# Batch Effects

Technical, non-biological factors that affect variation in data

Mary O'Neill, Ph.D.
Director, Single Cell Genomics
ONeillMB@uw.edu
@ONeillMB1

Anh Vo, M.S.
Bioinformatician
athuyvo@uw.edu

# How to deal with batch effects

- Best way to avoid batch effects is not to introduce them in the first place!
- If unavoidable, it is very important that study design does not confound batch and other variables
  - Nothing can salvage poor study design
- Make sure you have a batch effect
  - Sometimes (often?) batch correction introduces more artifacts than they alleviate
- Apply methods thoughtfully
  - Don't blindly trust methods
  - Know what they are doing, what to use them for, and where they can lead you astray
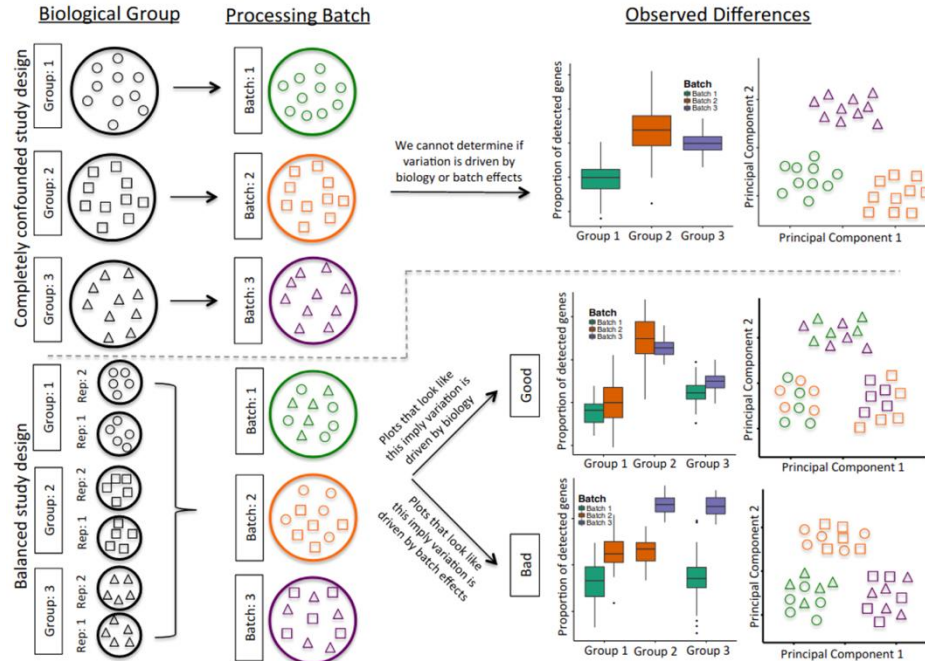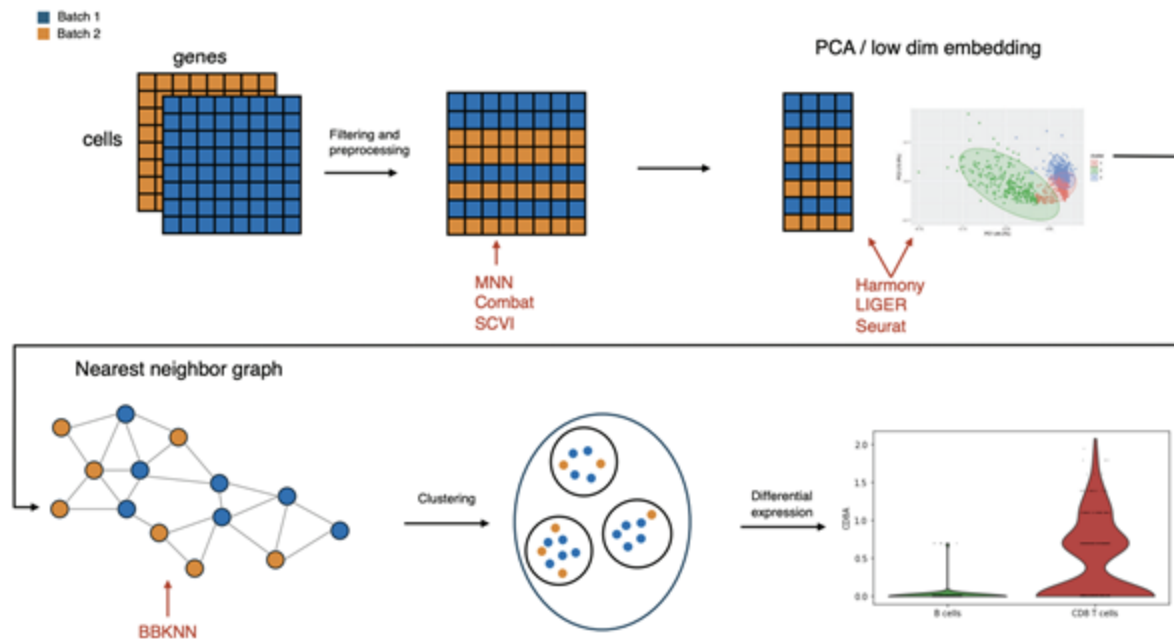
# How to deal with batch effects



Illustration of a confounded (top panels) and balanced (bottom panels) designs. Shapes denote different sample types (e.g. tissues or patients) and colours processing batches. In the confounded design it's impossible to disentangle biological variation from variation due to the processing batch. In the balanced design, by using tissue replicates and mixing them across batches, it is possible to distinguish between biological and batch-related variation. Figure from Hicks et al..
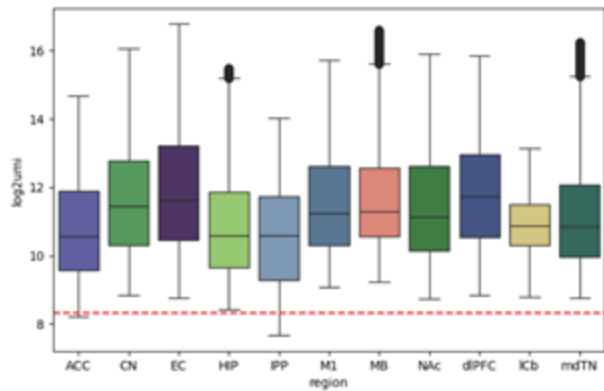
https://www.singlecellcourse.org/introduction-to-single-cell-rna-seq.html

# Batch Correction Methods

| | BBKNN | Combat | Harmony | LIGER | MNN | Seurat | SCVI |
|---|---|---|---|---|---|---|---|
| **Input** | KNN graph | Normalized count matrix | Normalized count matrix | Normalized count matrix | Normalized count matrix | Normalized count matrix | Raw count matrix |
| **Custom embedding** | None | None | Corrected embedding | Metagene / factor loadings | None | CCA | Learned lower dimensional latent space |
| **Correction object** | KNN graph | Count matrix | Embedding | Embedding | Count matrix | Embedding | Embedding |
| **Correction method** | Umap on merged neighborhood graph | Empirical bayes - linear correction method on the count values | Soft k-means - linear batch correction within small clusters in the embedded space | Quantile alignment of factor loadings | Mutual nearest neighbors - linear correction | Aligning canonical basis vectors to correct the embedding - lift the correction of the embedding to count space | Variational autoencoder - models the batch effect in a low dimensional space using a deep learning model, a new count matrix is imputed from the model. |
| **Returns** | Corrected KNN graph | Corrected count matrix | Corrected embedding | Corrected embedding | Corrected count matrix | Corrected count matrix | Corrected count matrix and corrected embedding |
| **Changes Count matrix** | No | Yes | No | No | Yes | Yes | Yes / Imputes new values |

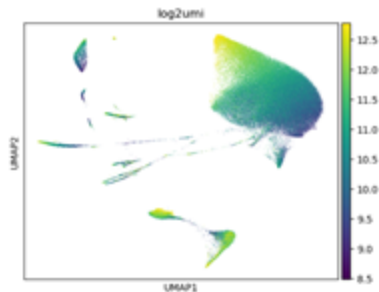Antonsson & Melsted, 2024 (BioXriv)

# Batch Correction Methods



Antonsson & Melsted, 2024 (BioXriv)

# Sidenote: we see minimal batch effects in sci-RNA-seq!



**Meidan UMI = 1,918 (~ 6x times higher than BICCN)**

Wei Yang, Unpublished

# Lessons from processing >1M PBMCs



Central African (80)
West European (80)
East Asian (62)

PBMCs

SARS-CoV-2
IAV

scRNA-seq

>10⁶ high-quality single-cell transcriptomes

AIMS

1. Characterize variability of the immune response to SARS-CoV-2 across human populations at single cell resolution (scRNA-seq)

2. Map genetic bases of immune variability in response to SARS-CoV-2 (eQTLs)

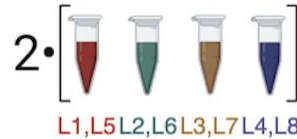3. Uncover natural selection and archaic introgression signals associated to virus-induced immune reponses

# Library design:



2 runs per week (16 runs total)
16 individuals/8 libraries
per experimental run

Each library contains 12 samples
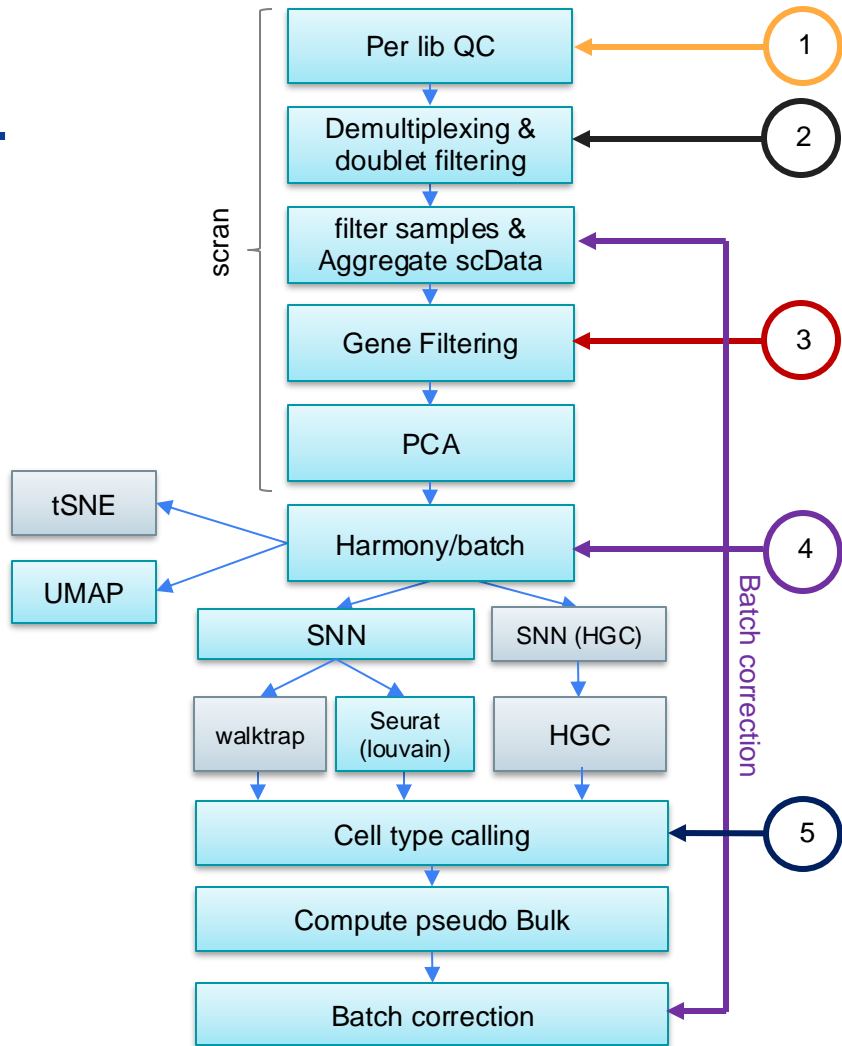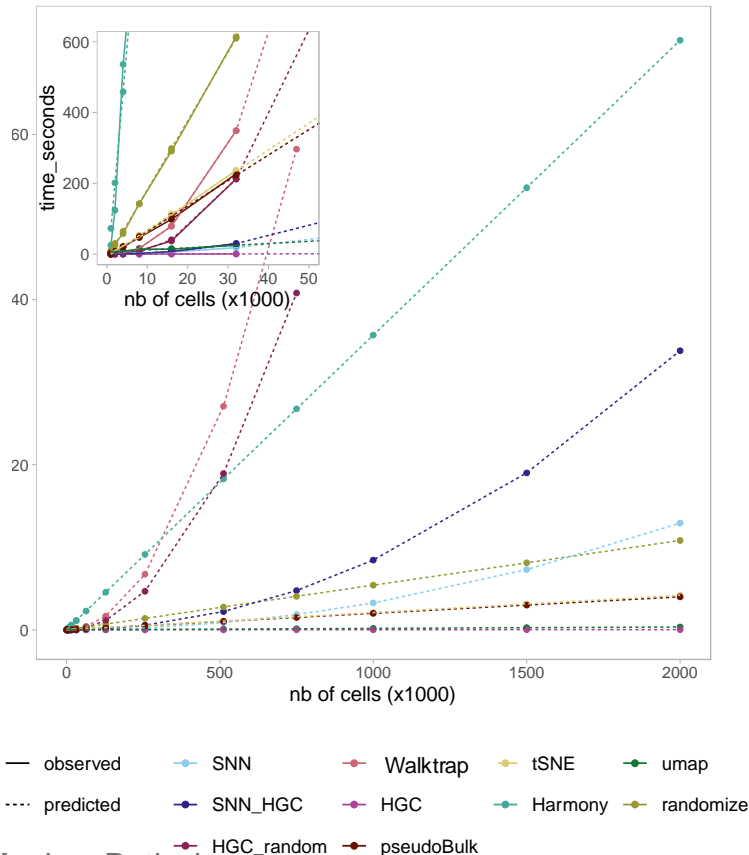(4 from each condition)

Each sample is done on two separate
libraries

target:
1,667 cells per
individual/condition

After QC:
~1500 cells on average
(median; IQR=558 cells)

# General pipeline overview



Slide from Maxime Rotival

# Size factor normalization & Batch effect removal:

Use `MultibatchNorm` to allow for differences in mean size factor across libraries
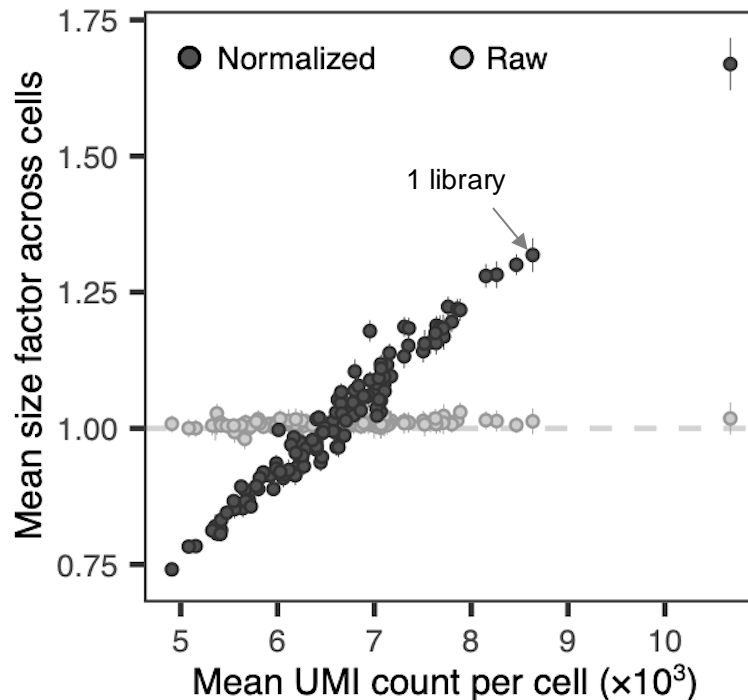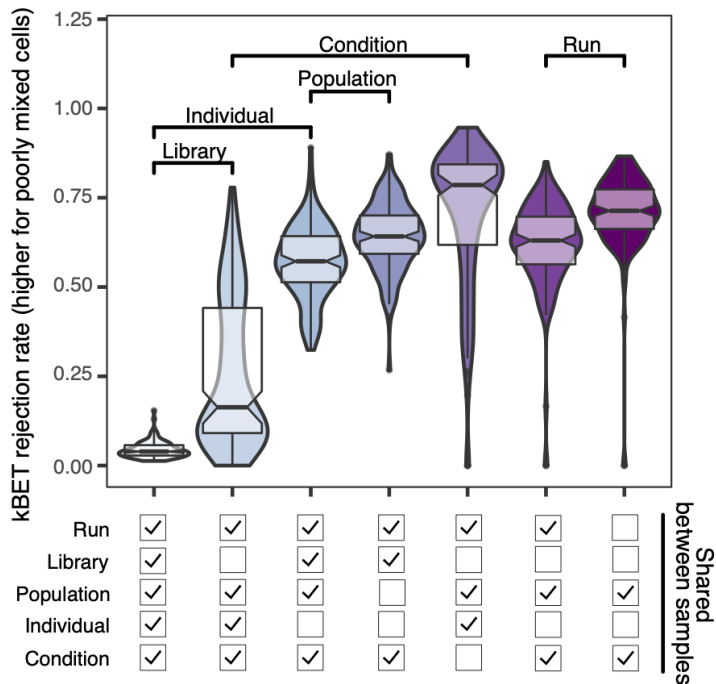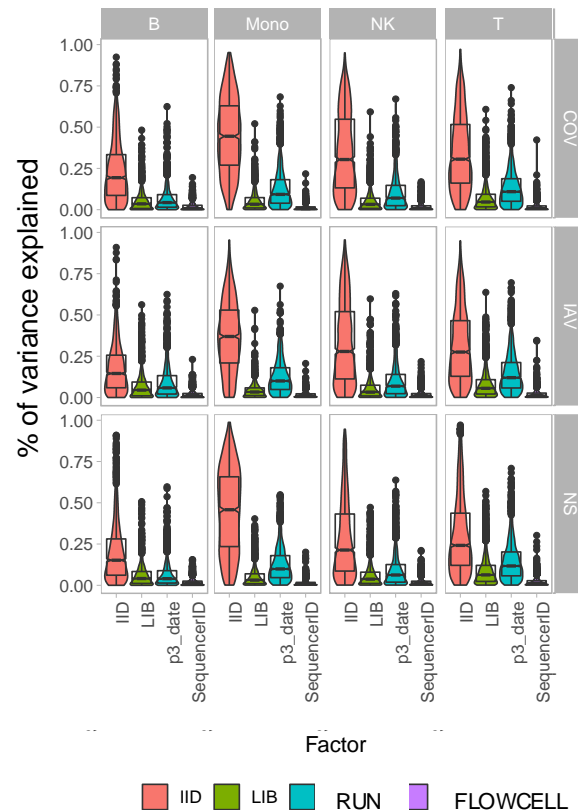
# Size factor normalization & Batch effect removal:

Use `MultibatchNorm` to allow for differences in mean size factor across libraries

Use `kBET` to estimate batch effects on cell mixing…

…and `Harmony` to correct for batch effects across experimental runs (for cell clustering purposes)

Aquino*, Bisiaux*, Li*, O'Neill*, *et al.* 2024 (Nature)

# Size factor normalization & Batch effect removal:

Use `MultibatchNorm` to allow for differences in mean size factor across libraries

Use `kBET` to estimate batch effects on cell mixing…

…and `Harmony` to correct for batch effects across experimental runs (for cell clustering purposes)

Use `linear mixed models (lme4)` to estimate & correct for batch effects at pseudobulk level
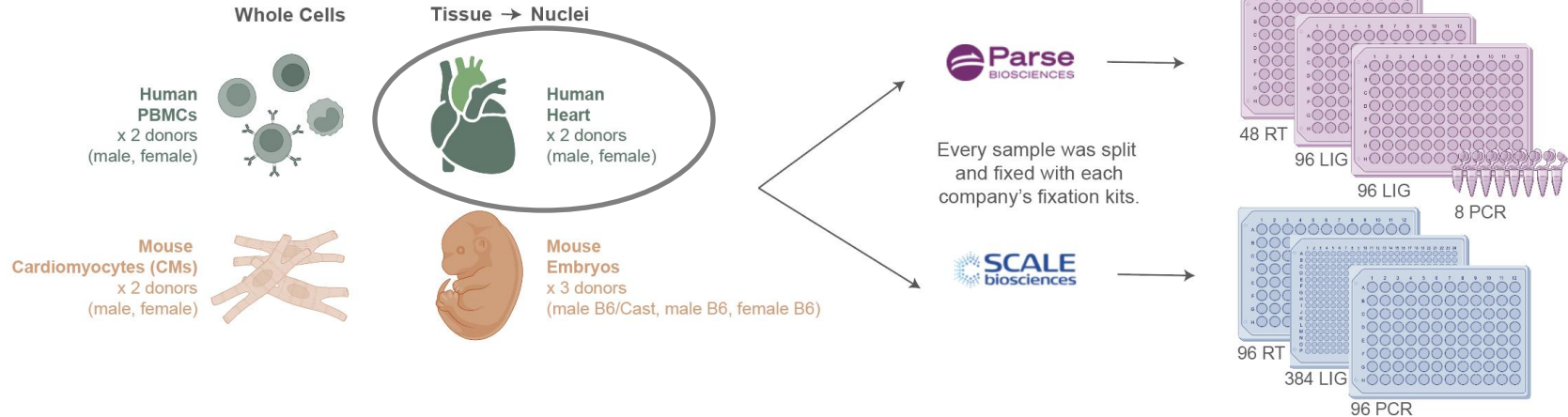(differential expression/ eQTL mapping)
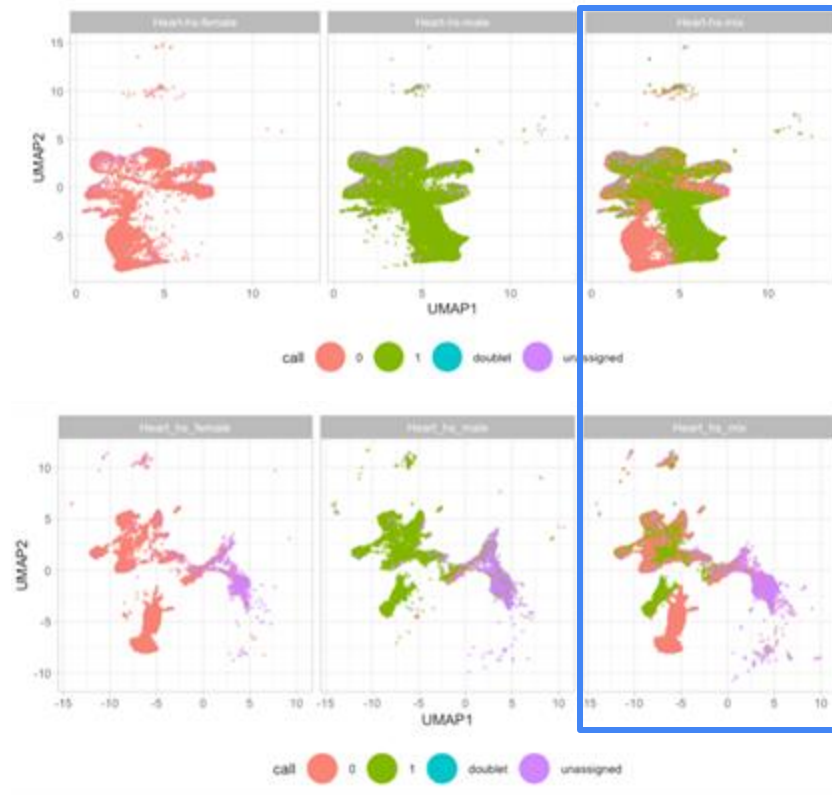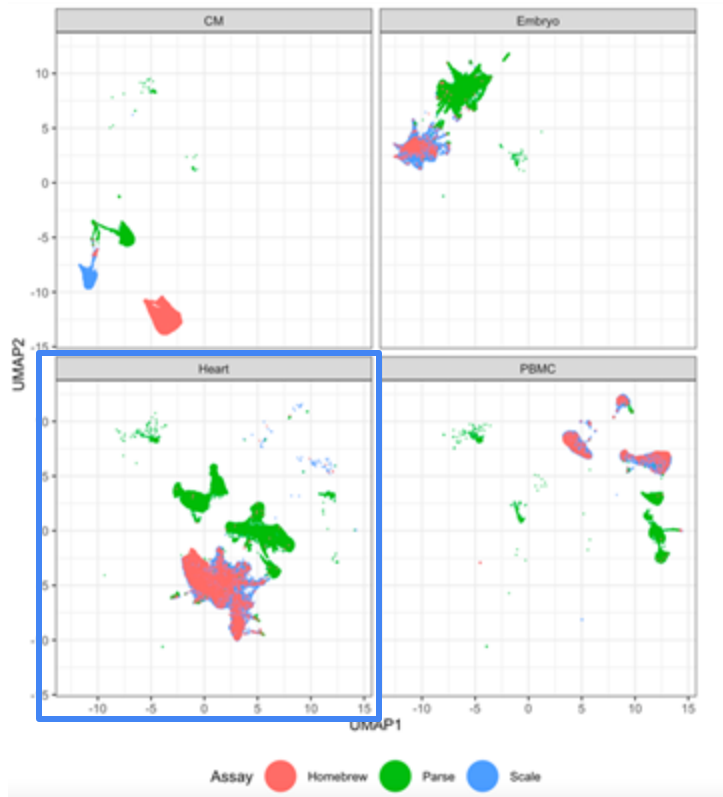
# Tutorial

Background on the dataset

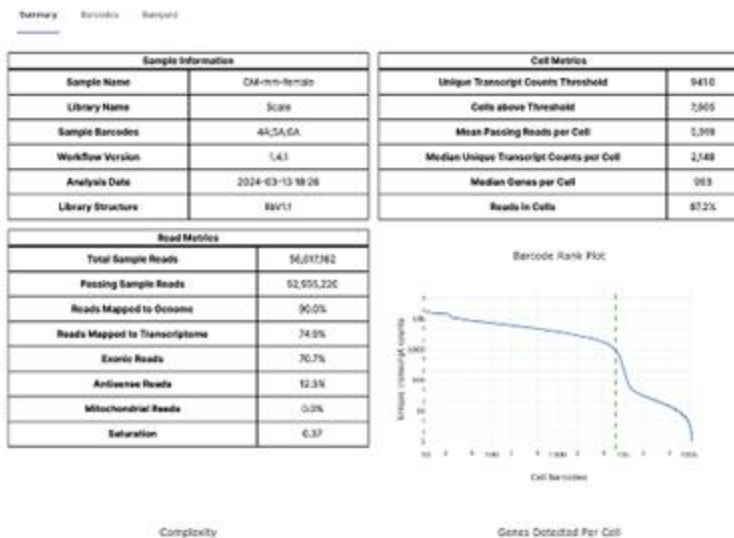# Combinatorial-indexing benchmarking (generated *in-house*)



Tutorial: random subsample of mixed sample of nuclei from two human hearts, processed with two different technologies

# Will play with human heart data today

# Each company has a pipeline



- ***Aggressively*** call cells – inflate UMIs?
- Both pipelines cut off almost an entire cell population (e.g. PBMCs) in barnyard samples
  - So, we start with every barcode >= 100 UMIs