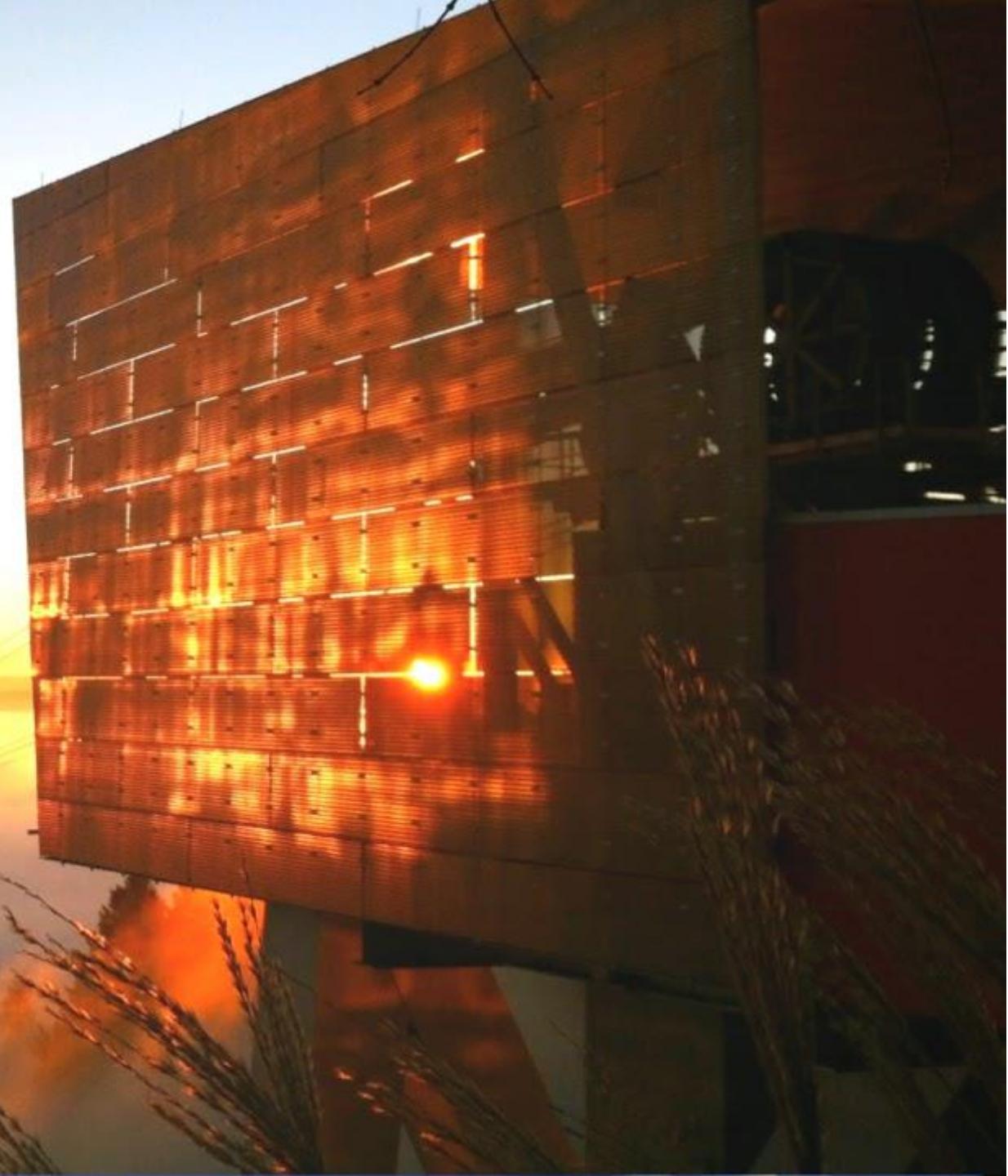


A functional overview of scRNA-seq technologies and methods

Nick Calistri

10/19/23

T32 Course



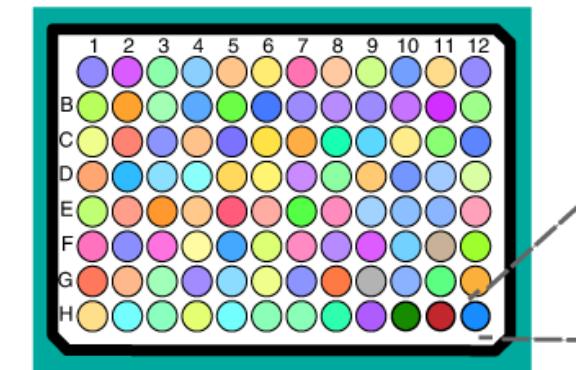
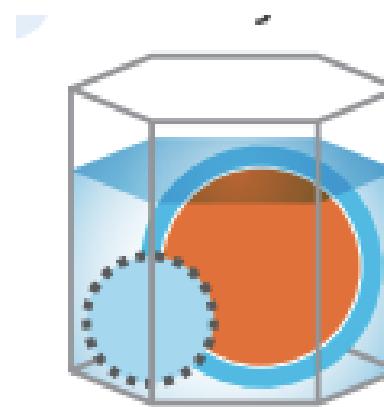
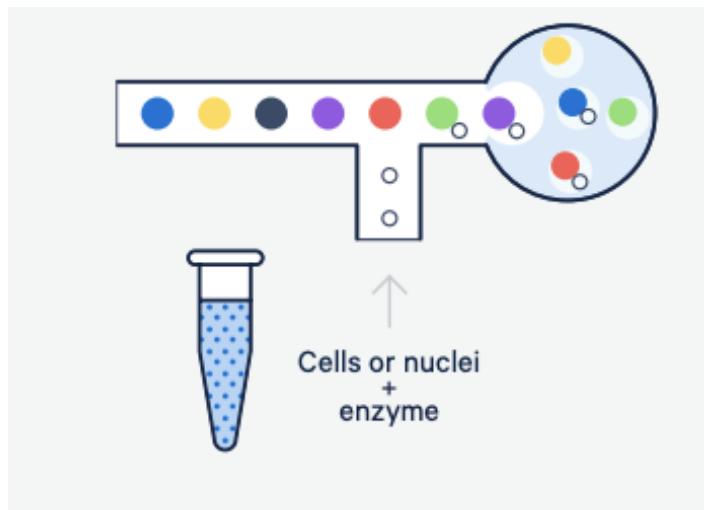
Outline

- Experimental process
 - Shared steps in creating scRNA-seq libraries for analysis
- Standard analysis steps
 - Normalization, dimensionality reduction and clustering
- Advanced techniques
 - Sample multiplexing, integration and cell-cell interaction analysis

scRNA-seq in 30 seconds

- Single cell
 - Need to isolate individual cells and index their contents for analysis
- RNA
 - Pros: ‘simple’ molecule that can be easily quantified
 - Cons: Proteins are typically the functional units of biology, RNA can miss factors such as post translational modification which greatly alters impact
- Sequencing
 - Quantification of order
 - Fundamentally, this is reading the order of nucleotides

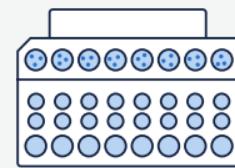
Experimental process of scRNA-seq



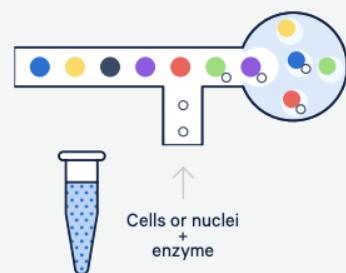
Isolation methods: Droplets

How it works: From cells and nuclei to sequencing-ready libraries

- 1 Gel Beads, cells or nuclei, enzymes, and partitioning oil are loaded onto a Next GEM chip.

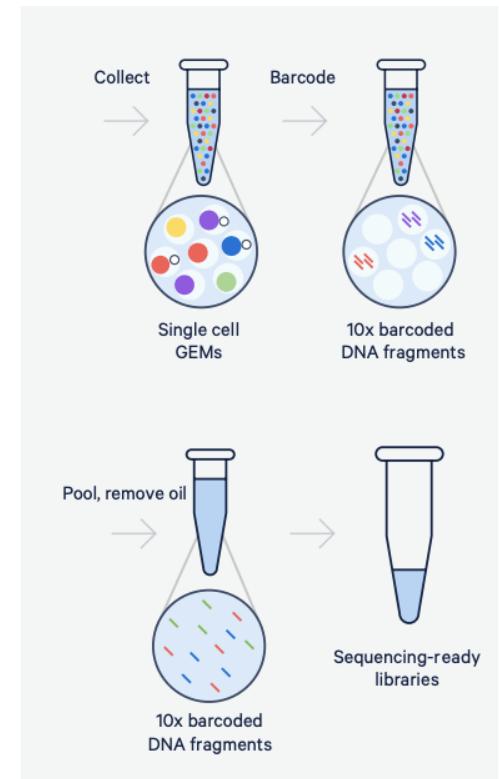


- 2 Within the Chromium instrument, barcoded Gel Beads are mixed with the cells or nuclei, enzymes, and partitioning oil to form tens of thousands of GEMs.

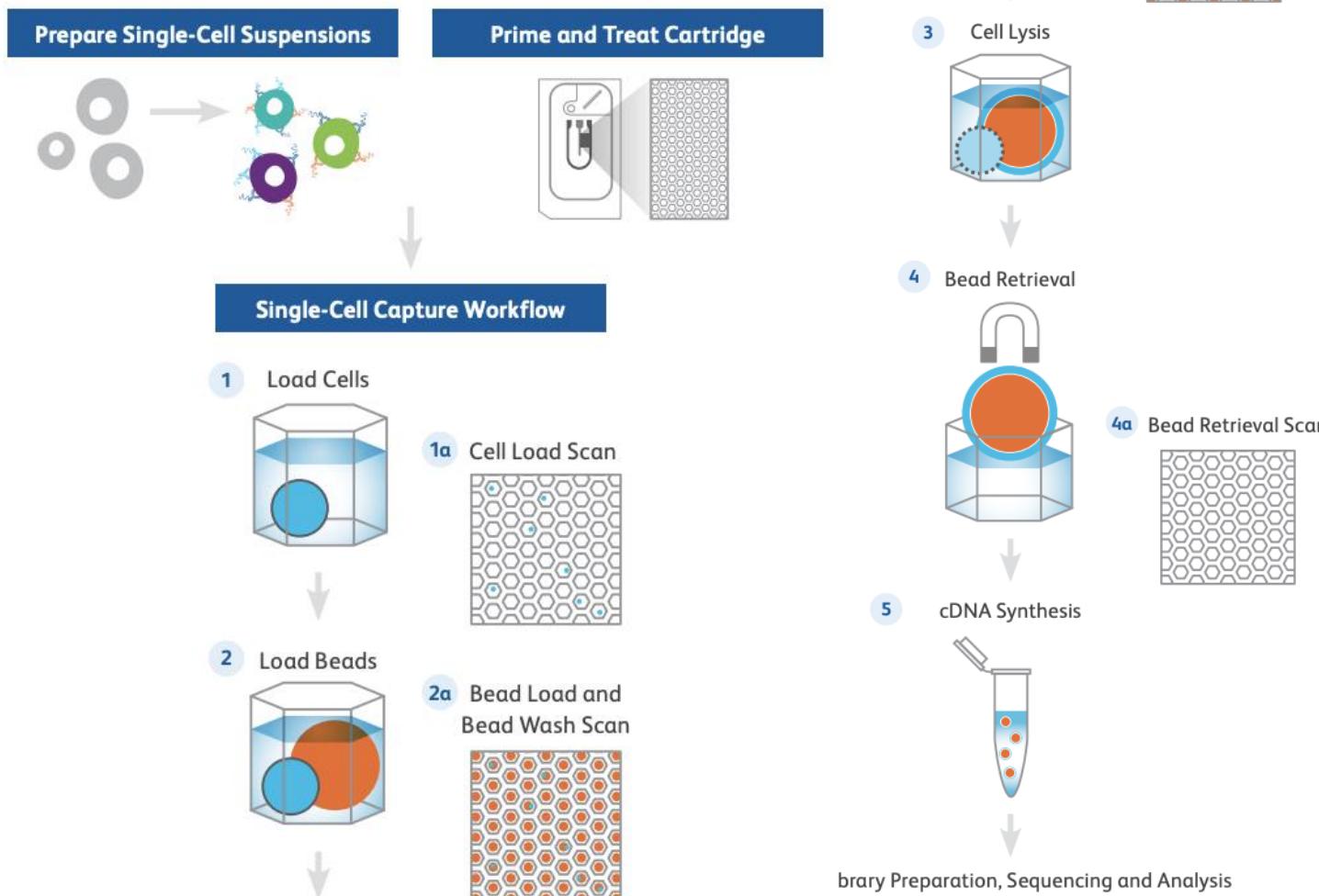


- 3 Each GEM acts as an individual reaction droplet in which the Gel Beads are dissolved and molecules of interest from each cell are captured and barcoded.

- 4 After barcoding, all fragments from the same cell or nucleus share a common 10x Barcode. Barcoded fragments for hundreds to tens of thousands of cells are pooled for downstream reactions to create short-read sequencer-compatible libraries.

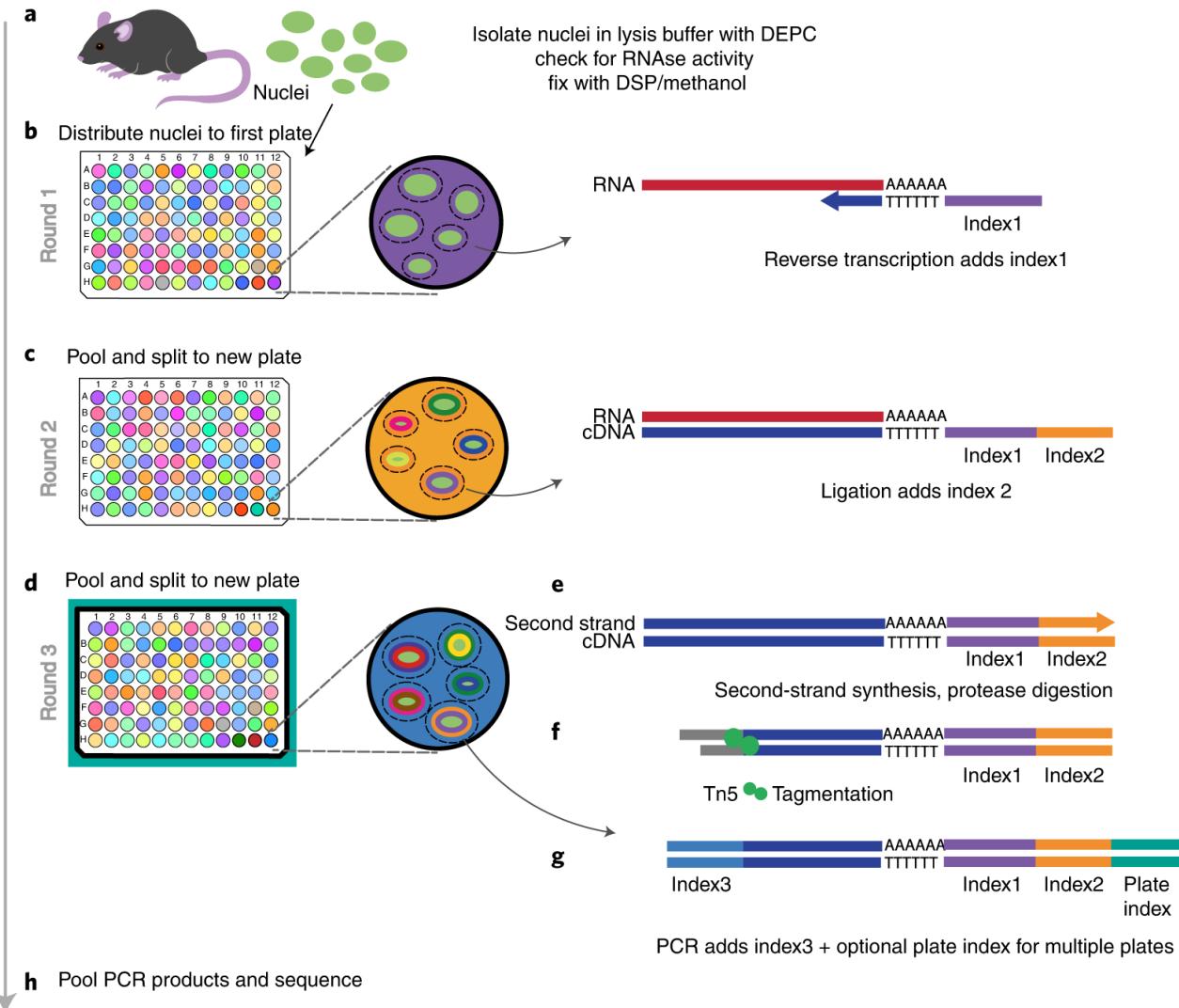


Isolation methods: Microwell



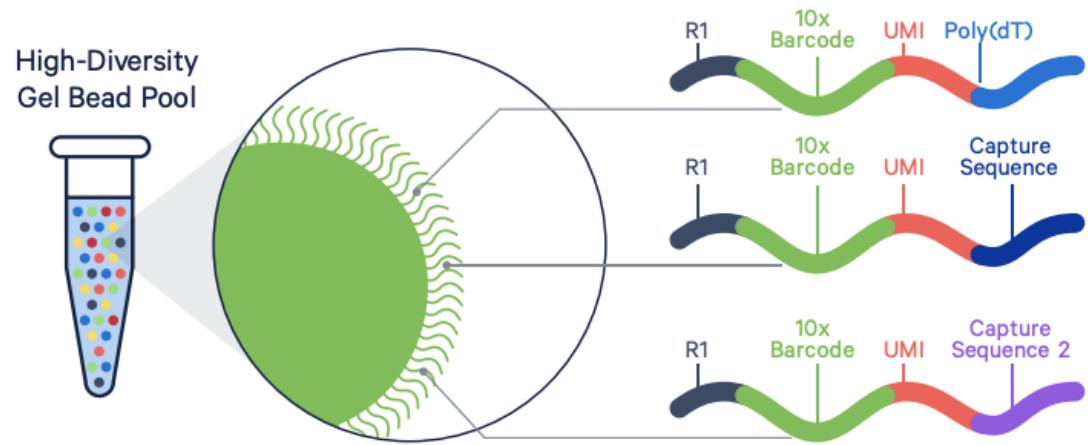
- **Pros**
 - Can visually inspect prior to step 3 to ensure good isolation (fraction of singlets)
- **Cons**
 - Custom hardware needed
 - Potentially higher ambient mRNA contamination

Isolation methods: Combinatorial Indexing



- Pros
 - Doesn't require specific hardware
- Cons
 - More index reads = fewer mRNA reads

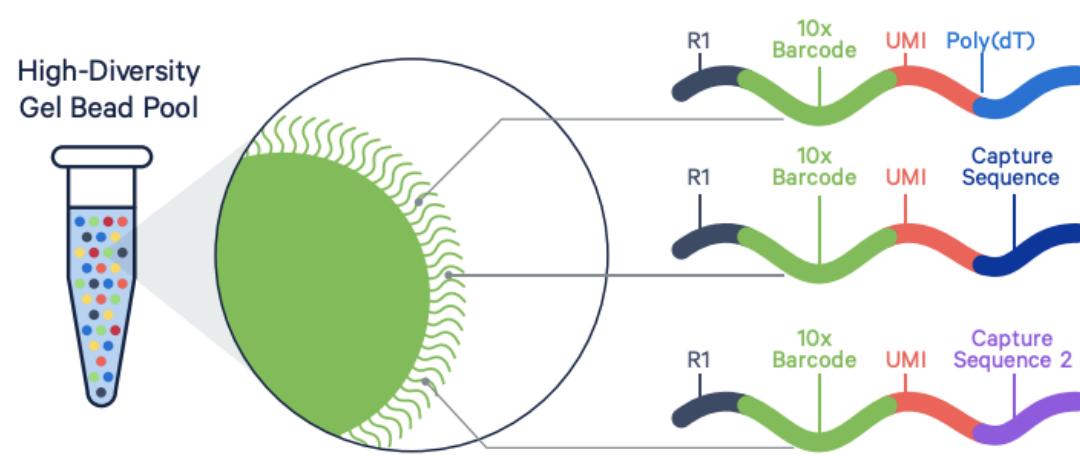
Unique Molecular Identifiers: The solution to amplification bias



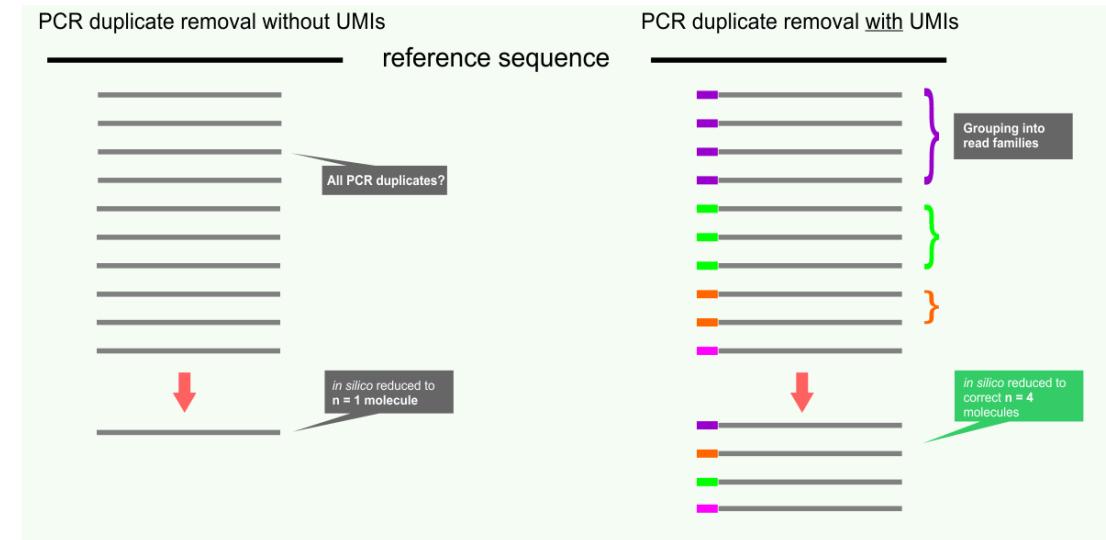
From 10X Genomics

- During amplification, specific mRNA are more likely to be replicated (nucleotide content, overall length etc)
- Each UMI on the 10X GEM bead represents a unique nucleotide sequence
- The unique sequence enables quantification of input mRNA strands

Unique Molecular Identifiers: The solution to amplification bias



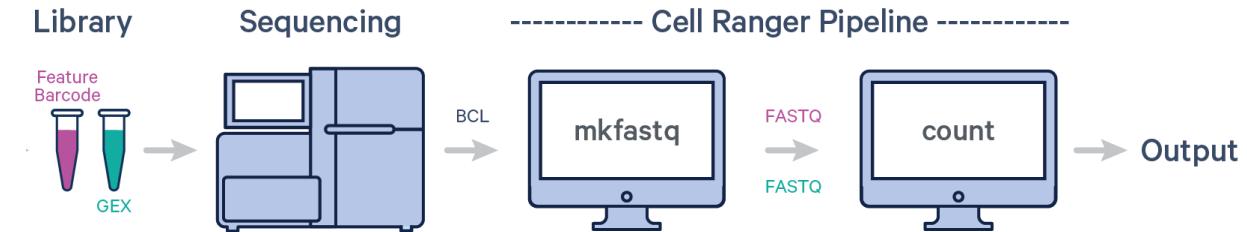
From 10X Genomics



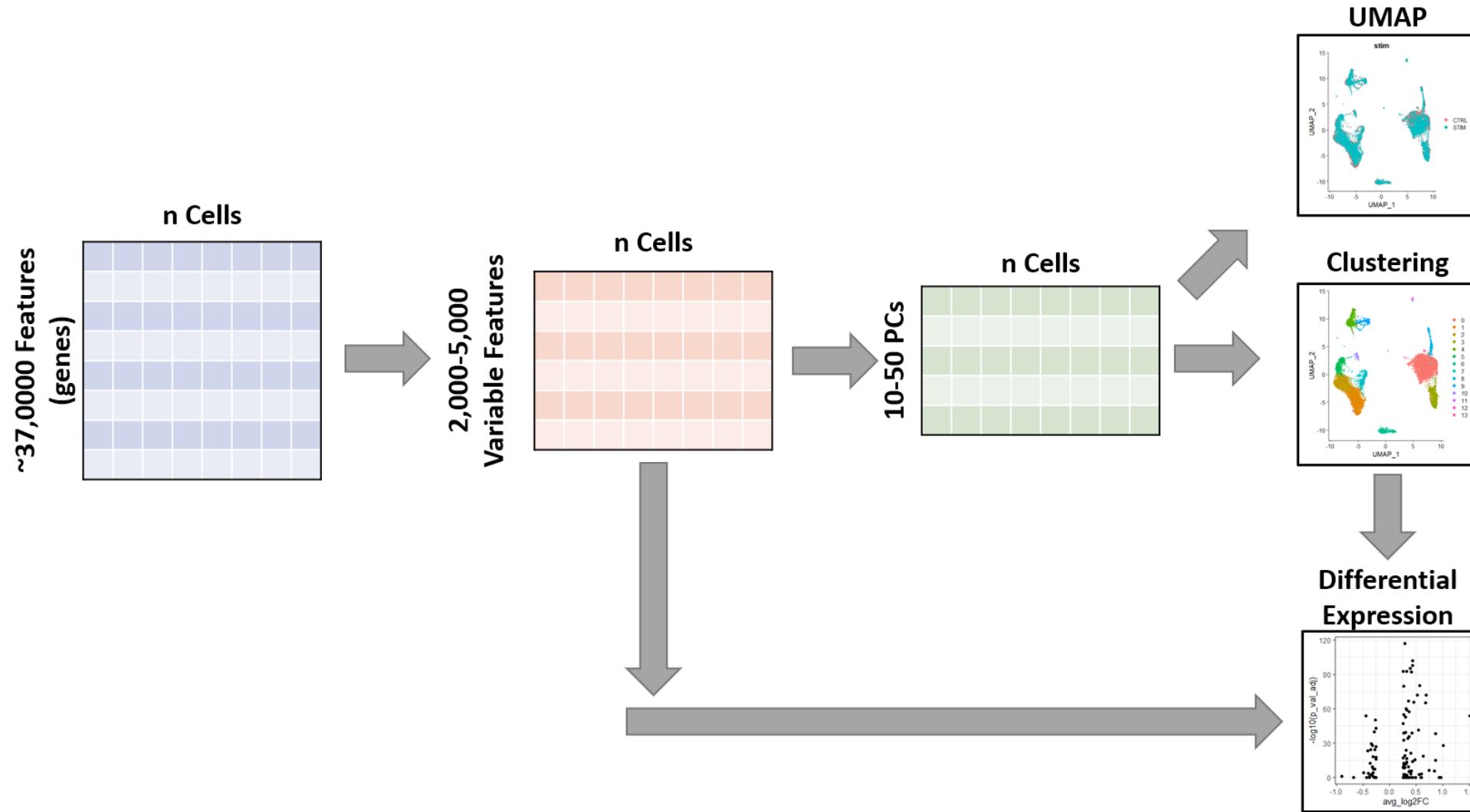
From UC Davis Genome Center

Alignment – reference matters!

- Raw reads returned as ‘.bcl’
 - Essentially the raw sequence reads
- Convert to ‘.fastq’ files
 - ‘Bcl2fastq’ or ‘mkfastq’ are two popular tools for this
 - MUCH smaller file size
 - Data typically organized with four files per ‘sample’ per flow cell
 - Two index files (I1, I2) and read files (R1, R2).
- Align to reference genome/transcriptome
 - ‘CellRanger’ is the 10X chromium alignment/counting tool
 - Aligns .fastq to reference – which gene is represented by each sequence?
 - Should make sure your reference includes any transgene or edited feature of interest!
 - Organizes the data in a ‘count matrix’. Each column is a cell, each row is a gene



Standard Analysis Steps

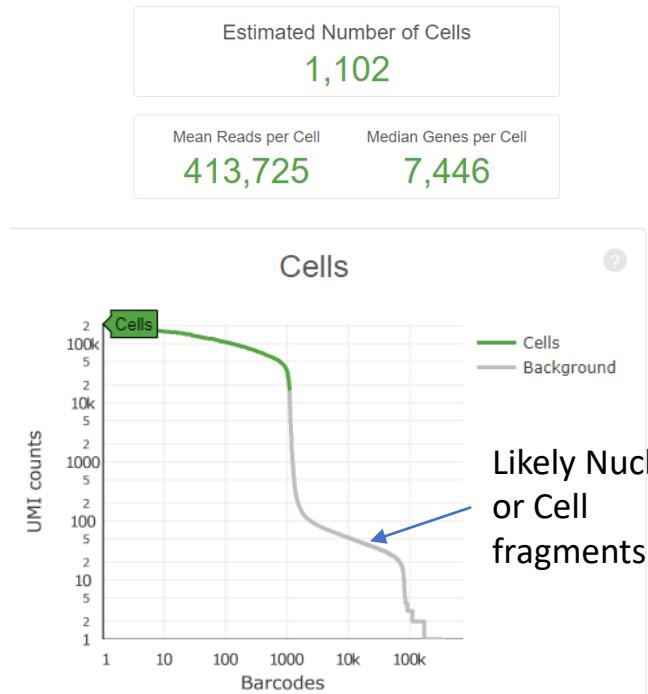


Comp background: “black box” understanding of algorithm is often enough!

- Aspire to know these three things of any algorithm/tool you’re using:
 - Input
 - Format, meaning, **units**
 - Output
 - Format, meaning, **units**
 - Goal
 - What is this algorithm trying to do?
- Bonus: Optimization criteria*
 - When does the algorithm decide to stop learning?
 - Can we evaluate the algorithm’s performance?
 - How well an algorithm performs is *usually* more important than how it works



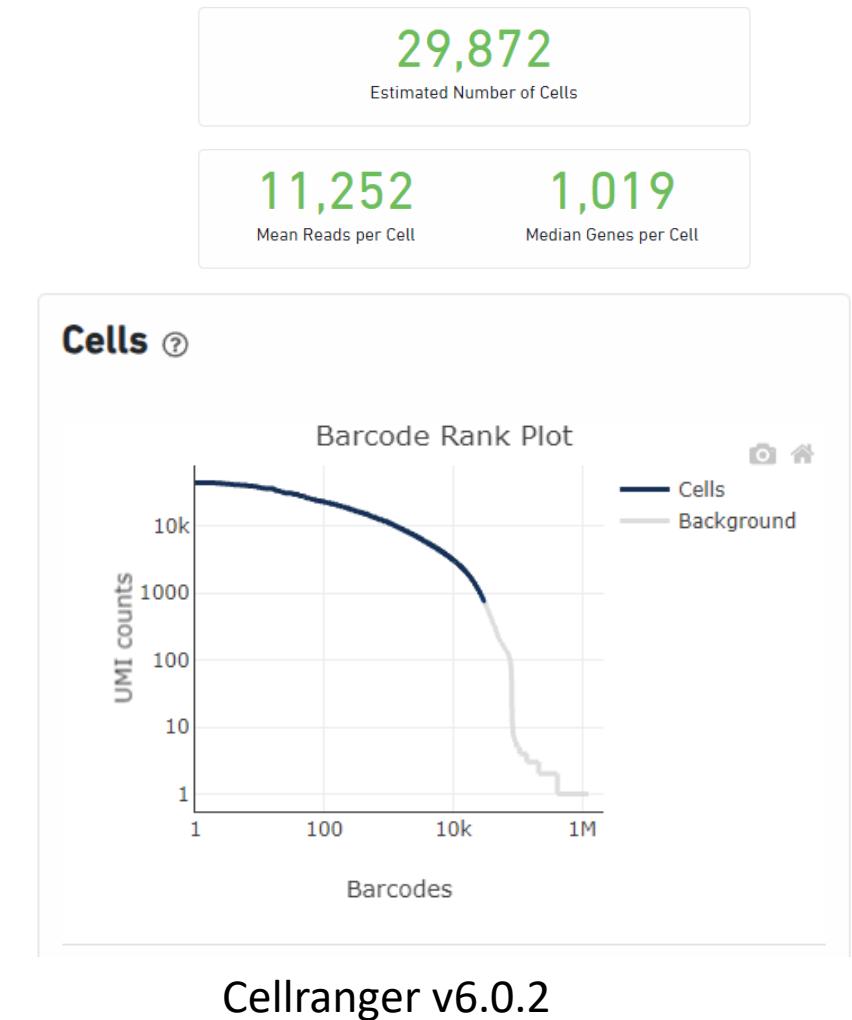
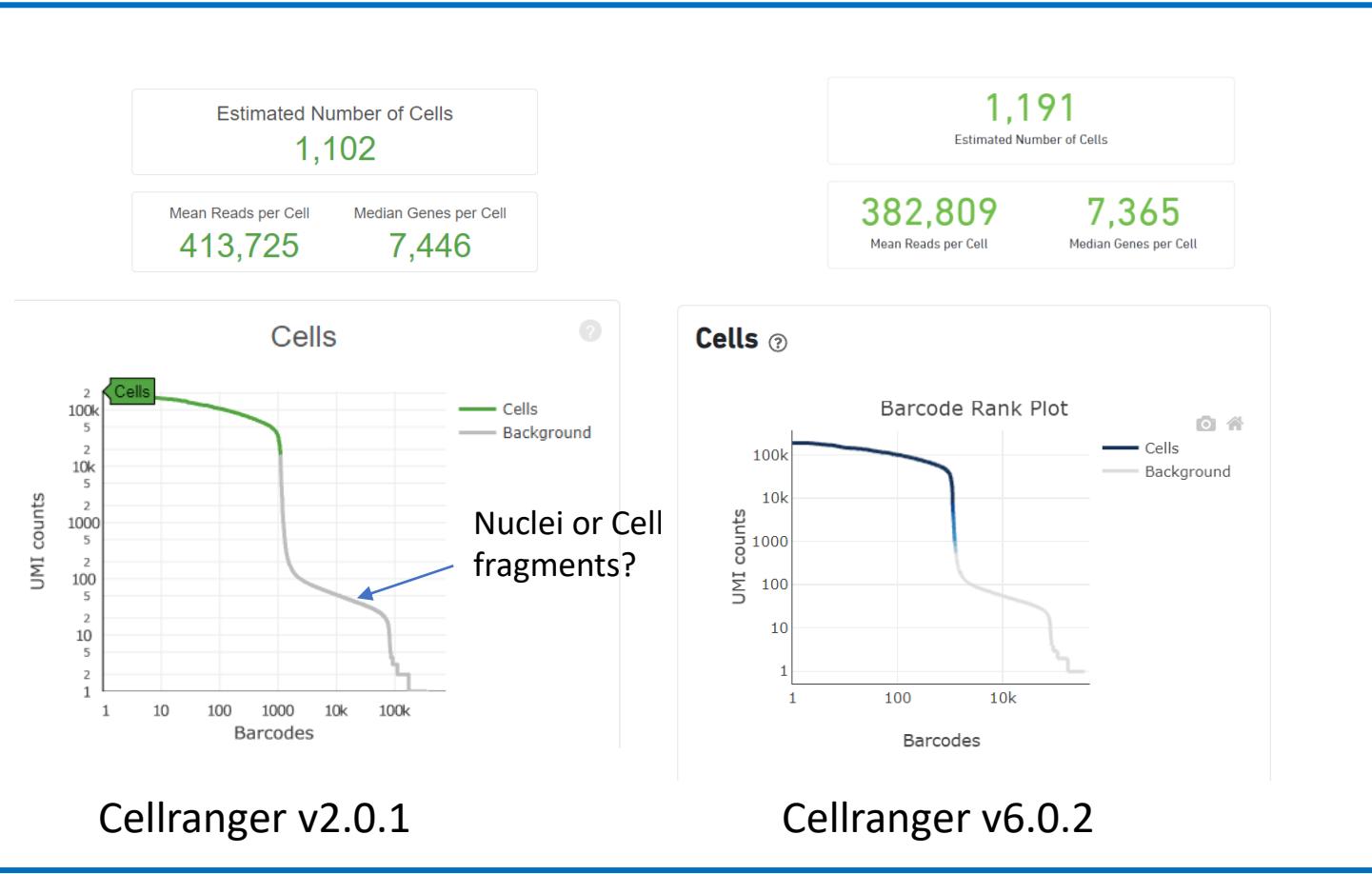
Cells versus droplets – how can we tell?



Cellranger v2.0.1

Cells versus droplets – how can we tell?

Same Data



scRNA-seq is an extremely sparse data set

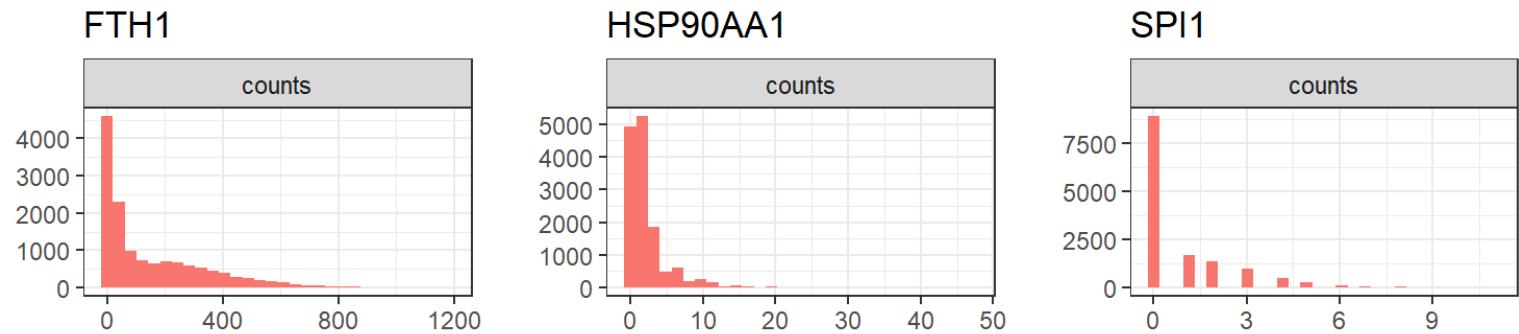
Data Set	Counts / cell	Zero count rate
Ifnb	2,037	95%
Pbmc3k	2,365	93.8%
Panc8	147,920	89.9%

Datasets accessed through package SeuratData



Image with 92.5% missing values

Normalization and scaling transforms the count data

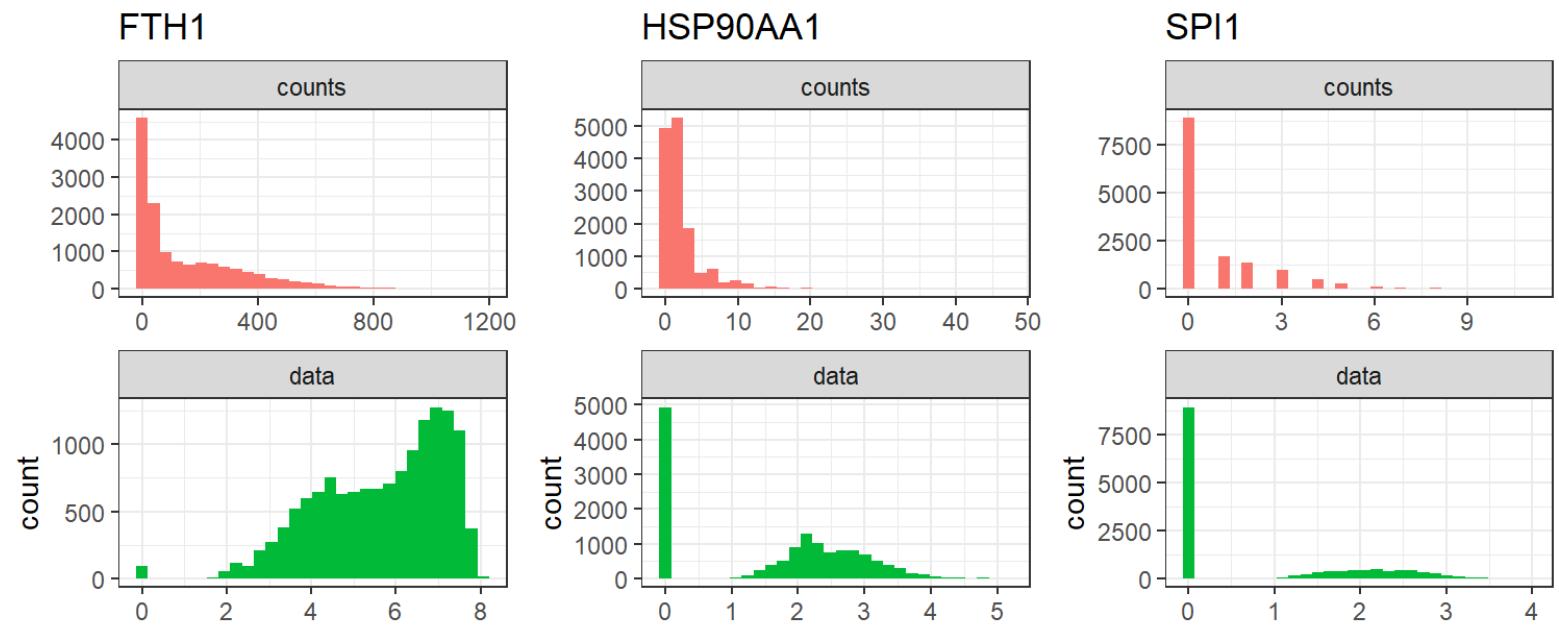


*Seurat by default caps scaled data at +/-10

Normalization and scaling transforms the count data

$$\text{NormalizedData}_{i,j} = \ln(\text{ScaleFactor} * \left(\frac{\text{count}_{i,j}}{\sum(\text{count}_j)} \right) + 1)$$

Effectively converting to relative, log scaled counts. This is normalizing each cell to itself (total number of counts)

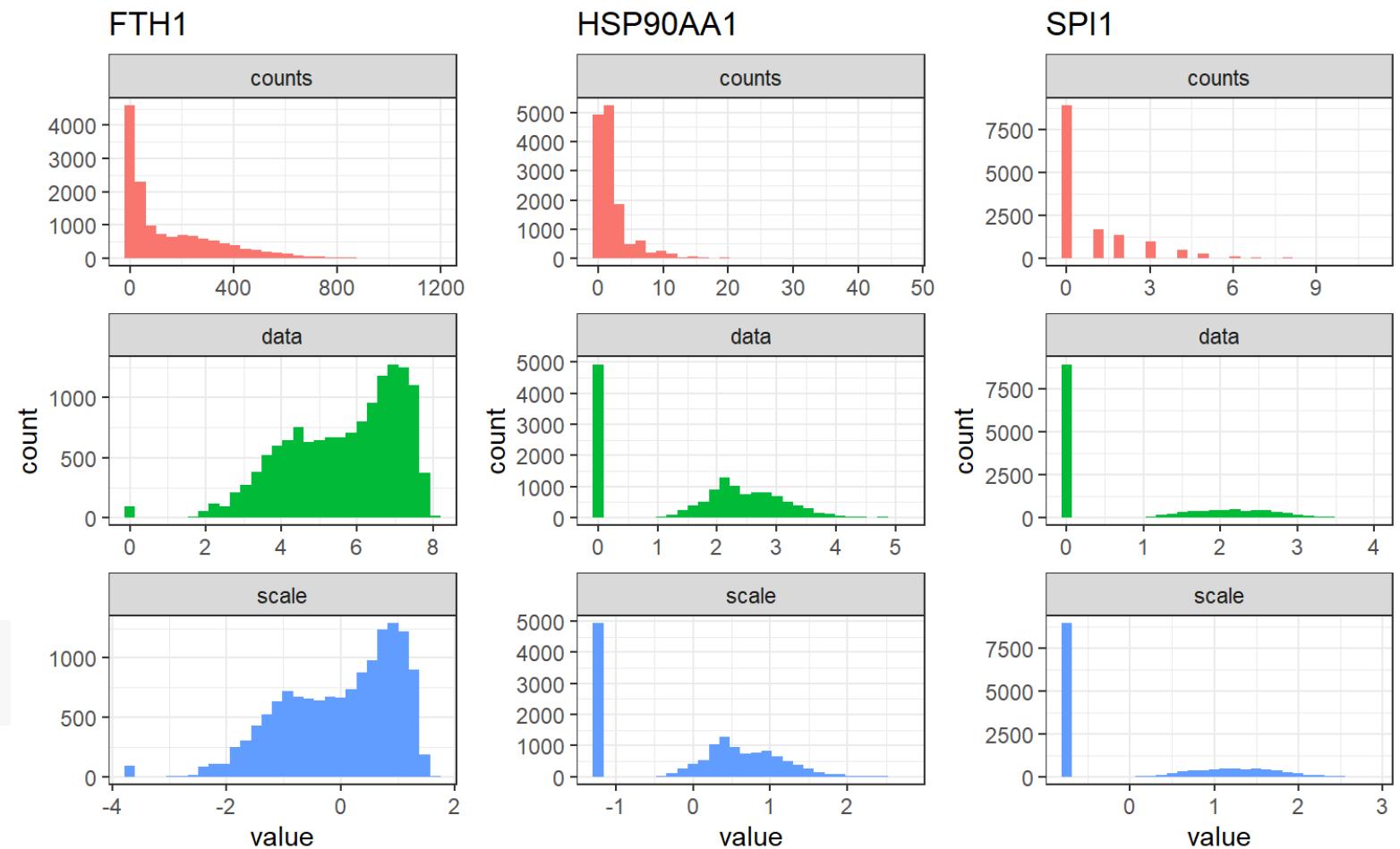


Normalization and scaling transforms the count data

$$\text{NormalizedData}_{i,j} = \ln(\text{ScaleFactor} * \left(\frac{\text{count}_{i,j}}{\sum(\text{count}_j)} \right) + 1)$$

$$\text{ScaledData} = \frac{\text{NormalizedData}_{i,j} - \text{mean}(\text{NormalizedData}_i)}{\text{std. dev}(\text{NormalizedData}_i)}$$

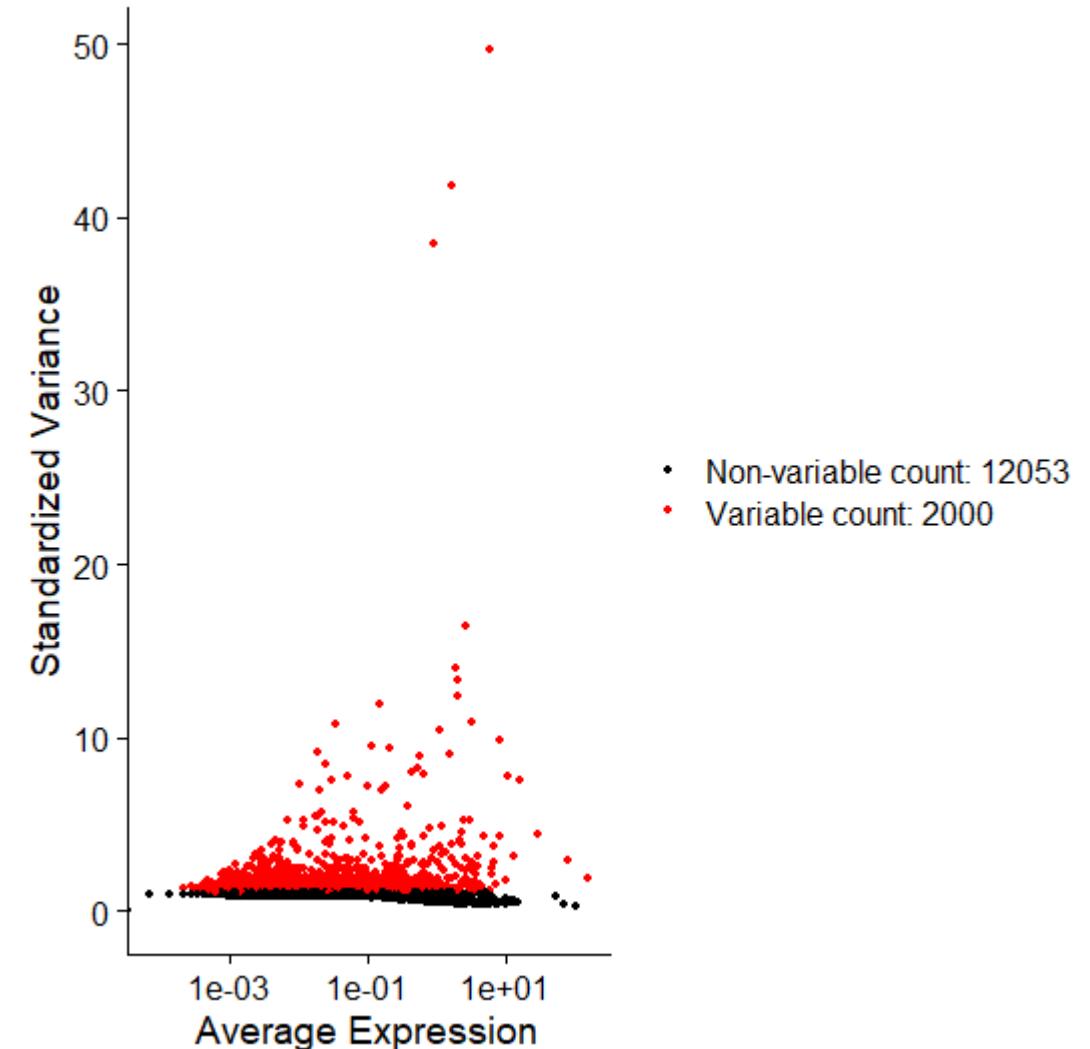
Z score of normalized data across all cells



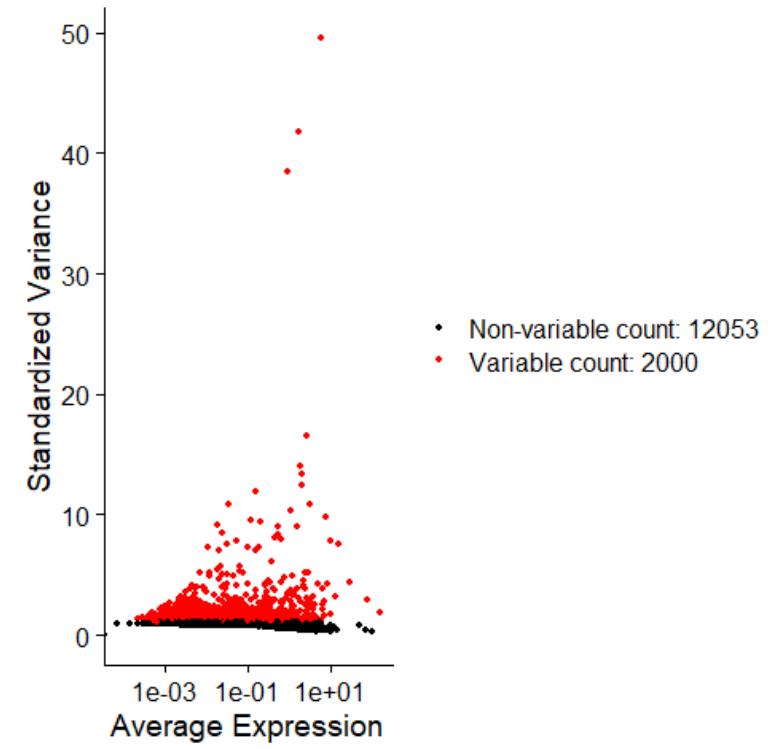
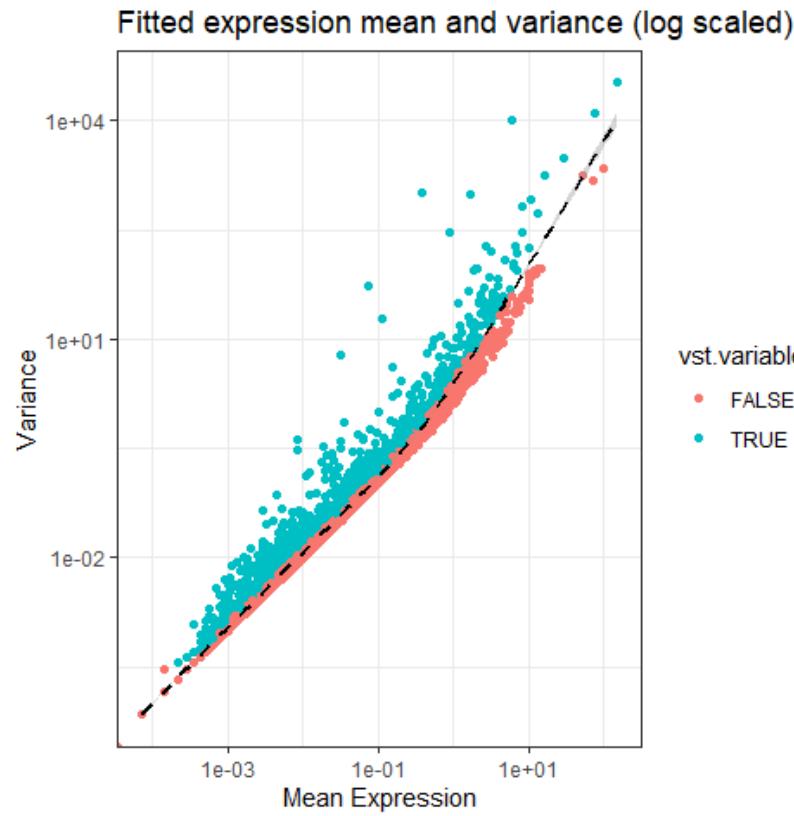
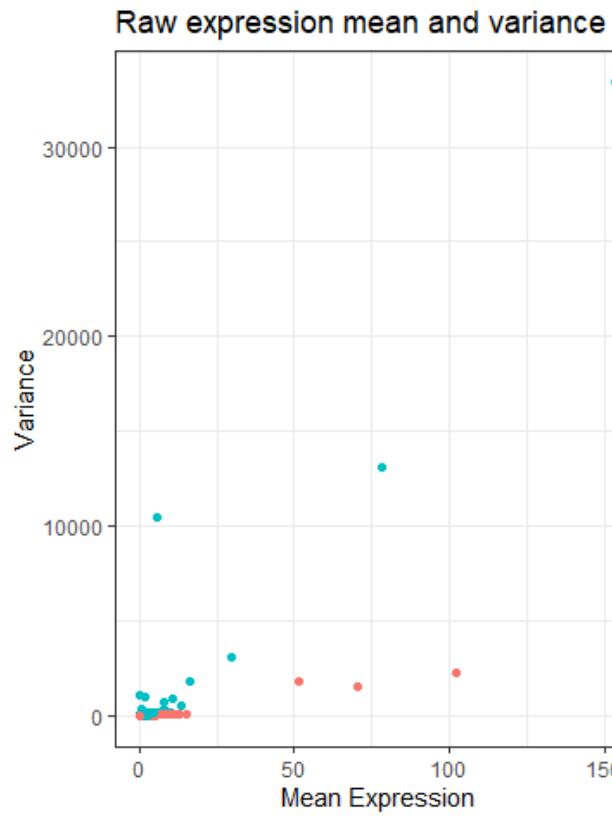
*Seurat by default caps scaled data at +/-10

What features are important?

- Variance stabilizing transformation
 - Regresses variance by mean expression
 - Orders features (genes) by variance beyond expected for mean expression
- Blackbox understanding:
 - Input: Mean gene expression, gene variance
 - Output: Top N most variable genes
 - Goal: Select most variable genes when accounting for expression distribution

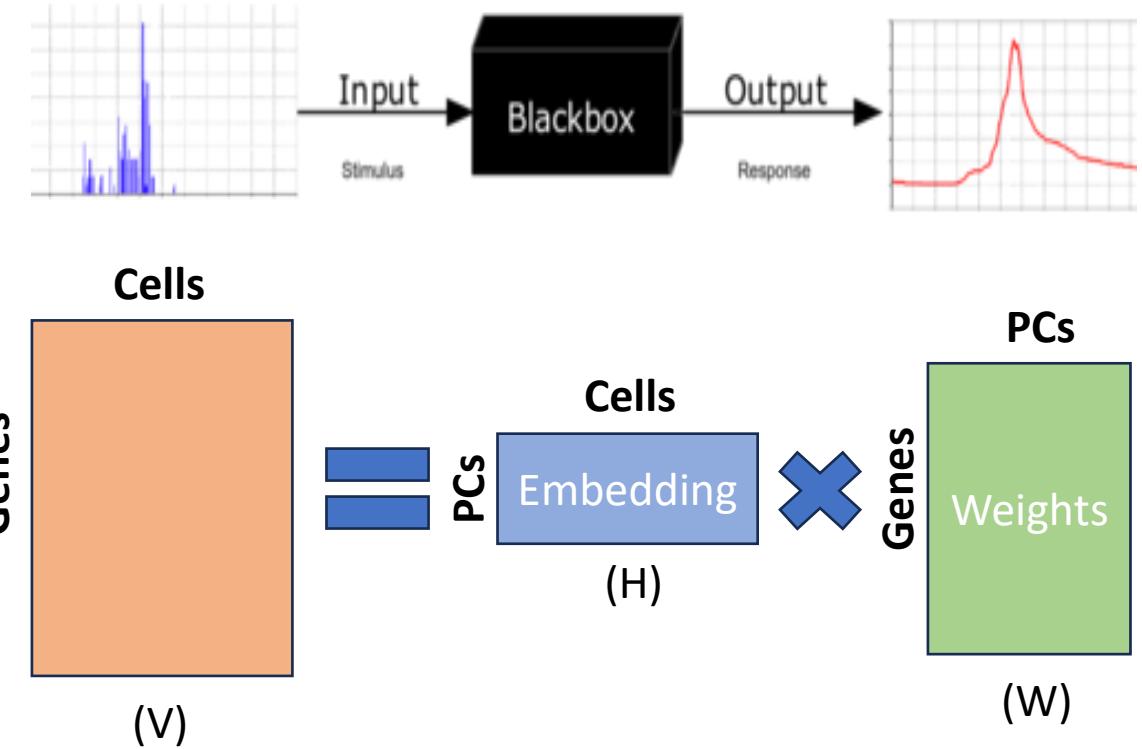


VST regresses expression from variance



Linear dimensionality – reduces noise and improves computational efficiency

- Blackbox understanding of PCA:
 - Input: features by observation matrix
 - For scRNA-seq: Genes by Cell
 - Normalized data
 - Filtered to just include variable features
 - Output:
 - Cell embedding
 - PC score by cell matrix
 - Weight matrix
 - Importance of each gene per PC
 - Goal:
 - Capture as much variation in gene x cell matrix in as few features as possible



Principal component analysis (PCA) – efficient recoding of data

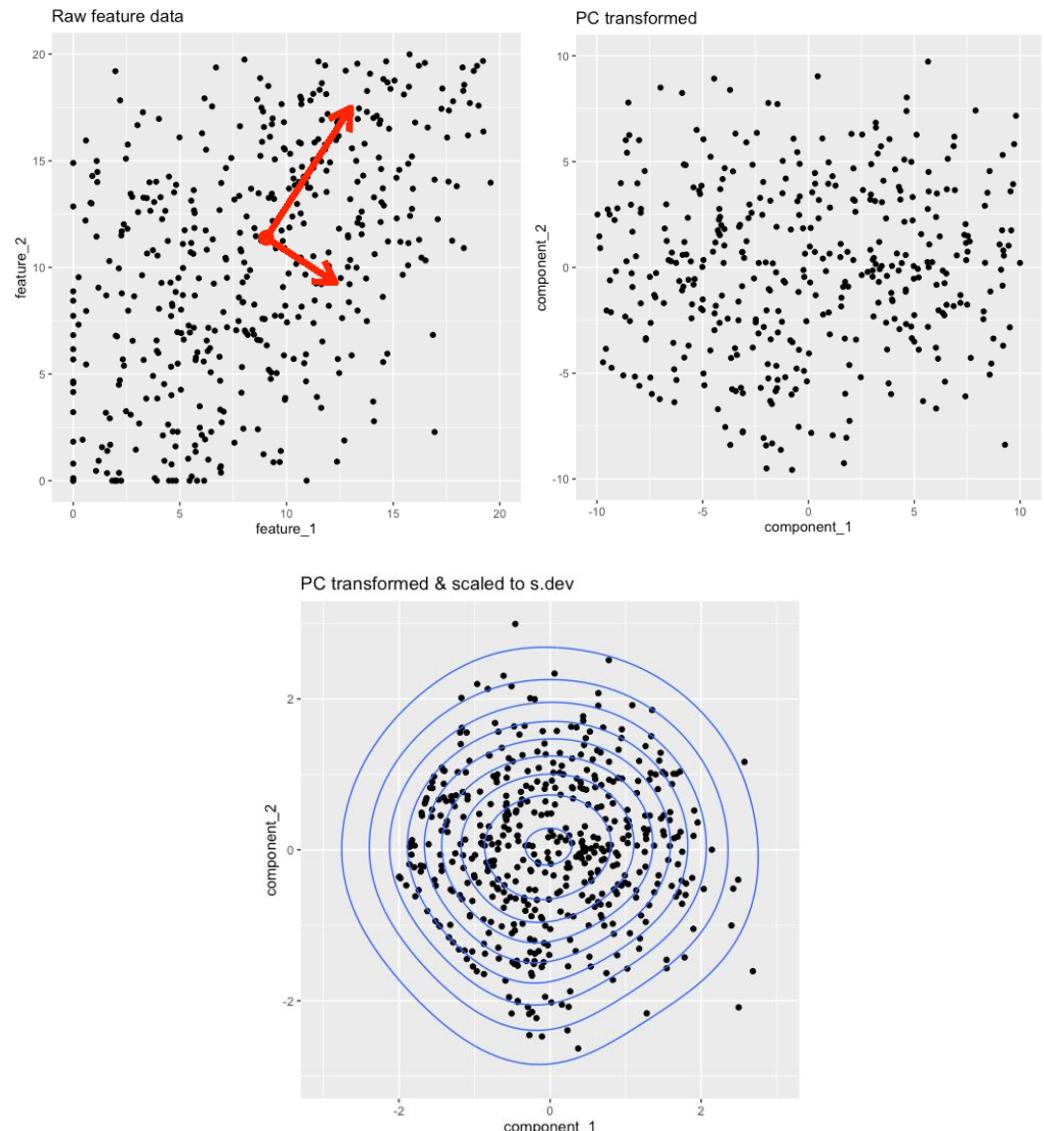
- Identifies perpendicular axes of high variation
- Recodes the data using these axes to describe the underlying data with these axes
- Number of possible components is the number of input features. Using X PCs when you have X input features just results in a rotation of the data and ordering of features by variance

Pros:

- Very efficient
- Deterministic* (same result every time)

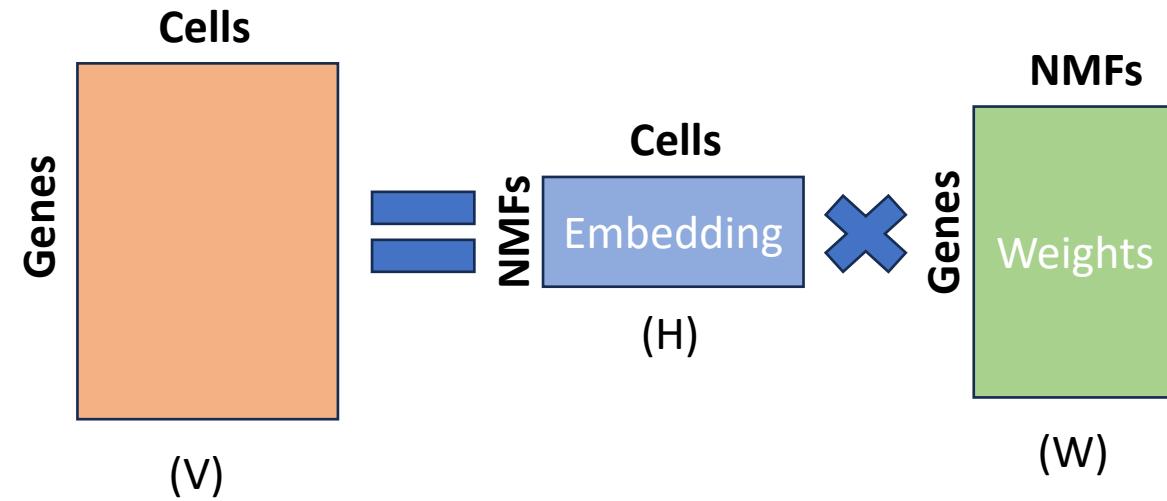
Cons:

- Not particularly interpretable



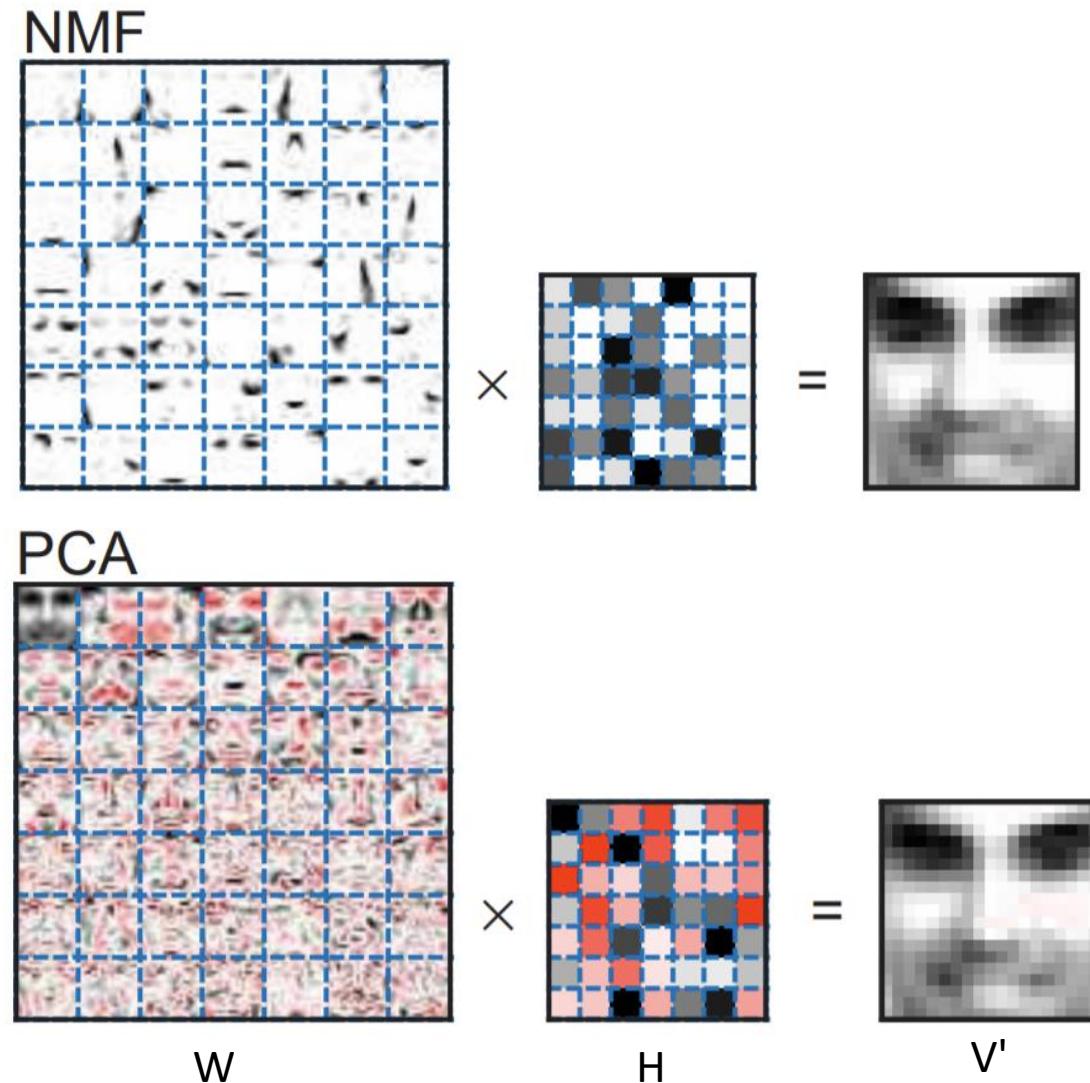
Non-negative Matrix Factorization (NMF) – Interpretable recoding of data

- Similar to PCA



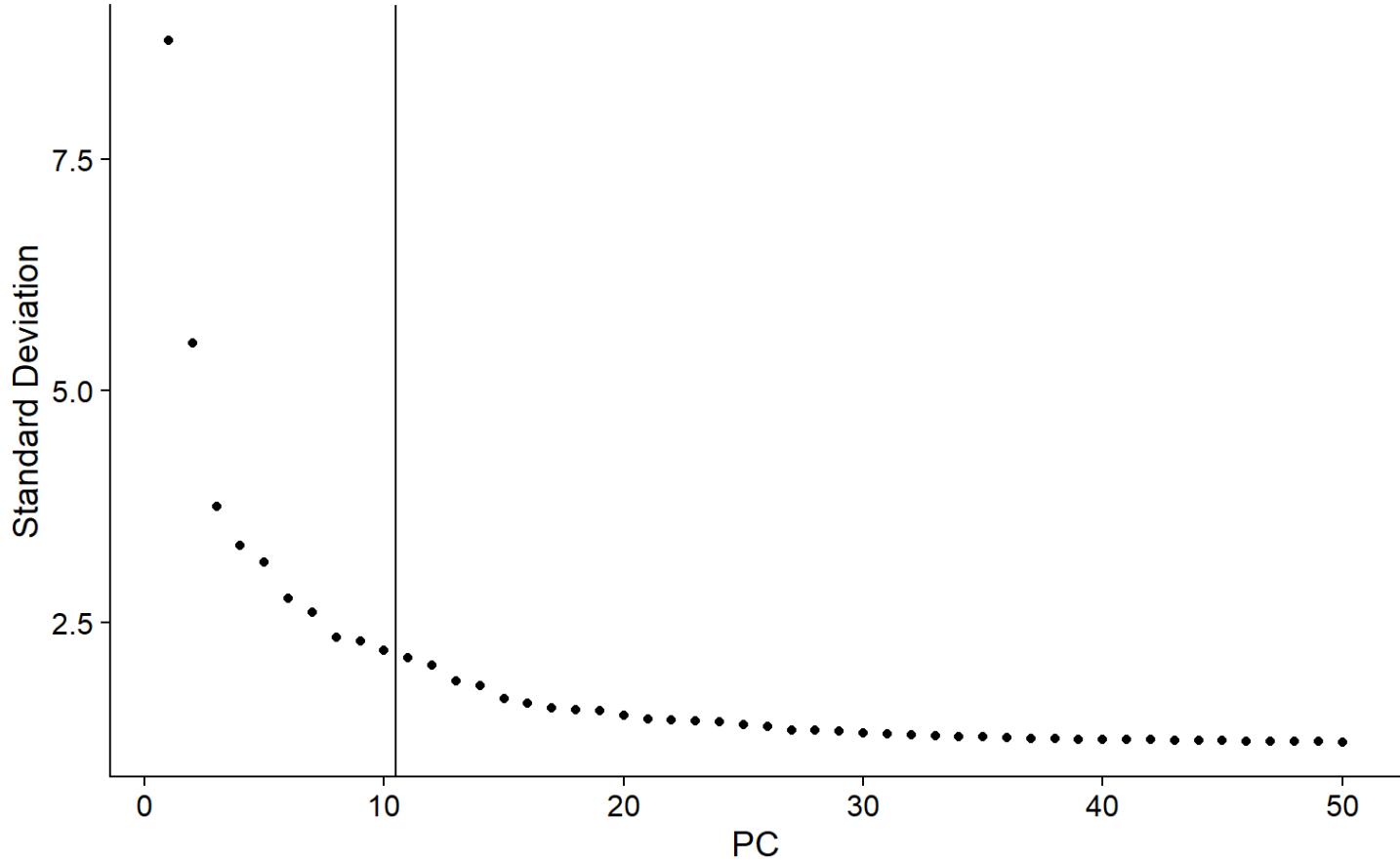
- Factor weights have strictly positive weights
 - This means factors are additive to reconstruct the original
- Stochastic (random seed, 'random' result) approach for most implementations
- Tends to require more latent features than PCA to reconstruct the original

Non-negative Matrix Factorization (NMF) – Interpretable recoding of data



How many LinDimRed features should we use?

- In practice, most people visualize the elbow plot and select some number of PCs where SD diminishes per additional PC
- Best practice is to sweep a number of PCs and compare clustering output to ensure your clusterings (and downstream analysis) is robust across a few values adjacent to your selection



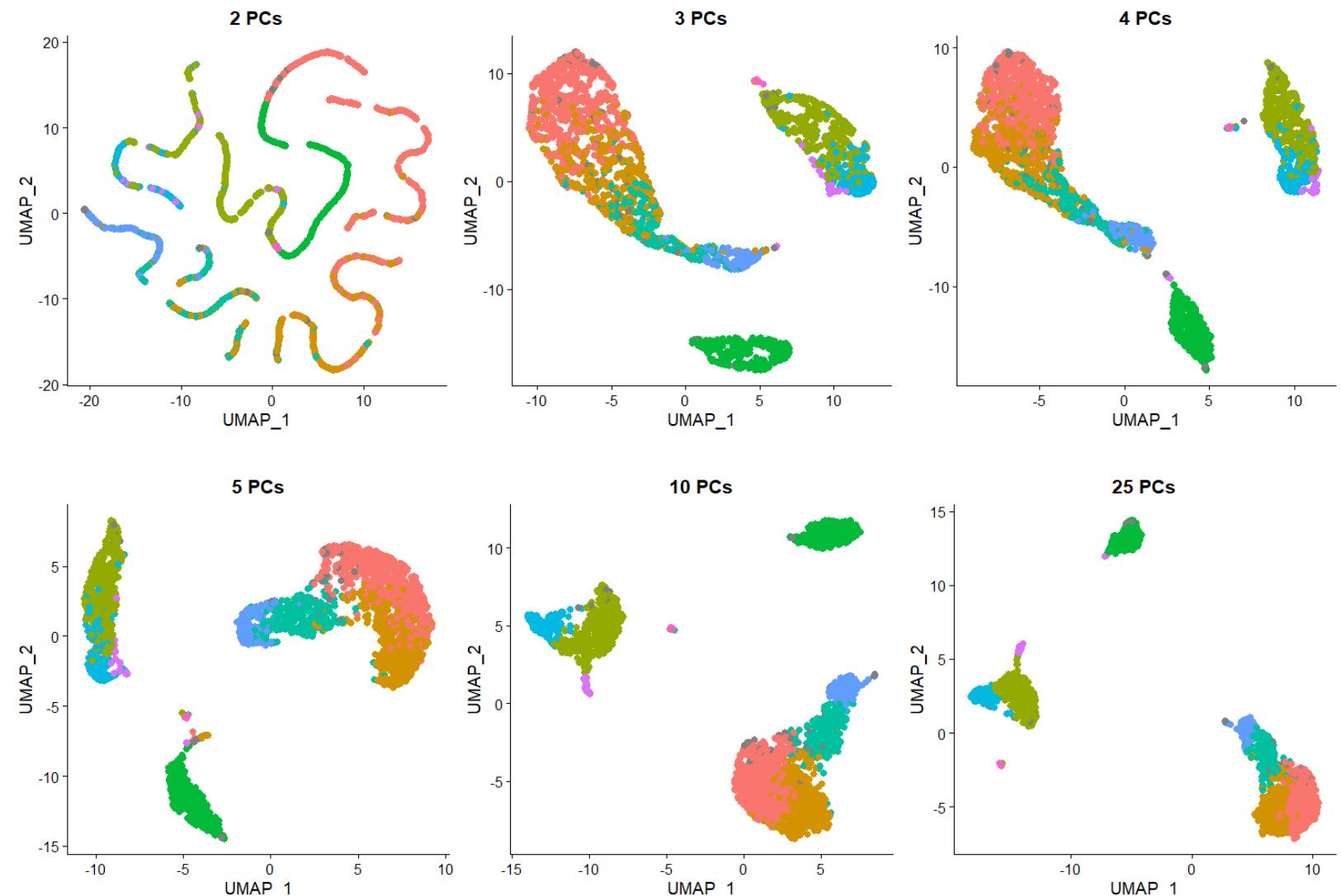
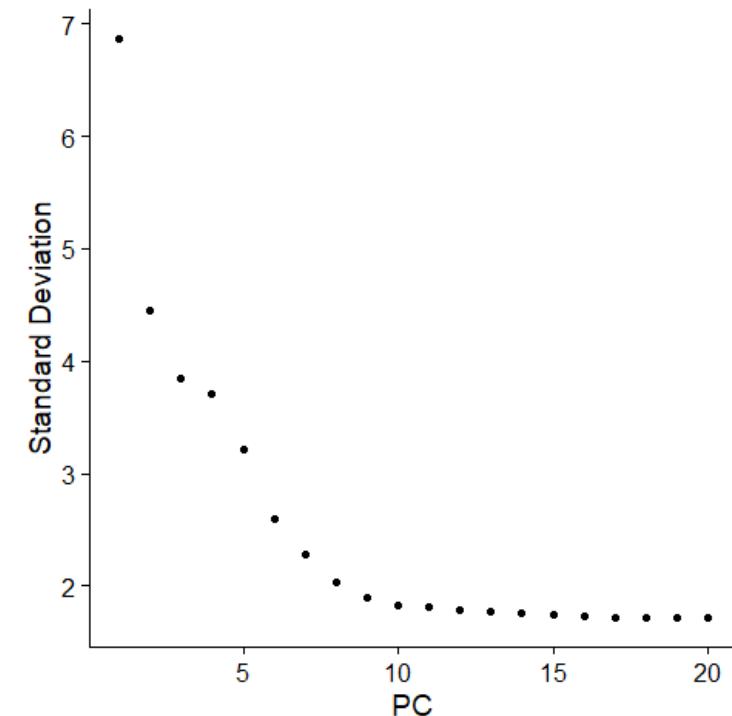
Non Linear Dimensionality reduction – because we still want to publish a manuscript

- Two popular methods: tSNE and UMAP
 - Both are primarily used to project into two dimensions
 - UMAP is newer/more popular.
 - Scales better with number of cells
 - Tends to better balance local versus global architecture
- It is the responsibility of the analyst to ensure UMAP/tSNE faithfully depicts underlying data, NOT to interpret from the visualization
- Output is HIGHLY dependent on parameters
 - <https://pair-code.github.io/understanding-umap/>
- Many ‘aberrations’ are diminished in 3d instead of 2d – consider computing 3d UMAP and identifying a representative 2d slice instead! (Credit: Adey lab)

Non linear DimRed ‘blackbox’ understanding

- Input:
 - Linear latent embeddings (PCA scores per cell)
- Output:
 - 2 dimensional coordinates
- Goal:
 - Position cells in low dimensions (2d umap) in such a way that relative distance from high dimensions (PCA) are conserved
 - Clearly from the Pachter paper we reviewed, the low dimensional visualization is a compromise

Sweeping through principal components



Clustering: Leveraging similarity to circumvent sparsity



Image with 92.5% missing values

Clustering: Leveraging similarity to circumvent sparsity



Image with 92.5% missing values



'clustered' image

Clustering: Leveraging similarity to circumvent sparsity



Image with 92.5% missing values



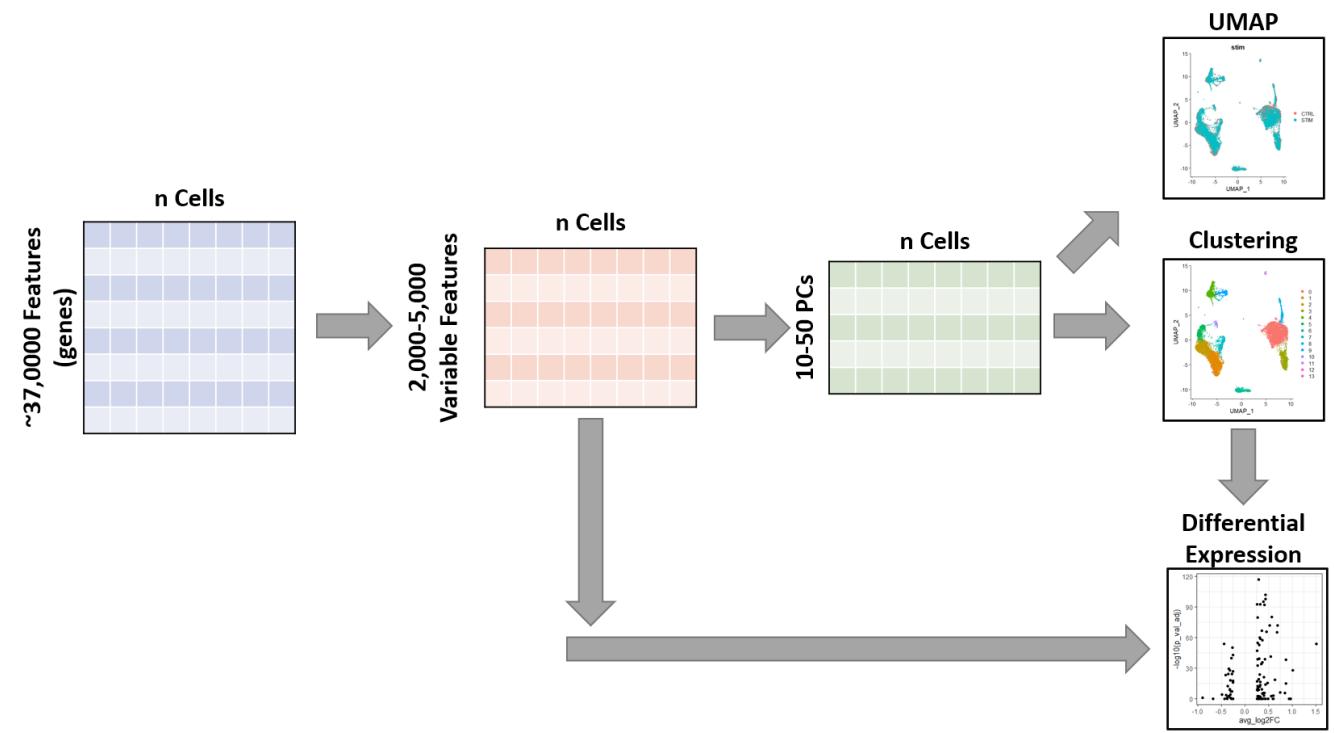
'clustered' image



Original

Clustering – black box understanding

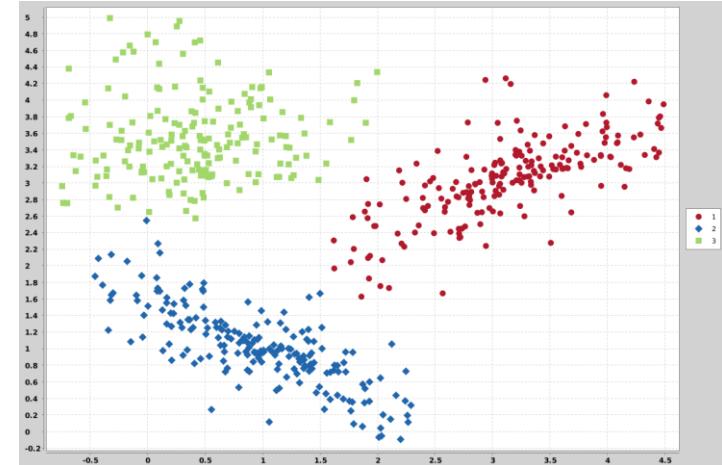
- Input:
 - Cell embedding in linear latent space (Principal component embedding)
- Output:
 - Label (cluster) for each cell
- Goal:
 - Cluster labels represent communities which are similar within feature space



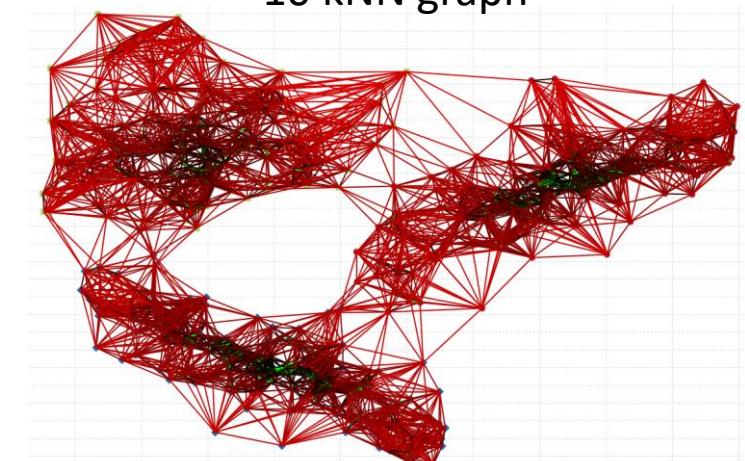
Clustering – let the data tell us about itself

- Two popular methods for scRNA-seq:
Louvain and Leiden algorithms
 - Both are ‘community detection’ methods that optimize for **modularity**
 - Modularity is the ratio of edges within a community compared to outside the community

2 Dimension plot of data



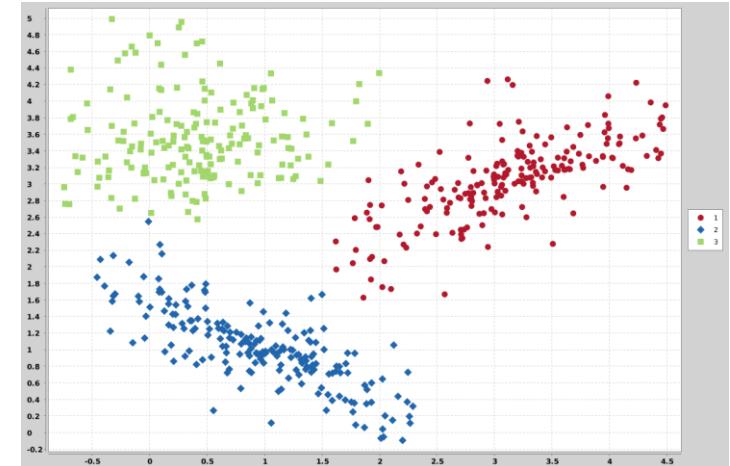
10 kNN graph



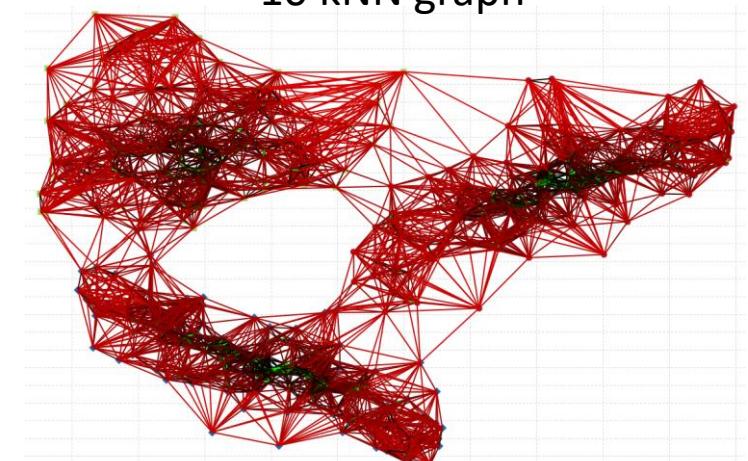
Clustering – let the data tell us about itself

- Two popular methods for scRNA-seq:
Louvain and Leiden algorithms
 - Both are ‘community detection’ methods that optimize for **modularity**
 - Modularity is the ratio of edges within a community compared to outside the community
 - Some methods (Seurat default) use the Shared Nearest Neighbor (jaccard index of KNN) for edge weights

2 Dimension plot of data



10 kNN graph

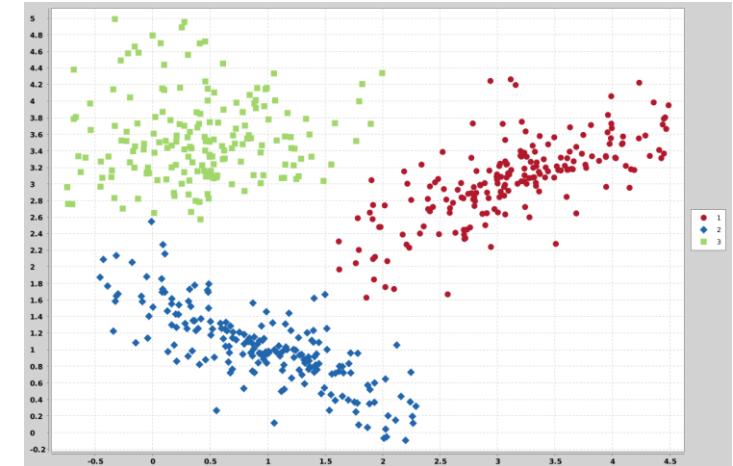


Clustering – let the data tell us about itself

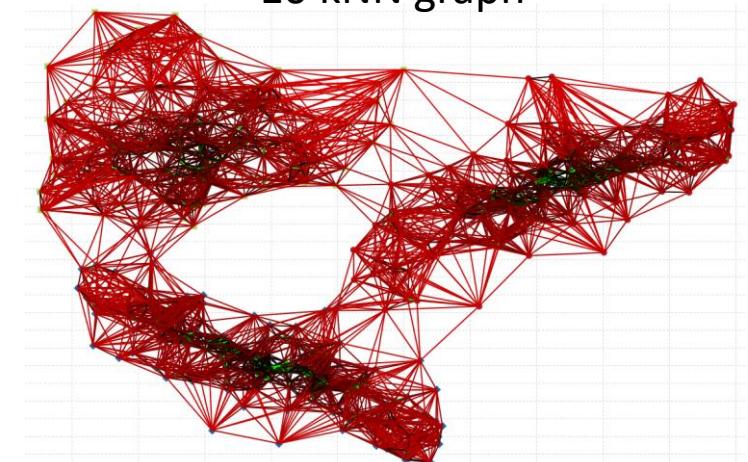
- Two popular methods for scRNA-seq:
Louvain and Leiden algorithms
 - Both are ‘community detection’ methods that optimize for **modularity**
 - Modularity is the ratio of edges within a community compared to outside the community
 - Some methods (Seurat default) use the Shared Nearest Neighbor (jaccard index of KNN) for edge weights
 - **Resolution** parameter determines when to separate communities (high = more communities, low = fewer)

(Leiden approach is theoretically ‘better’ as it strictly prevents disconnected communities)

2 Dimension plot of data



10 kNN graph



How do we know the clustering is correct?

How do we know the clustering is correct?

- We don't!

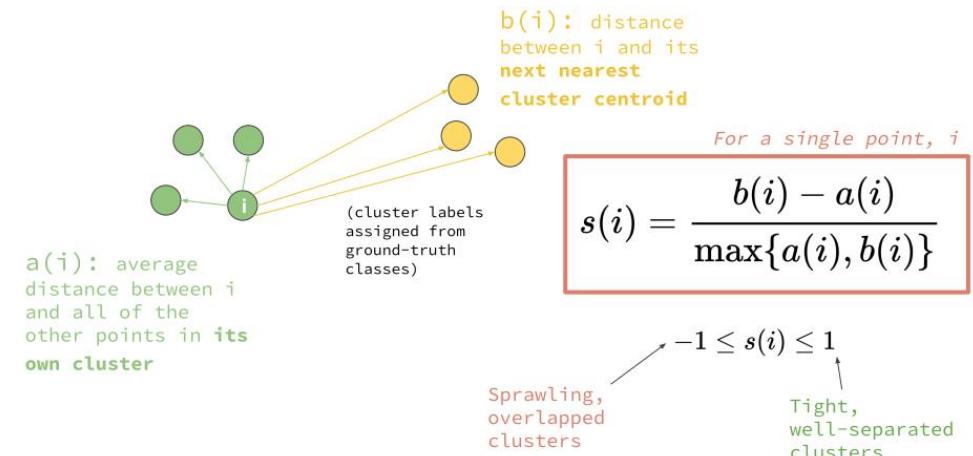
How do we know the clustering is correct?

- We don't!
 - Clustering is inherently 'unsupervised', meaning we don't have a known answer to check against
 - Clustering output is entirely dependent on prior processing steps. Change any upstream parameters and you can have drastically different results
 - No quantification method will capture every possible feature of a cell
 - Example: Biological 'celltypes' were historically defined by surface proteins, NOT transcriptional profile
 - Stable surface protein probably means either:
 - High turn over of protein => high expression level
 - Low turn over of protein => low expression level

How do we know the clustering is correct?

- We don't!
- But we can evaluate similarity within clusters:
 - (medoid) silhouette width – evaluates the mean distance to same-cluster medoid versus nearest non-same cluster medoid
 - Root Mean Squared Deviation – ‘noise’/dispersion within cluster
 - Gap statistic – ratio of assigned within-cluster dispersion versus null reference distribution

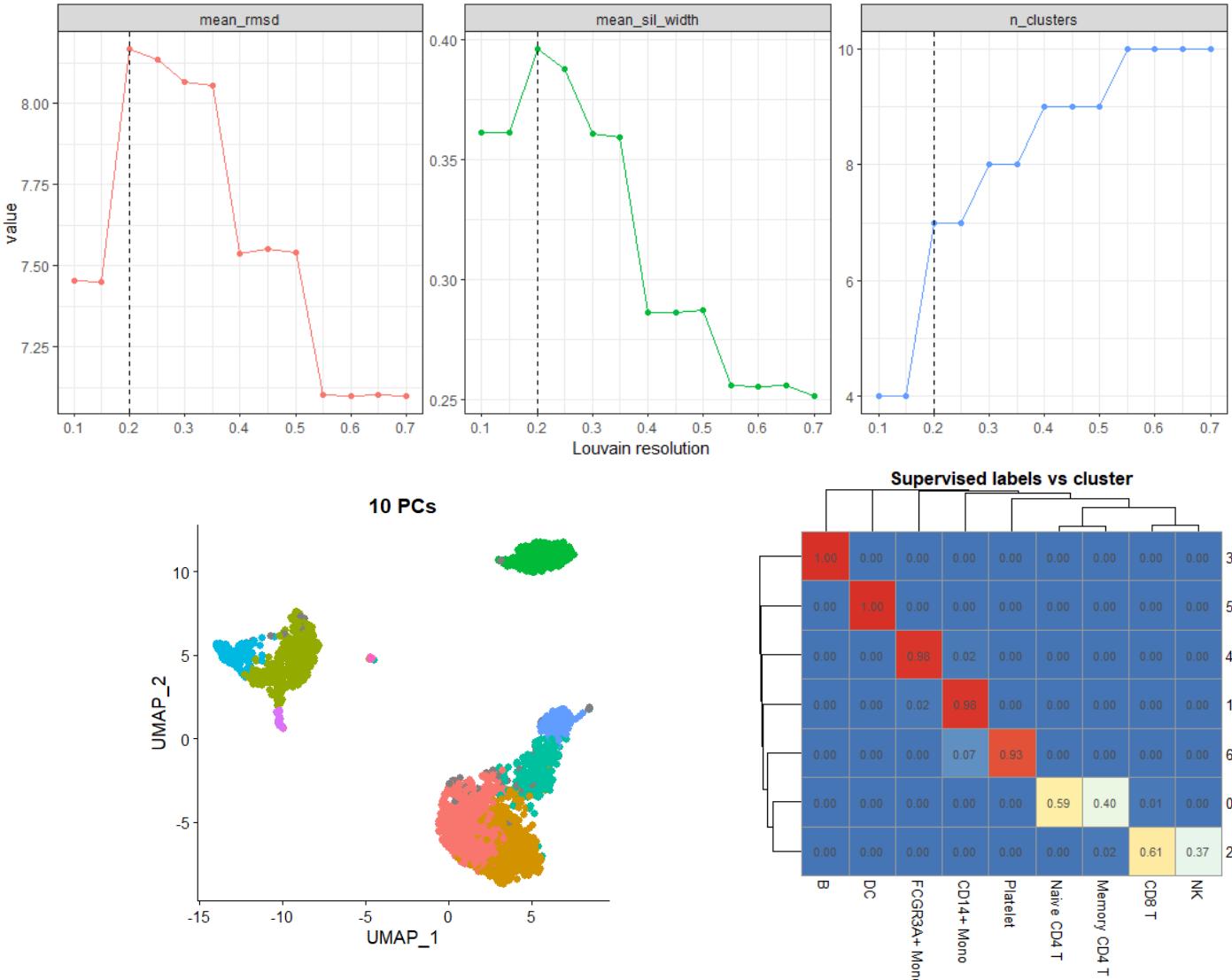
Silhouette width example



<https://www.platform.ai/post/the-silhouette-loss-function-metric-learning-with-a-cluster-validity-index>

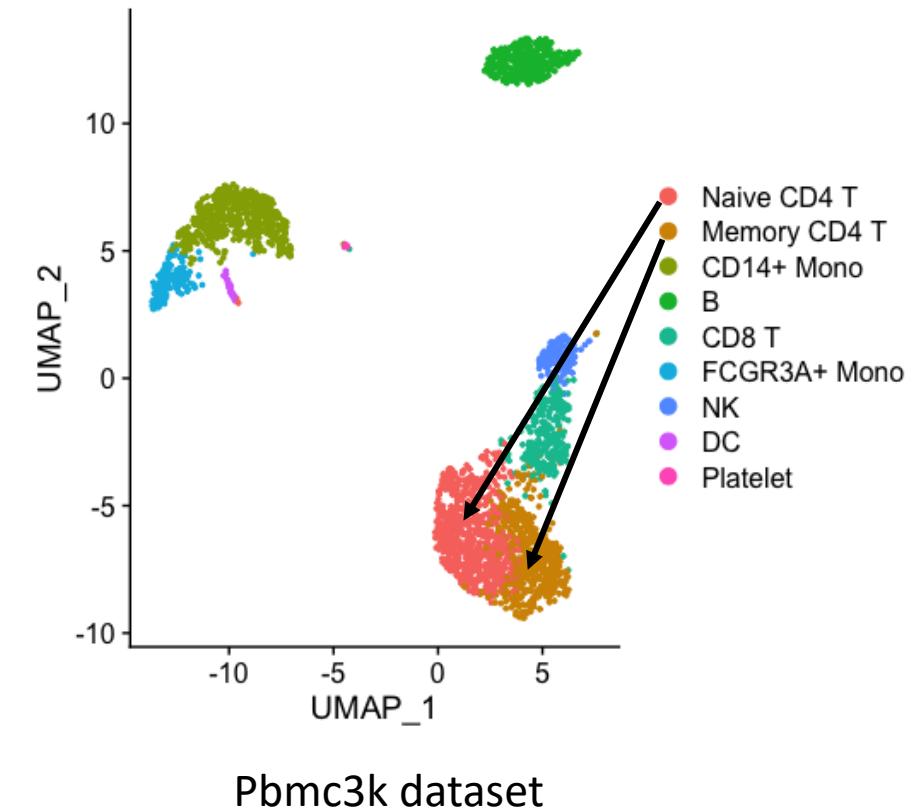
Clustering metric example: Seurat pbmc3k data

1. Normalize/scale
data/variable/features/PC
A already run
 2. Recluster at multiple
resolutions
 3. Evaluate RMSD and
medoid silhouette width
at each clustering
resolution



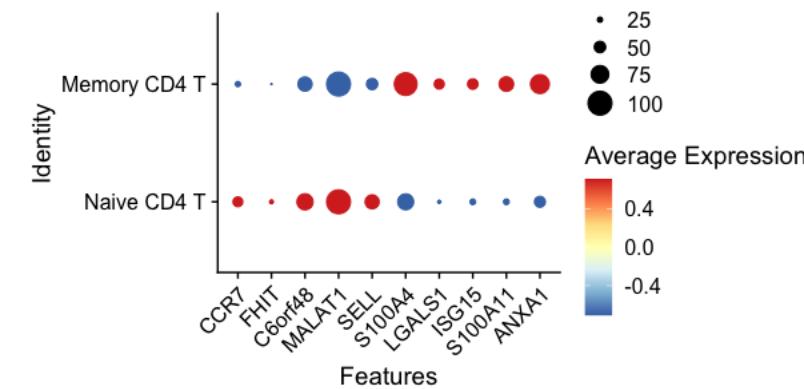
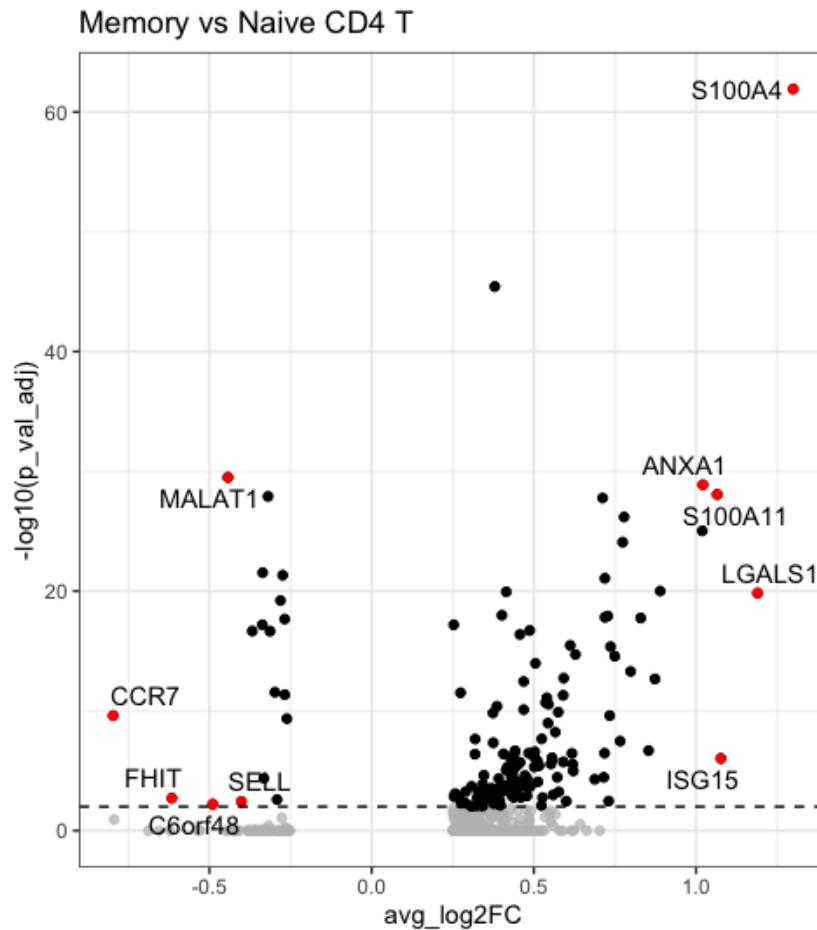
Differential expression analysis – what's different?

- Input
 - Gene expression matrix (typically at the normalized level)
 - Query label (cluster/celltype you're testing)
 - **Reference label** (cluster(s)/celltype(s) you're testing against)
- Output
 - Gene level log-fold change
 - Significance of difference
 - Fraction of query/reference expressing(>0 UMI)

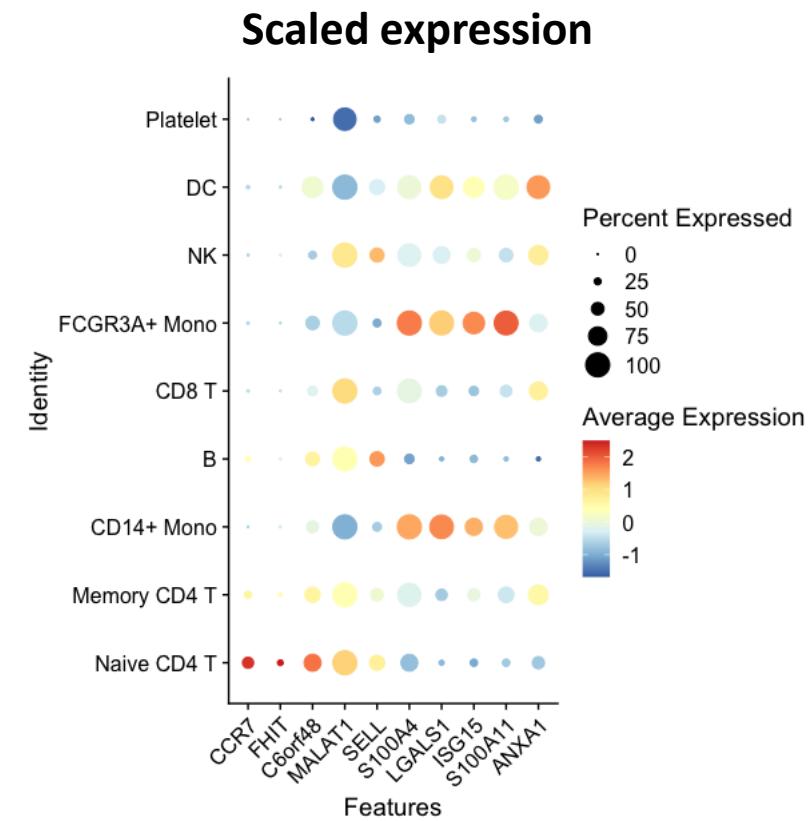
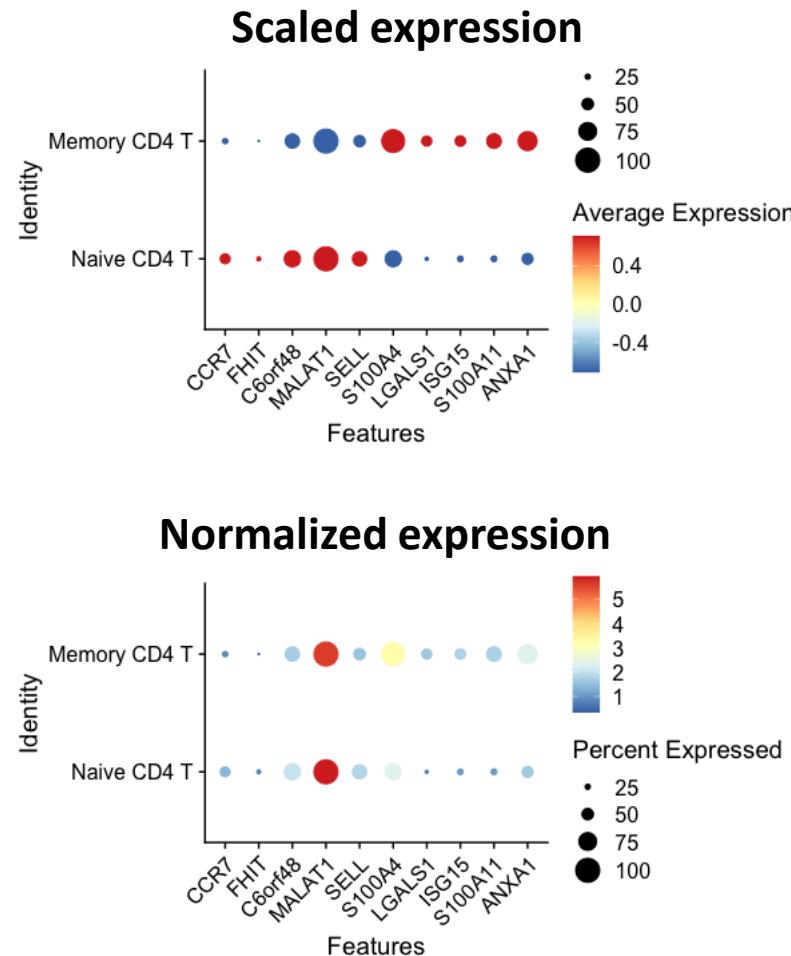


Differential expression analysis

- Common simple version:
 - Wilcoxon test
- How to visualize:
 - Volcano plot
 - Heatmap/dotplot

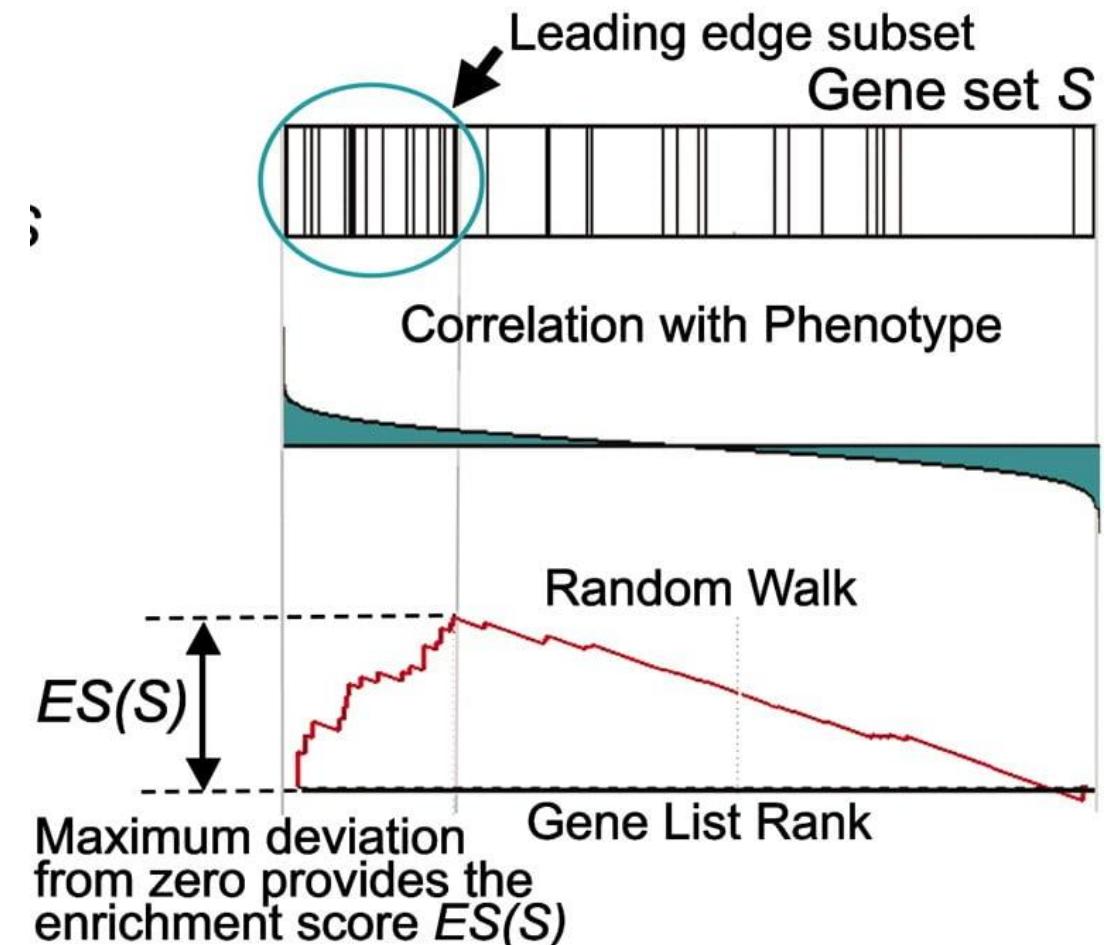


'Absolute' gene expression versus relative



Pathway analysis – leverage prior knowledge

- Example: Gene Set Enrichment Analysis
 - Input:
 - Query: List of genes ordered by (relative) expression
 - Reference: List of unordered genes, all involved in some biological process
 - Output:
 - Enrichment score
 - Significance
 - Goal:
 - Identify how overrepresented the query genes are in the query data set



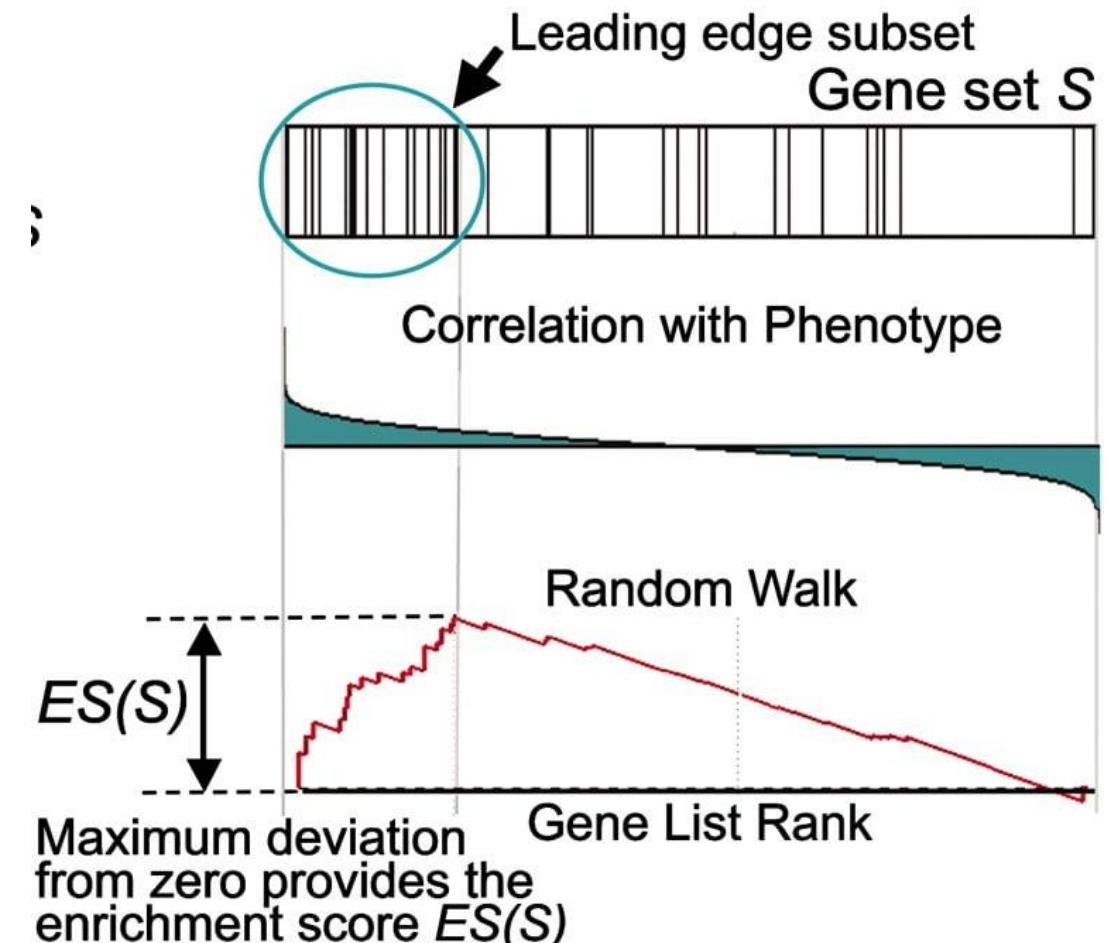
GSEA approach

1. Order your query gene list by (relative) expression
2. Annotate whether each gene is in the reference gene list
3. Compute running enrichment score
 - A. Sweep high-> low expression
 - B. If a gene is in the reference gene list then:

$$E.Score_{x+1} = E.Score_x + \frac{1}{n_{reference}}$$

Otherwise if gene is not in reference:

$$E.Score_{x+1} = E.Score_x + \frac{-1}{n_{query} - n_{reference}}$$



Celltype Identification – helps communicate and interpret findings

Marker gene database

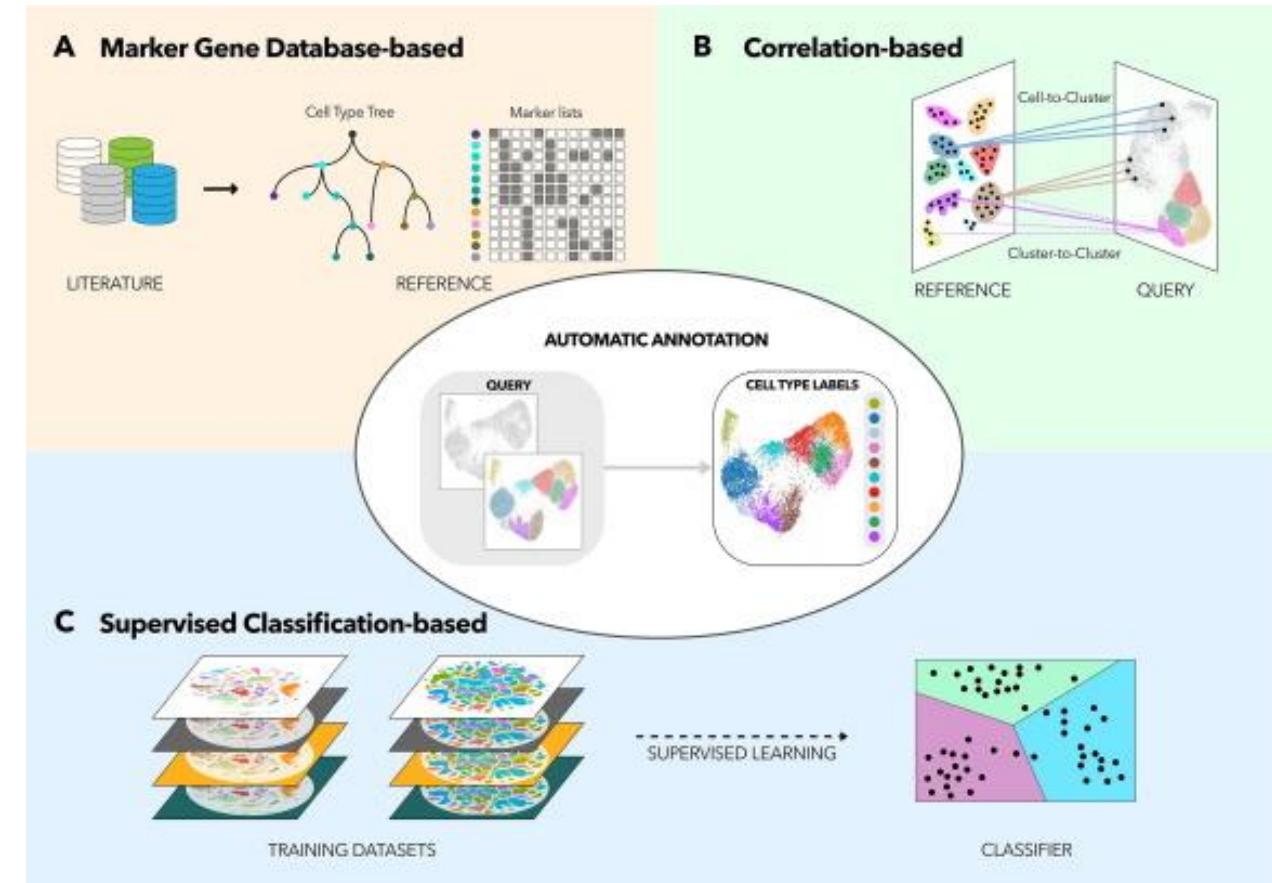
- Match expression via some scoring system

Correlation-based

- Use multiple correlation or mapping metrics to transfer labels from a reference data set to the query data set

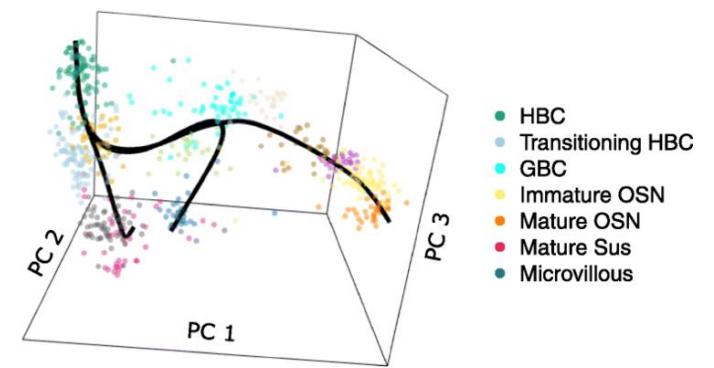
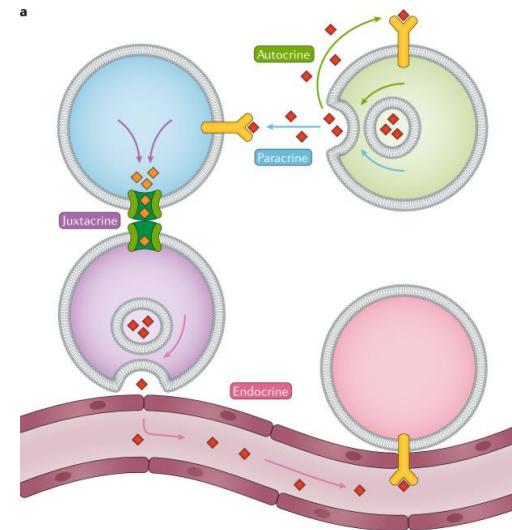
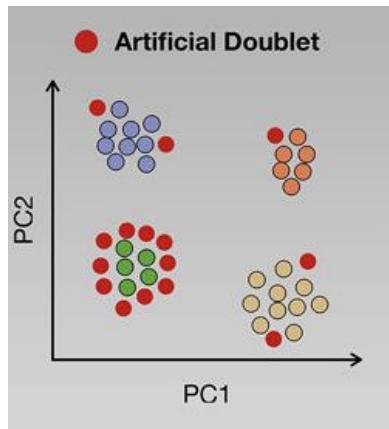
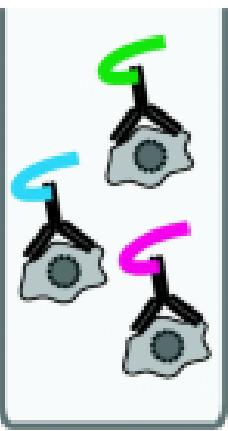
Supervised Classification-based

- Some model is trained to identify cell type, and then used on your query data to apply a celltype label



Automated methods for cell type annotation on scRNA-seq data
10.1016/j.csbj.2021.01.015

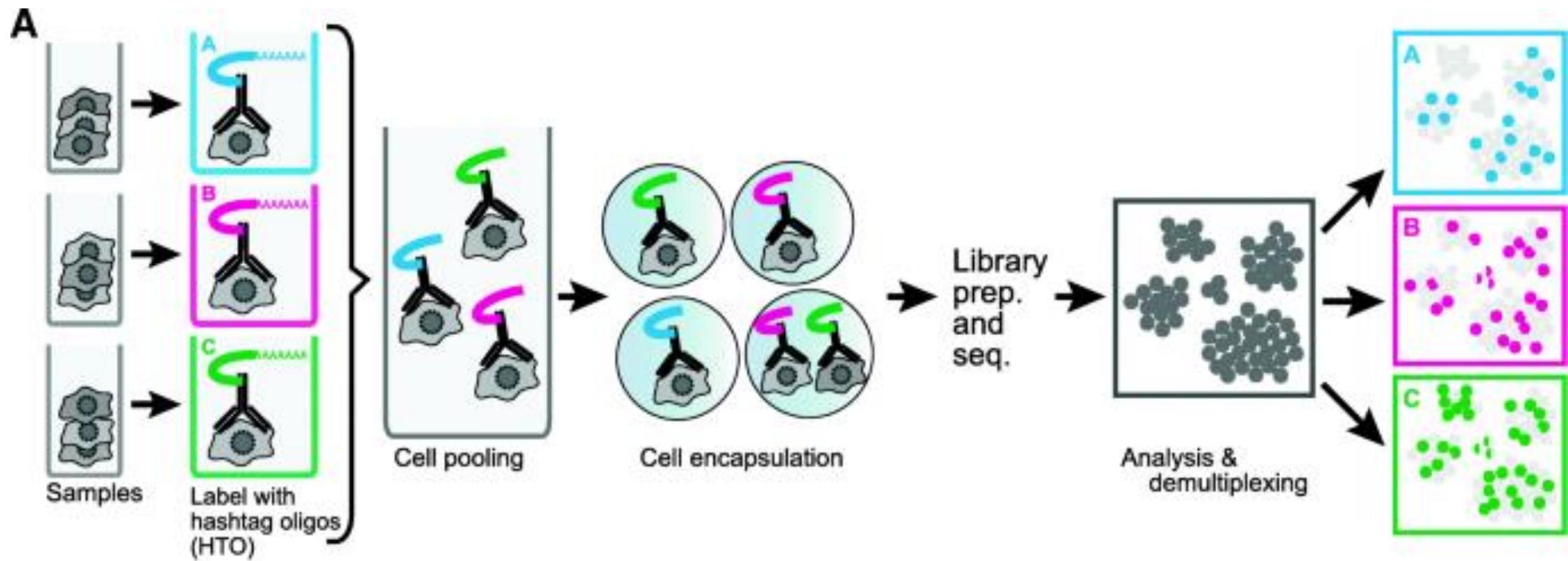
Advanced methods in scRNAseq Analysis



Topics:

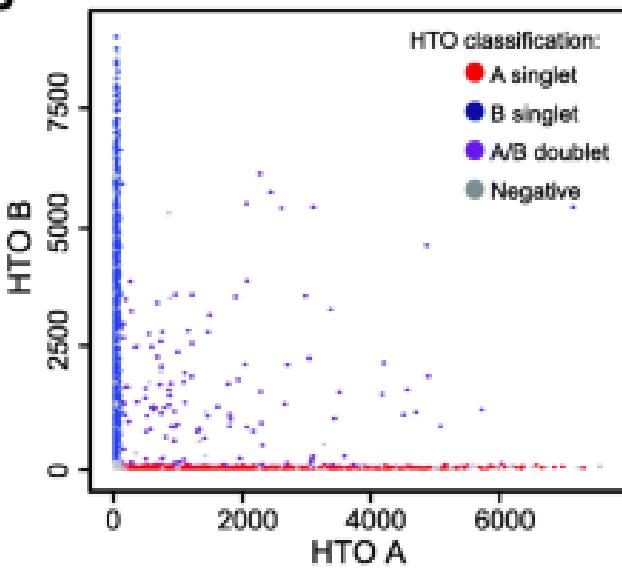
- Hashtag demultiplexing
- In silico doublet identification
- Interaction analysis
- Trajectory Inference
- Cross species analysis vignette (classifier & integration)

Sample Multiplexing with Hashtag Oligonucleotides

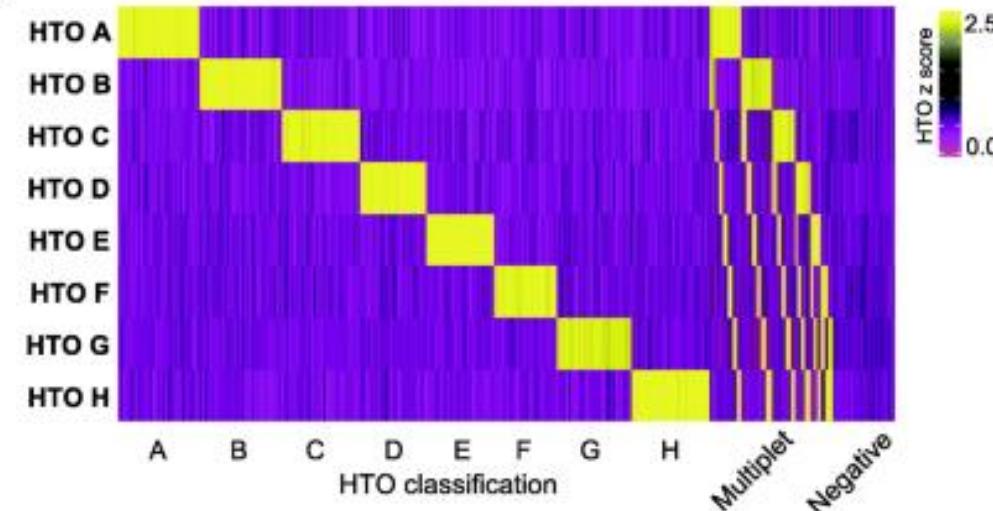


Demultiplexing QC plots

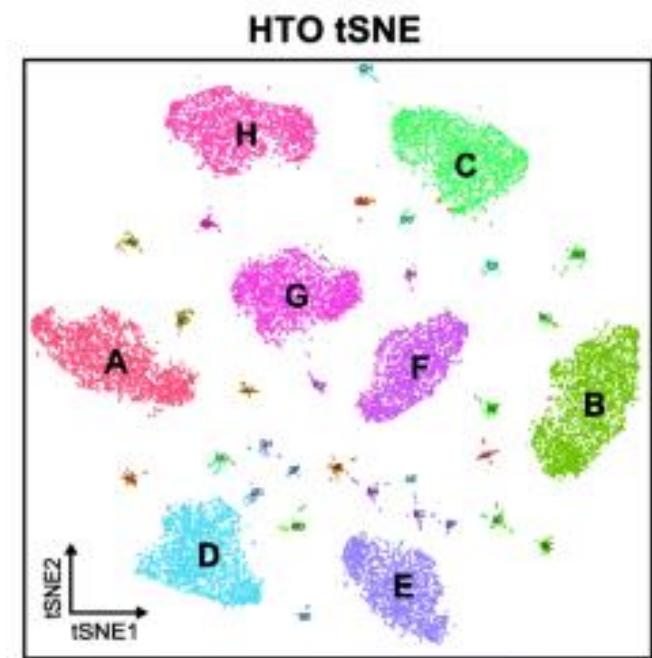
B



C



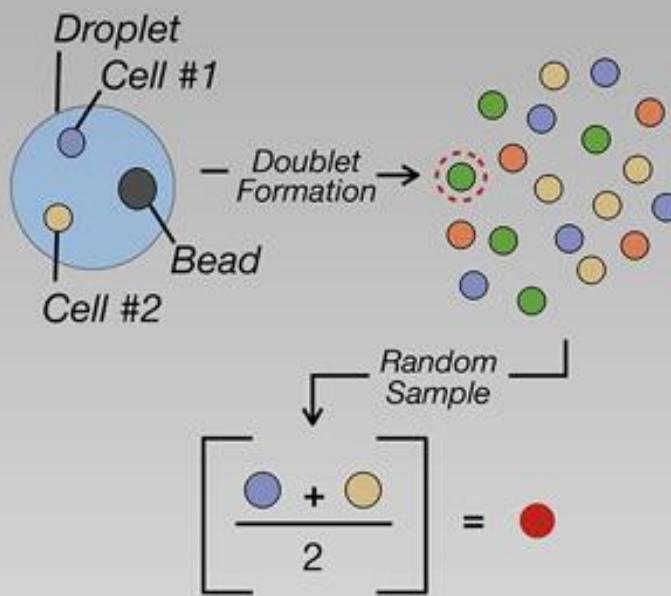
D



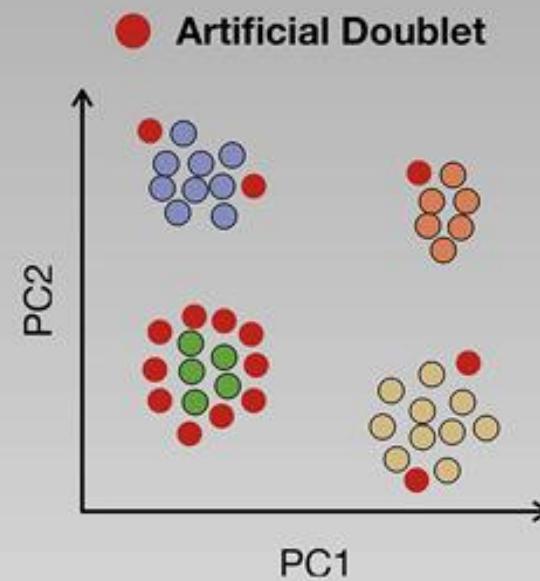
In-silico doublet identification

DoubletFinder Overview

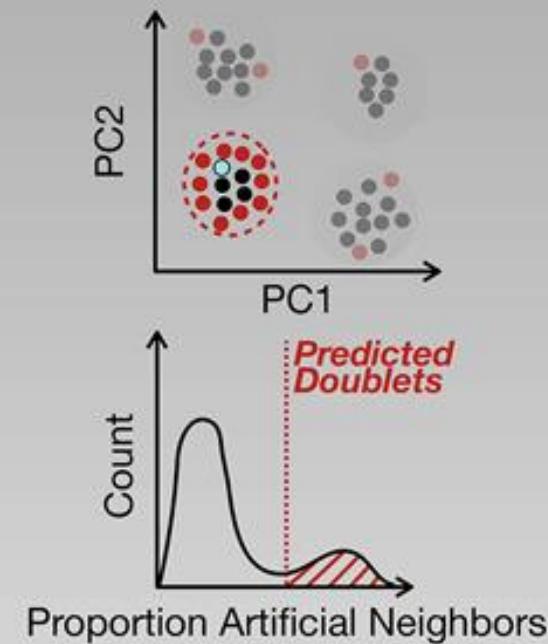
(1) Simulate Doublets



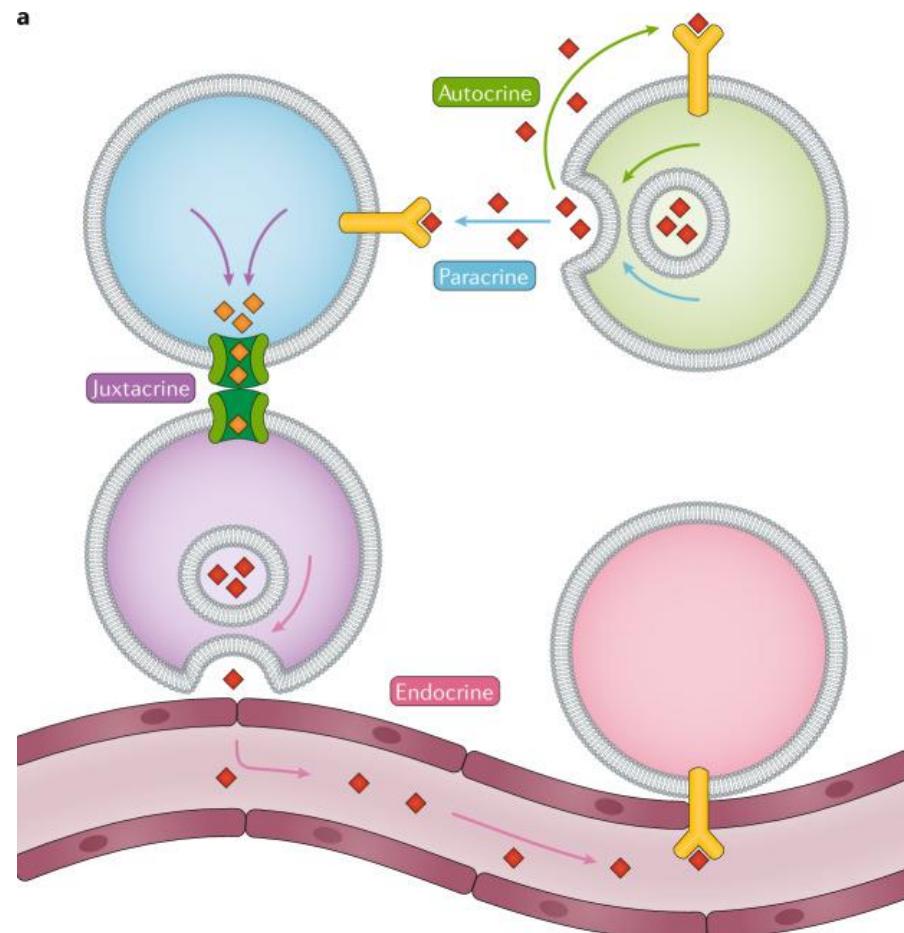
(2) Dimensionality Reduction



(3) Doublet Identification



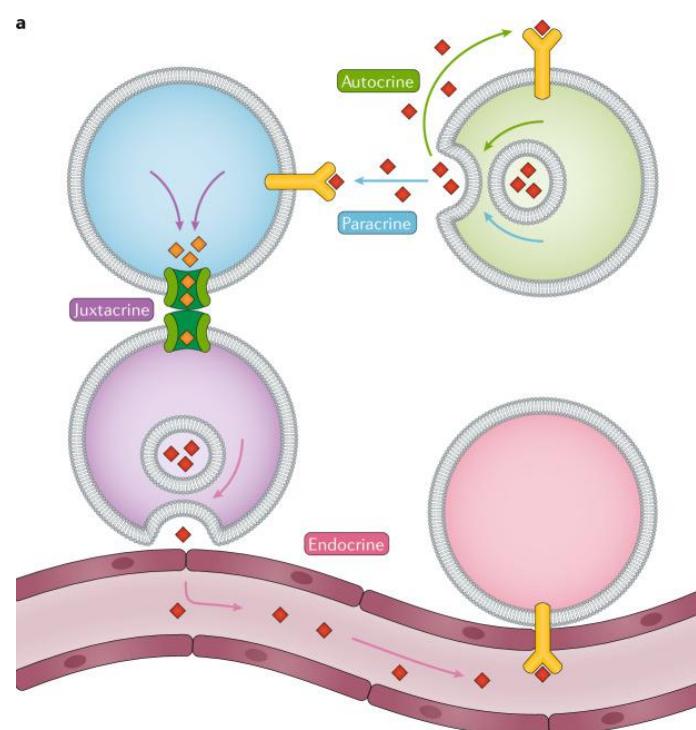
Interpreting cell-cell interactions with scRNA-seq



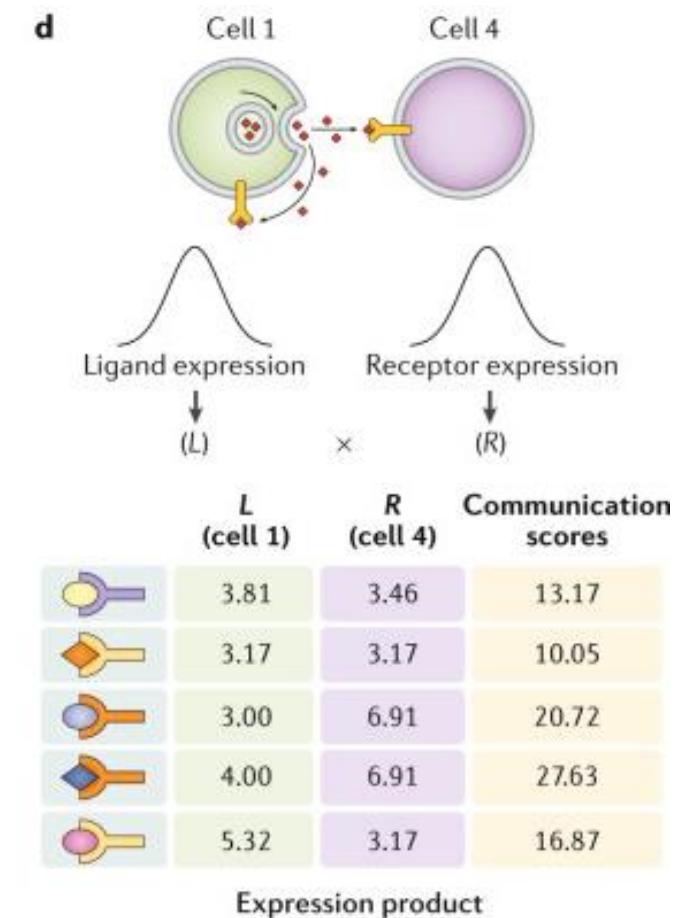
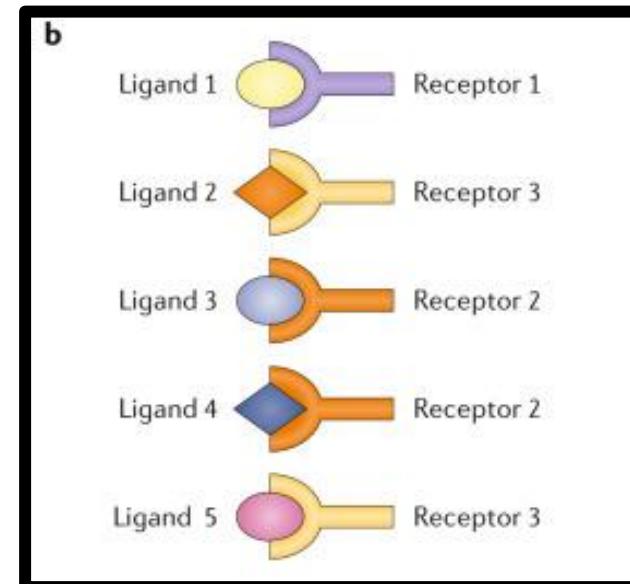
Ligand – signaling molecule

Receptor - target of ligand binding,
initiates signaling cascade

Interpreting cell-cell interactions with scRNA-seq



Dictionary of ligand-receptor interactions

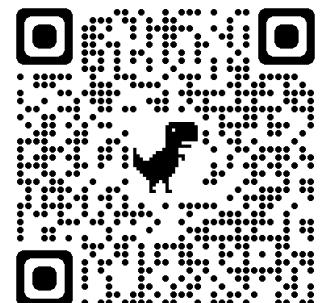


Trajectory inference – identifying paths of continuity through similarity

- [prior trajectory inference talk](#)

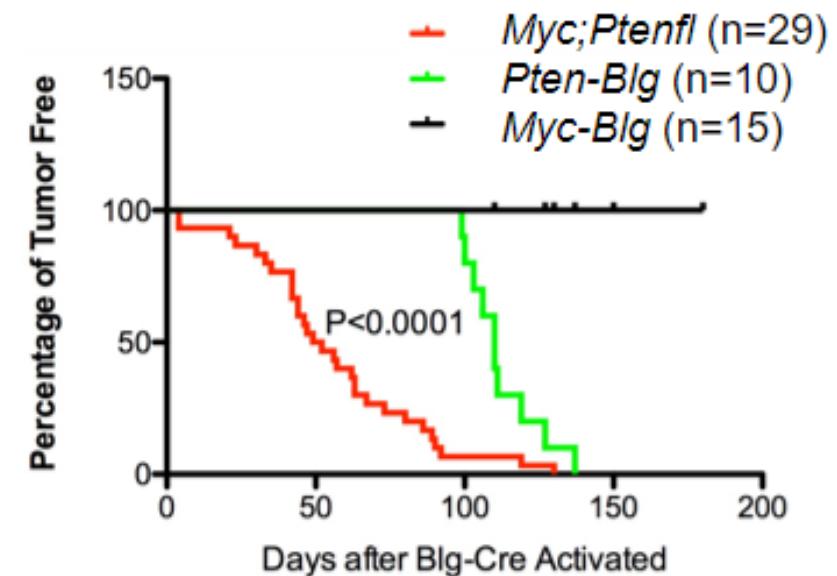
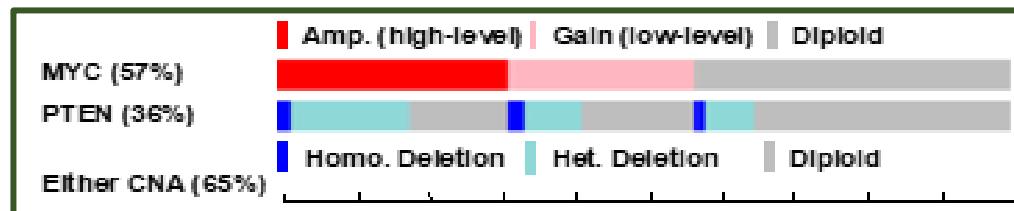
Cross species analysis vignette

Following data from: Doha, Wang, Calistri et al. MYC Deregulation and PTEN Loss Model Tumor and Stromal Heterogeneity of Aggressive Triple-Negative Breast Cancer. *Nat Commun* **14**, 5665 (2023). <https://doi.org/10.1038/s41467-023-40841-6>



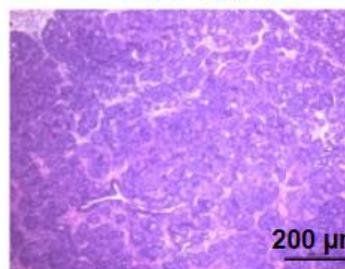
Myc;Ptenfl tumor phenotypes correspond to specific histologic subtypes

Cohort of 309 human TNBCs

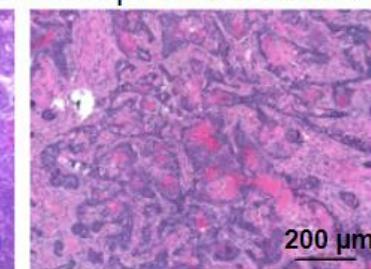


Myc;Ptenfl Stroma Poor

IDC-Solid

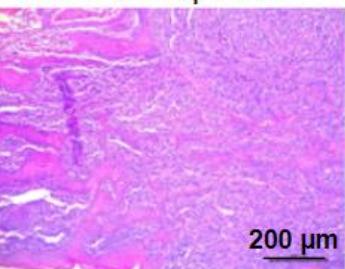


IDC- squamous features

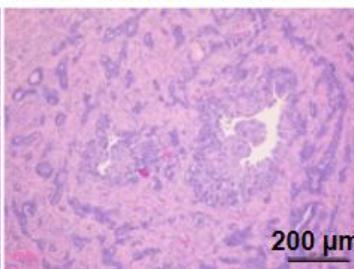


Myc;Ptenfl Stroma Rich

IDC- metaplastic



IDC- lobular features



200 μ m

200 μ m

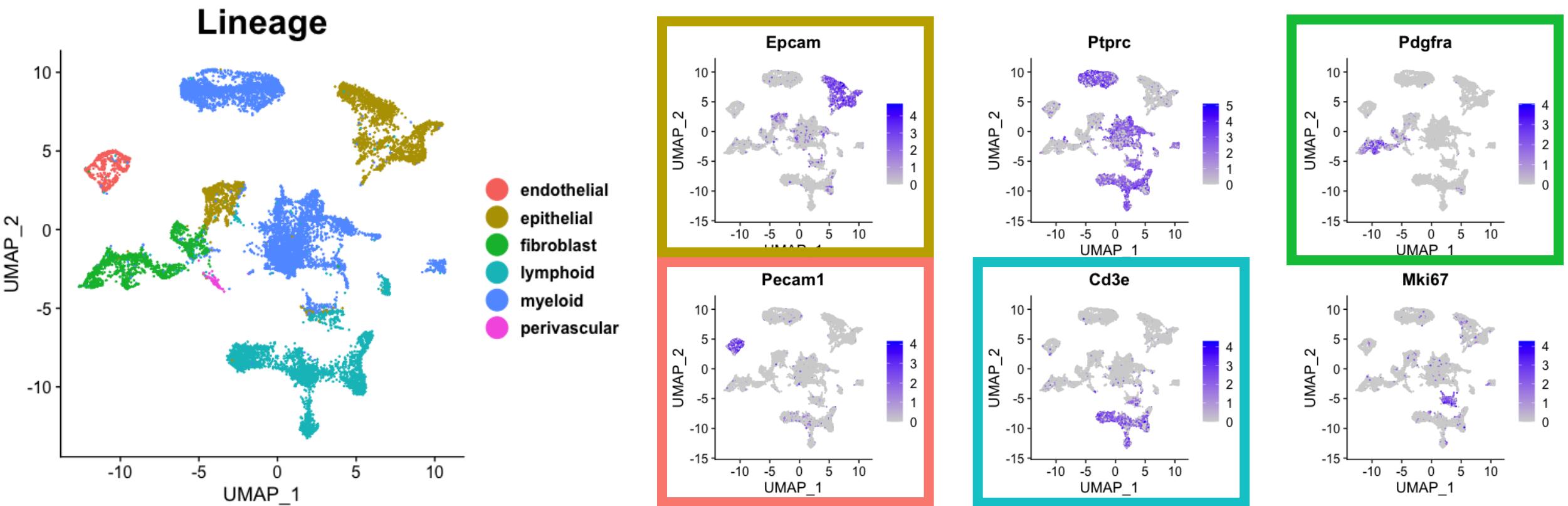
200 μ m

200 μ m

Myc;Ptenfl(cluster1)

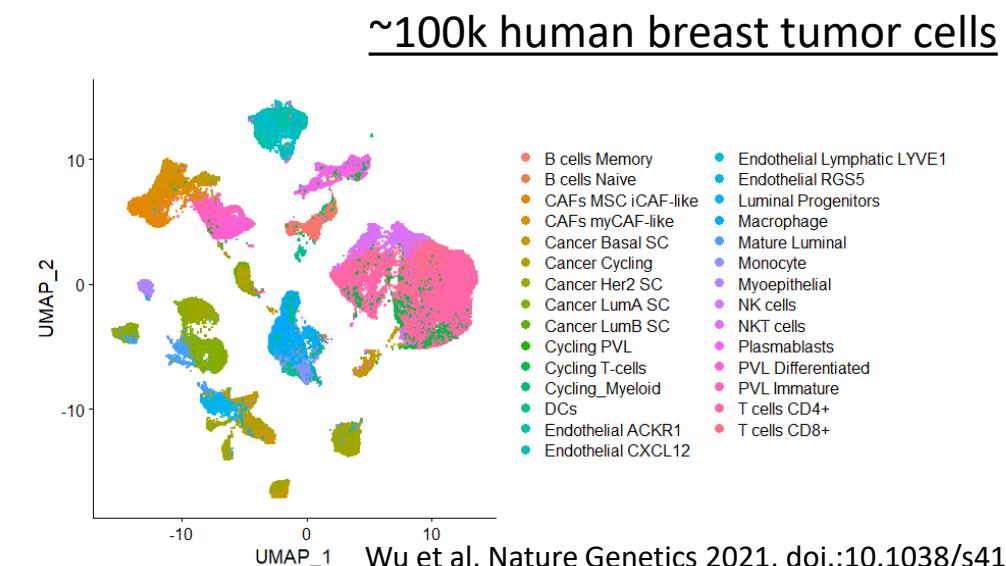
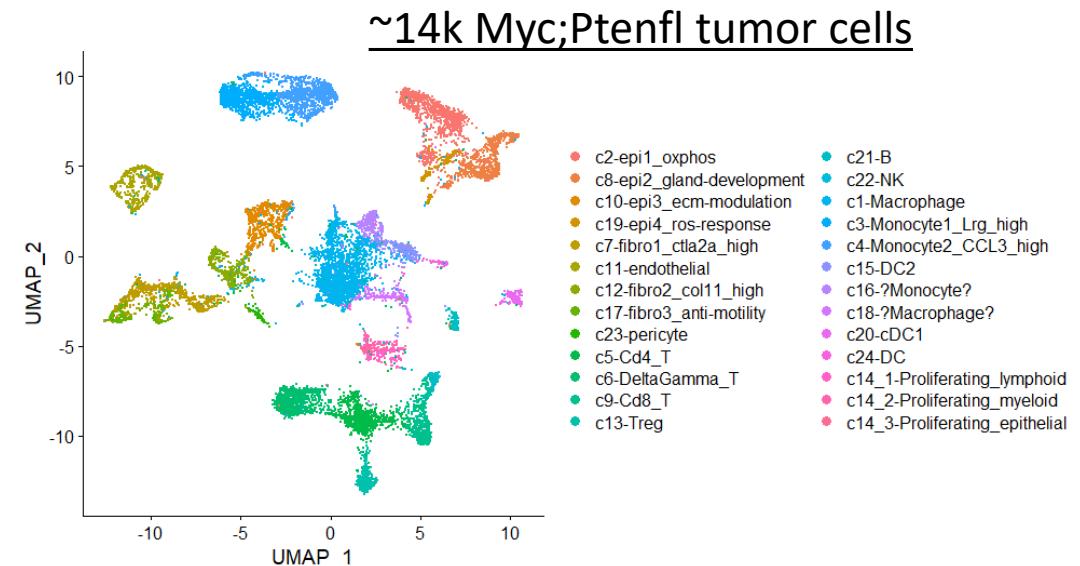
Myc;Ptenfl(cluster-2)

Clusters aligned with known lineages



That's neat.... But is it clinically relevant?

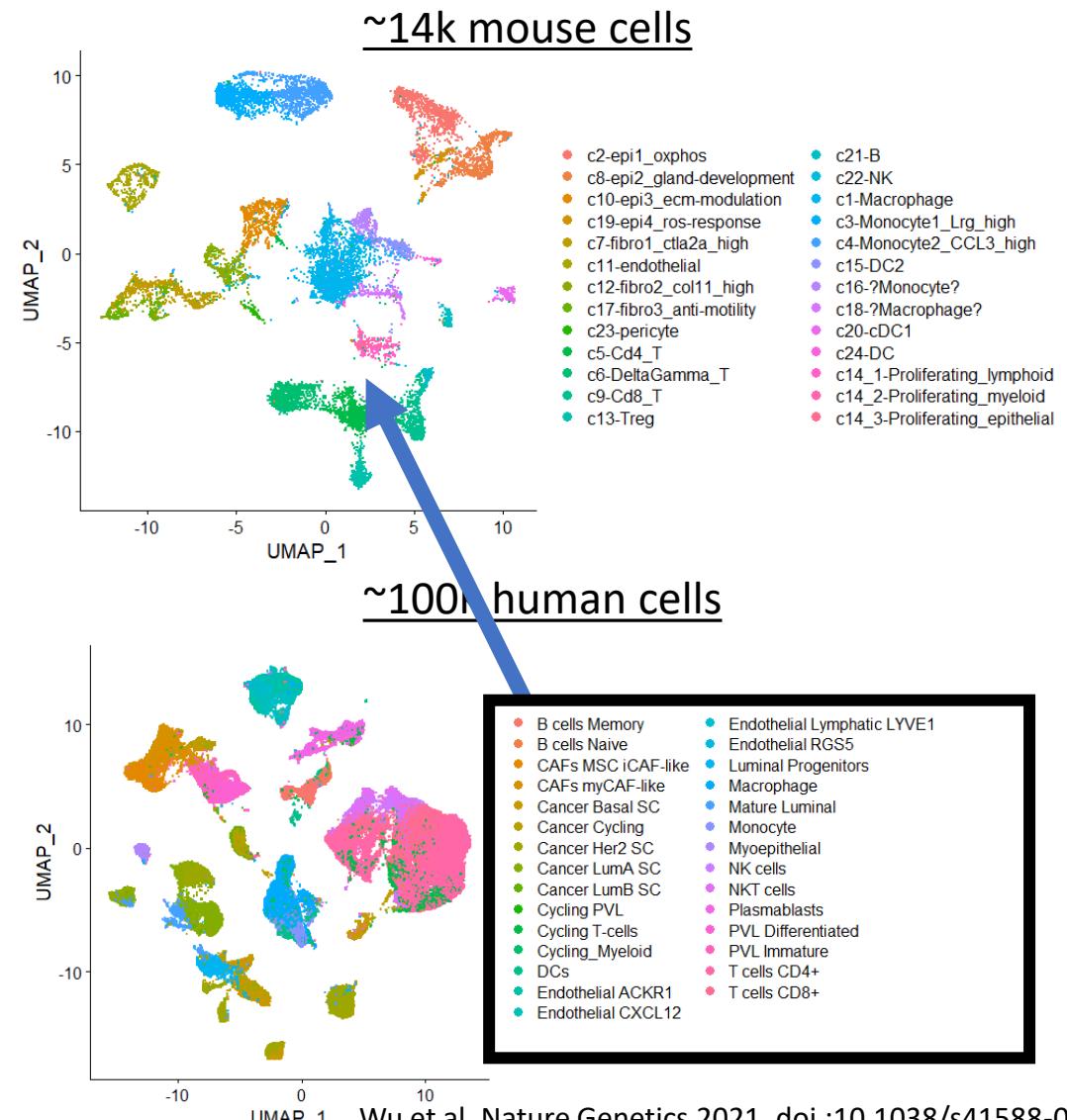
Are the mouse cell states representative of human disease?



Are the mouse cell states representative of human disease?

1. Classification (Supervised)

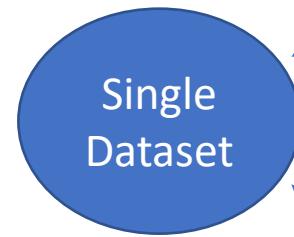
- Train a classifier on one species using a subset of shared features and apply it to the other



Are the mouse cell states representative of human disease?

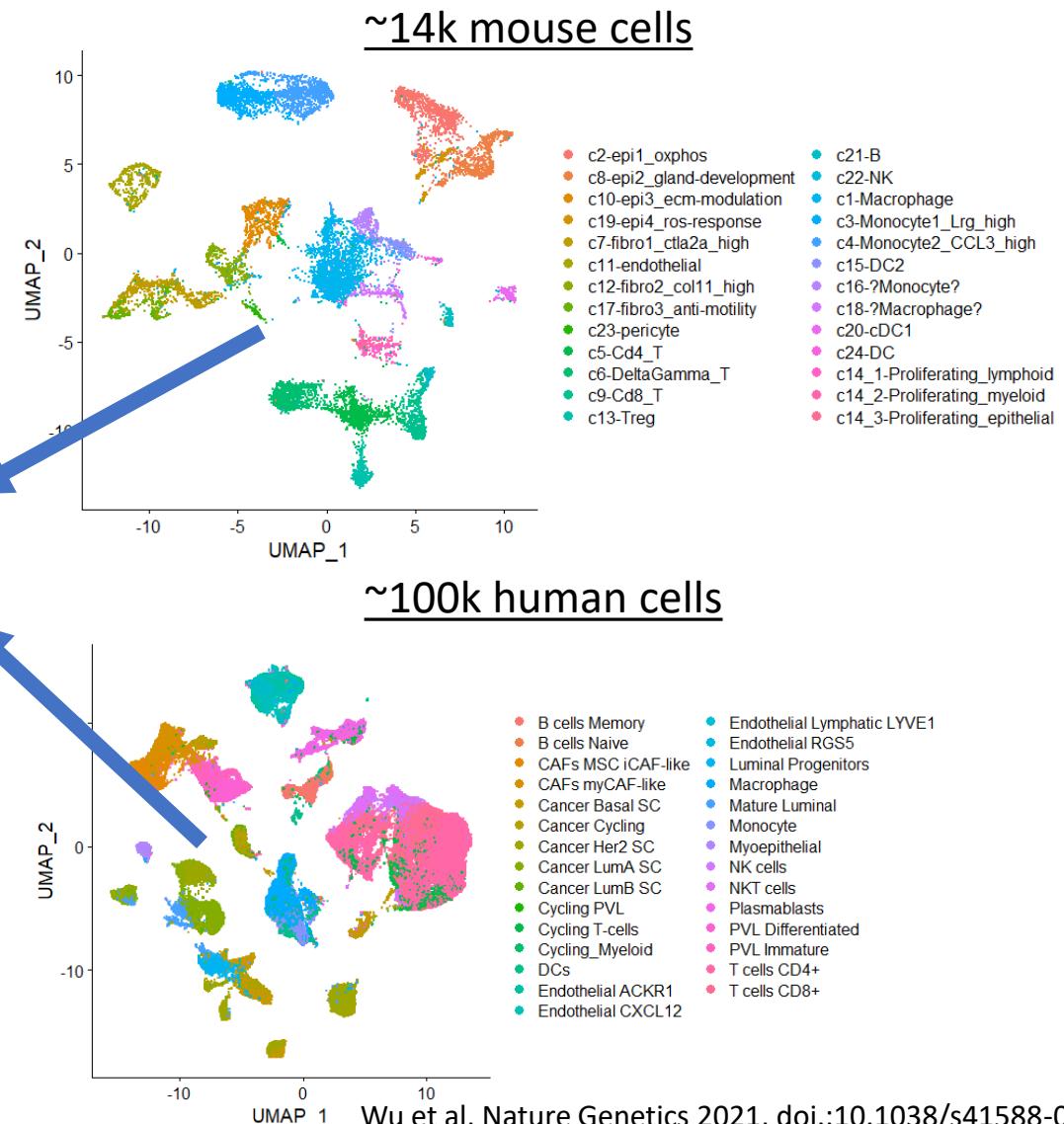
1. Classification (Supervised)

- Train a classifier on one species using a subset of shared features and apply it to the other

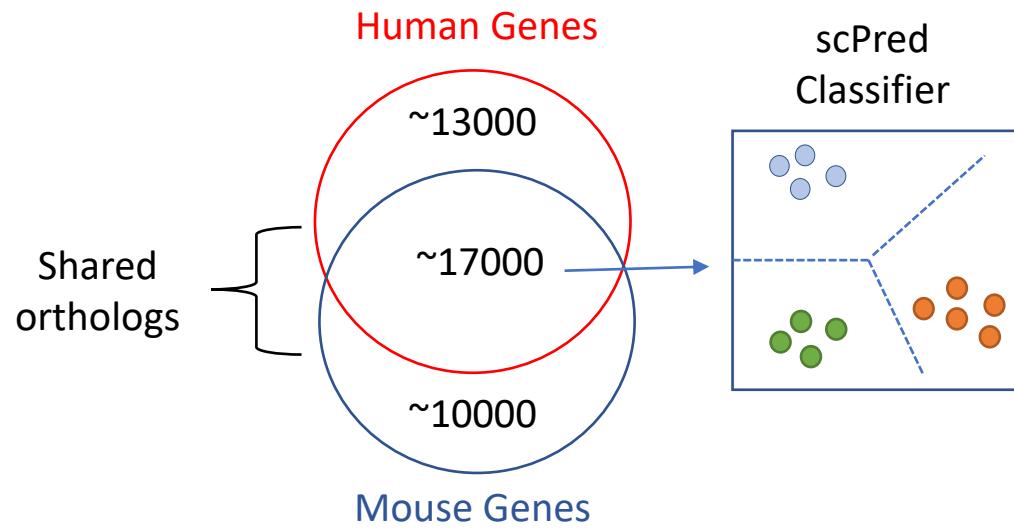


2. Data integration (Unsupervised)

- Species-normalize the data to enable simultaneous analysis within the unified feature space



Classification approach

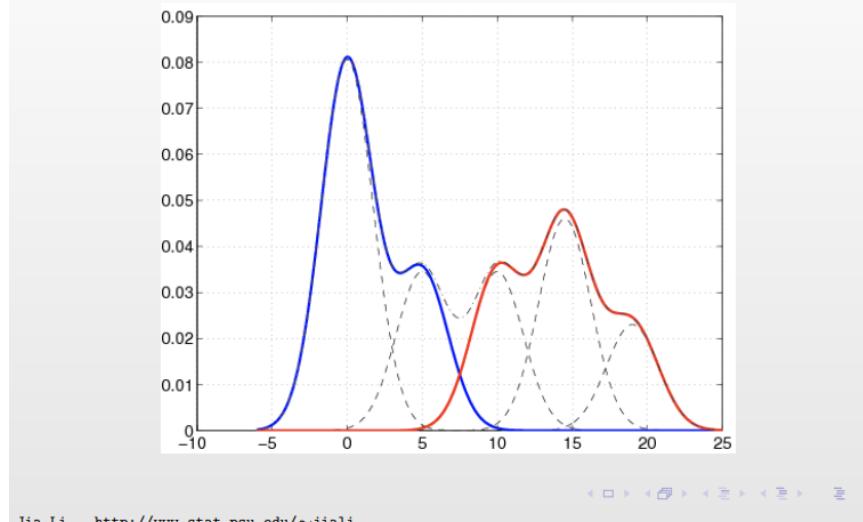


1. Scale and normalize each species data set individually
2. Identify orthologous features (convert mouse features to human where possible)
3. Identify human classes to train (provided by initial publication)
4. Train classifier with scPred to identify these classes using only the shared feature space
5. Apply the classifier to the second species using the same shared feature space

scPred: accurate supervised method for cell-type classification from single-cell RNA-seq data

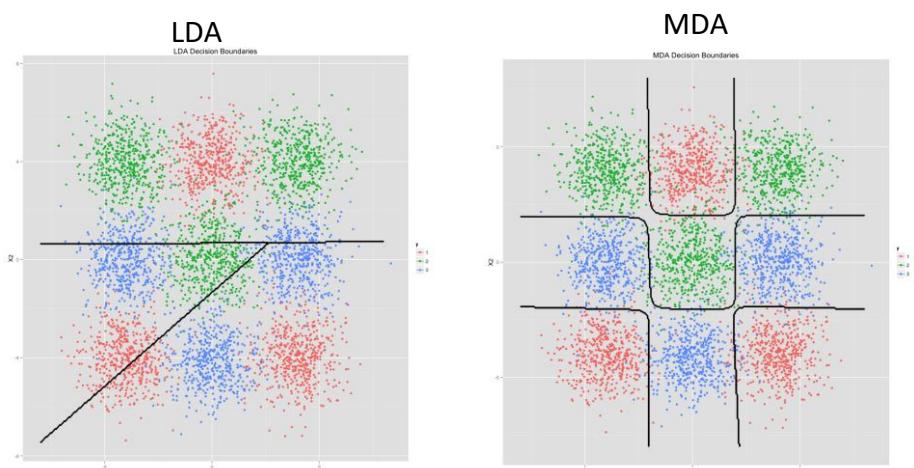
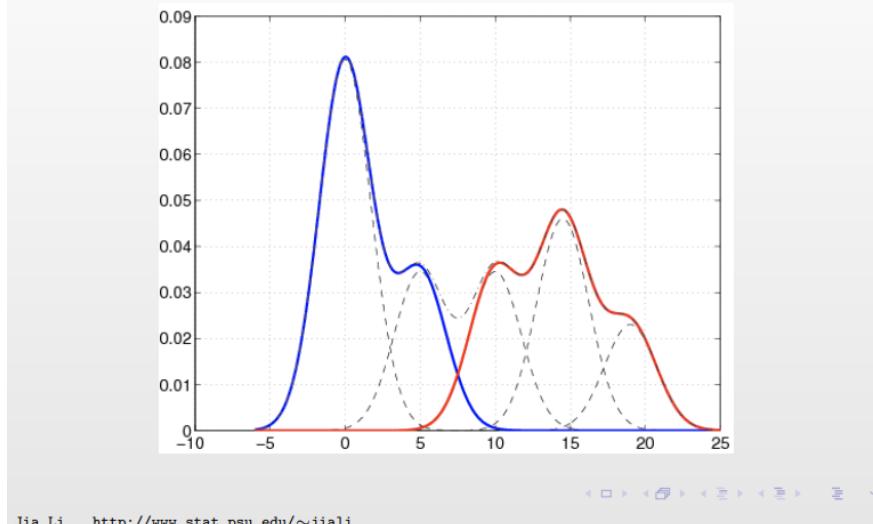
Mixture Discriminant Analysis for classification

A 2-classes example. Class 1 is a mixture of 3 normals and class 2 a mixture of 2 normals. The variances for all the normals are 3.0.



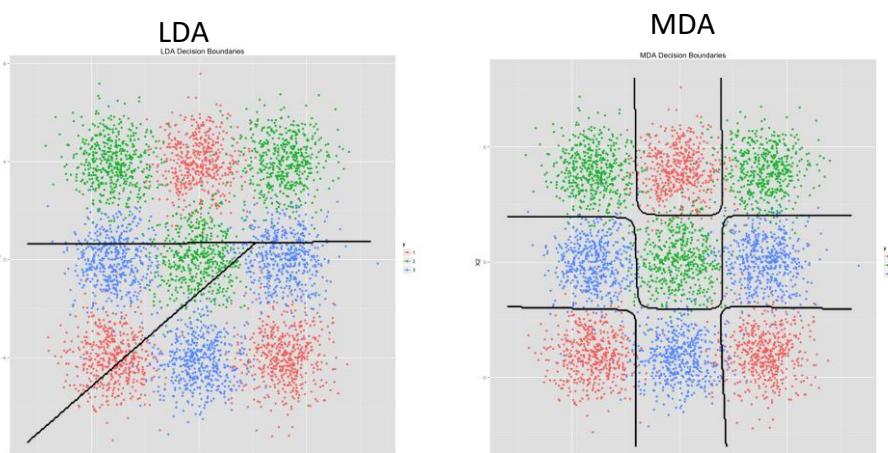
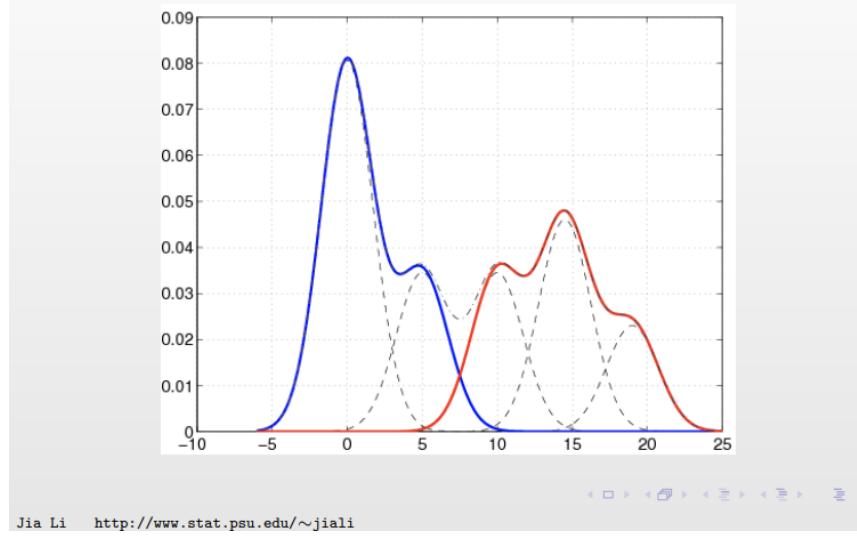
MDA vs LDA

A 2-classes example. Class 1 is a mixture of 3 normals and class 2 a mixture of 2 normals. The variances for all the normals are 3.0.

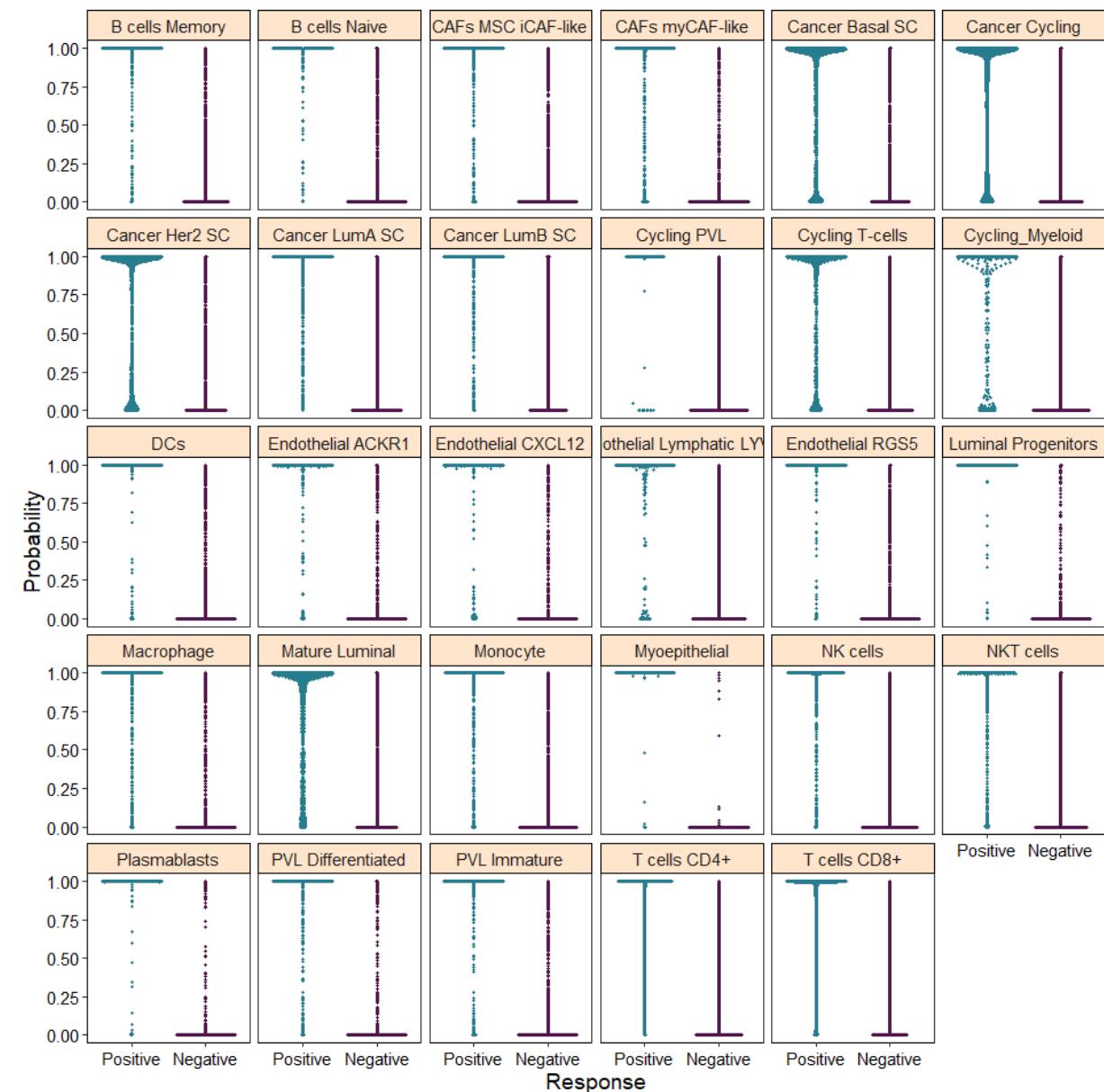


MDA worked well for training human classes

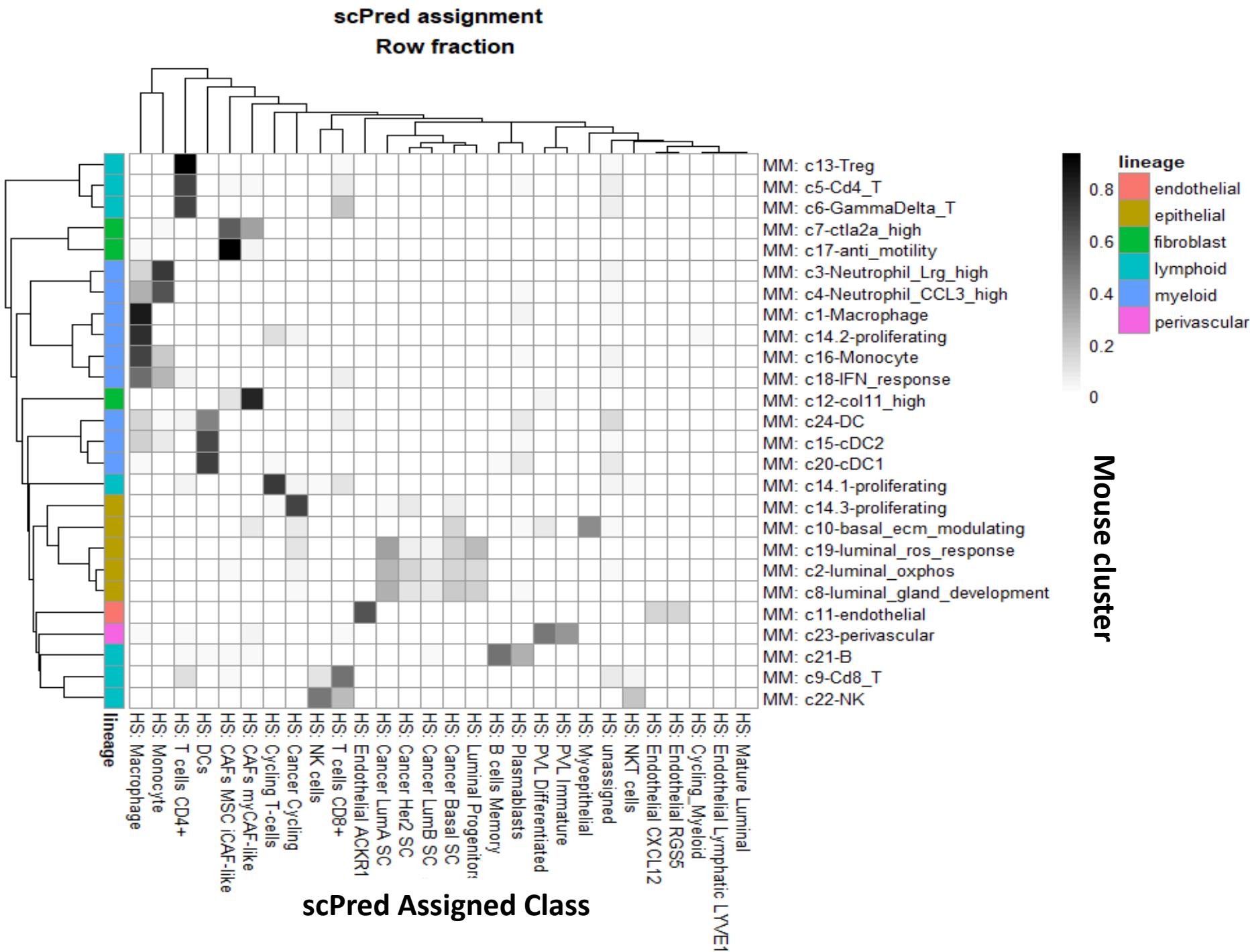
A 2-classes example. Class 1 is a mixture of 3 normals and class 2 a mixture of 2 normals. The variances for all the normals are 3.0.



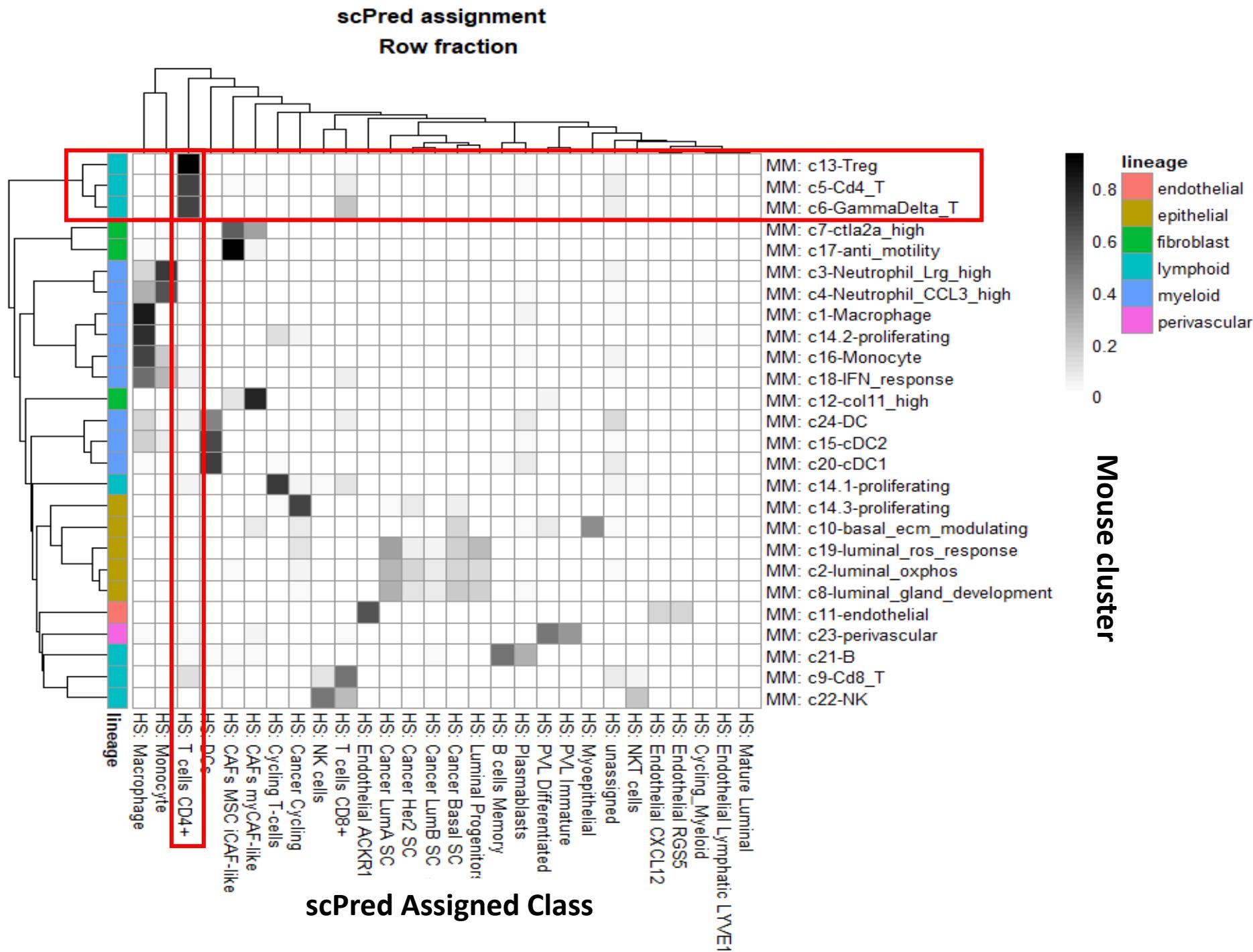
<https://www.r-bloggers.com/2013/07/a-brief-look-at-mixture-discriminant-analysis/>



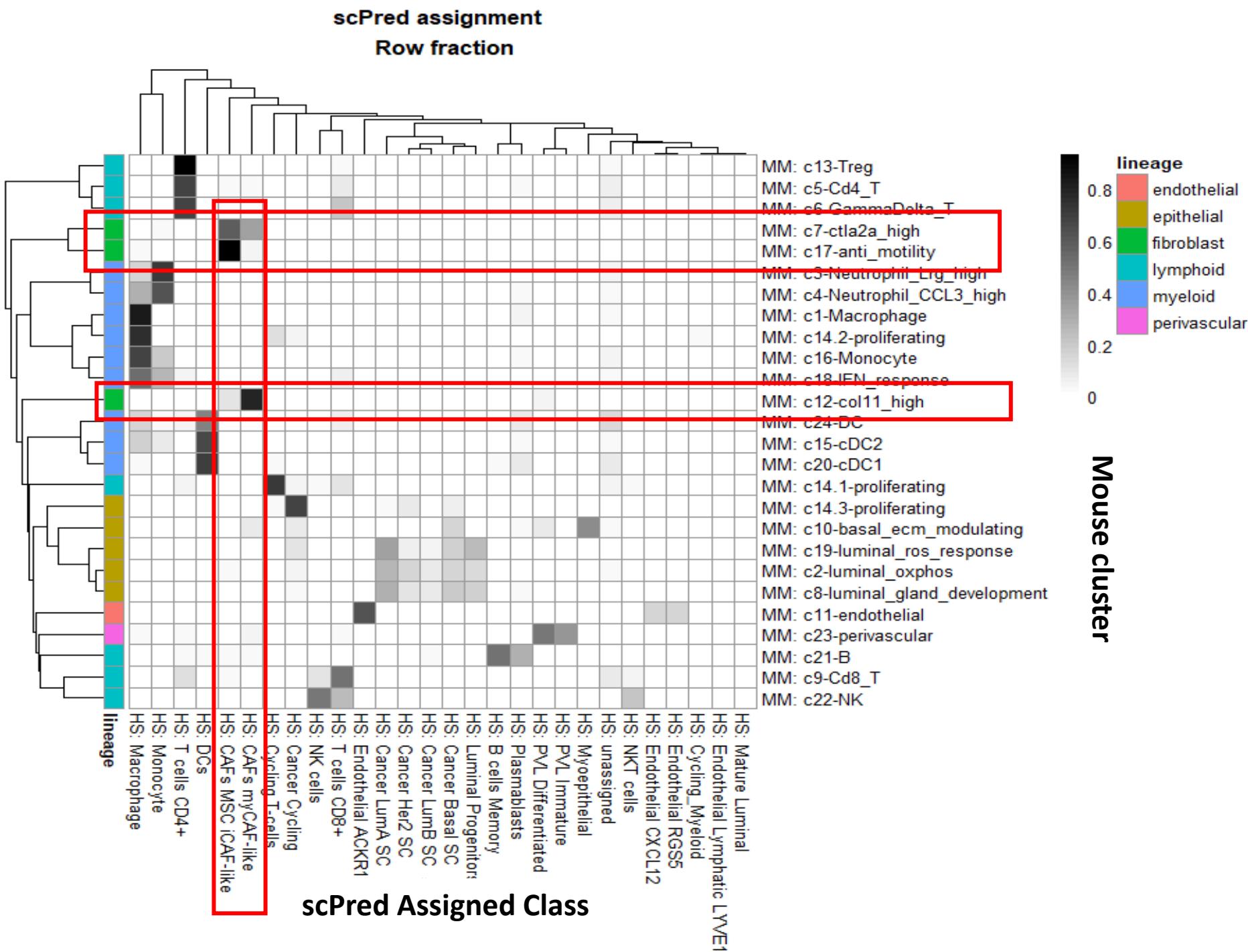
Human
trained
MDA
classifier
applied to
mouse



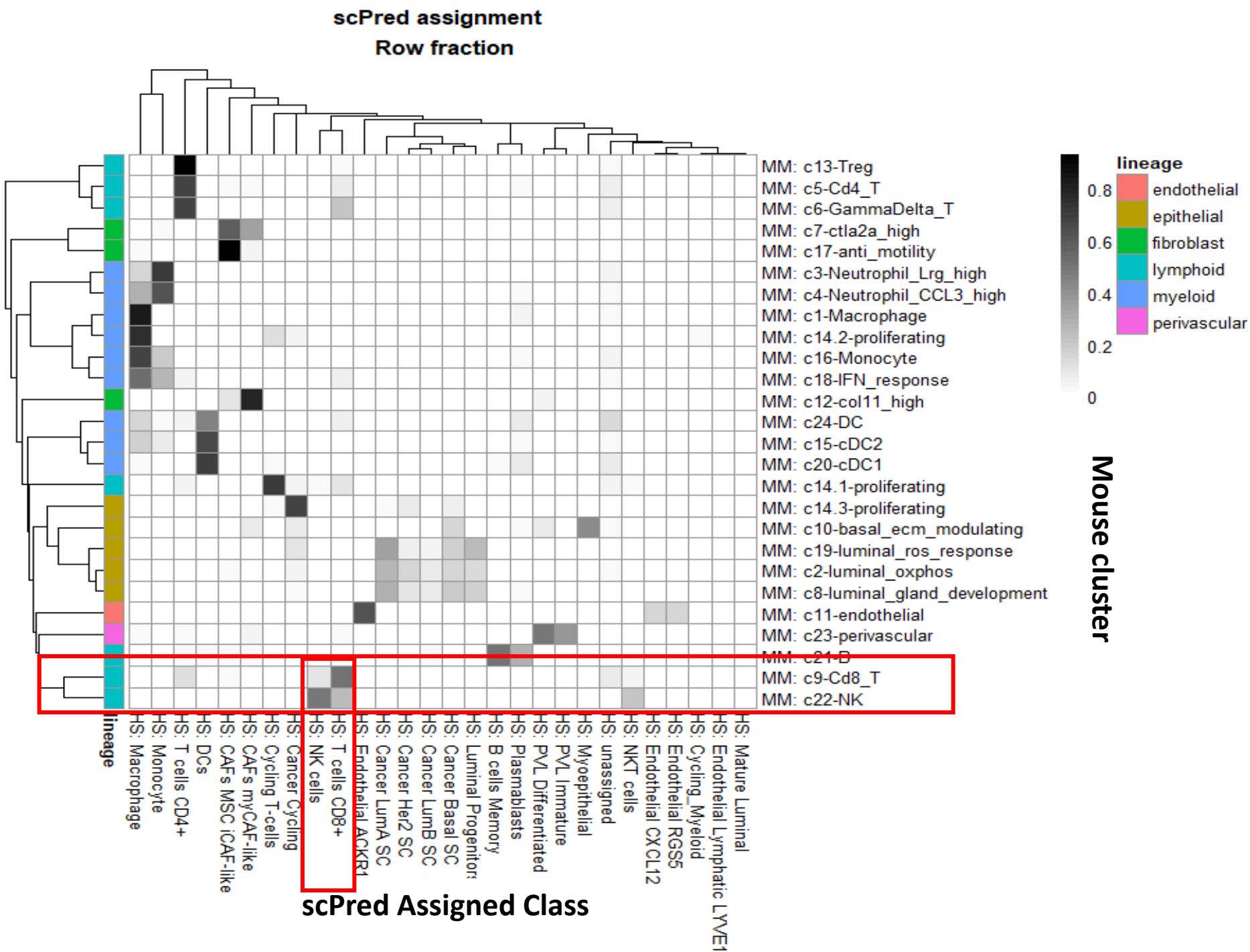
Human
trained
MDA
classifier
applied to
mouse



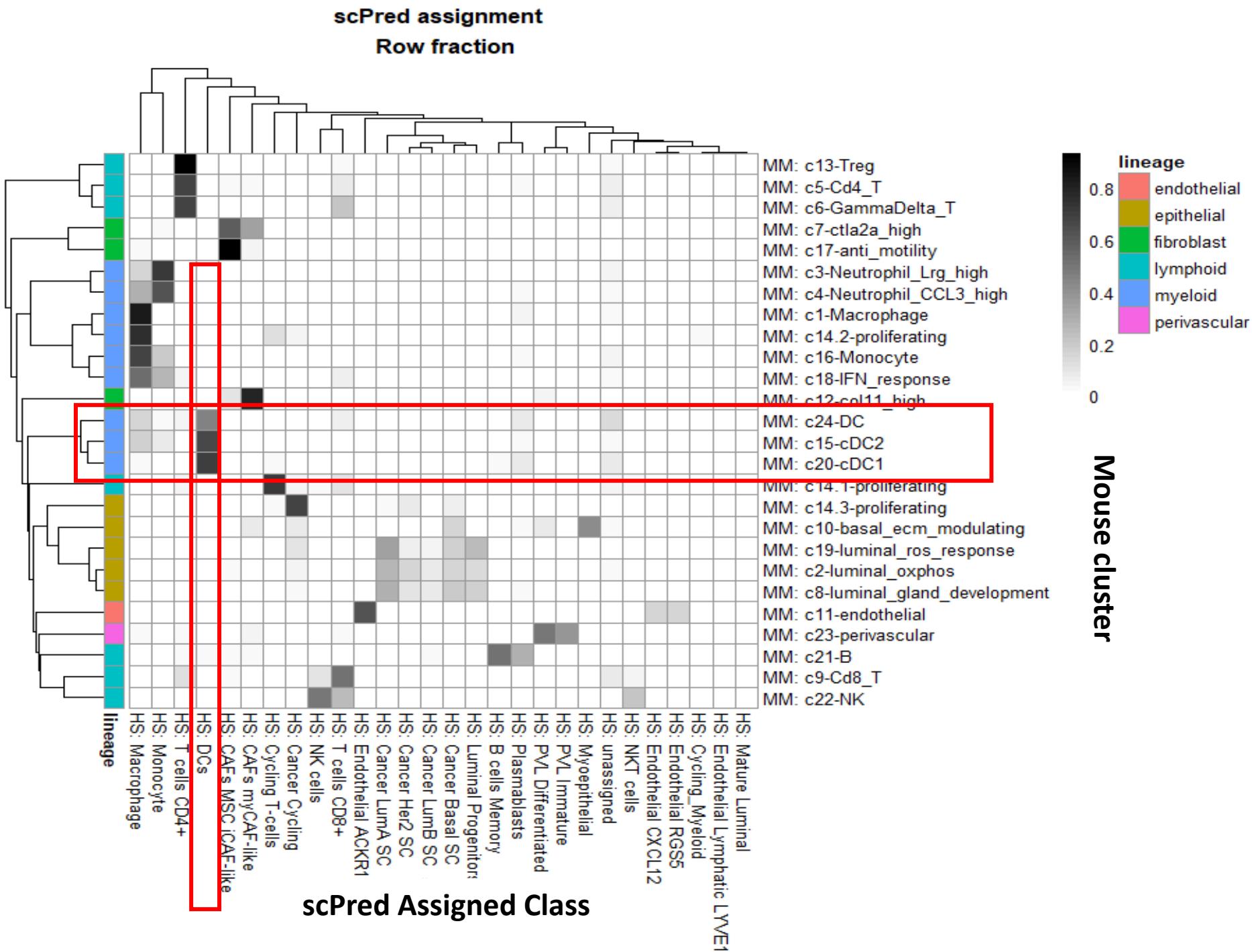
Human
trained
MDA
classifier
applied to
mouse



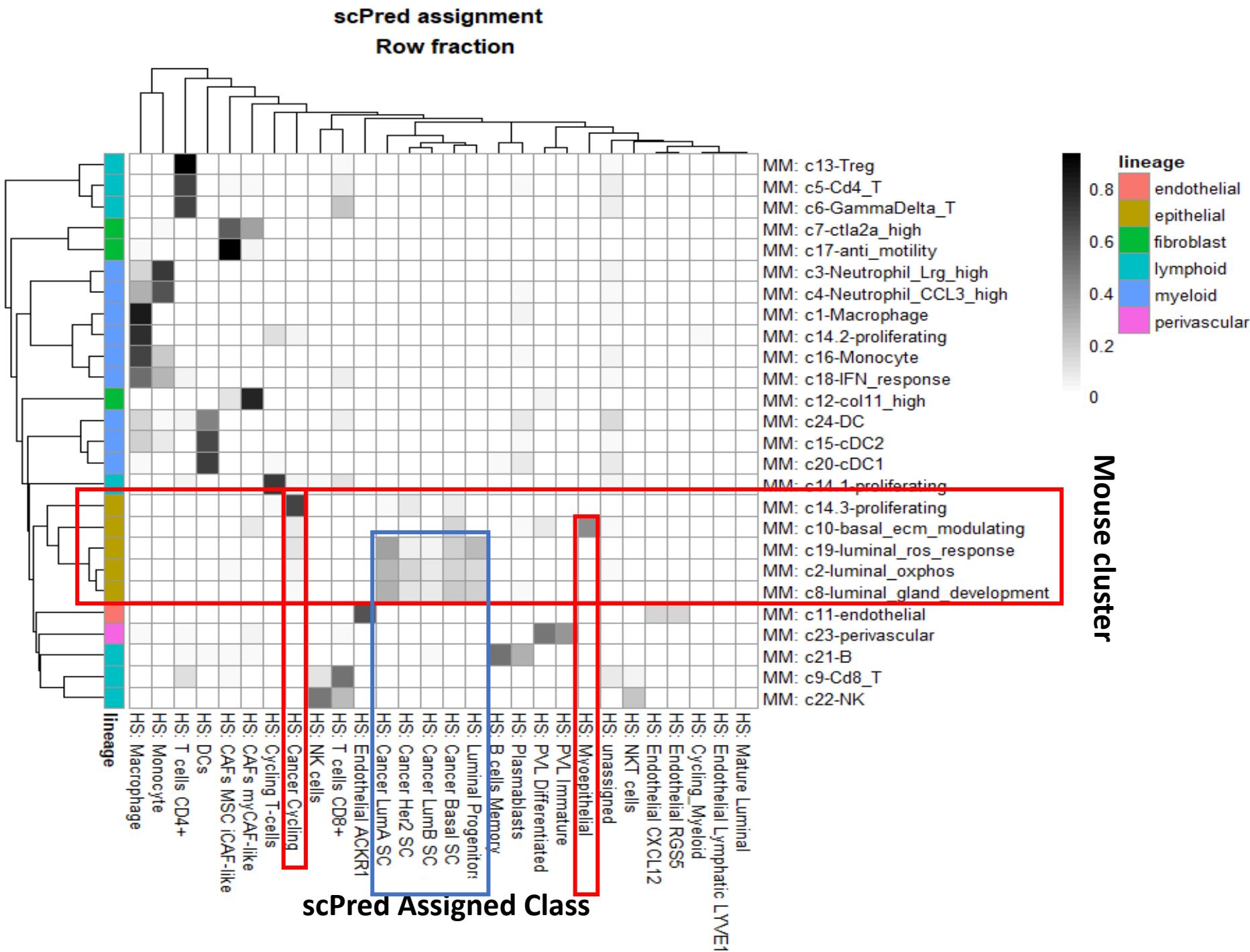
Human
trained
MDA
classifier
applied to
mouse



Human
trained
MDA
classifier
applied to
mouse



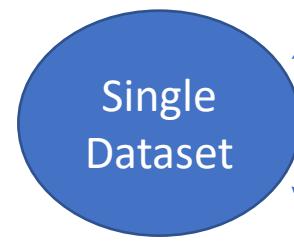
Human
trained
MDA
classifier
applied to
mouse



Are the mouse cell states representative of human disease?

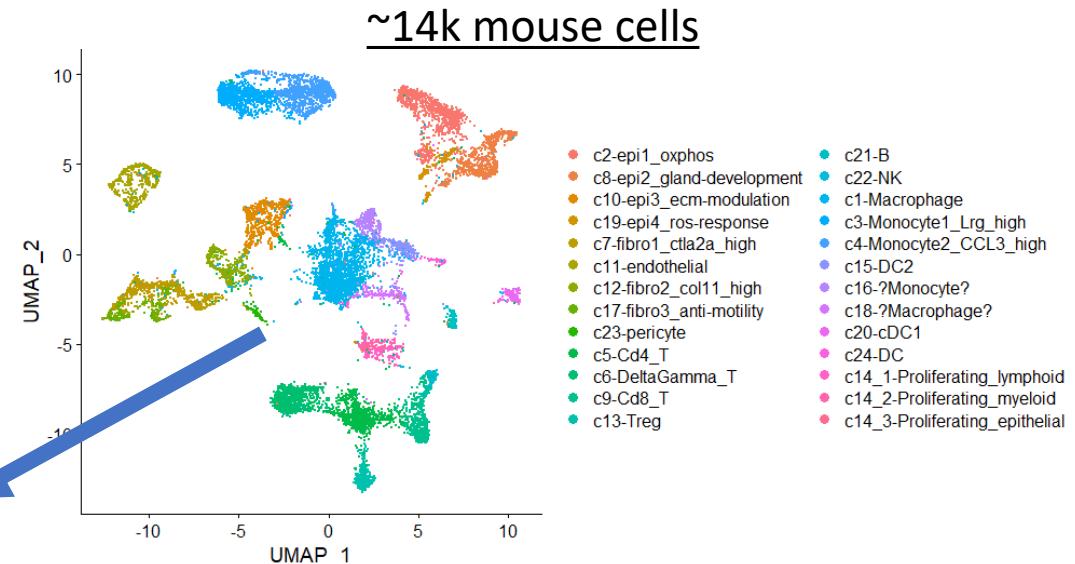
1. Classification (Supervised)

- Train a classifier on one species using a subset of shared features and apply it to the other

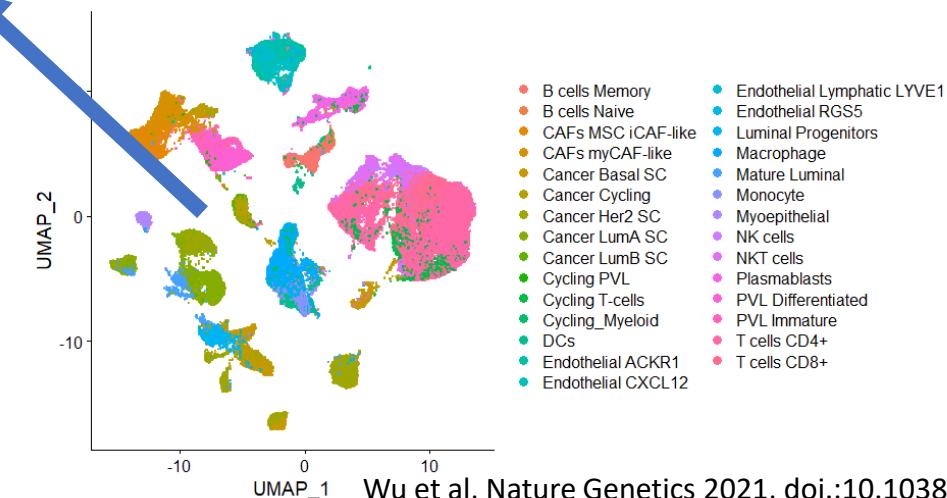


2. Data integration (Unsupervised)

- Species-normalize the data to enable simultaneous analysis within the unified feature space



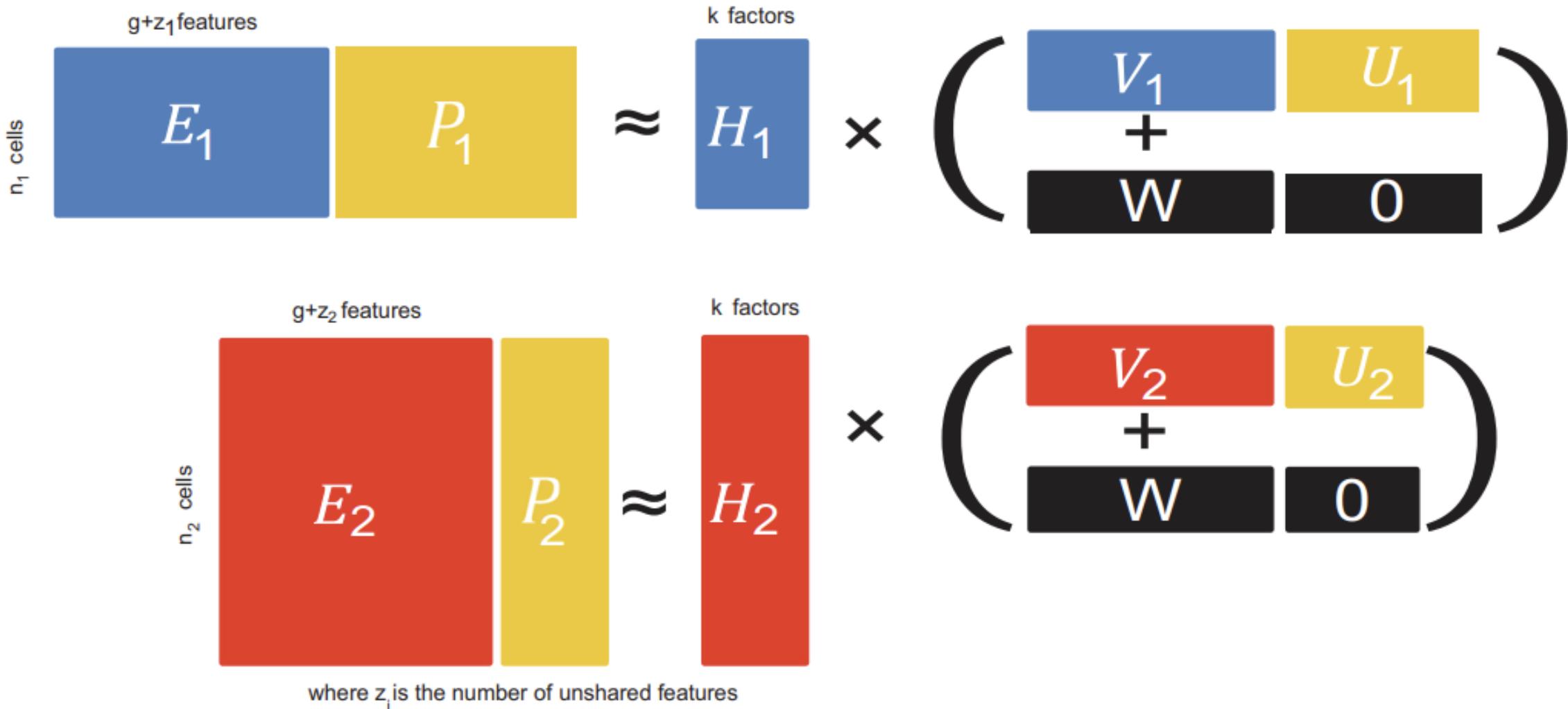
~100k human cells



UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative matrix factorization

(More tangible example in a few slides, don't worry!)

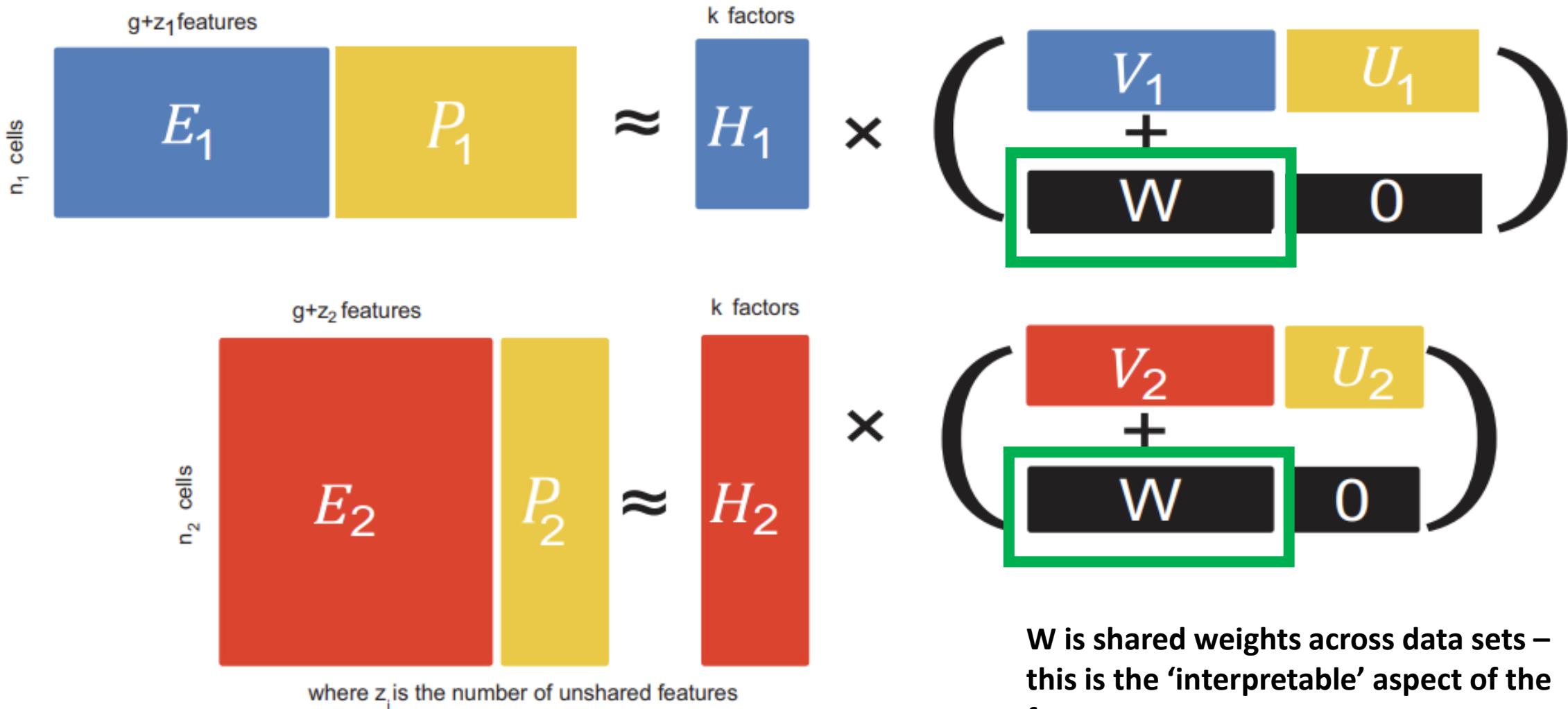
April R. Kriebel¹ & Joshua D. Welch^{1,2}✉



UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative matrix factorization

(More tangible example in a few slides, don't worry!)

April R. Kriebel¹ & Joshua D. Welch^{1,2}✉

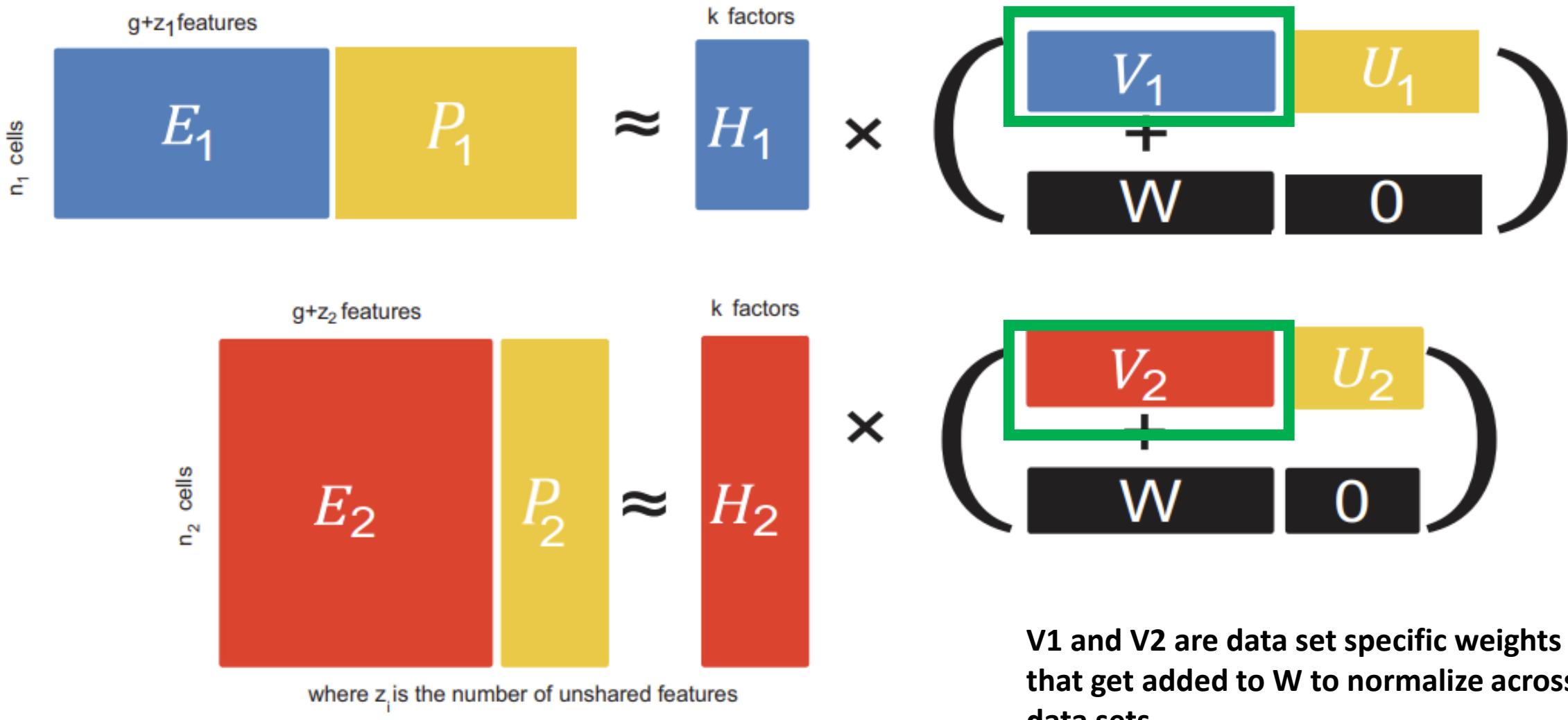


W is shared weights across data sets – this is the ‘interpretable’ aspect of the factors

UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative matrix factorization

(More tangible example in a few slides, don't worry!)

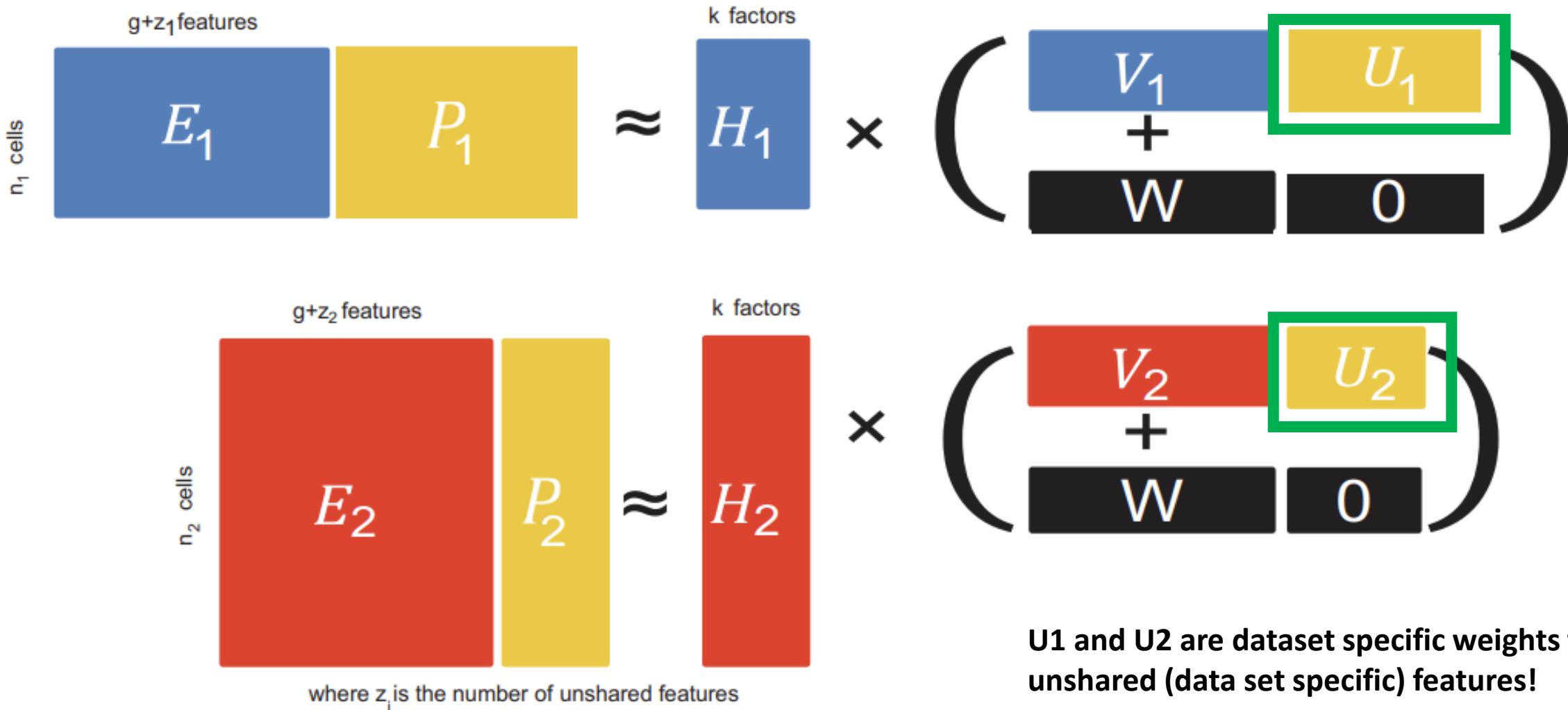
April R. Kriebel¹ & Joshua D. Welch^{1,2}✉



UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative matrix factorization

(More tangible example in a few slides, don't worry!)

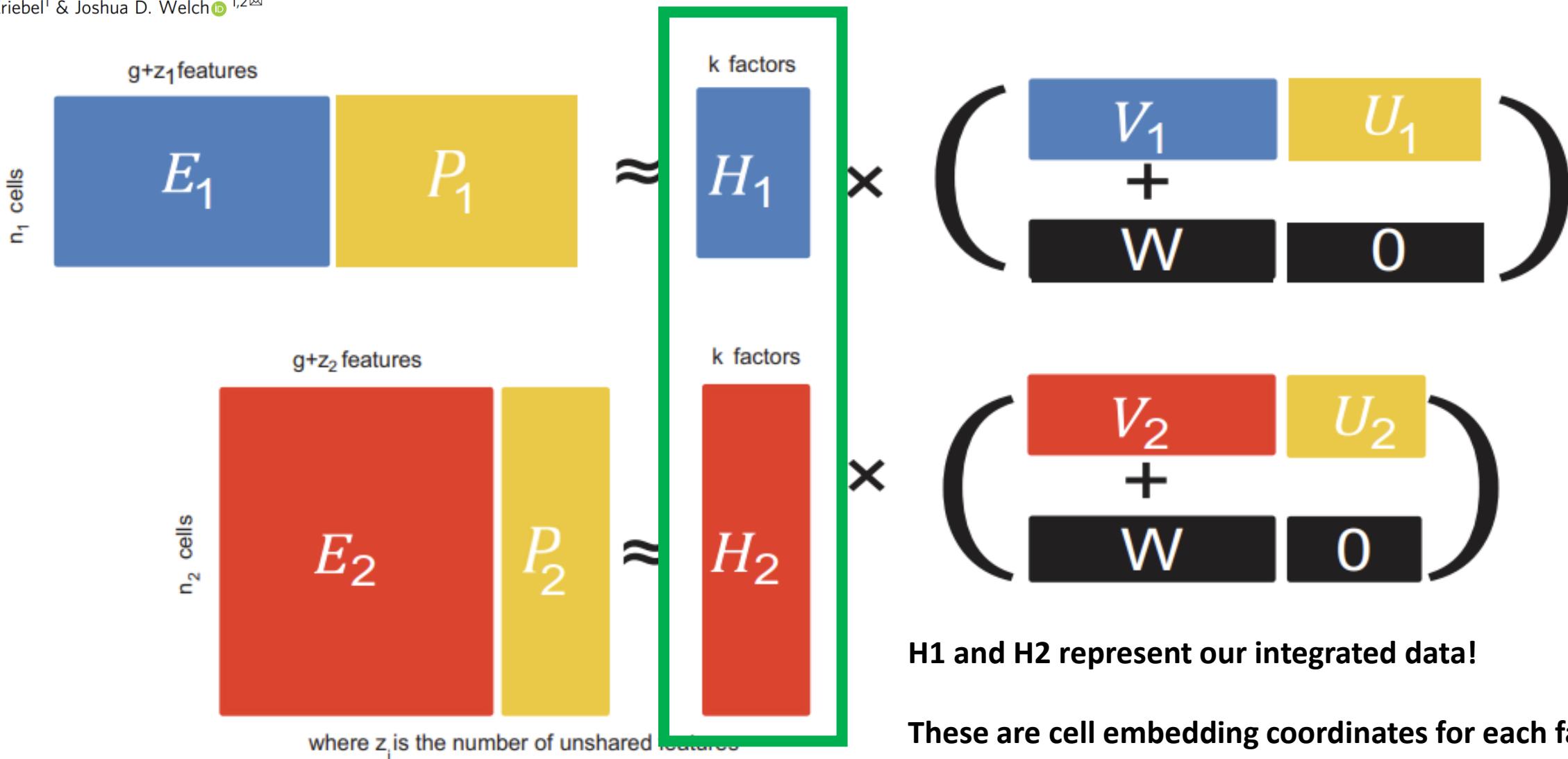
April R. Kriebel¹ & Joshua D. Welch^{1,2}✉



UINMF performs mosaic integration of single-cell multi-omic datasets using nonnegative matrix factorization

(More tangible example in a few slides, don't worry!)

April R. Kriebel¹ & Joshua D. Welch^{1,2}✉



Example of learning from unshared features

Data set 1: Dogs Looking forward



Data set 2: Dogs looking right



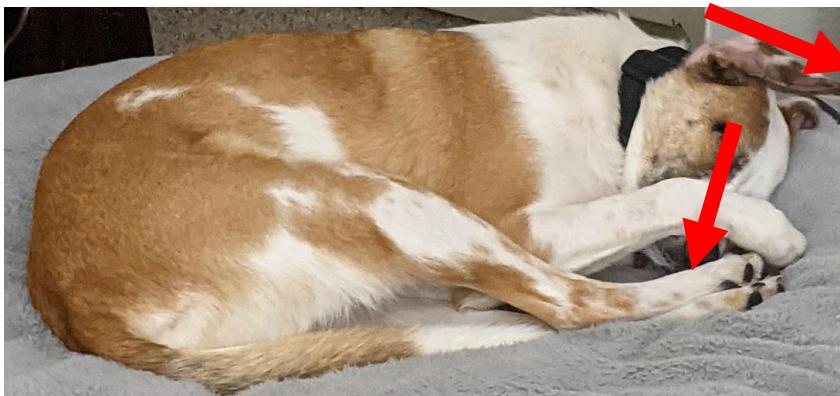
*Assume normalization occurs within each dataset so length measurements are accurate

Example of learning from unshared features

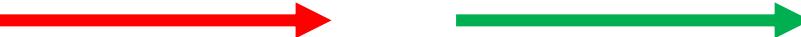
Data set 1: Dogs Looking forward



Data set 2: Dogs looking right



*Assume normalization occurs within each dataset so length measurements are accurate



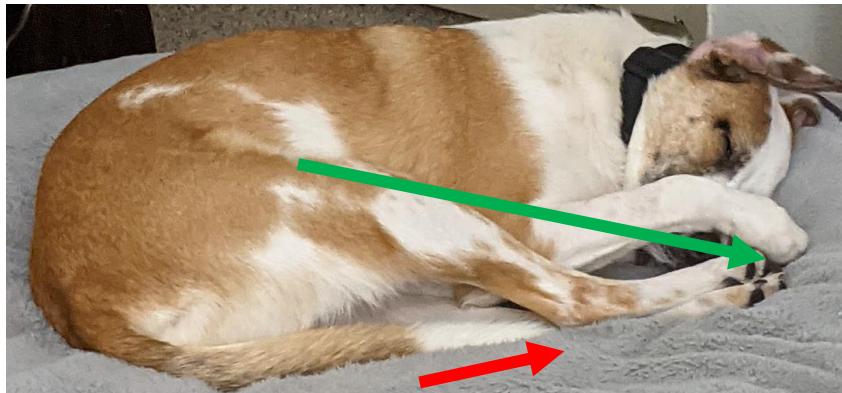
Topic (Factor)	Shared features	DS1 features	DS2 features
Head size	Right ear length Right eye-to-snout distance	Left ear length Left eye-to-snout distance	-

Example of learning from unshared features

Data set 1: Dogs Looking forward



Data set 2: Dogs looking right



*Assume normalization occurs within each dataset so length measurements are accurate



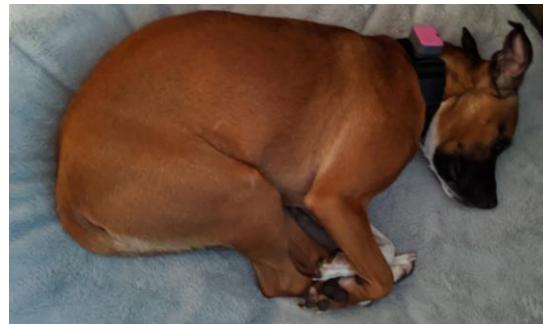
Topic (Factor)	Shared features	DS1 features	DS2 features
Head size	Right ear length Right eye-to-snout distance	Left ear length Left ear-to-snout distance	-
Tail length	Length of white tip	-	Length of rear right leg

Example of unshared features used for topic modeling

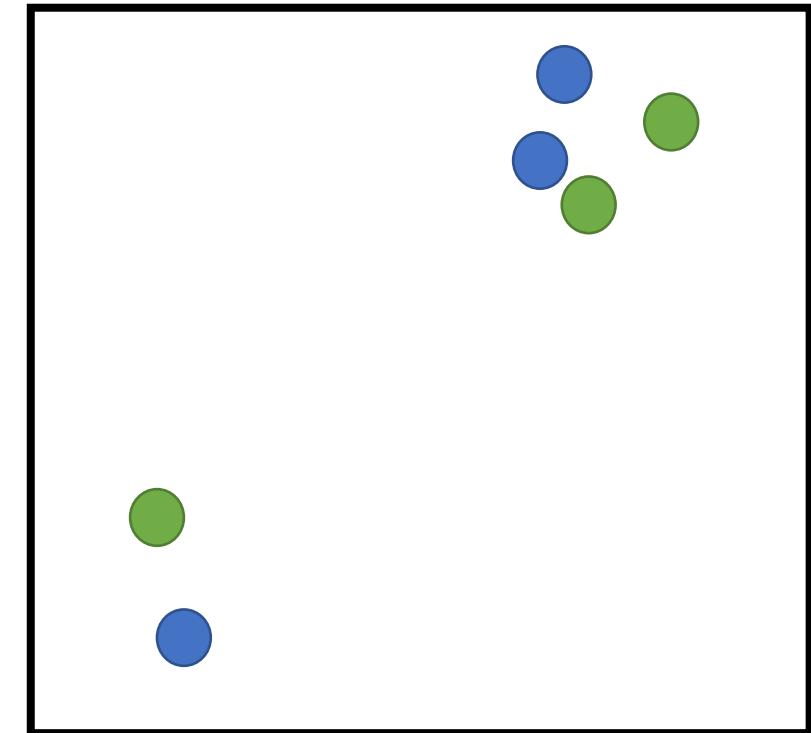
Dataset 1: Dogs looking forward



Dataset 2: Dogs looking right



Factor 2: Tail length



Factor 1: Head Size

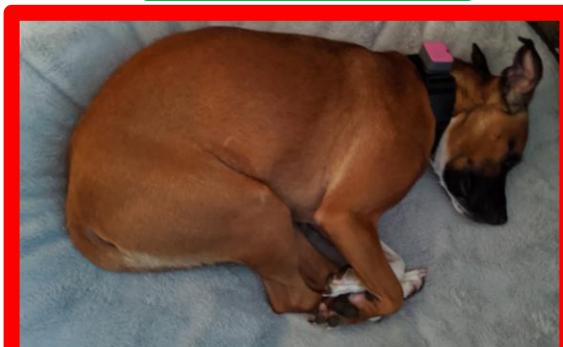
DS1
DS2

Example of unshared features used for topic modeling

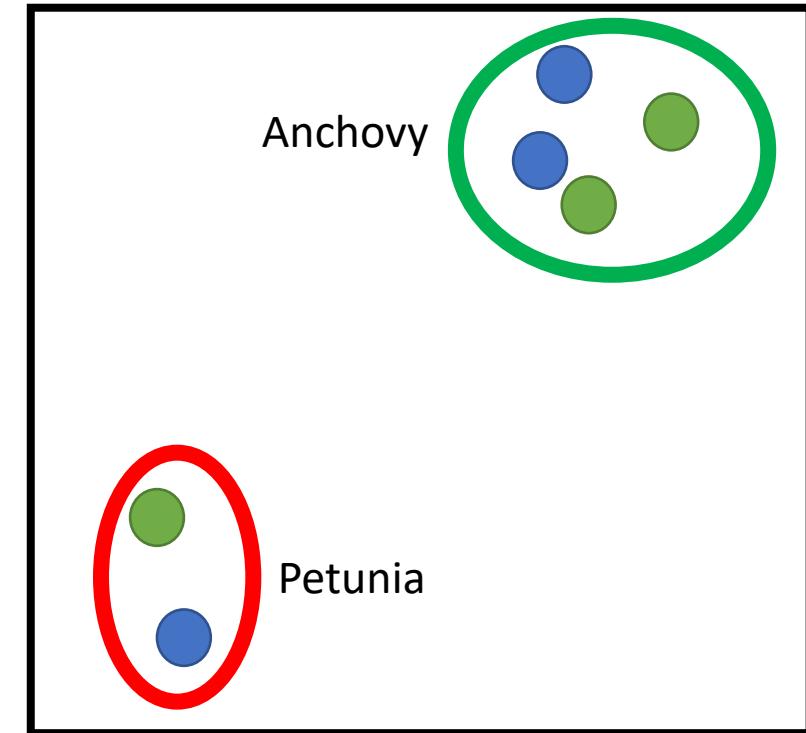
Dataset 1: Dogs looking forward



Dataset 2: Dogs looking right



Factor 2: Tail length



Factor 1: Head Size

Data integration plan

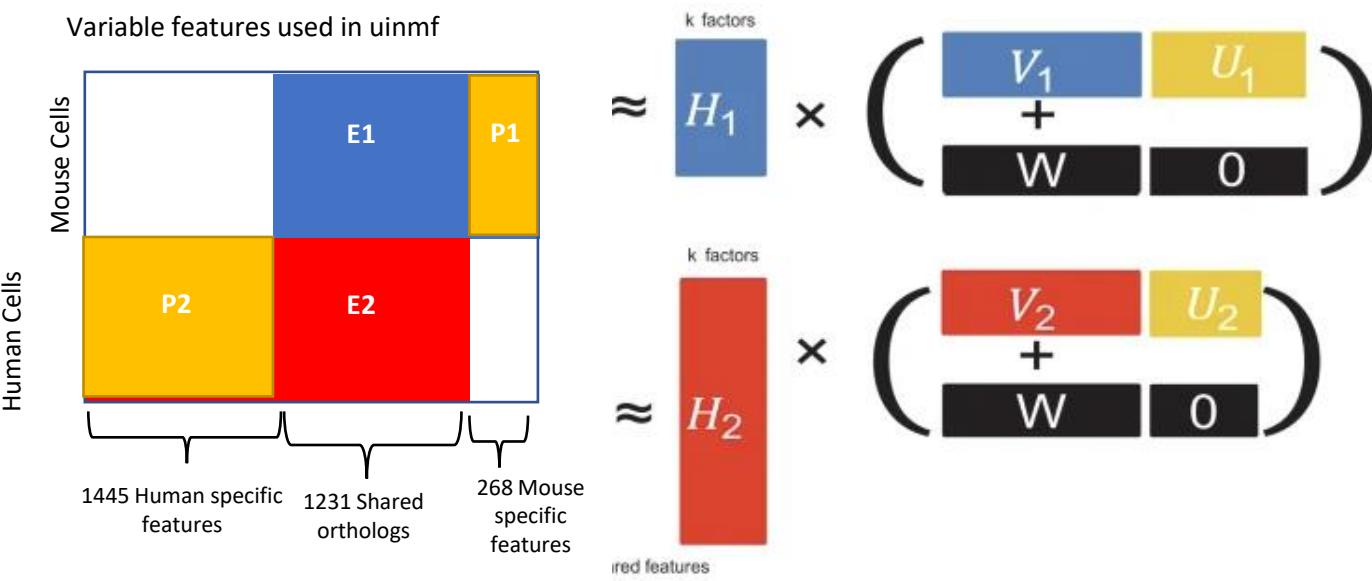
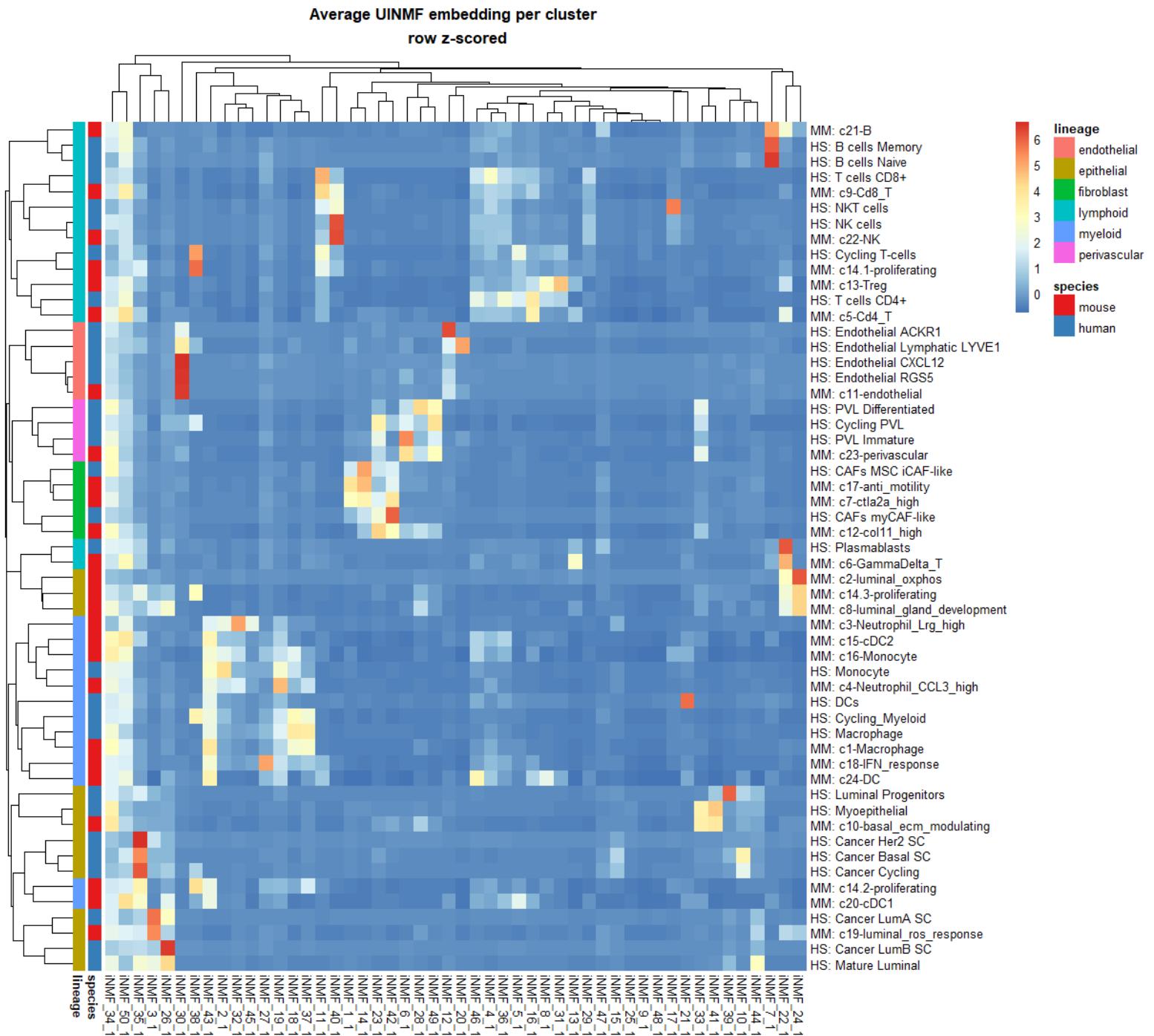


Figure adapted from welch et al

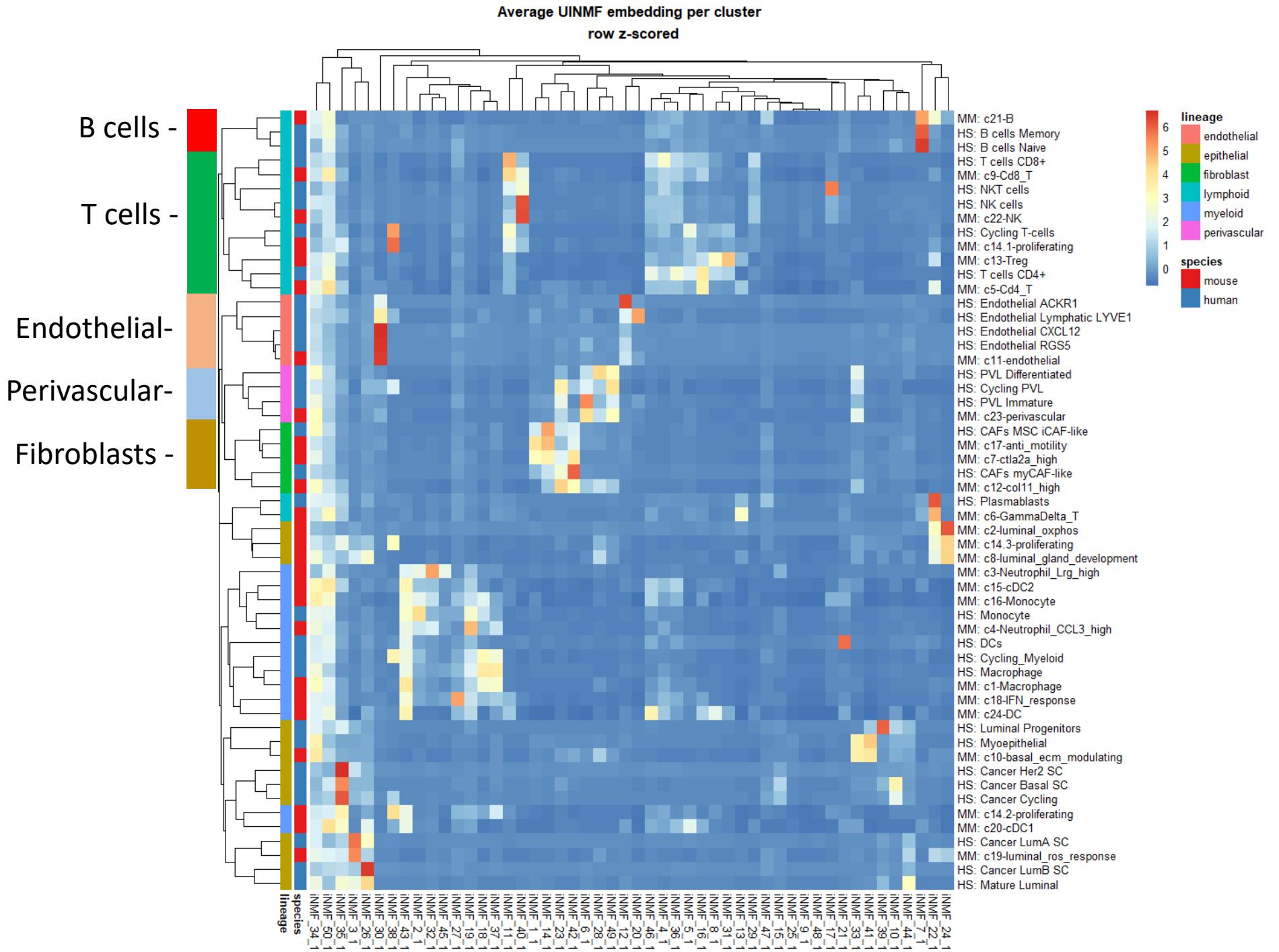
1. Assign celltypes for each data set using only that specie's data set
2. Scale and normalize each data set separately
3. Identify variable shared features
4. Identify variable unshared features
5. Perform UINMF to integrate across species
6. Cluster on UINMF loadings (H matrix)
7. Assess relative co-clustering of previously identified cell types

UINMF Factors
recapitulate
known biology

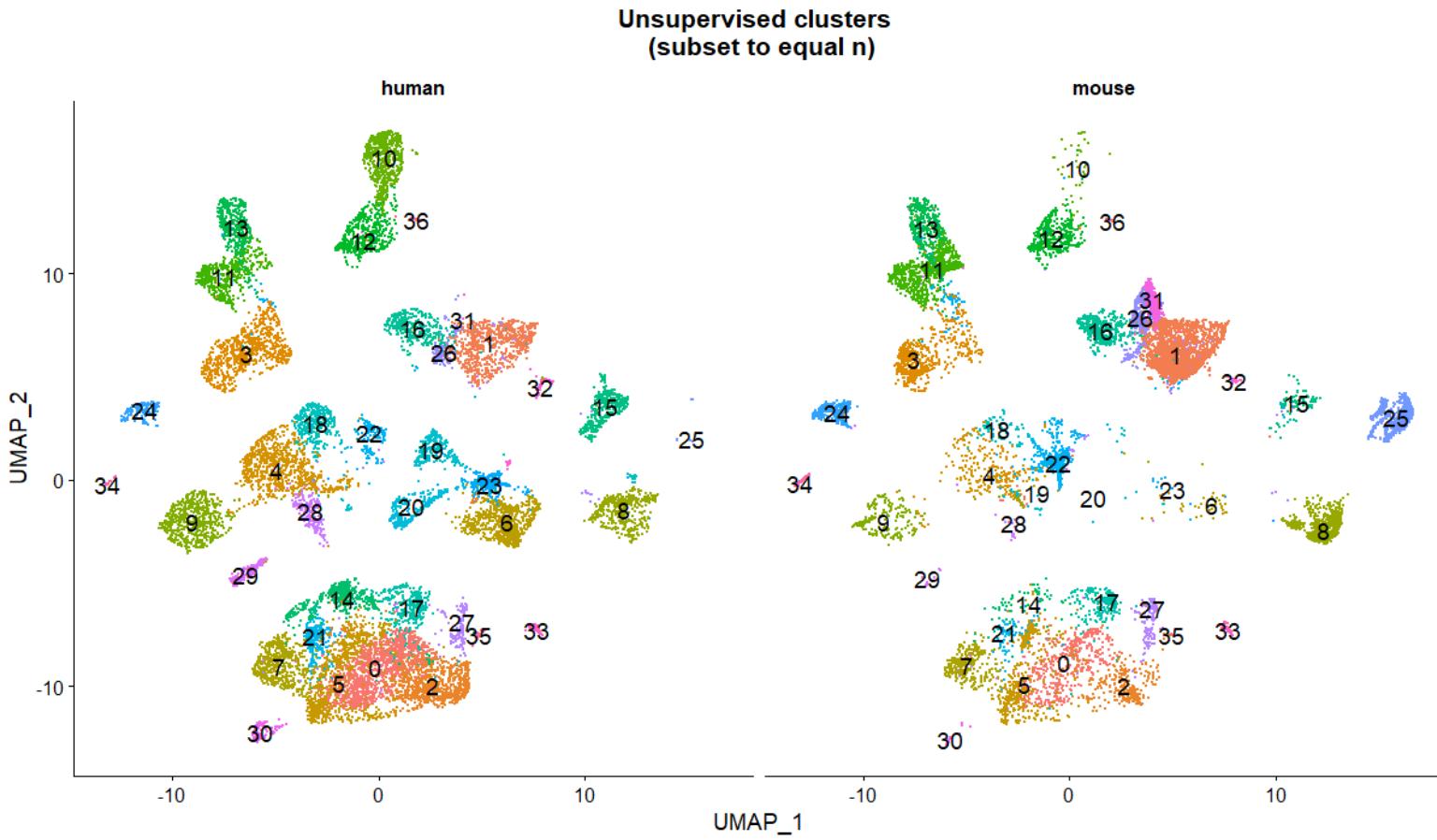
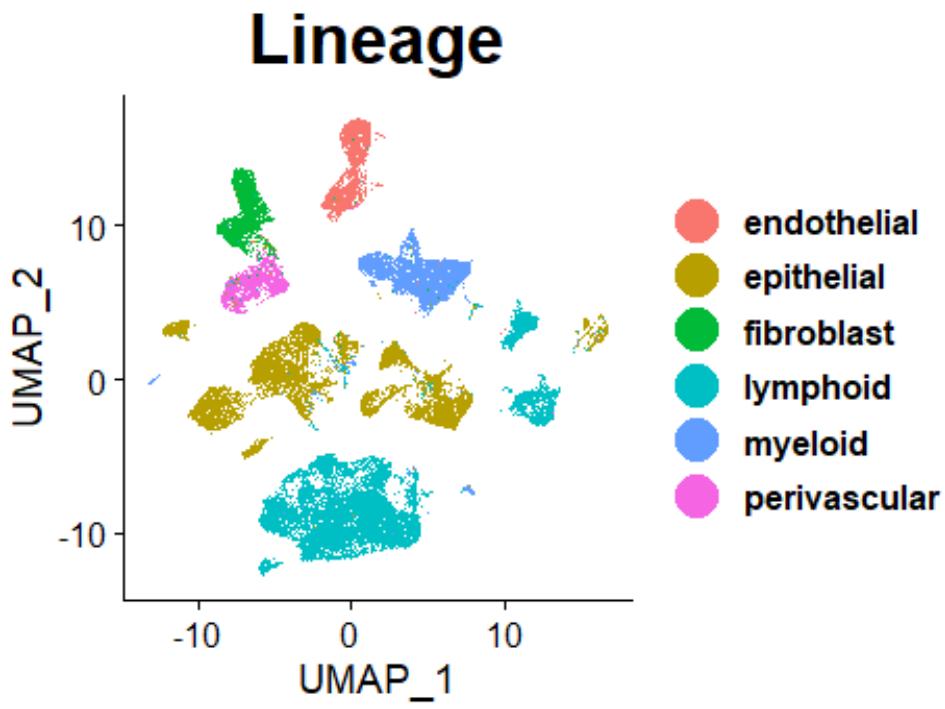
Original celltype assignment
(from uni-species analysis)



Original celltype assignment (from uni-species analysis)



UINMF maintains separation of cell lineages across species



Most important points

- **Background matters.** Almost every analysis tool/algorithm is highly dependent on the context of the data set. This needs to be considered when interpreting any results.
- Always attempt to have at least a ‘blackbox’ understanding of an algorithm:
 - Input data
 - Output data
 - Goal of algorithm
 - Optimization criteria*
 - Can we extract and evaluate how well the approach worked?

Databases to know

- Public data repositories
 - Gene expression omnibus [<https://www.ncbi.nlm.nih.gov/geo/>]
 - NCBI affiliated gene expression repository. Most studies get published here
 - Broad Single Cell Portal [https://singlecell.broadinstitute.org/single_cell]
 - Interactive tools let you explore the data before pulling it down
- Celltype references
 - CZI [<https://cellxgene.cziscience.com/>]
 - Largest(?) reference data base for gene expression
 - CellMarkerDB [<http://bio-bigdata.hrbmu.edu.cn/CellMarker/>]
 - Data supported marker genes for celltypes identified in human or mouse samples

Important tools

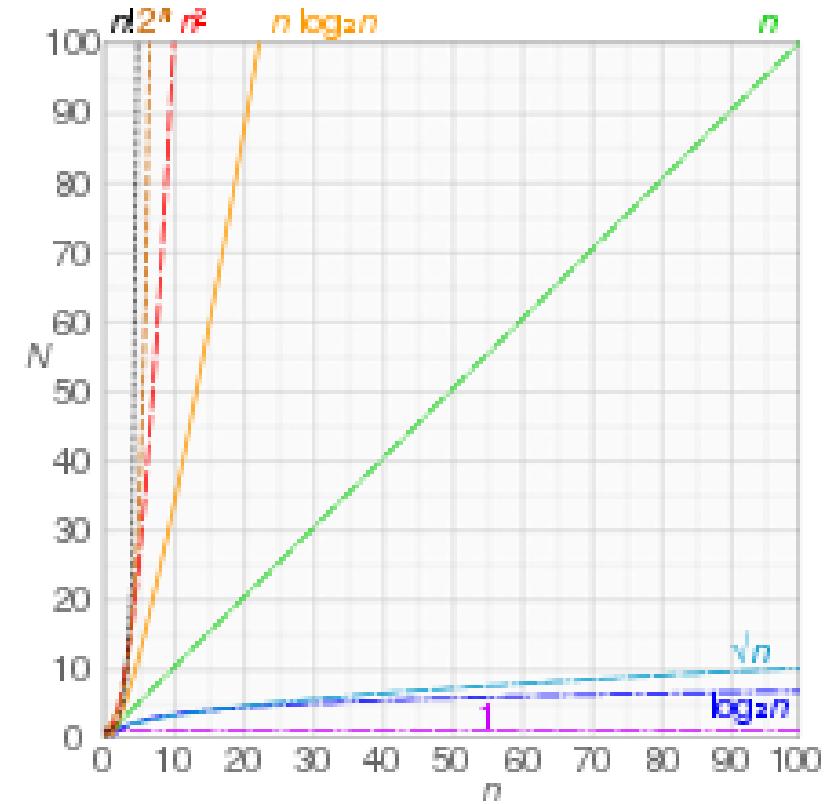
- General analysis:
 - Seurat [R] – Most popular, excellent vignettes. Nested data structure needs to be learned
 - SCRAN [R] – uses the ‘Single Cell Experiment/SCE’ data structure which is common to MANY biology tools
 - Scanpy [python]– well documented and feature rich
- Droplet / early data tools
 - DropletUtils [R] – Droplet level analysis, read-based downsampling and more
 - SoupX [R] – Remove ambient mRNA contamination
- Geneset enrichment
 - Msigdbr [R] – easily access Msigdb gene sets
 - ClusterProfiler [R] – Interfaces with tons of DBs and enables various enrichment analyses

Important tools 2

- Integration
 - Seurat [R] – Integrates at feature level, slow. Can incorrectly impute data!
 - Harmony [R] – integrates at principal component level, very fast
 - rLiger [R] – integrates at NMF level, slow. Can integrate unshared features
- Gene name conversion
 - BiomaRt [R] – AMAZING tool.
- Cell-cell interaction
 - Cellchat [R] – general analysis, great visualization tools
 - Nichenetr [R] – Focused target-receiver analysis with receptor-> downstream gene expression change filtering

Comp background: Big O notation describes the efficiency of an algorithm

- Big O \approx ‘Order of ...’
 - Describes how the time to compute something relates to the length of elements used to compute
 - Example:
 - $O(N)$
 - Normalizing data takes some set amount of time per cell
 - $O(N^2)$
 - Pairwise correlations – calculated between all possible pairs
 - Pairwise distance calculations – calculated between all possible pairs
- When the number of features (M) \approx the number of samples (N), then this can impact run time!
 - $O(M \cdot N) == O(N^2)$ when $M \approx N$
 - This means that for very large data sets it is sometimes necessary to use less precise methods that operate in quicker times



https://en.wikipedia.org/wiki/Big_O_notation