

# Batch effects in scRNA-seq data

Committee meeting

---

Almut Lütge

DMLS - University of Zürich

2021-02-26



Swiss Institute of  
Bioinformatics

# Batch effects

# Batch effects

**Systematic differences** between same type of cells [.. ][caused by technical sources] (Oslokov, 2019)

Differences between data sets [..][, that] occur due to **uncontrolled variability** in **experimental factors**, e.g., reagent quality, [.. and] can interfere with downstream analyses if not explicitly modelled. (Lun, 2019)



**Systematic unwanted variation**

New Results

 [Comment on this paper](#)

## **CellMixS: quantifying and visualizing batch effects in single cell RNA-seq data**

 Almut Lütge,  Joanna Zyprych-Walczak, Urszula Brykczynska Kunzmann,  HelenaL Crowell, Daniela Calini,  Dheeraj Malhotra,  Charlotte Soneson,  Mark D Robinson

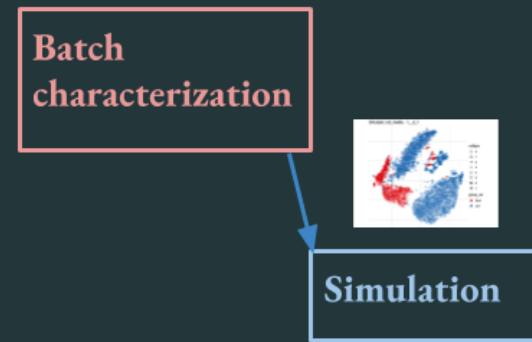
**doi:** <https://doi.org/10.1101/2020.12.11.420885>

This article is a preprint and has not been certified by peer review [[what does this mean?](#)].

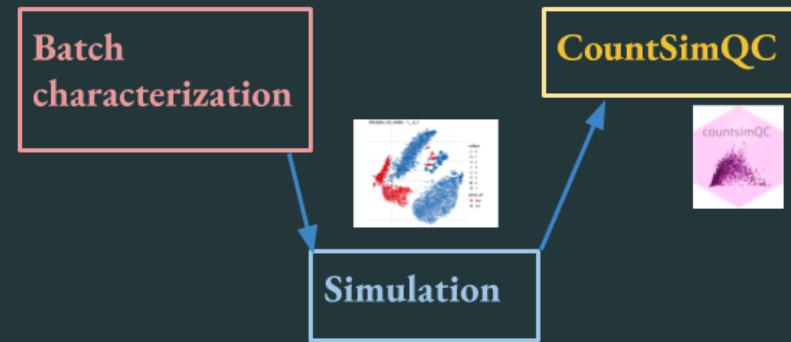
# Project overview

Batch  
characterization

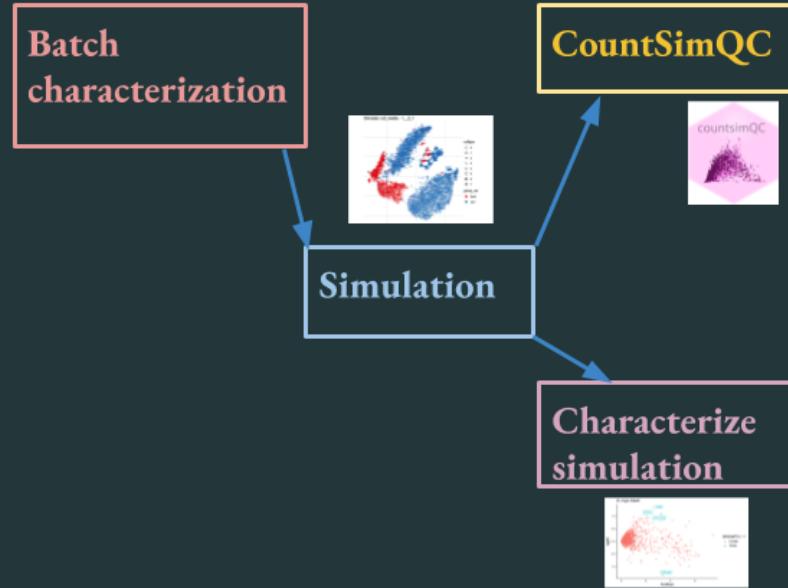
# Project overview



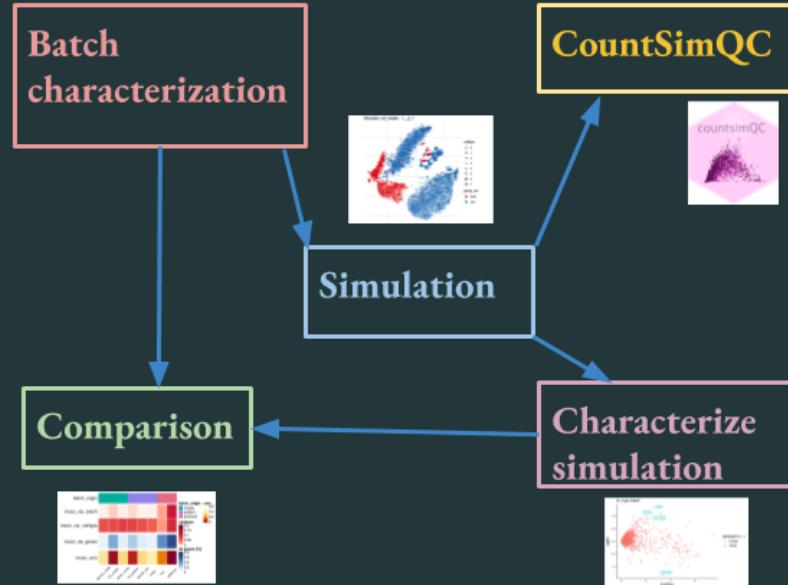
# Project overview



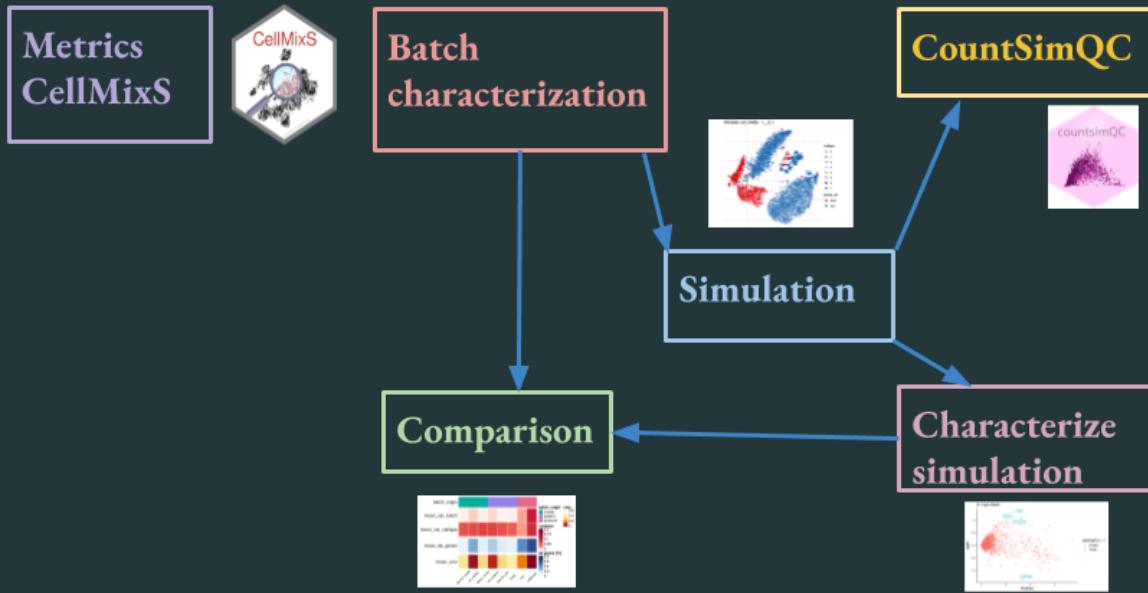
# Project overview



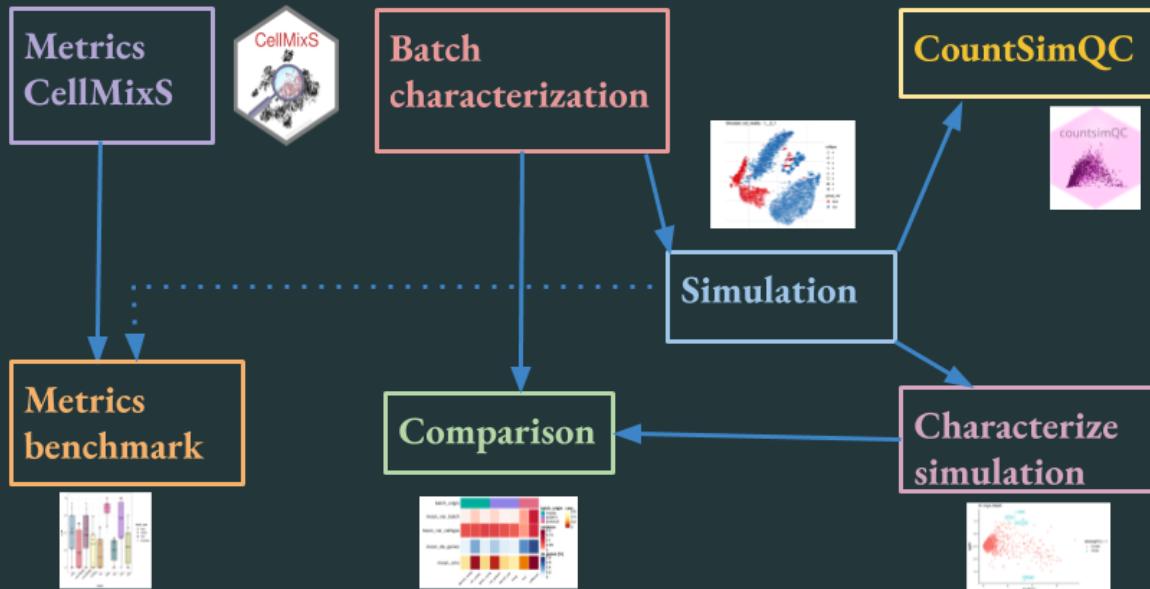
# Project overview



# Project overview



# Project overview



Batch characterization

# Batch characterization

- 7 datasets
- 9 batch effects
- *patient, protocol, storage*
- **Variance partitioning**
- **logFC** distribution and correlation
- DE genes/overlap

## Variance partition

$$Y_g = \mu + X_p \alpha_{pg} + X_b \beta_{bg} + X_{p:b} \gamma_{(p:b)g} + \epsilon_g$$

$Y_g$ : normalized expression gene g

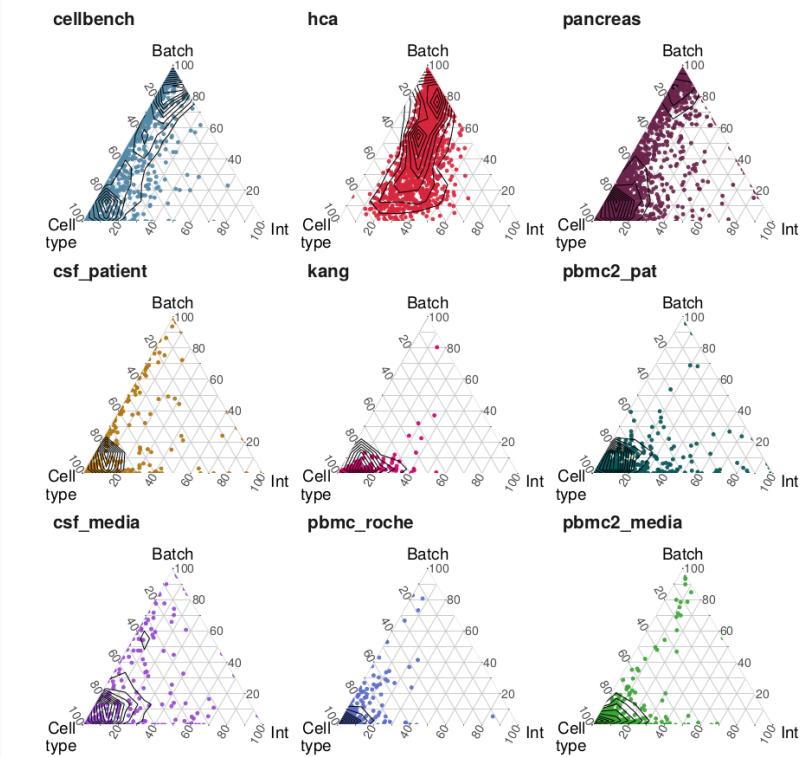
$\mu$ : baseline expression

$X_p, X_b, X_{p:b}$ : design matrices for the (random) cell types, batches and interactions effect

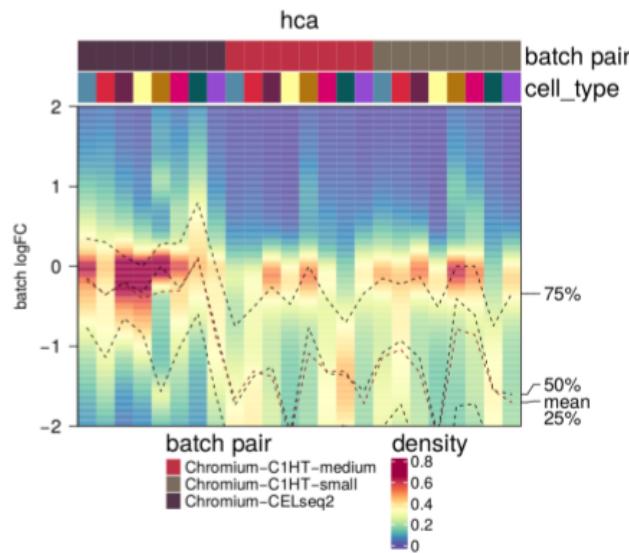
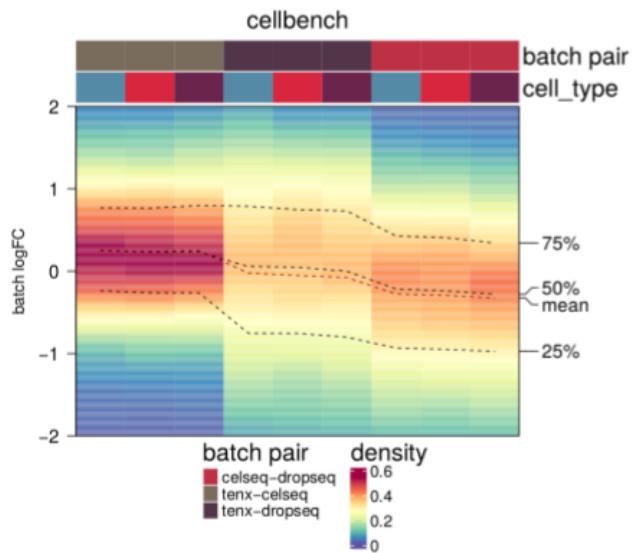
$\alpha_{pg} \sim N(0, \sigma_{pg}^2)$ , ...: corresponding random effects

$\epsilon_g \sim N(0, \sigma_g^2)$  corresponding errors

# Percent variance explained by ...:



# Log fold change distributions



# Metrics

# Metrics

## Cell level:

- Cellspecific Mixing Score (**cms**)
- Local Inverse Simpson Index (**lisi**)
- Mixing Metric (**mm**)
- Shannon's **entropy**

# Metrics

## Cell level:

- Cellspecific Mixing Score (**cms**)
- Local Inverse Simpson Index (**lisi**)
- Mixing Metric (**mm**)
- Shannon's **entropy**

## Celltype level:

- k-nearest neighbour Batch Effect Test (**kBet**)
- Average silhouette width (**asw**)
- Graph connectivity (**graph**)

# Metrics

## Cell level:

- Cellspecific Mixing Score (**cms**)
- Local Inverse Simpson Index (**lisI**)
- Mixing Metric (**mm**)
- Shannon's **entropy**

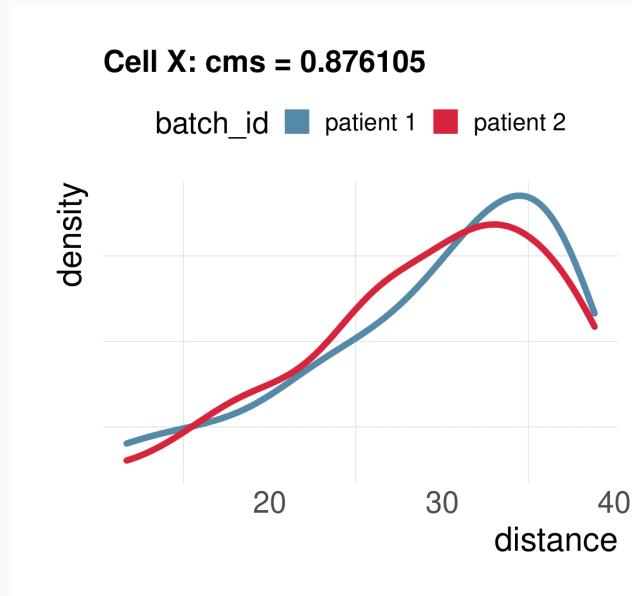
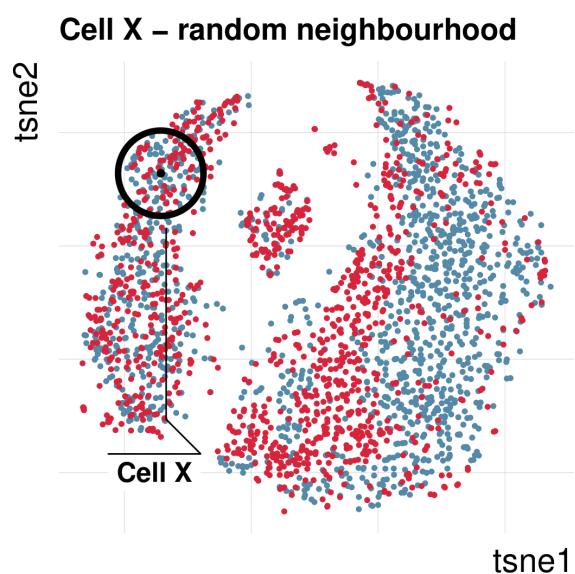
## Global level:

- Principal component regression (**pcr**)

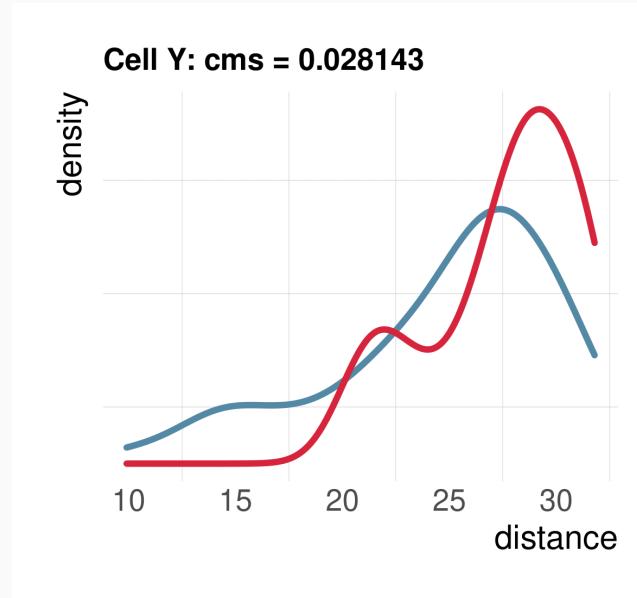
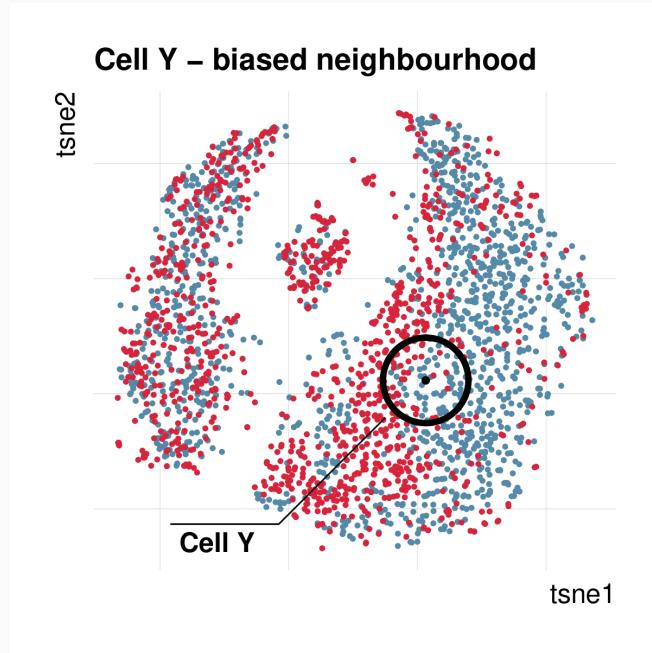
## Celltype level:

- k-nearest neighbour Batch Effect Test (**kBet**)
- Average silhouette width (**asw**)
- Graph connectivity (**graph**)

# Cell-specific Mixing Score (cms)



# Cell-specific Mixing Score (cms)



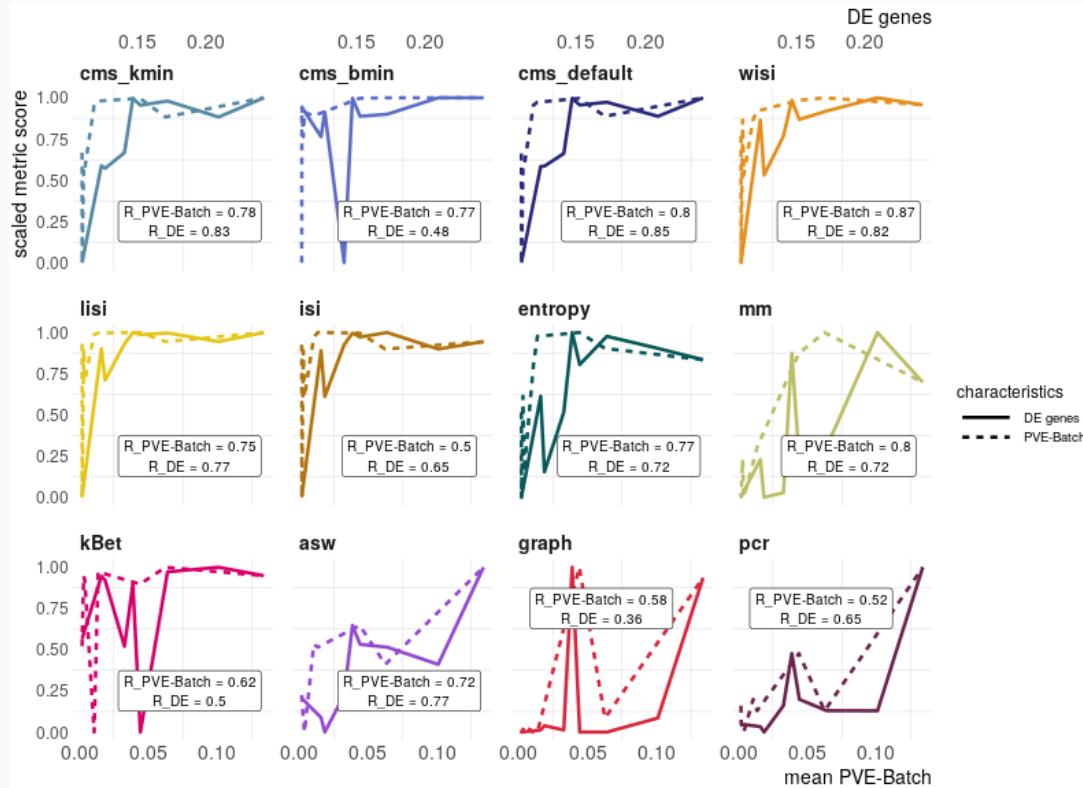
Benchmark

## Task 1: Batch characteristics

*Aim: Test whether metrics reflect batch strength across datasets*

Spearman correlation of metrics with surrogates of batch strength (e.g., PVE-Batch and proportion of DE genes between batches) across datasets

# Batch characteristics



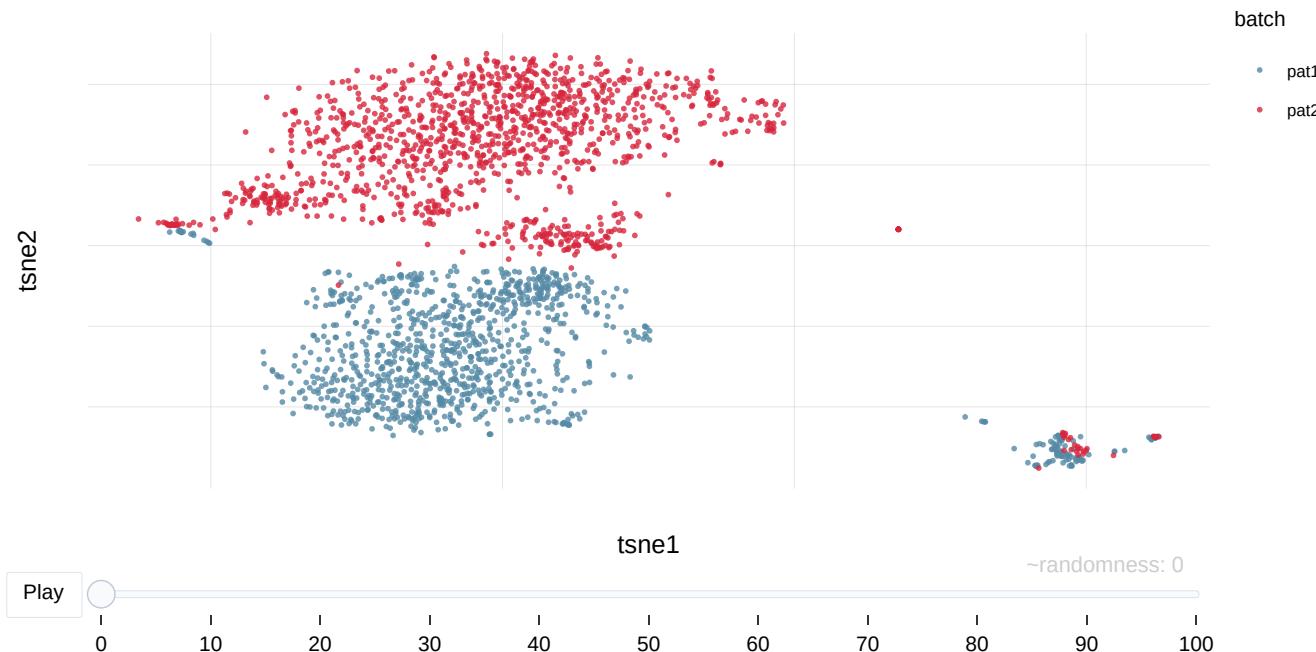
## Task 2: Batch label permutation

*Aim: Negative control and test whether metrics scale with randomness*

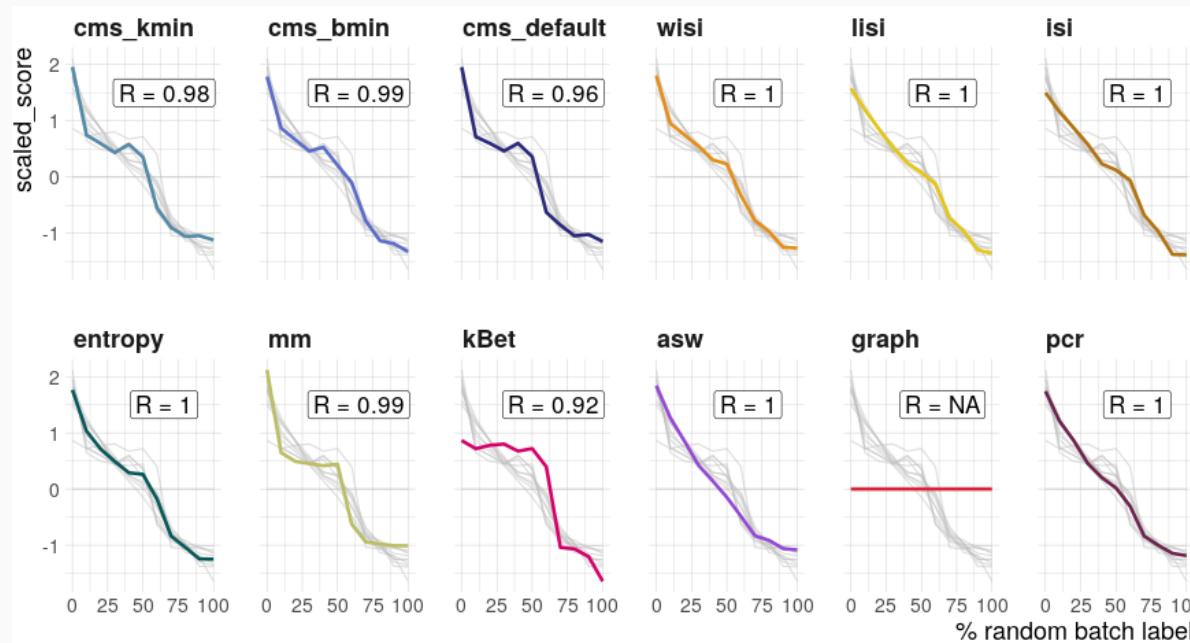
Spearman correlation of metrics with the percentage of randomly permuted batch label

# Batch label permutation

Batch label permutation



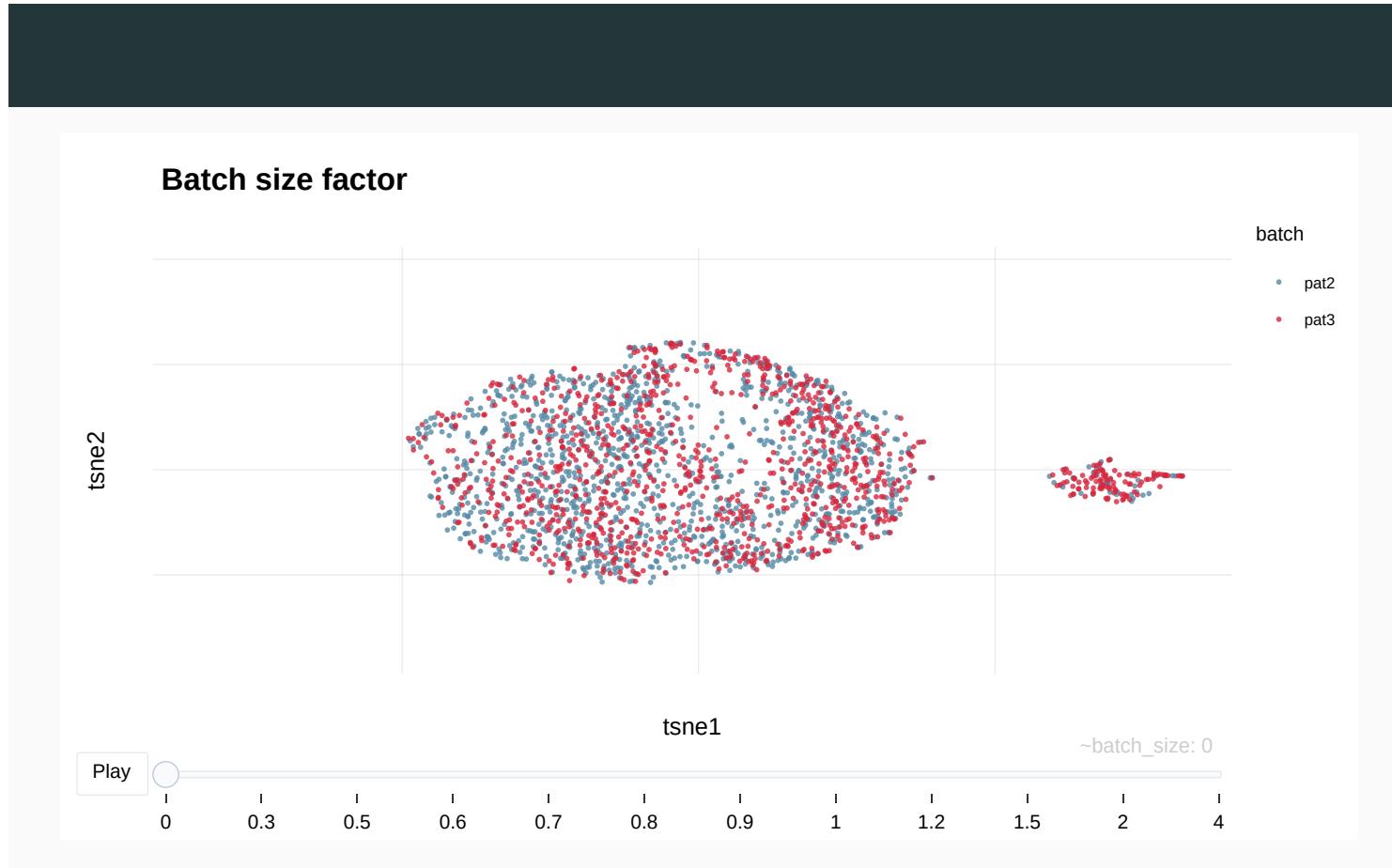
# Batch label permutation



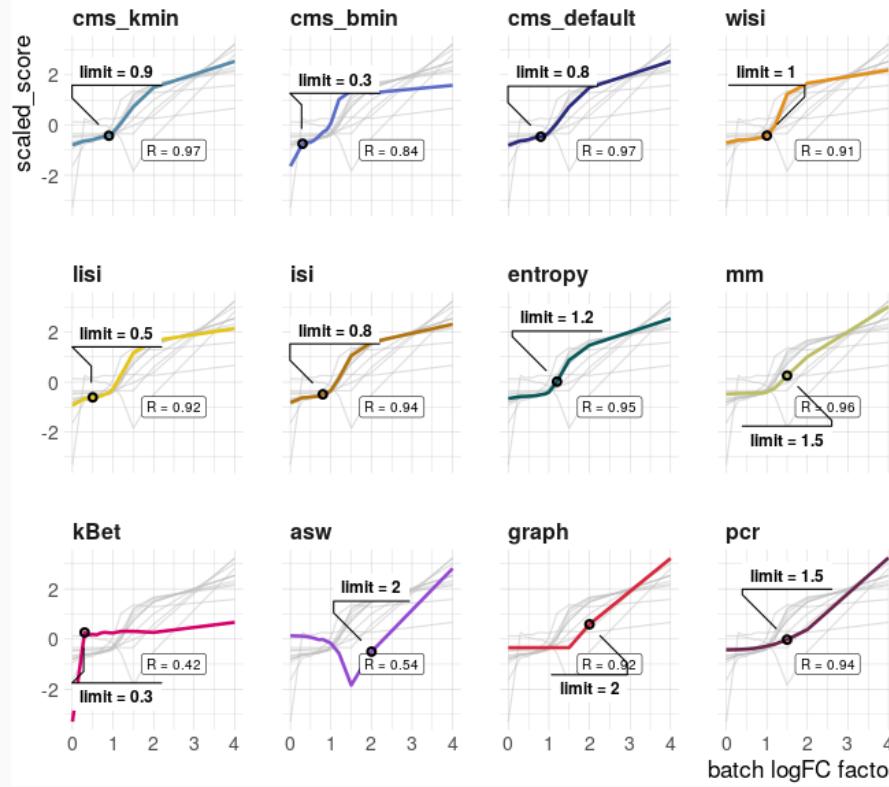
## Task 3: Scaling and detection limits

*Aim: Test whether metrics scale with (synthetic) batch strength; Estimate lower limit of batch detection*

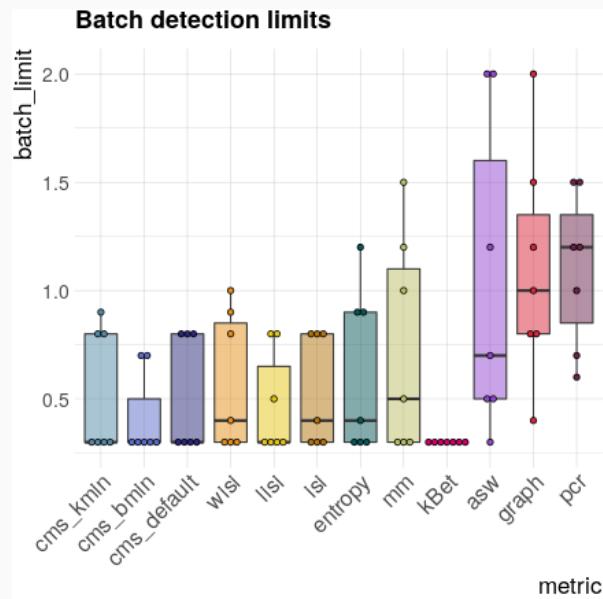
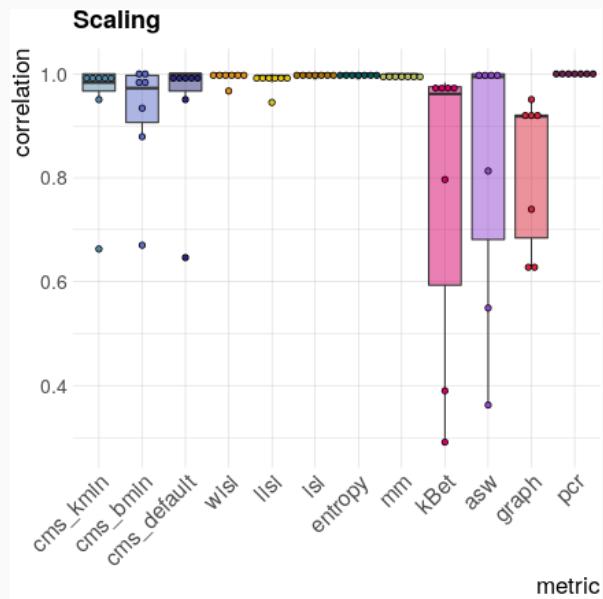
Spearman correlation of metrics with the batch logFC in simulation series on the same dataset; Minimal batch logFC that is recognized from the metrics as batch effect



# Scaling and Sensitivity



# Scaling and Sensitivity

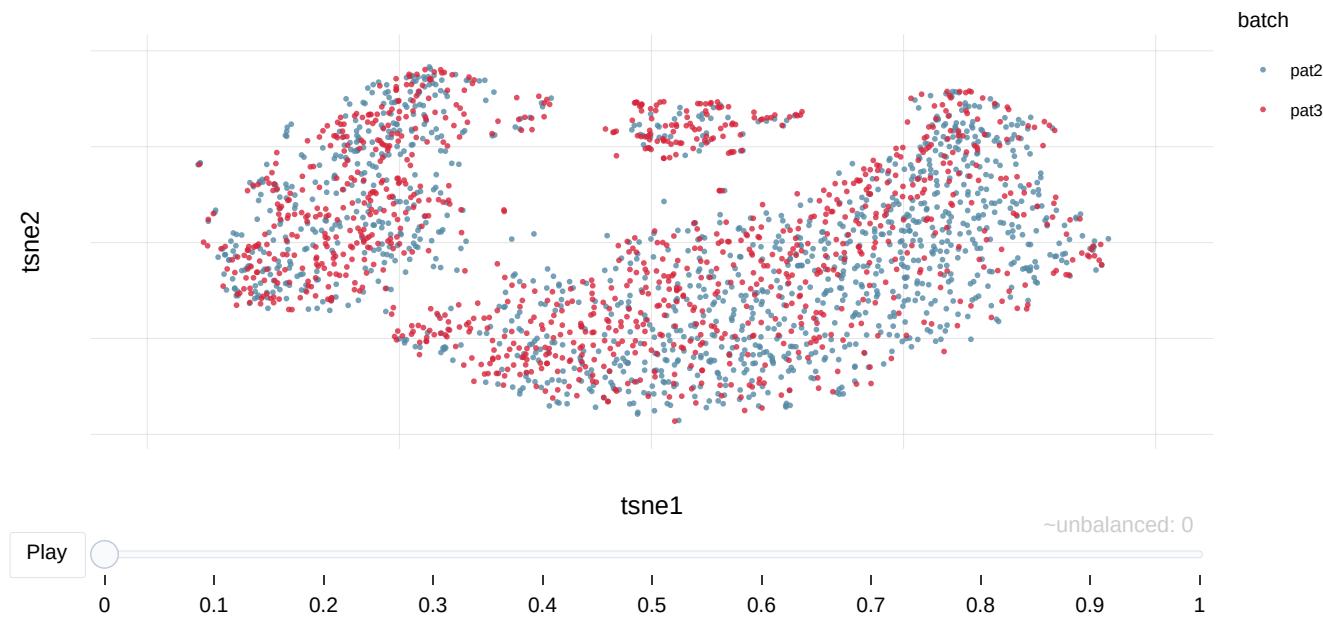


## Task 4: Imbalanced batches

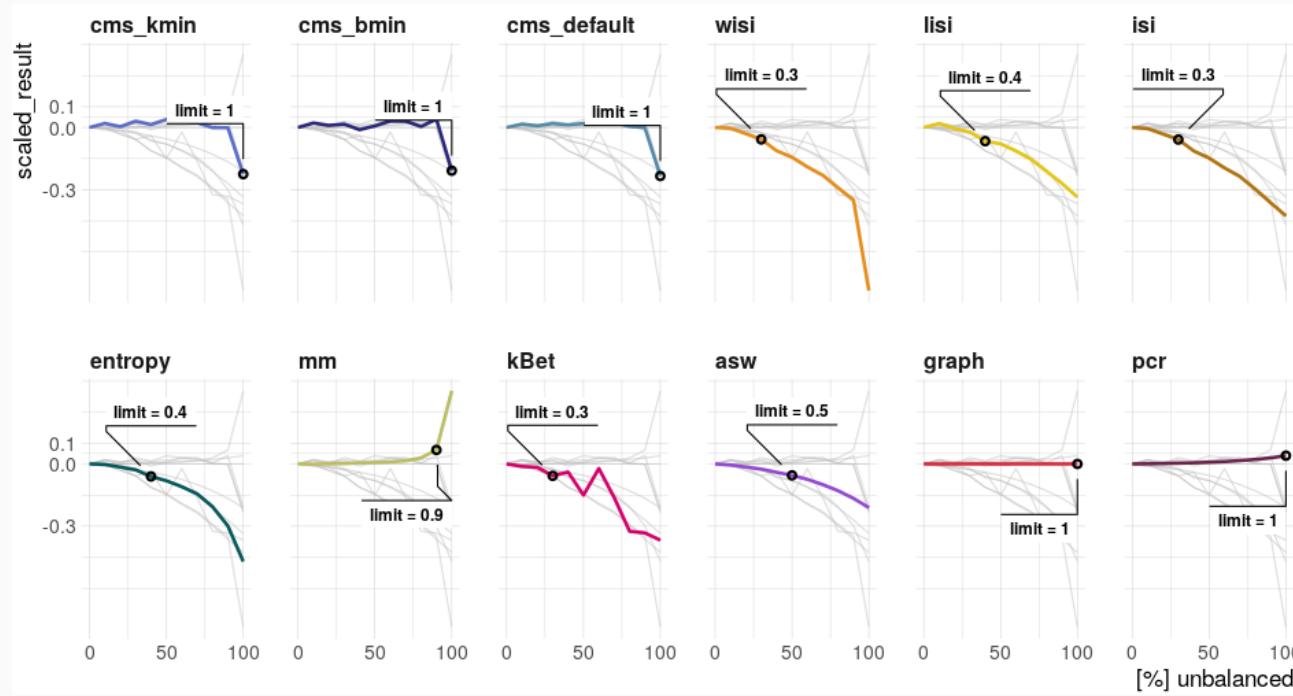
*Aim: Reaction of metrics to imbalance cell type abundance within the same dataset*

Test sensitivity towards imbalance of cell type abundance

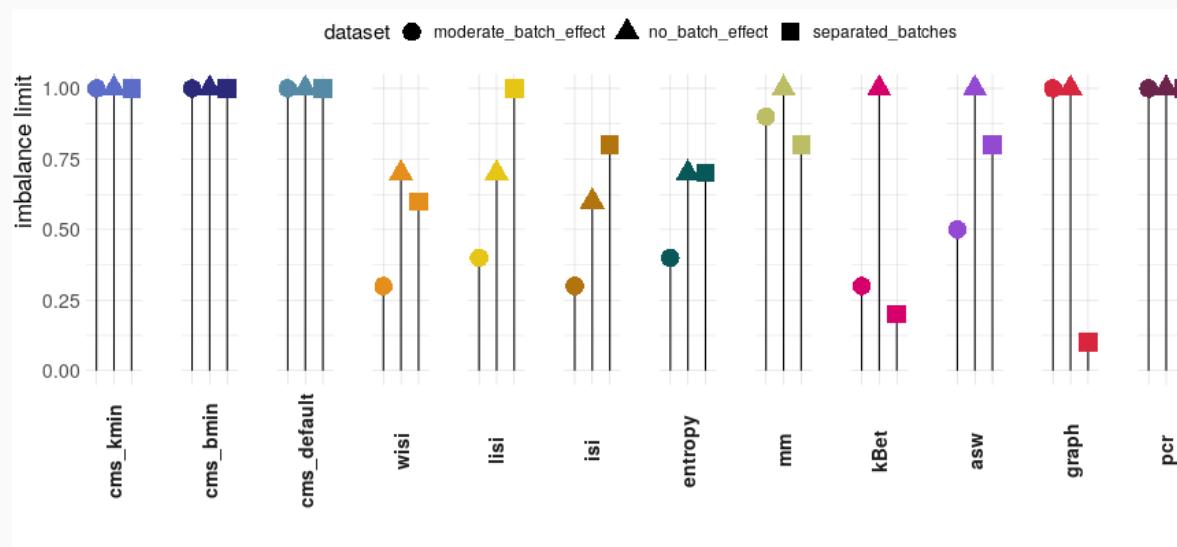
## Imbalanced batch effects



# Imbalanced batch effects

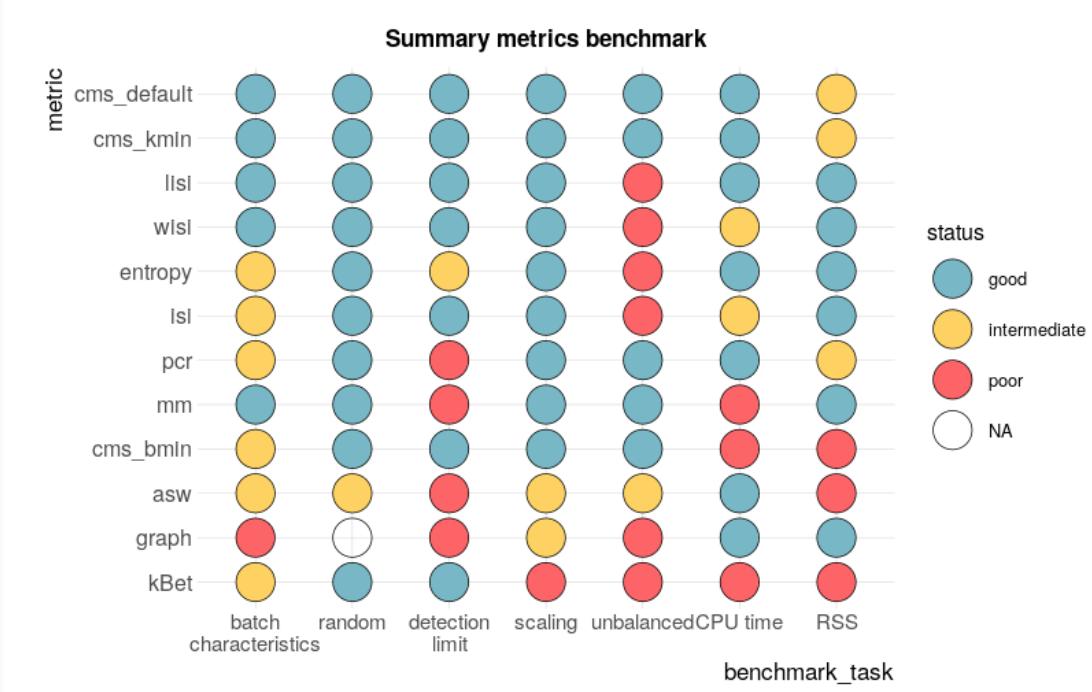


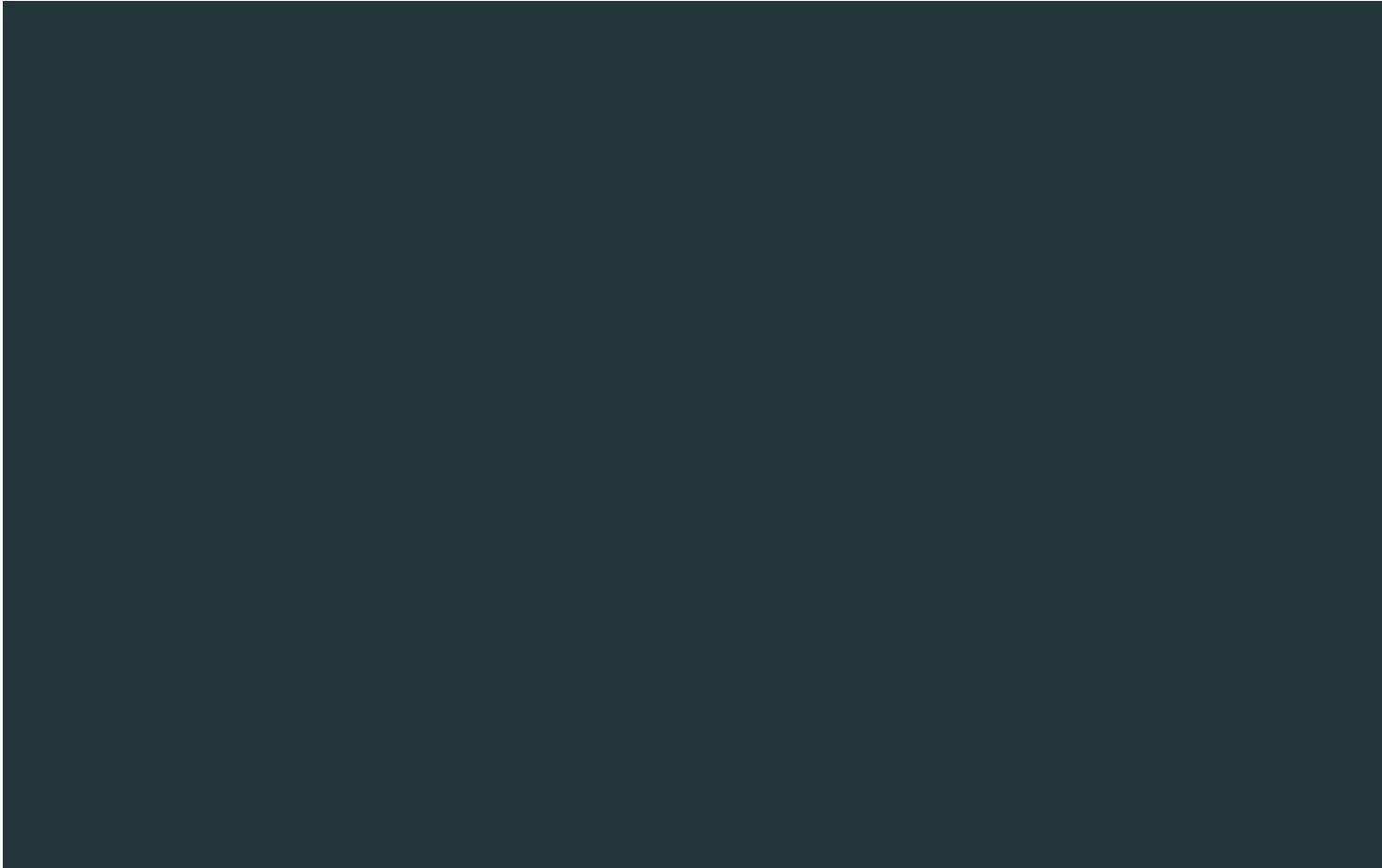
# Imbalanced batch effects



# Summary

# Summary





2. *Omni\_batch*: open and continuous  
benchmarking of single cell batch correction  
methods

# State of the art

Research | Open Access | Published: 16 January 2020

## A benchmark of batch-effect correction methods for single-cell RNA sequencing data

Hoa Thi Nhu Tran, Kok Siong Ang, Marion Chevrier, Xiaomeng Zhang, Nicole Yee Shin Lee, Michelle Goh & Jinmiao Chen 

*Genome Biology* 21, Article number: 12 (2020) | [Cite this article](#)

## Flexible comparison of batch correction methods for single-cell RNA-seq using BatchBench

Ruben Chazarra-Gil , Stijn van Dongen, Vladimir Yu Kiselev , Martin Hemberg 

*Nucleic Acids Research*, gkab004, <https://doi.org/10.1093/nar/gkab004>

Published: 01 February 2021 Article history ▾

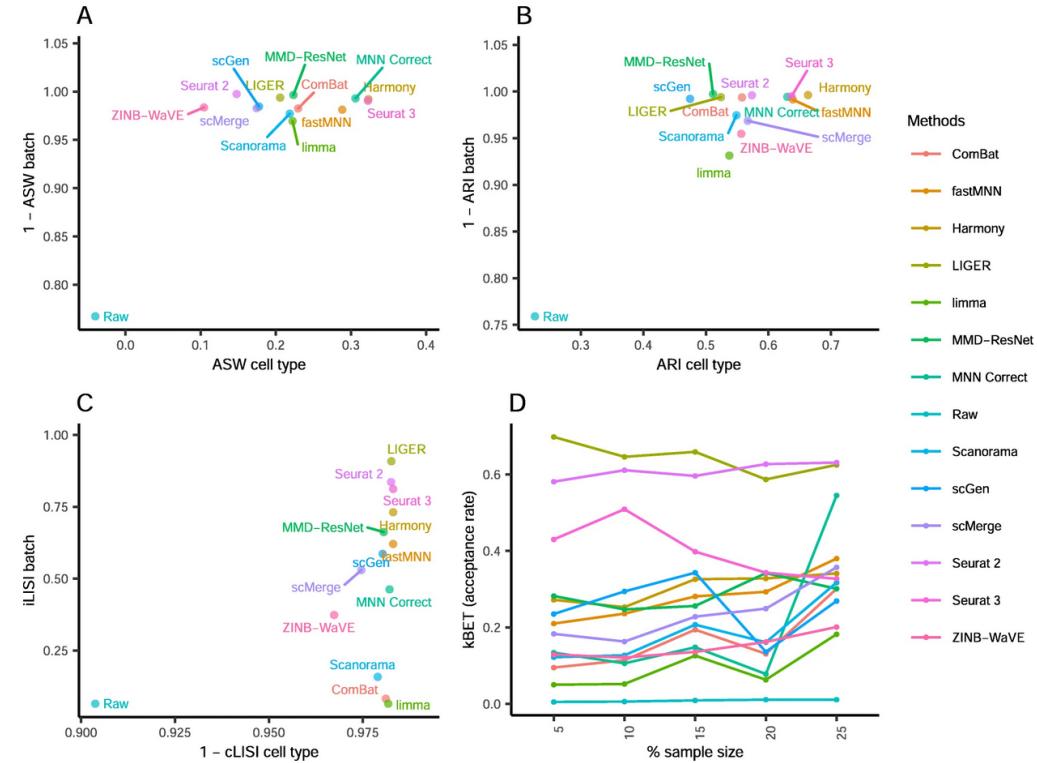
## Benchmarking atlas-level data integration in single-cell genomics

 MD Luecken,  M Büttner,  K Chaichoompu, A Danese,  M Interlandi,  MF Mueller,  DC Strobl,  L Zappia,  M Dugas,  M Colomé-Tatché,  FJ Theis

doi: <https://doi.org/10.1101/2020.05.22.111161>

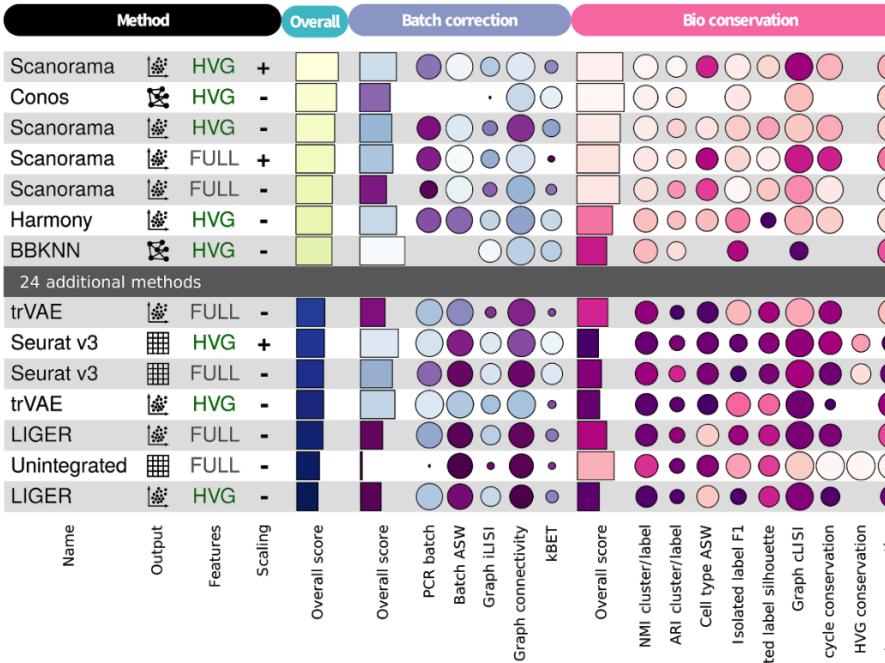
This article is a preprint and has not been certified by peer review [what does this mean?].

# Benchmark A



# Benchmark B

a



Output

- gene
- embed
- graph

Scaling

- + scaled
- unscaled

Ranking

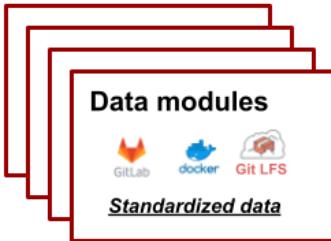


Score

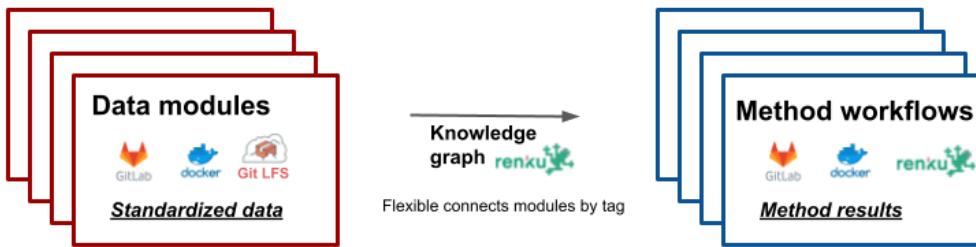


b

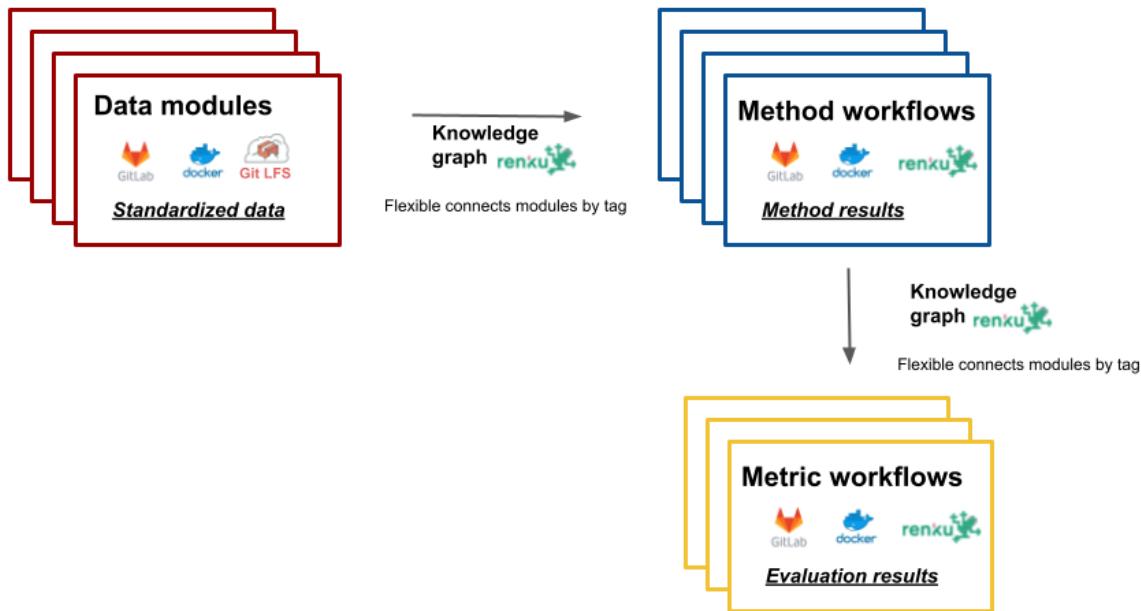
# Omni-benchmark



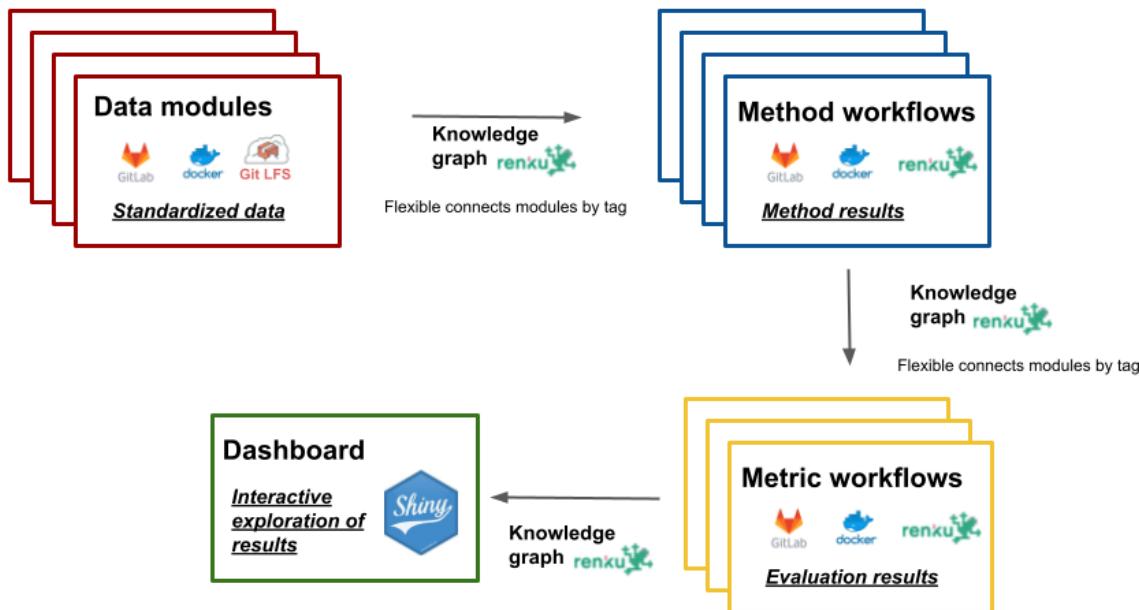
# Omni-benchmark



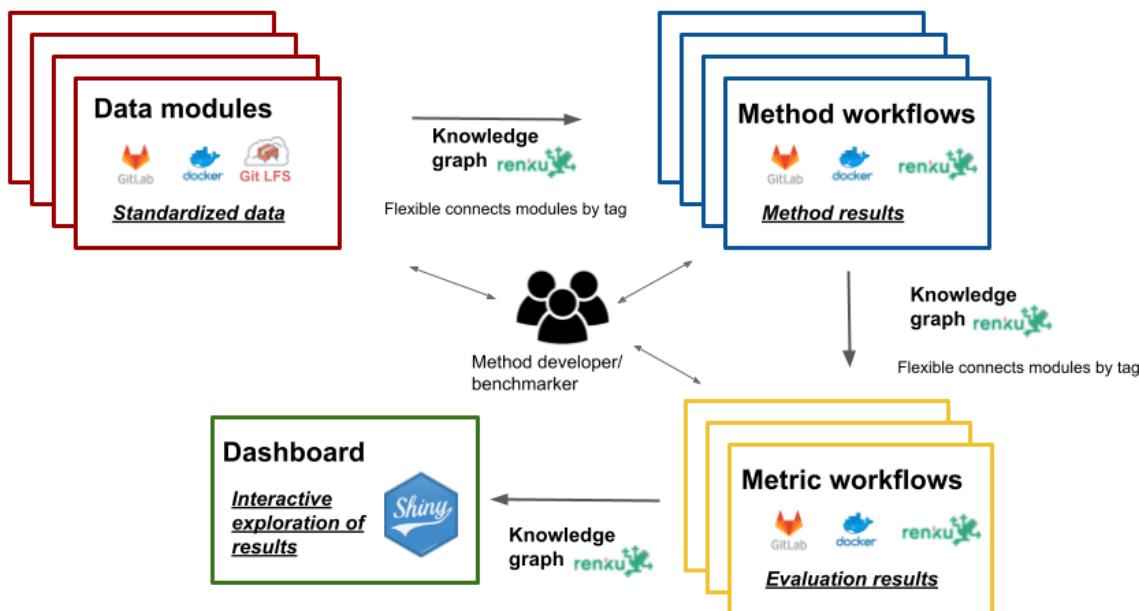
# Omni-benchmark



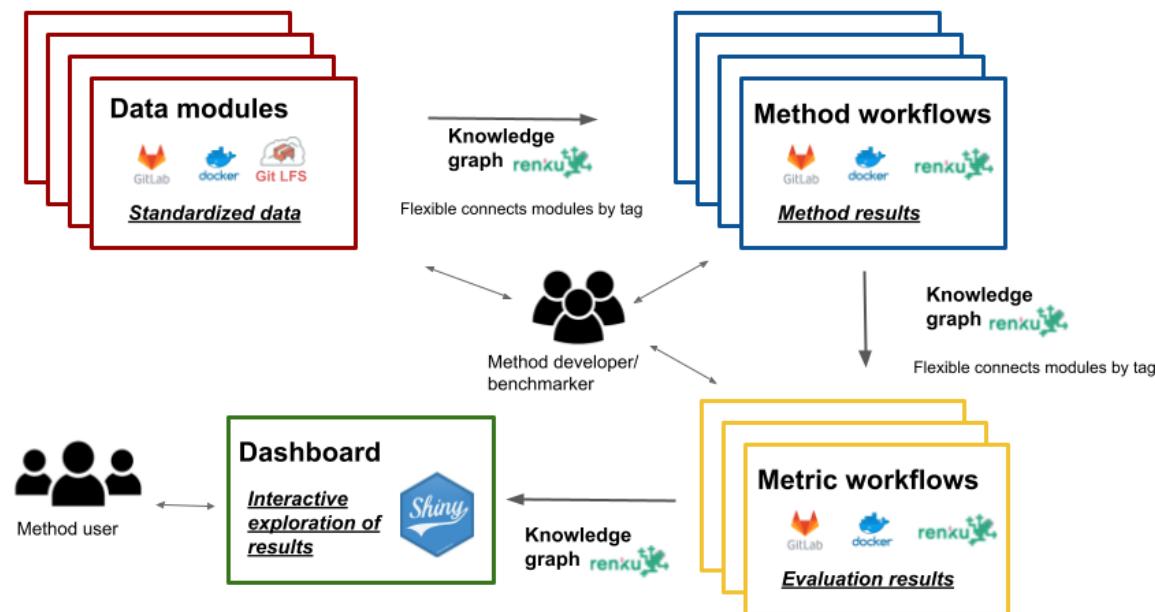
# Omni-benchmark



# Omni-benchmark



# Omni-benchmark



### **3. Teaching and course work**

## Teaching and course work

**Teaching:** BIO 134 (170 hours), R course 

**Course work:** 8/12 credits

**Missing:** Method seminar (1 credits)  
Transferable skills courses (2 credits)  
Pre-doctoral exam (1 credit)



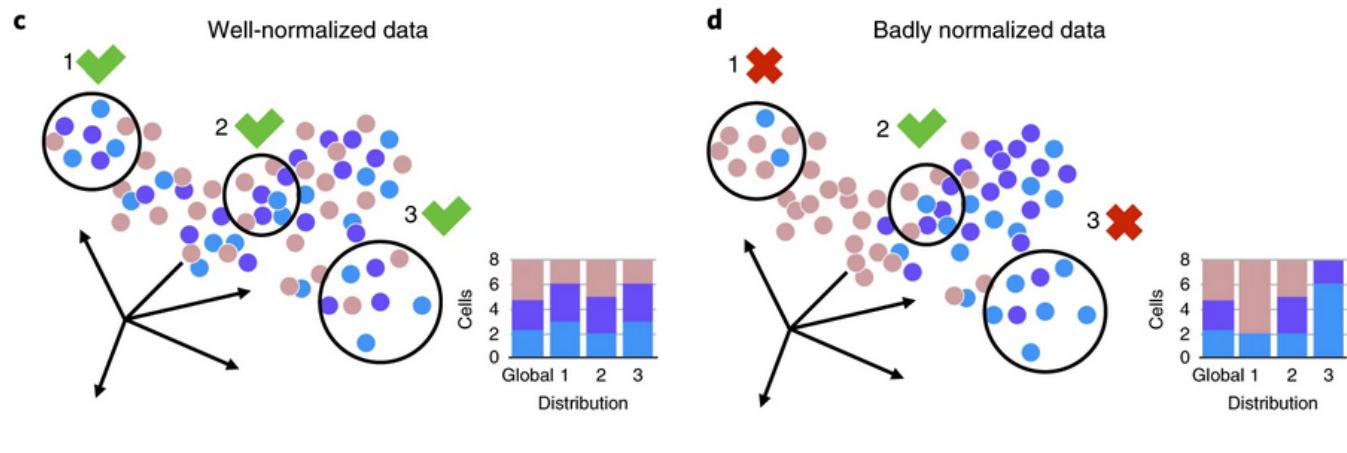
Extra slides

# Local Inverse Simpson Index (lisi)

- neighborhood diversity
- effective number of batches
- neighbor weighting:
  - euclidean distance --> wisi
  - no weighting --> isi
  - Gaussian kernel based weighting --> lisi

$$\frac{1}{\sum_{b=1}^B p(b)}$$

# k-nearest neighbour batch effect test (kBET)



(Buttner et. al., 2019)

# Mixing metric (mm)

