

# Batch Effects

Technical, non-biological factors that affect variation in data

Mary O'Neill, Ph.D.

Director, Single Cell Genomics

[ONeillMB@uw.edu](mailto:ONeillMB@uw.edu)

@ONeillMB1

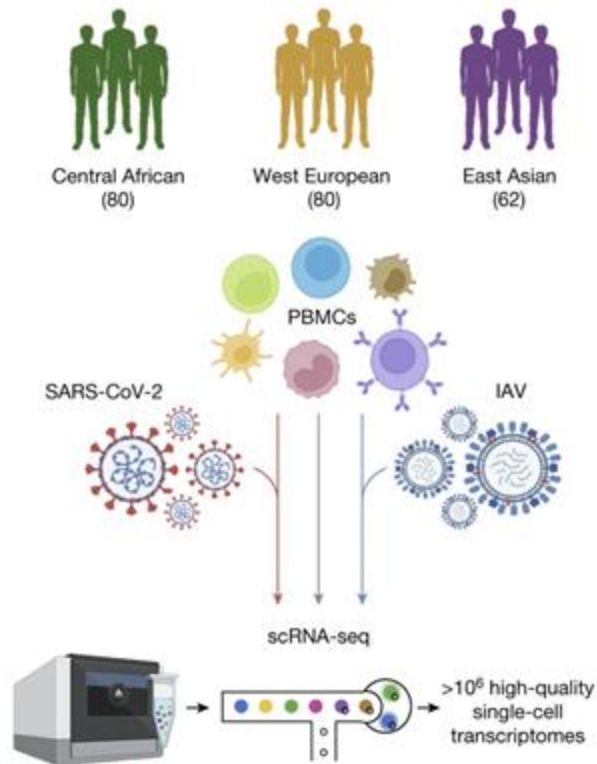
Anh Vo, M.S.

Bioinformatician

[athuyvo@uw.edu](mailto:athuyvo@uw.edu)

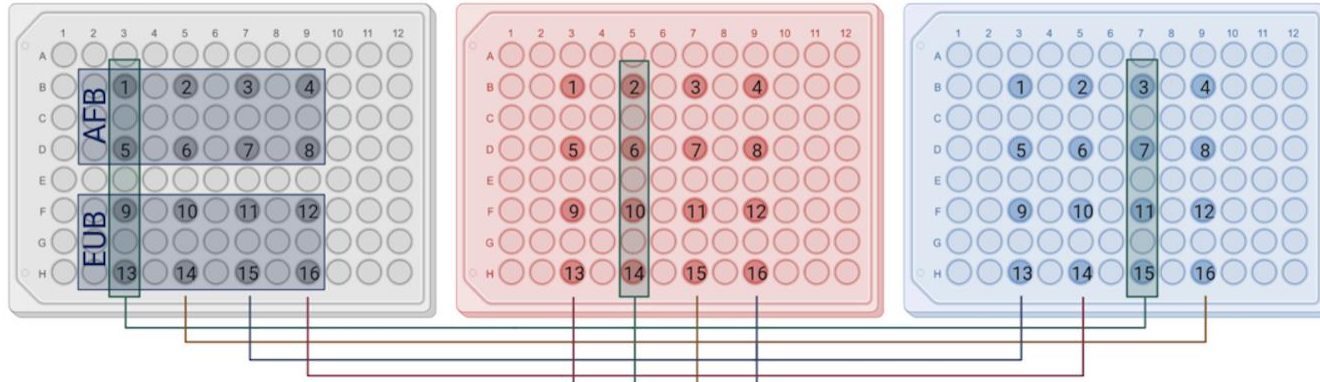


# Lessons from processing >1M PBMCs



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 responses

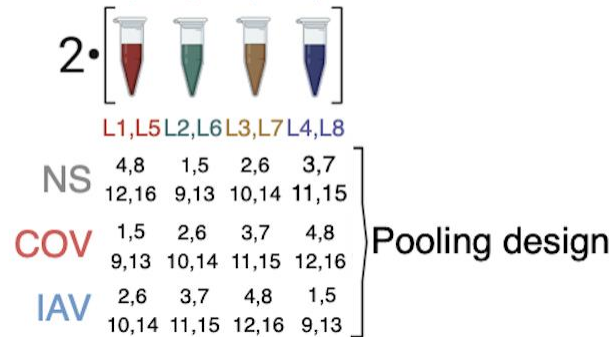
# Library design:



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

Each library contains 12 samples  
(4 from each condition)

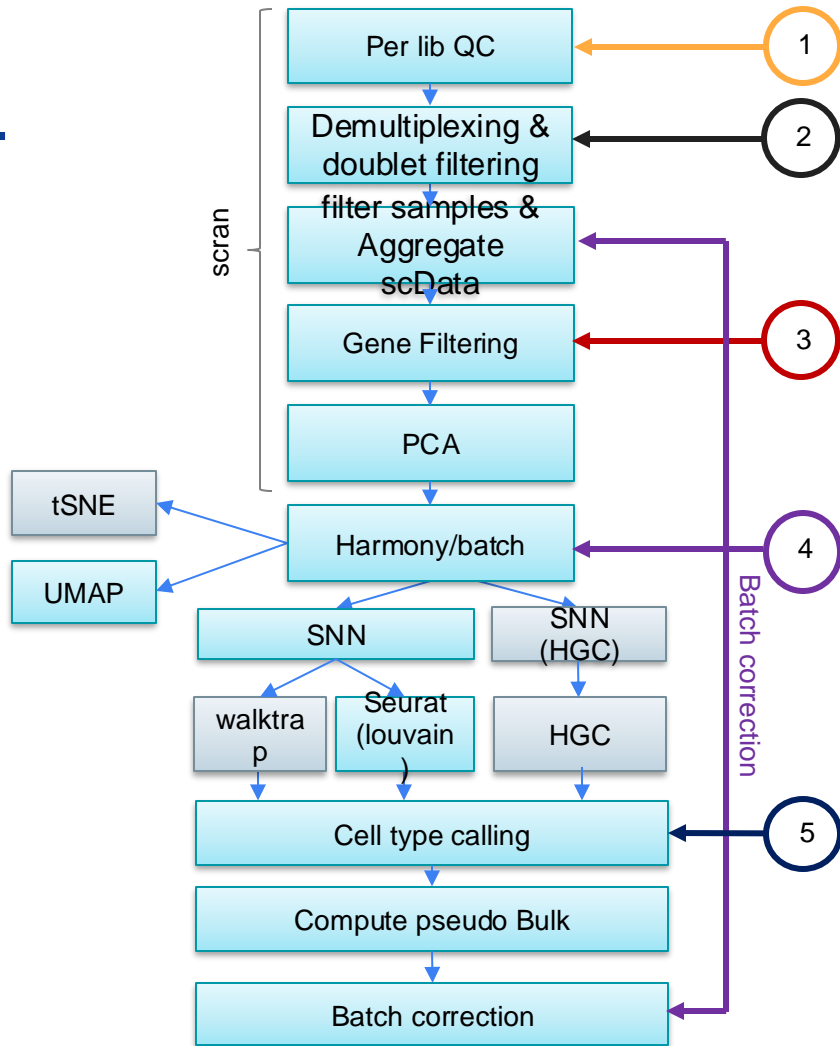
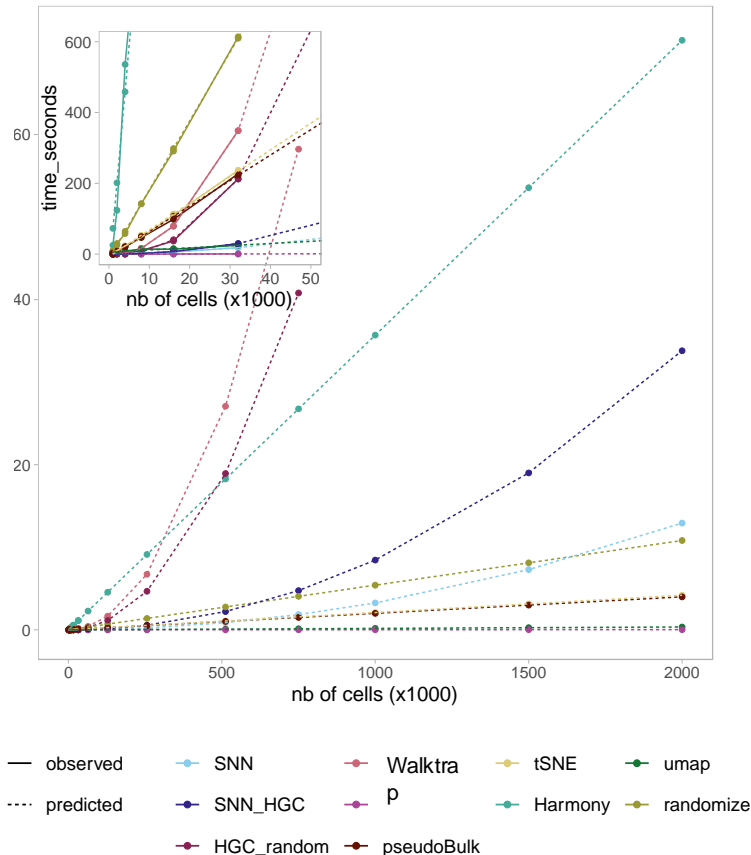
Each sample is done on two separate  
libraries



**target:**  
1,667 cells per  
individual/condition

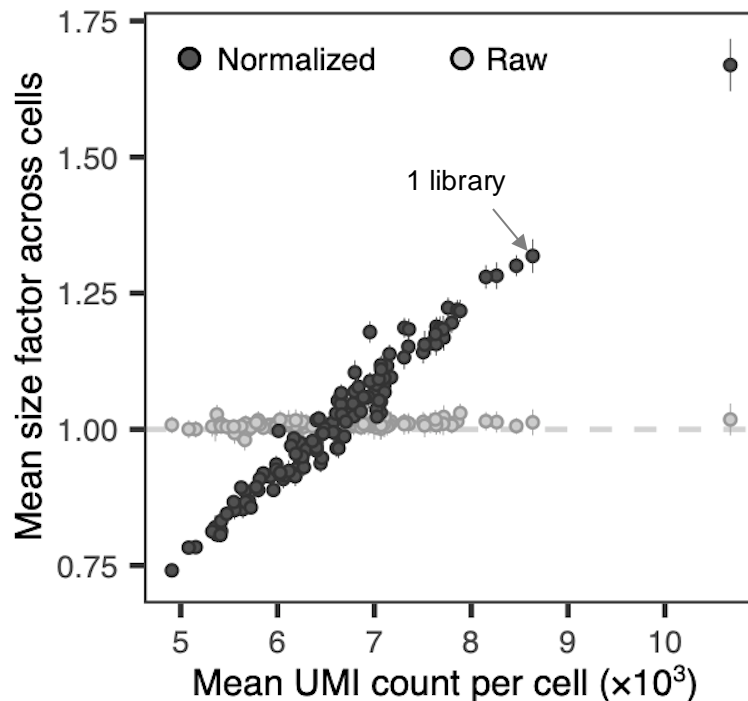
**After QC:**  
~1500 cells on average  
(median; IQR=558 cells)

# General pipeline overview



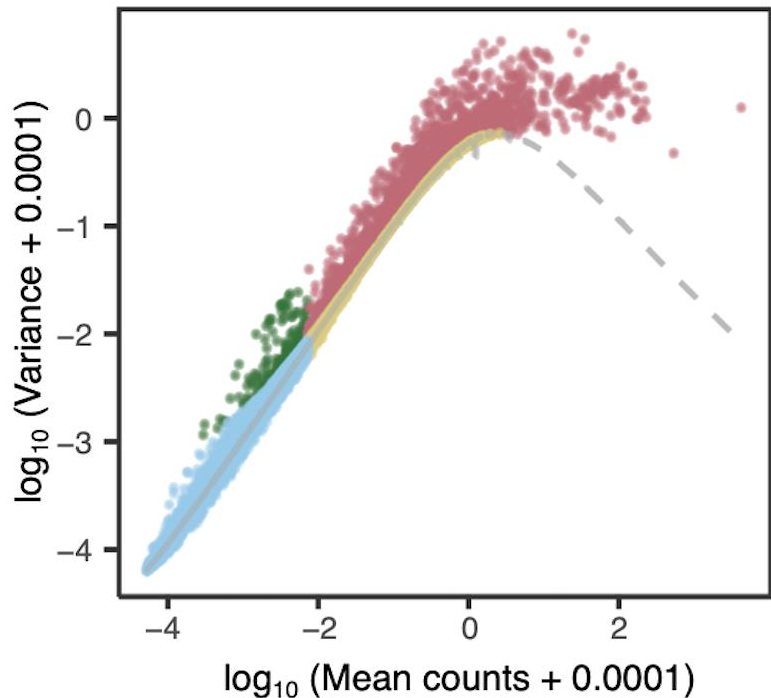
# Size factor normalization & Batch effect removal:

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

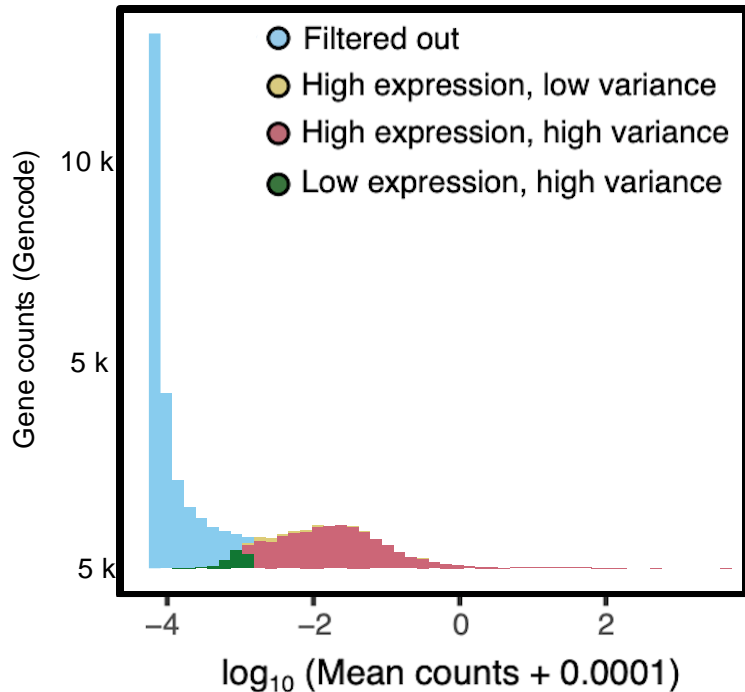


# Gene Filtering :

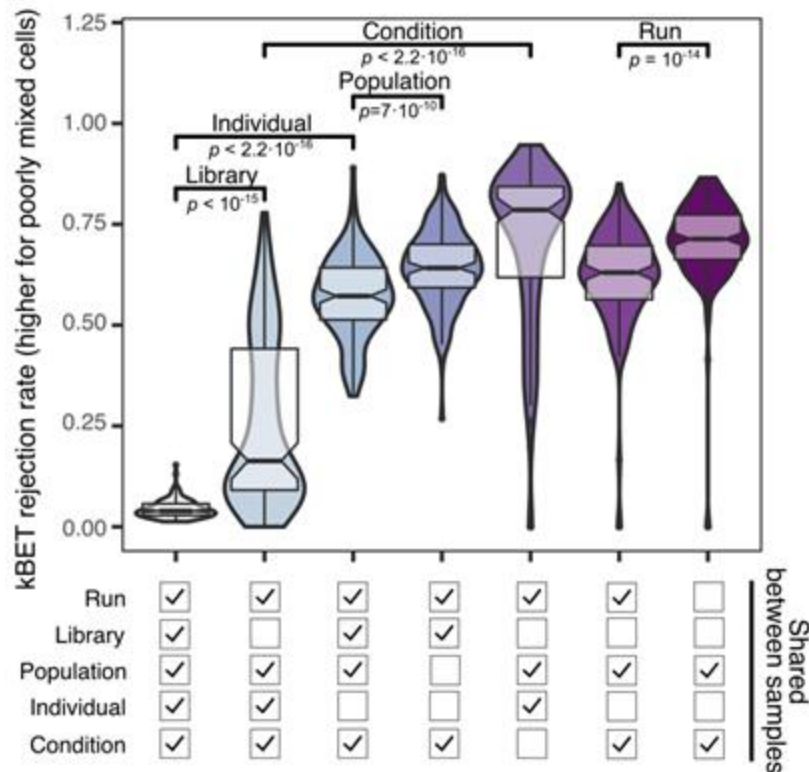
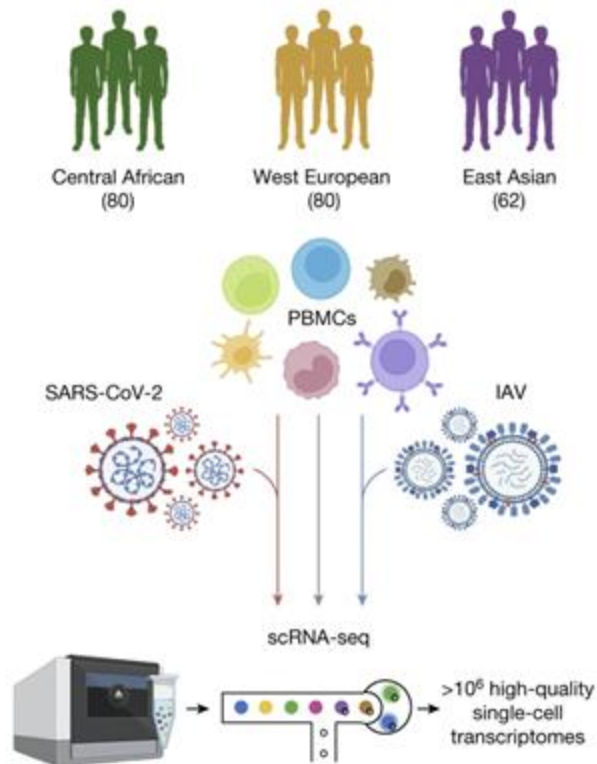
Use `modelGeneVarByPoisson` to  
estimate biological/technical variance



Keep genes with either high  
expression/high variance



# Sources of variation





# How to deal with batch effects

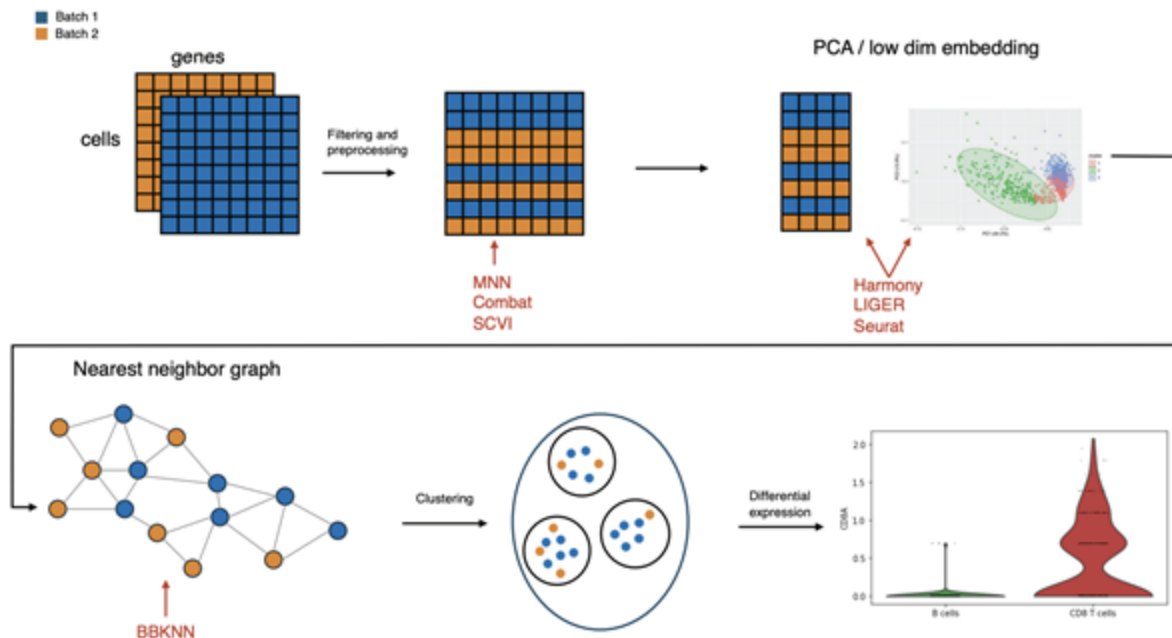
- Best way to avoid batch effects is not to introduce them in the first place!
- If unavoidable, it is very important that study design does not confound batch and other variables
  - Nothing can salvage poor study design
- Make sure you have a batch effect
  - Sometimes (often?) batch correction introduces 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



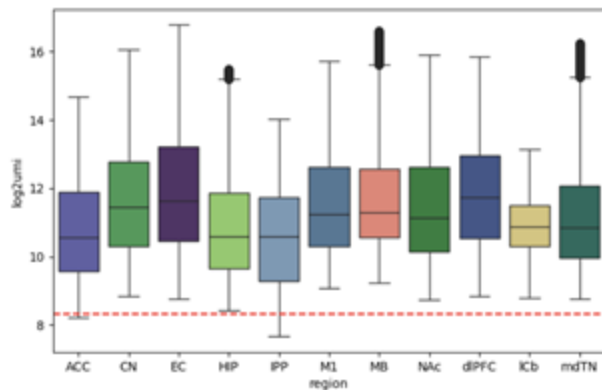
# Batch Correction Methods

	BBKNN	Combat	Harmony	LIGER	MNN	Seurat	SCVI
<b>Input</b>	KNN graph	Normalized count matrix	Normalized count matrix	Normalized count matrix	Normalized count matrix	Normalized count matrix	Raw count matrix
<b>Custom embedding</b>	None	None	Corrected embedding	Metagene / factor loadings	None	CCA	Learned lower dimensional latent space
<b>Correction object</b>	KNN graph	Count matrix	Embedding	Embedding	Count matrix	Embedding	Embedding
<b>Correction method</b>	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.
<b>Returns</b>	Corrected KNN graph	Corrected count matrix	Corrected embedding	Corrected embedding	Corrected count matrix	Corrected count matrix	Corrected count matrix and corrected embedding
<b>Changes Count matrix</b>	No	Yes	No	No	Yes	Yes	Yes / Imputes new values

# Batch Correction Methods

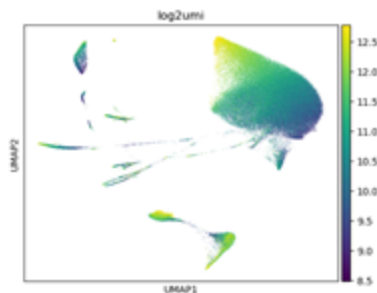


# Minimal batch effect in sci-RNA-seq!



BICCN atlas

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



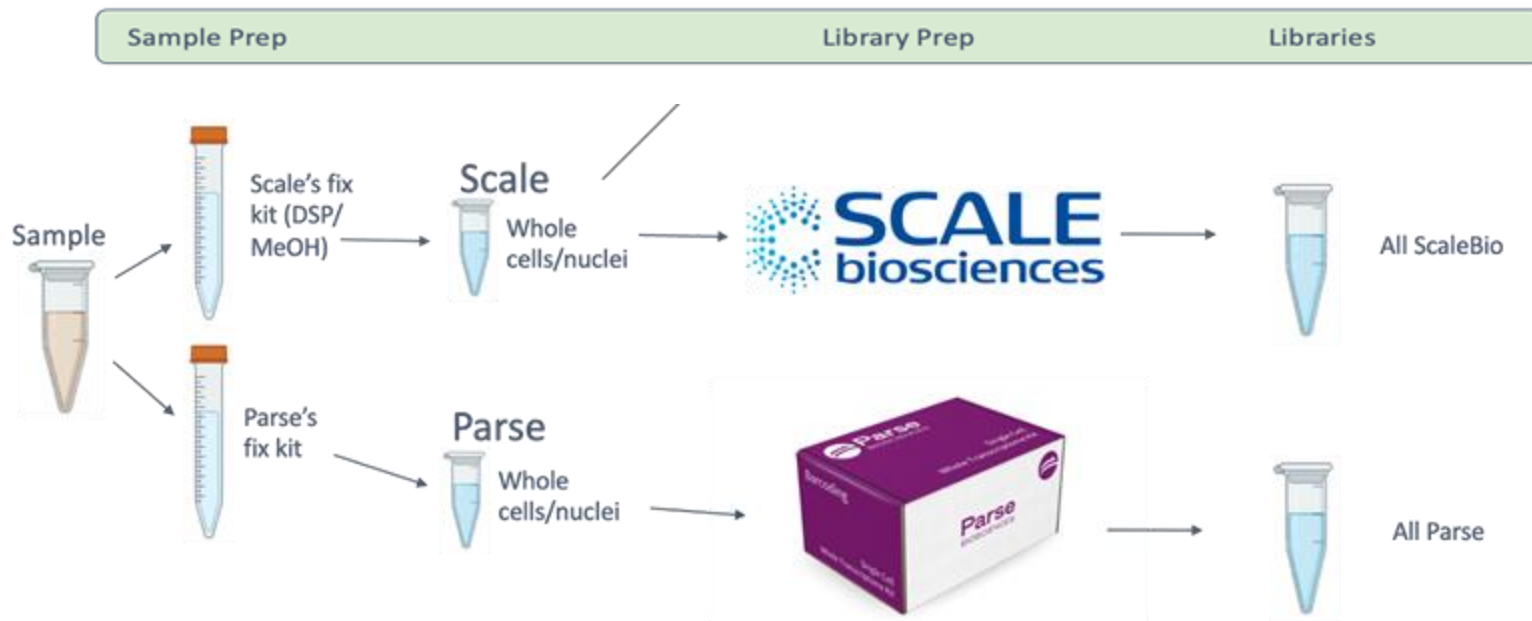
- Snyder-Mackler\_RNA3-051\_nova\_data
- Snyder-Mackler\_RNA3-056\_057\_nova\_data
- Snyder-Mackler\_RNA3-057\_franken\_novaseq\_data
- Snyder-Mackler\_RNA3-058\_059\_nova\_data
- Snyder-Mackler\_RNA3-059\_franken\_novaseq\_data
- Snyder-Mackler\_RNA3-060\_061\_nova\_data
- Snyder-Mackler\_RNA3-061\_franken\_novaseq\_data
- Snyder-Mackler\_RNA3-062\_franken\_novaseq\_data
- Snyder-Mackler\_RNA3\_053\_054\_nova\_data

# Tutorial

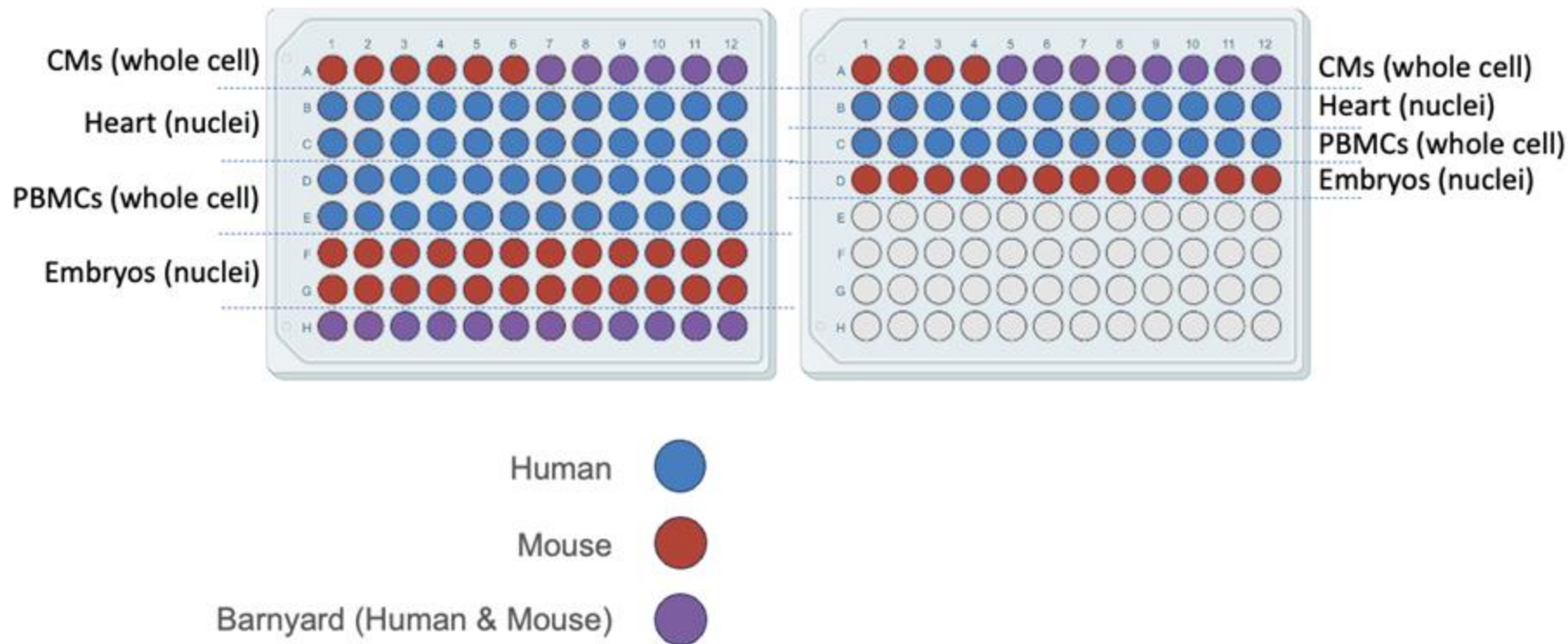
Background on the dataset



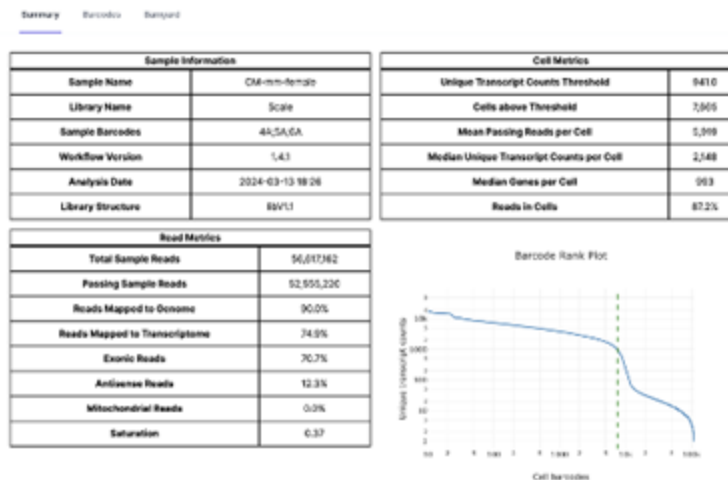
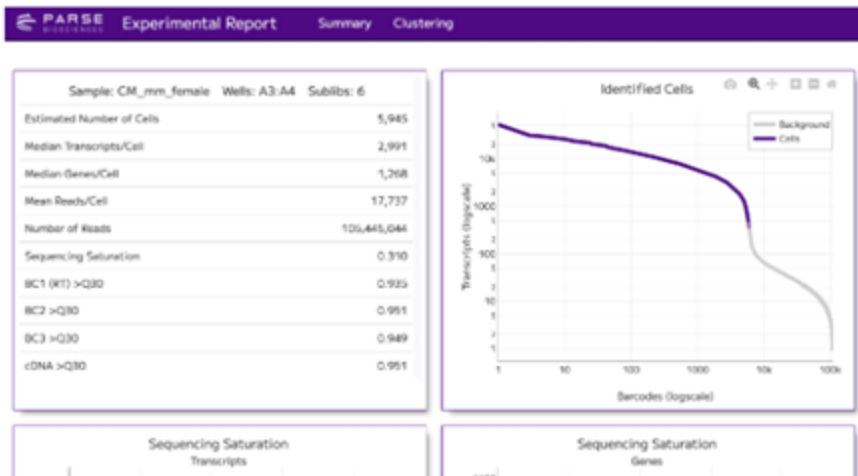
# Study Design



# Study Design



# Each company has a pipeline

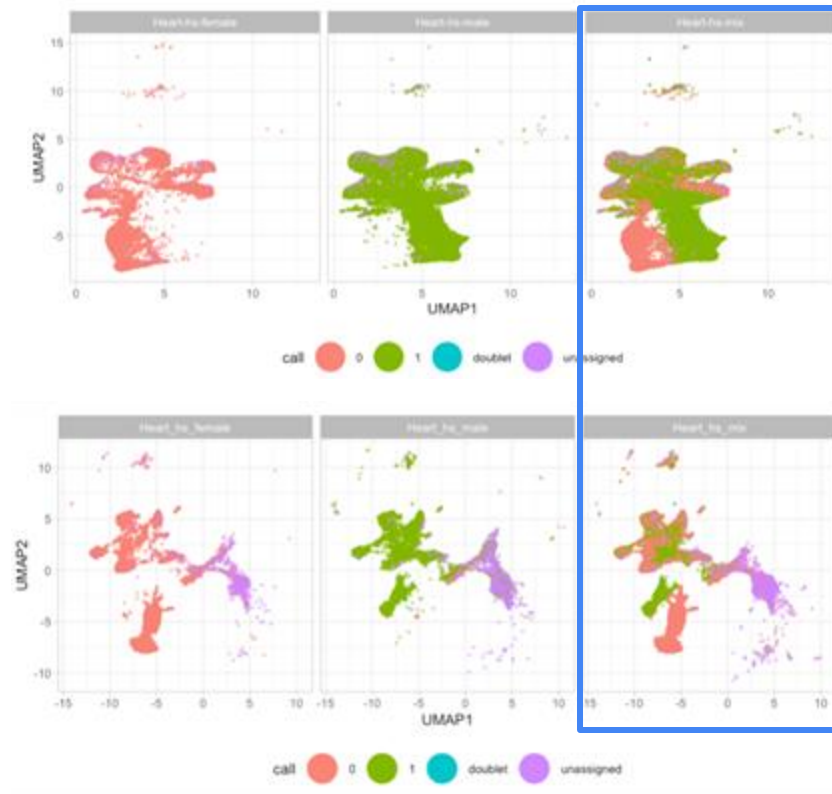
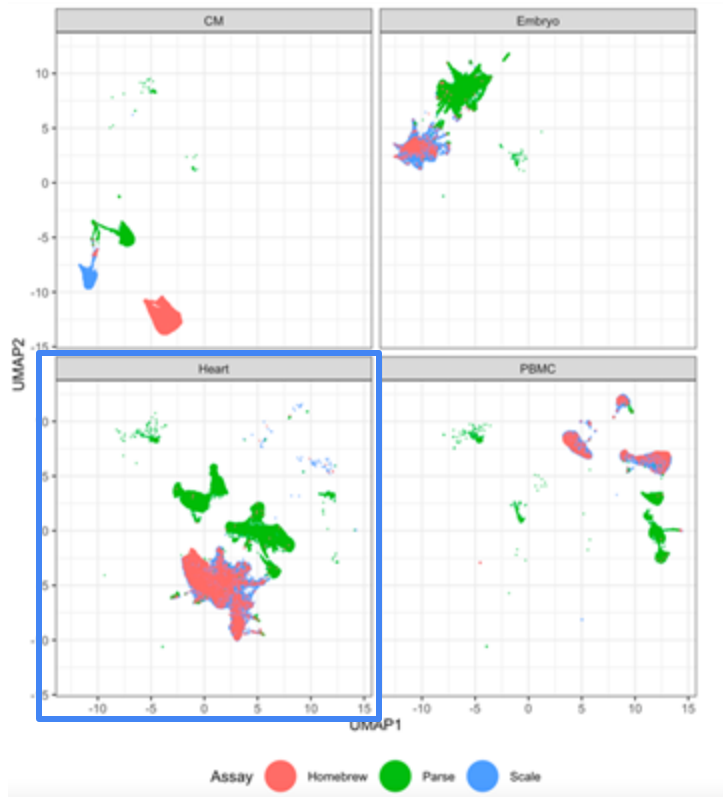


Complexity

Genes Detected Per Cell

- **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  $\geq 100$  UMIs
- Eventually, will try to use a near universal pipeline. For now, using respective pipelines.

# Will play with human heart data today





Thank you!

