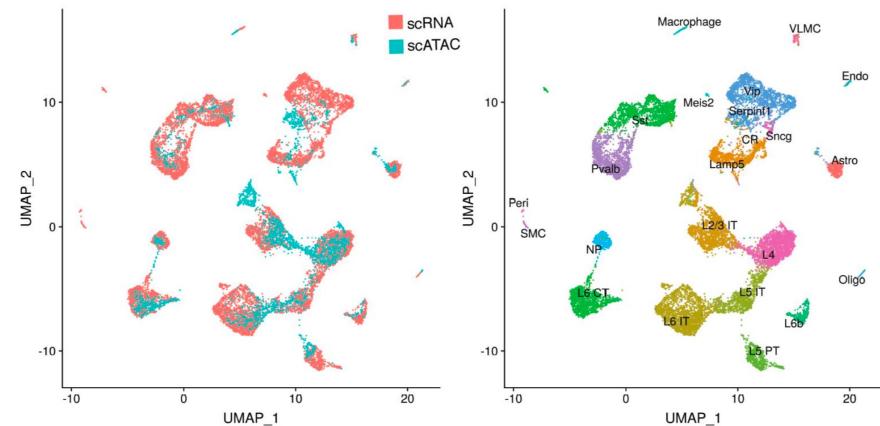


# Sources of variation in single-cell: OMICs

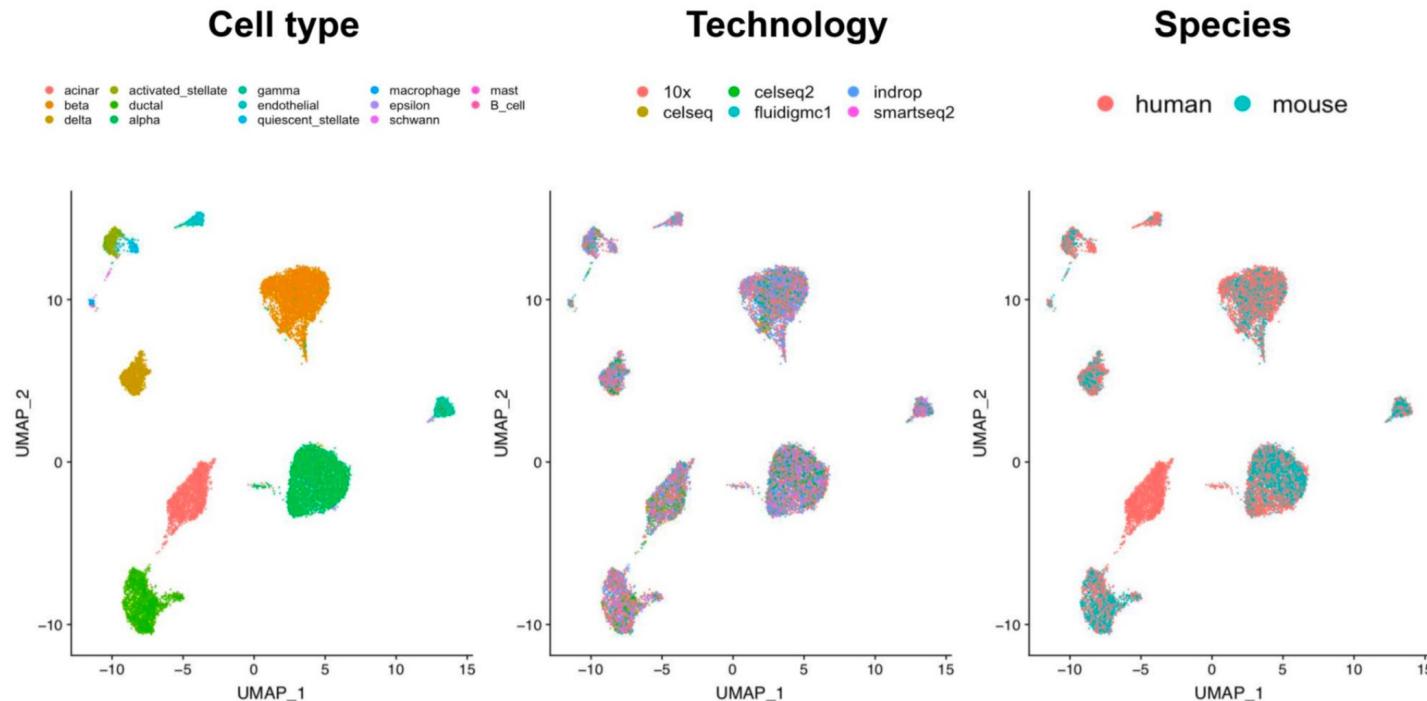
snRNAseq + snNMTseq



scRNAseq + scATACseq



# Sources of variation in single-cell: Species

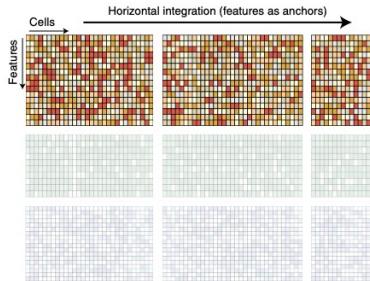


Retinal bipolar datasets: 51K cells, 6 technologies, 2 Species

# Types of data integration

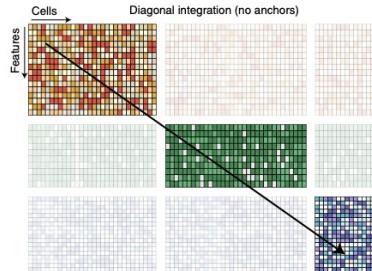
## Horizontal

- Shared features
- Different samples



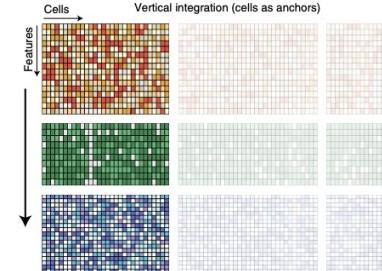
## Diagonal

- Different features
- Different samples



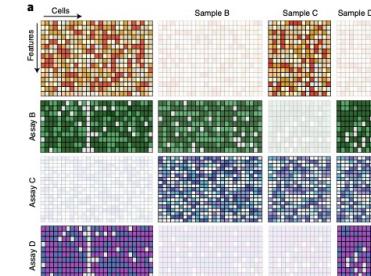
## Vertical

- Different features
- Shared samples



## Mosaic

- Some features shared
- Some samples shared

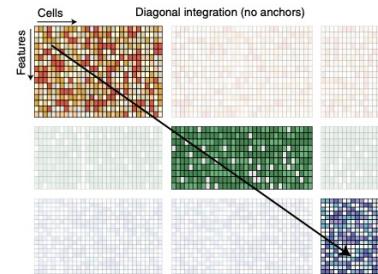
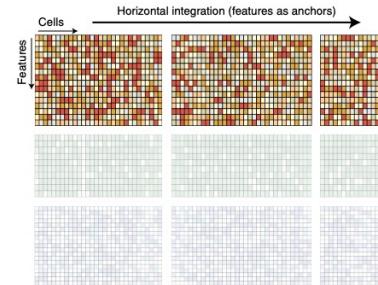


# Types of data integration

**Table 1 | Overview of common data integration methods classified according to their anchor choice**

Integration task	Method	Ref.
Vertical (global)	CCA	112
Vertical (global)	JIVE	70
Vertical (global)	PLS	71
Vertical (global)	MCIA	113
→ Vertical (global)	MOFA+	65
→ Vertical (global)	scAI	114
→ Vertical (global)	iNMF	38
Vertical (global)	Seurat v4	11
Vertical (local)	Spearman's rank correlation coefficient	50
Vertical (local)	LMM	51
→ Horizontal	MNN	21
→ Horizontal	Seurat v3	22
Horizontal	LIGER	23
→ Horizontal	Harmony	24
Horizontal	Scanorama	29
Horizontal	BBKNN	25
Horizontal	scVI	26
Horizontal	scmap	28
Horizontal	conos	27
Diagonal	MATCHER	77
Diagonal	MMD-MMA	78
Diagonal	SCIM	115
Diagonal	UnionCom	116
→ Diagonal	coupledNMF	117

# Horizontal and Diagonal integration

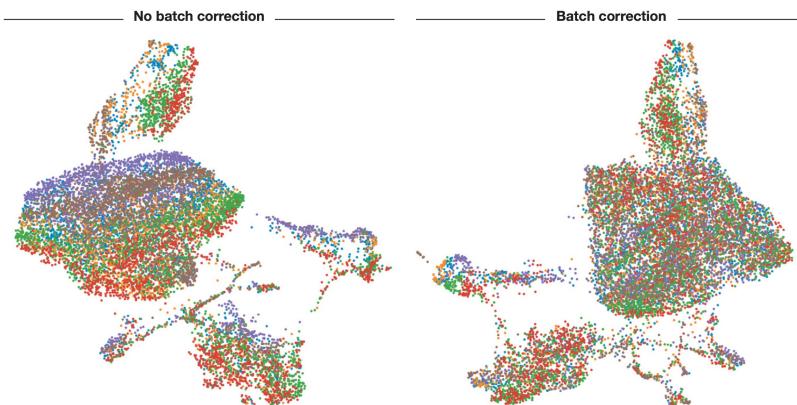


# ComBat

Uses empirical Bayes regression on shared gene factors

Works well on simpler small-medium datasets

All datasets need to be similar in cell type composition



- limma::removeBatchEffect()
- seurat::ScaleData()
- sva::combat()
- batchelor::rescaleBatches()

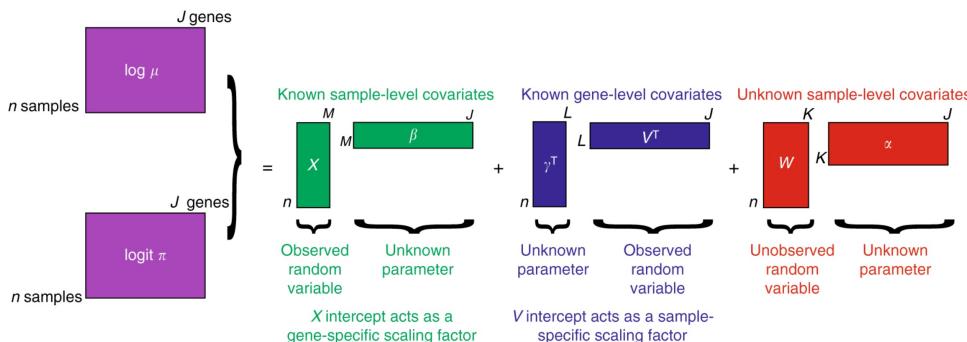
## Known issues with ComBat:

1. Do not account for differences in population composition
2. Assume batch effect is additive
3. Prone to overcorrection (in cases of partial confounding)
4. Will fail in large datasets with complex mixture of cell type

# ZIMB-WaVE (factor analysis)

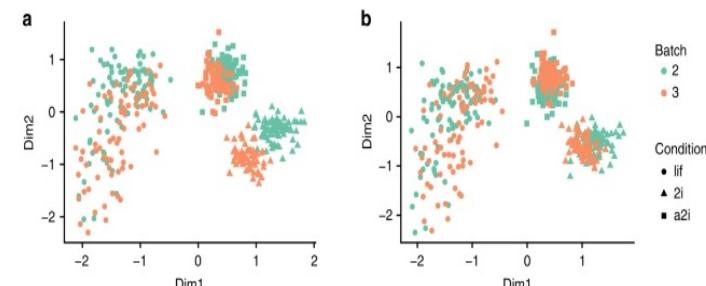
Zero-Inflated Negative Binomial-based Wanted Variation Extraction

Can accommodate both gene- and cell-level covariates



## Known issues with ZIMB-WaVE:

1. Do not account for differences in population composition
2. Assume batch effect is additive
3. It will be slow on large datasets
4. Will fail in large datasets with complex mixture of cell type



# Mutual Nearest Neighbours (MNN)

Dimensionality reduction via multibatch PCA with all datasets

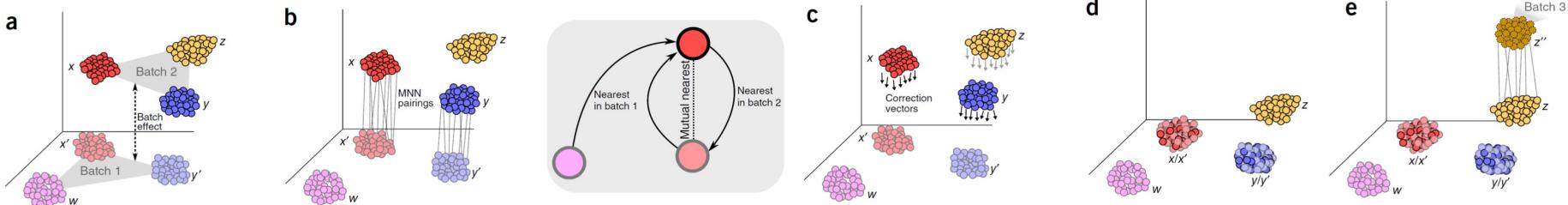
Find kNN within datasets and kMNN across datasets

Compute merging vectors for each cell

It scales well on large datasets

## Known issues with MNN:

1. There is at least one cell population that is present in both batches
2. The batch effect is almost orthogonal to the biological subspace
3. Batch effect variation is much smaller than the biological effect variation between different cell types



# Scanorama

Python

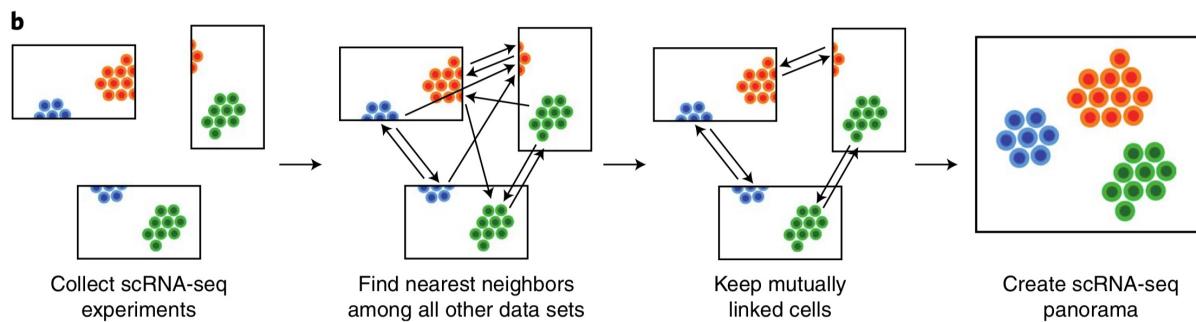
Dimensionality reduction via SVD with all datasets

Find K-NN within and across datasets

Compute merging vectors

## Known issues with Scanorama:

1. There is at least one cell population that is present in both batches
2. The batch effect is almost orthogonal to the biological subspace
3. Batch effect variation is much smaller than the biological effect variation between different cell types



# Harmony

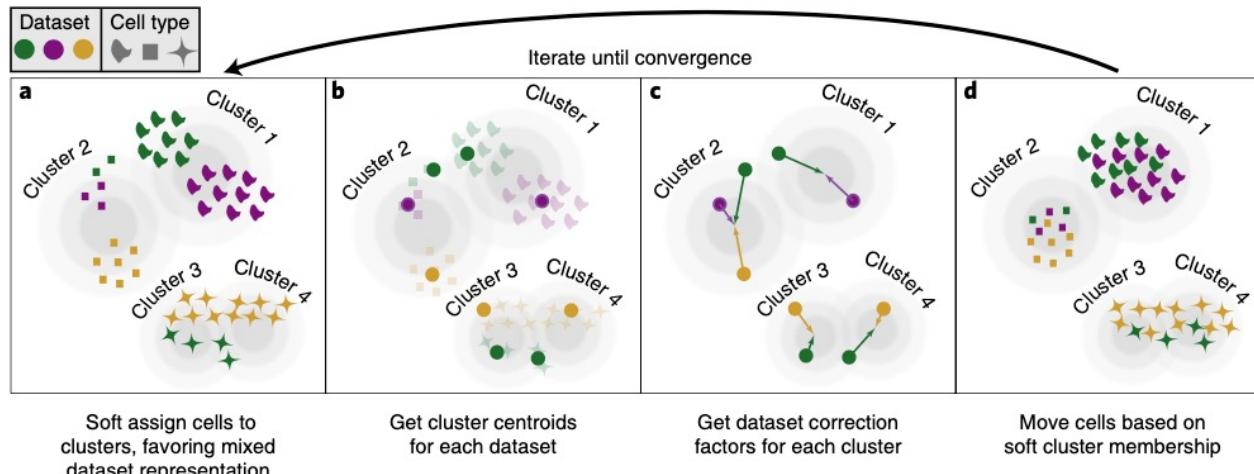
Employs fuzzy k-means clustering to speed-up computations

Find K-NN to cluster centroids across datasets

Compute merging vectors for each cell

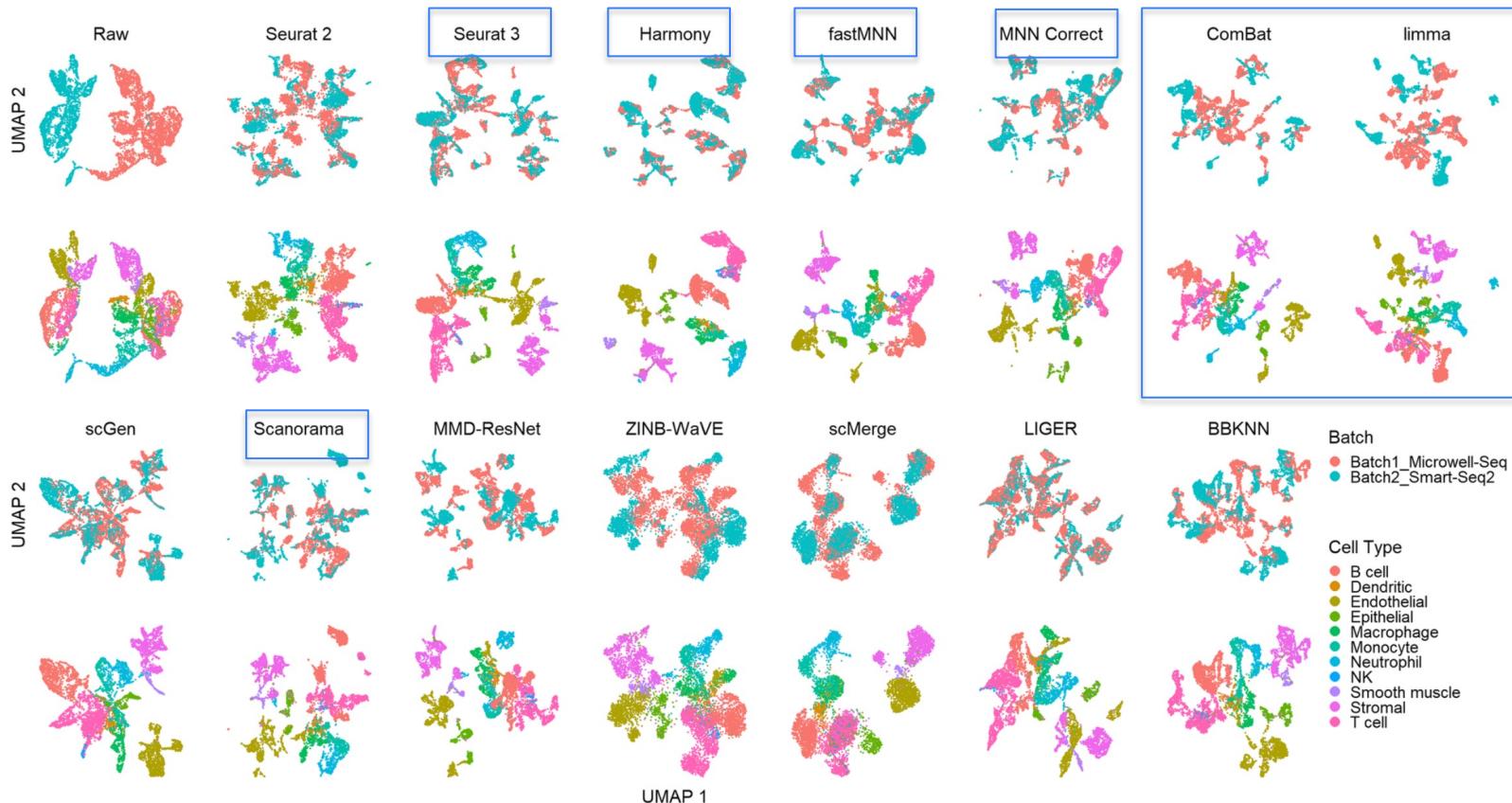
## Known issues with Harmony:

1. Need to find optimal number of clusters for your dataset
2. Need to define convergence criteria / thresholds

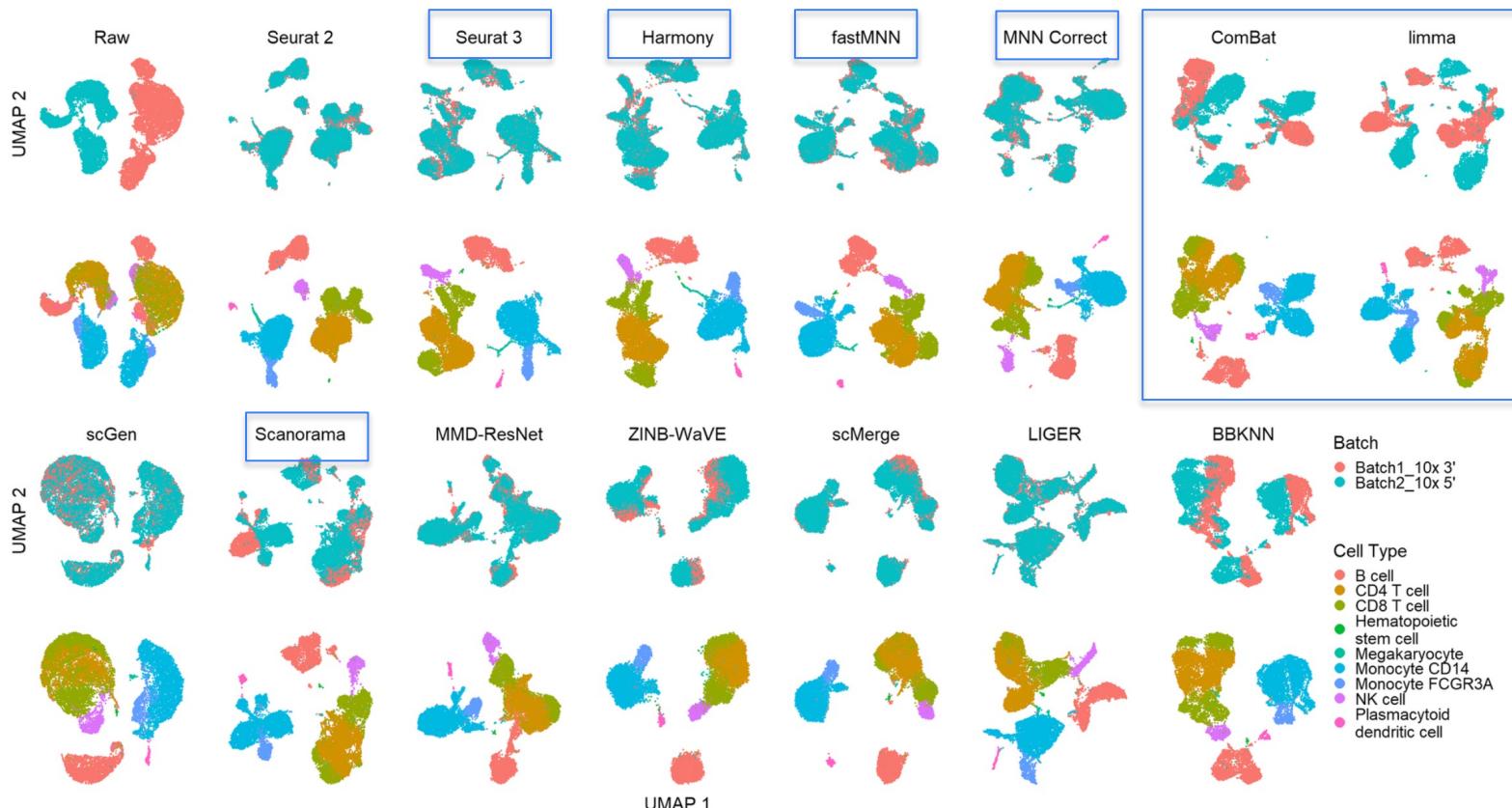


# Integration performance assessment

# Integration performance assessment



# Integration performance assessment



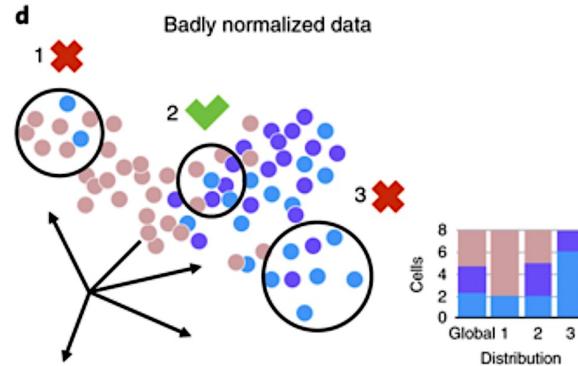
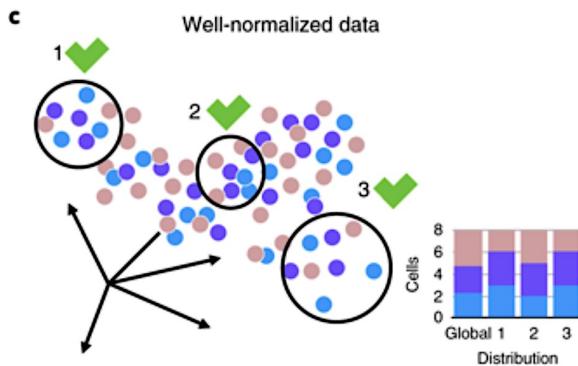
# Integration performance assessment

Goal:

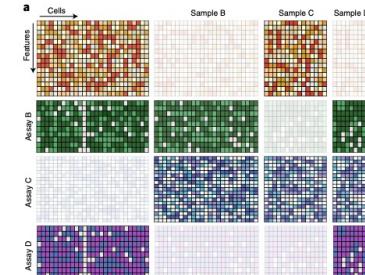
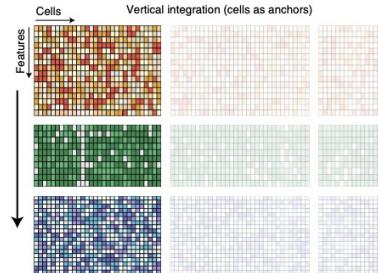
- The batch-originating variance is erased
- Meaningful heterogeneity is preserved
- No artefactual variance is introduced

What it practically means:

- Similar cell types are intermixed across batches
- We are not mixing distinct cell types (across or within batches)
- We do not separate similar cells within batches



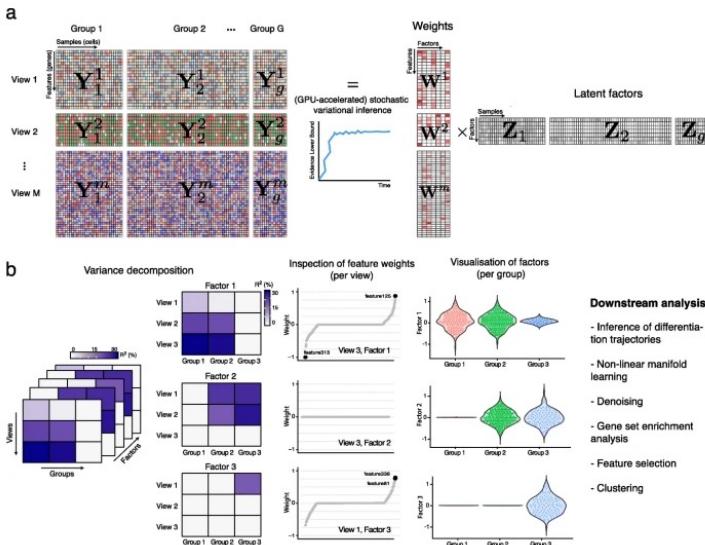
# Vertical and Mosaic integration



# Factor analysis via NNMF

Accurate, flexible and robust methods

## MOFA+

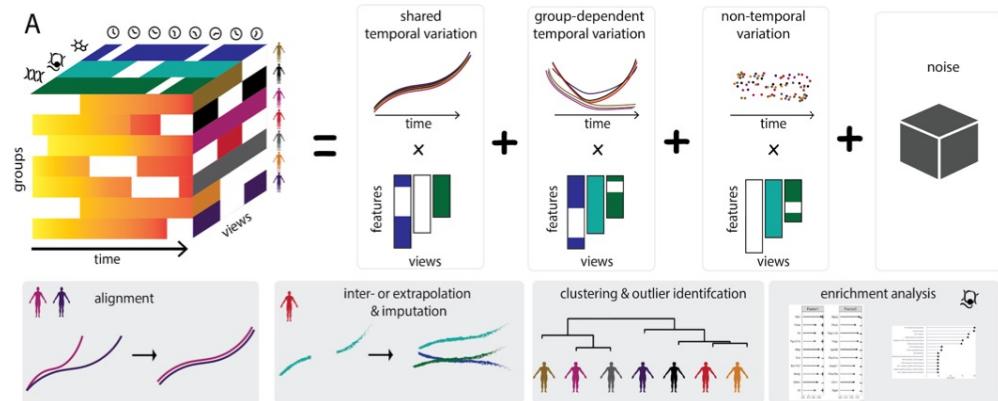


[Argelaguet et al \(2020\) Genome Biology](#)

Known issues with NNMF for single-cell:

1. Will fail in large datasets with complex mixture of cell type

## MEPHISTO



[Velten et al \(2020\) BioRxiv](#)

# Weighted Nearest Neighbours (WNN)

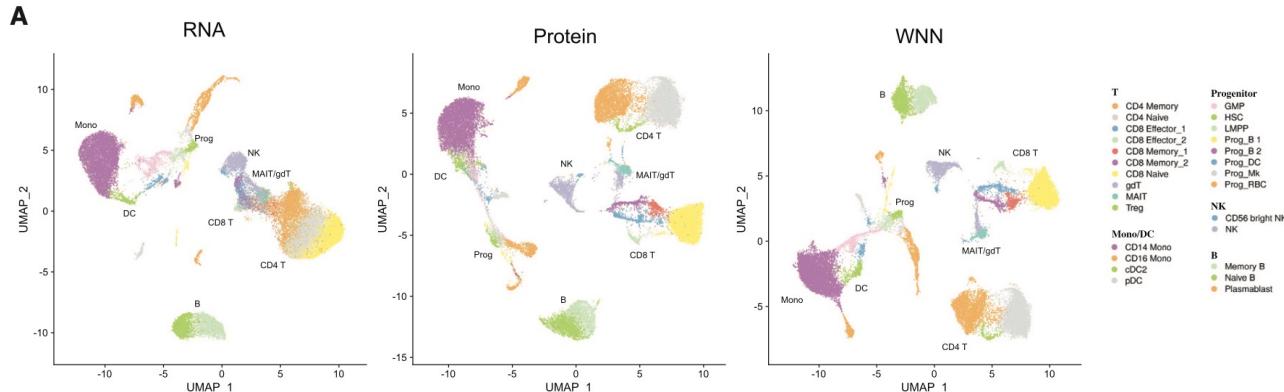
Computes KNN graphs within modalities

Graphs are then weighted and merged (just like in similarity network fusion: SNF)

Can use different similarities for each OMIC

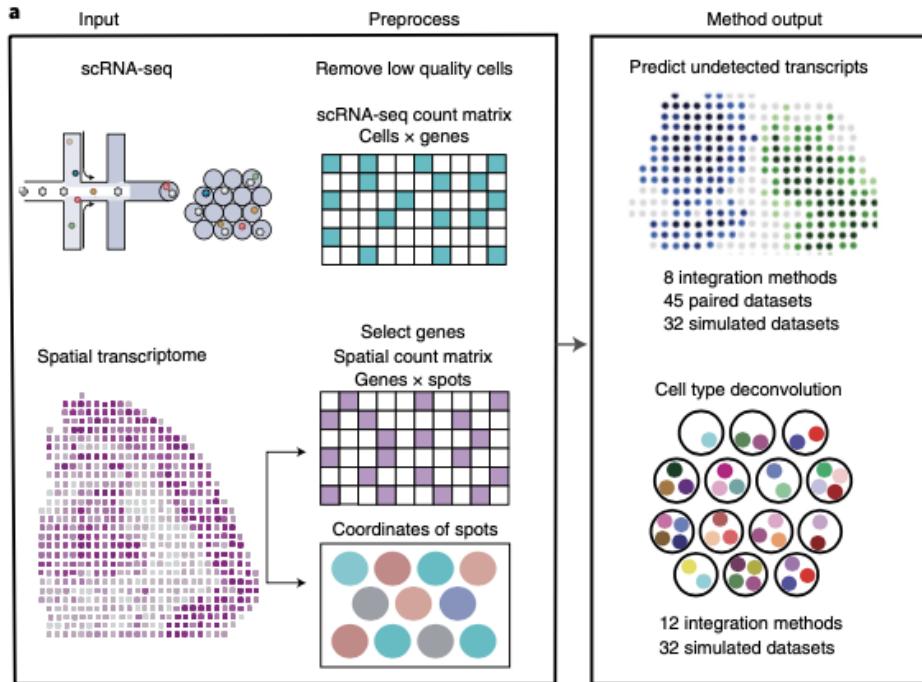
**Known issues with WNN:**

1. Graph weights are somewhat arbitrary



# Data integration between single-cell and spatial data

# Types of integration tasks



Gene imputation

Cell type deconvolution

# Types of integration tasks

Marker gene based	Anchor based	Probabilistic Modelling	Optimization based
<p>Extract marker genes (MG) for each cell type from SC data</p> <p>Compute enrichment score for each set of MGs in spatial locations</p> <p>Normalize to make scores sum to 1</p> <p><b>Ex:</b> Itai et.al</p>	<p>Find anchors between modalities (MNNs). Create correction vector based on differences in expression.</p> <p>Use correction vectors to remove platform effects. Integrated data sets.</p> <p>Transfer labels of single cells to spatial data points.</p> <p><b>Ex:</b> Seurat</p>	<p>Assume gene expression follows certain statistical distributions.</p> <p>Joint model for SC and spatial data. Learn cell type parameters from SC data, use to deconvolve spatial data (when mixed).</p> <p>Correct for eventual platform differences</p> <p><b>Ex:</b> stereoscope, RCTD, cell2location</p>	<p>Find spatial location where each cell is most likely to reside.</p> <p>Tries to simultaneously optimize terms such as:</p> <ul style="list-style-type: none"> <li>• Cell density</li> <li>• UMI distribution across genes within spots</li> <li>• gene distribution across spots</li> </ul> <p><b>Ex:</b> Tangram</p>
GSEA	Cell type classification MNN	Deconvolution WNNLS	KNN imputation

# spaGE: imputation of genes using single-cell data

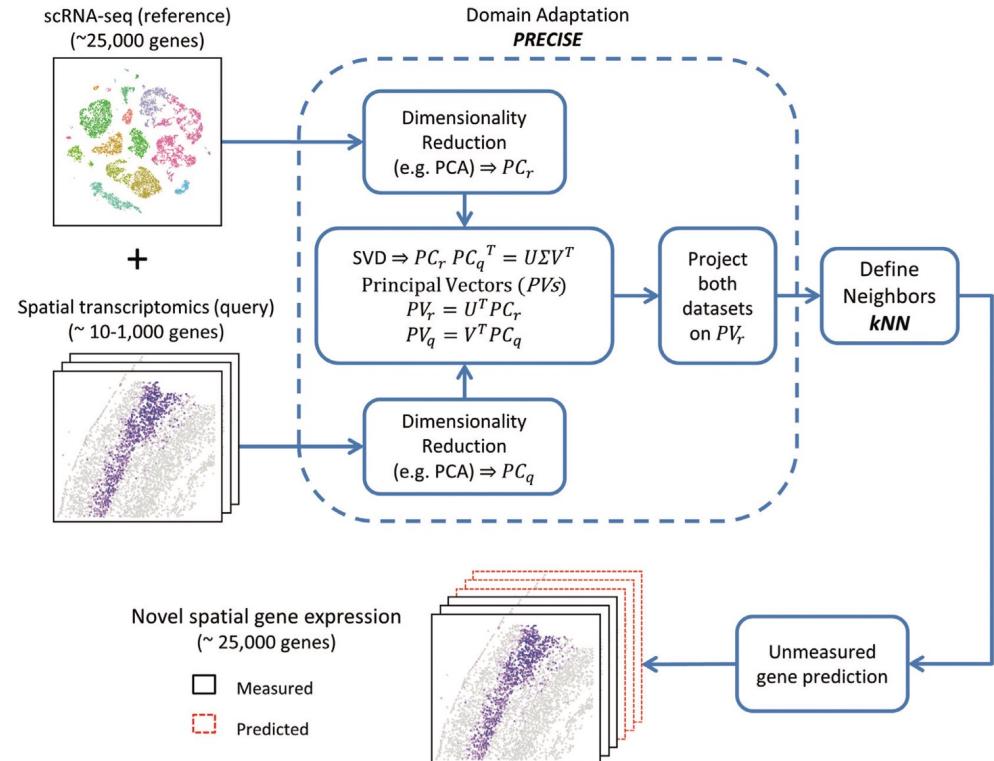
Single-cell genome-wide dataset (20K genes)

In situ probe-based spatial data (10-1K genes)

Integration of single-cell and spatial data  
(diagonal integration)

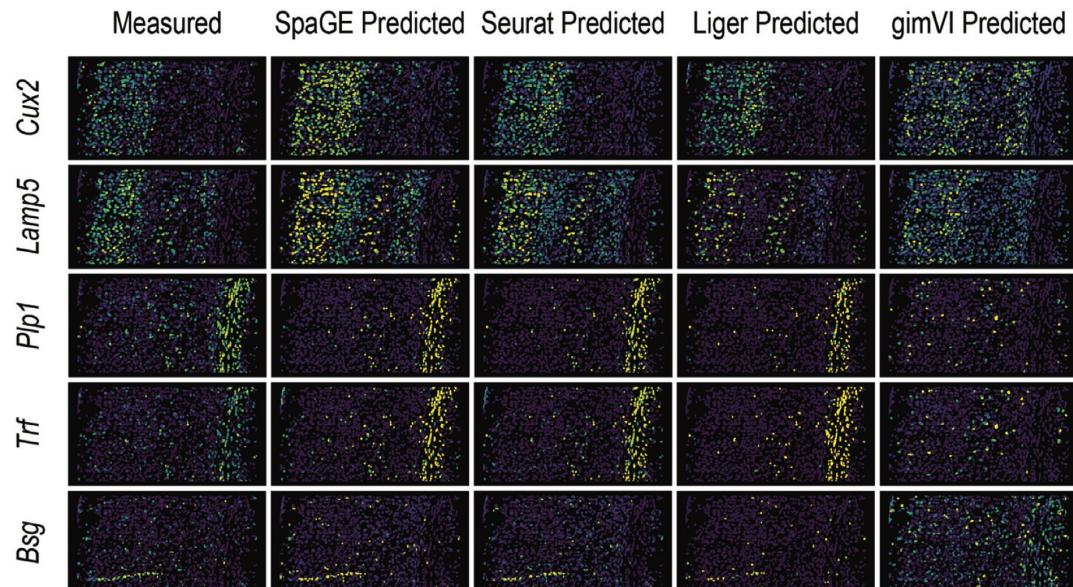
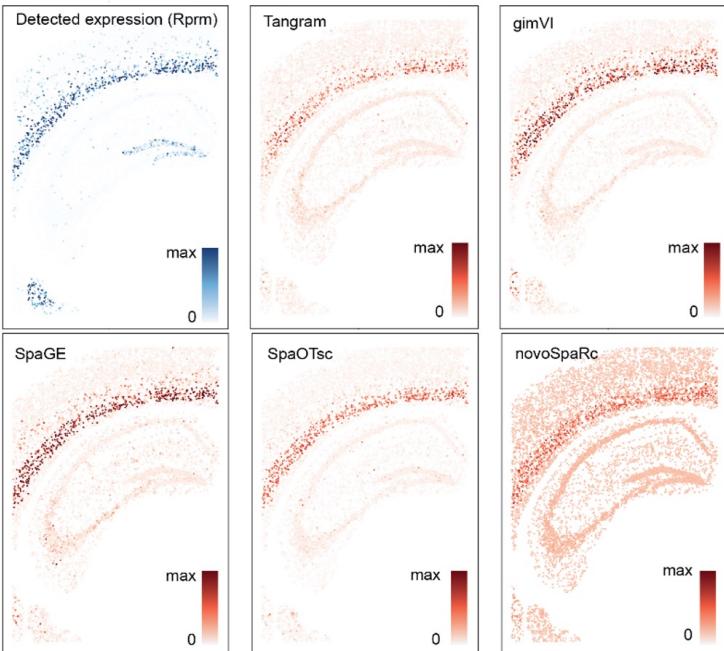
Define KNN in the integrated space across  
technologies (single cell vs spatial)

Predict gene expression value based on the  
expression of the single-cell



# Gene imputation method comparison

A



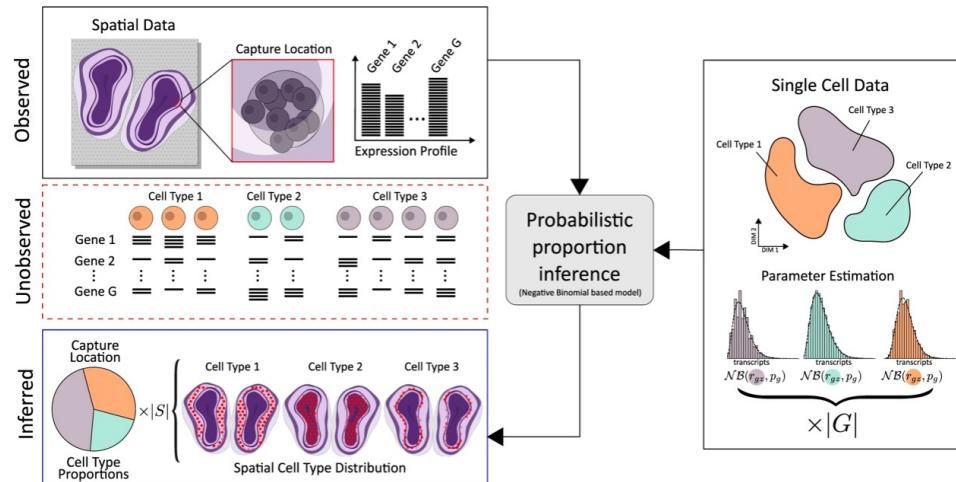
# Deconvolution: Stereoscope

Creates a NB glm model on the single cell dataset

Creates a NB glm model on the bulk dataset, including a variable accounting for cell type proportions ( $w$ ), given the gene-wise priors estimated from the single cell data

$$y_{gc} \sim \mathcal{NB}(s_c r_{gz_c}, p_g), \quad s_c = \sum_{g \in G} y_{gc}.$$

$$x_{sg} \sim \mathcal{NB}(\beta_g v_s^T r_g + \gamma_s \epsilon_g, p_g).$$

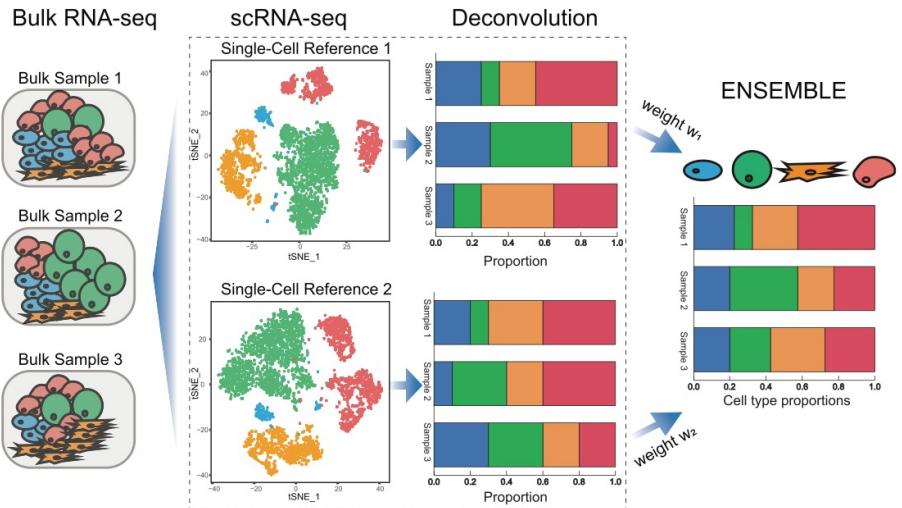


# Deconvolution: SCDC

Based on weighted non-negative least squares (WNNLS).

First, we create a basis-matrix contain the weighted gene expression across cell types in the single-cell dataset. Where P is the proportion of cell mixture (known).

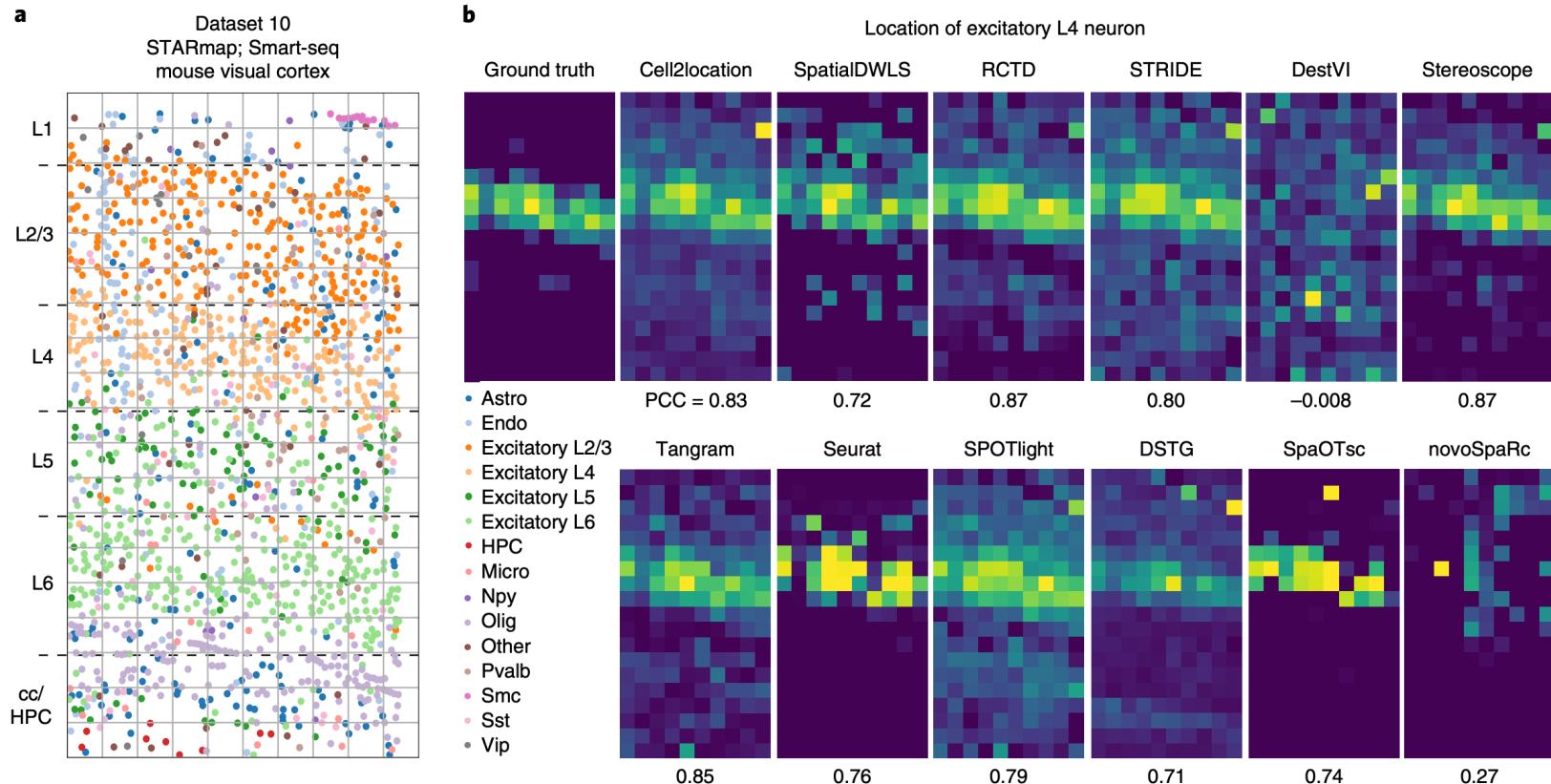
Then, on the bulk data, SCDC try to minimize the distance between the basis matrix from single-cell and the bulk.



$$\hat{\mathbf{Y}} = \hat{\mathbf{B}}\hat{\mathbf{P}},$$

$$\mathbf{Y} \approx \mathbf{B}\mathbf{P},$$

# Deconvolution method comparison



# Conclusions

Graph-based integration methods work well on very large datasets ( > 100K samples/cells ).

Matrix factorization are very flexible and accurate, but might not work in very large datasets.

Different integration techniques exist for different purposes, so please read up on benchmarking papers

Investigate which method is best on your case, and try a few (*most of them will agree to a certain extent*).

# Thank you!

**Paulo Czarnewski**

*ELIXIR Single-Cell Omics community co-lead  
National Bioinformatics Infrastructure Sweden (ELIXIR-SE)*