

1. Interactive scRNA-Seq analysis with the Single Cell Toolkit (SCTK)

2. Decontamination of ambient RNA with DecontX

W. Evan Johnson, Ph.D. and Joshua Campbell, Ph.D.
Boston University School of Medicine

CSBC scRNA-seq working group

Single Cell RNA-seq

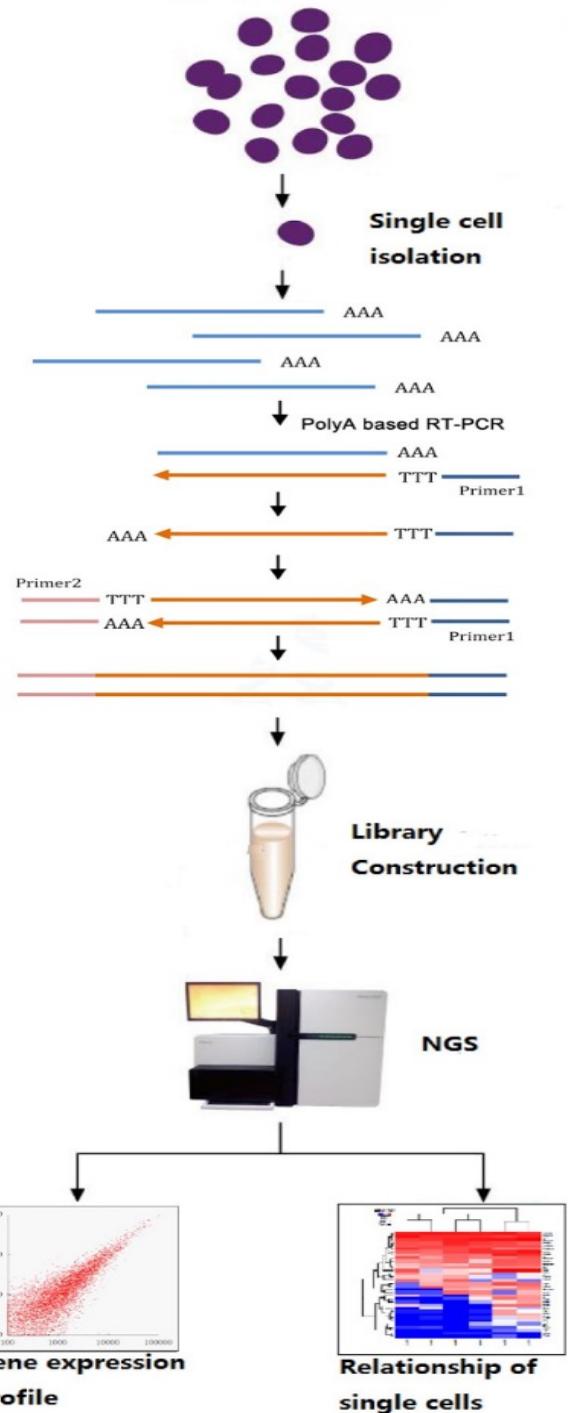
(The Future of RNA-sequencing)

Single Cell Profiling

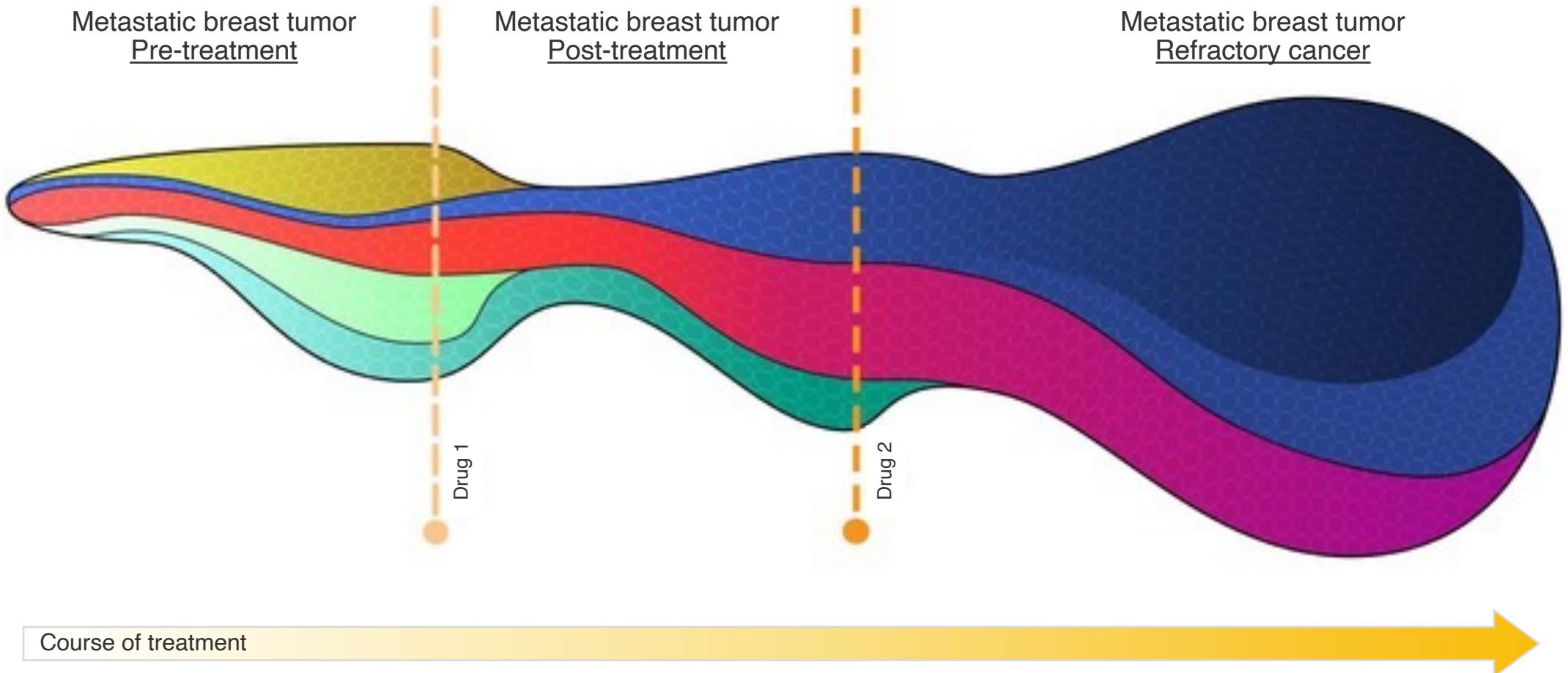
Applications for Single Cell Sequencing:

- **Cancer:** intra-tumoral heterogeneity
- **Development:** characterize every cell in blastocyst
- **Other applications:** novel/rare cell type discovery
- **Infectious Diseases:** combine with metagenomics to explore host response (e.g. TB)
- And more!

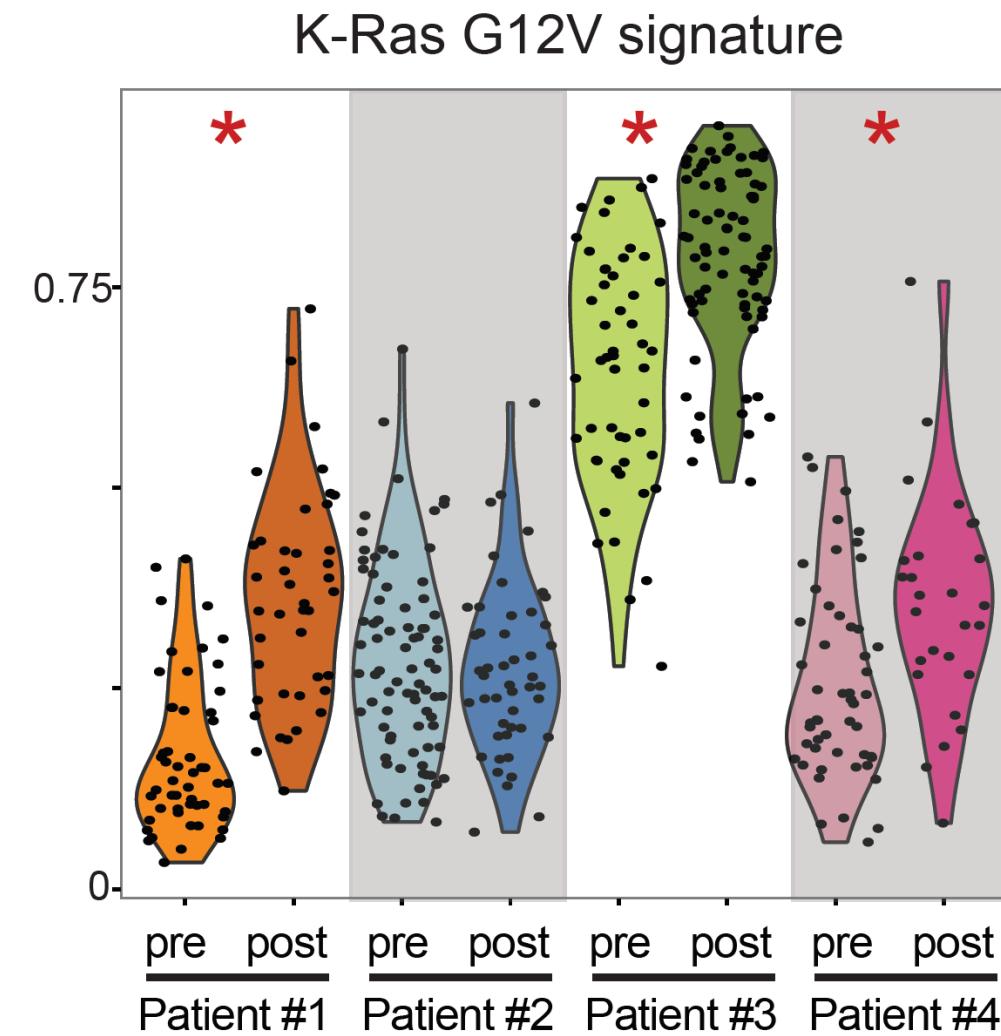
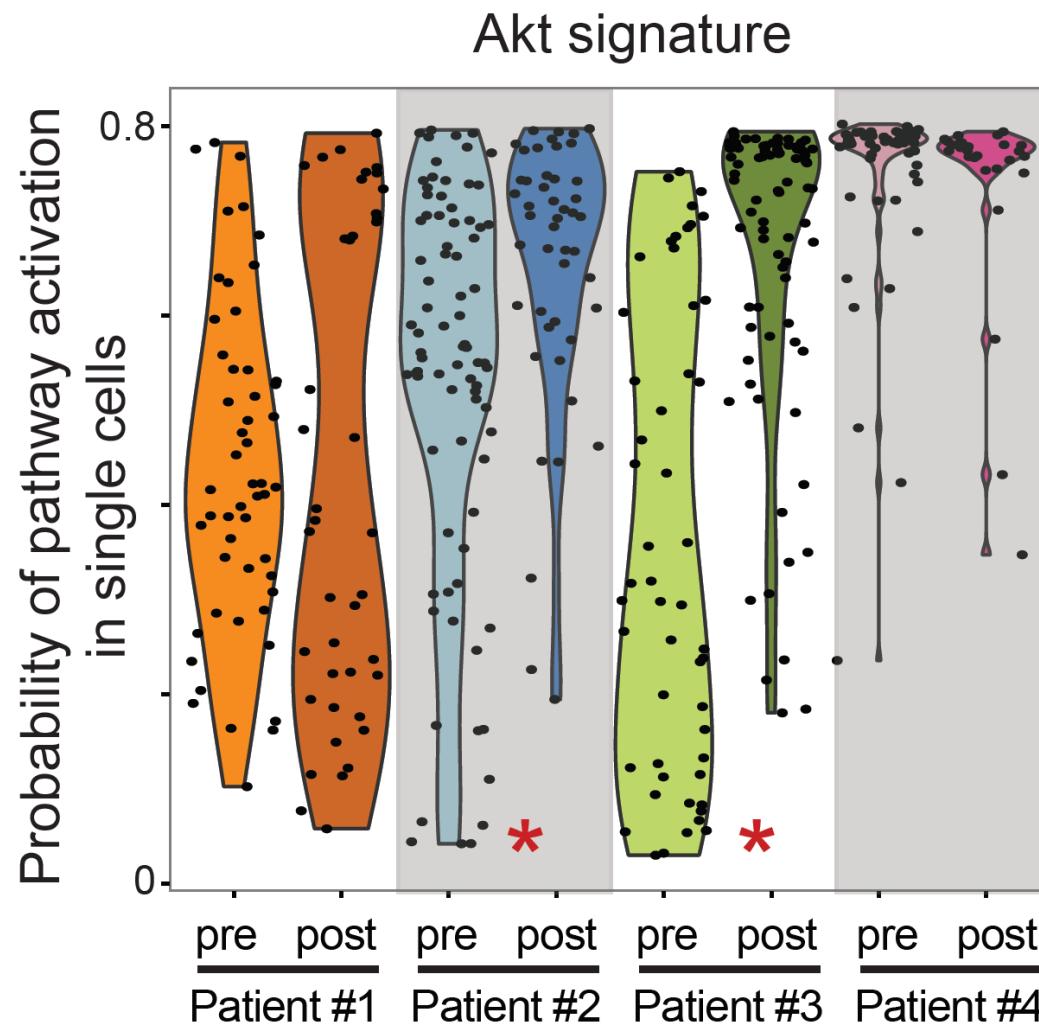
The Future of RNA-sequencing



Tumor Heterogeneity

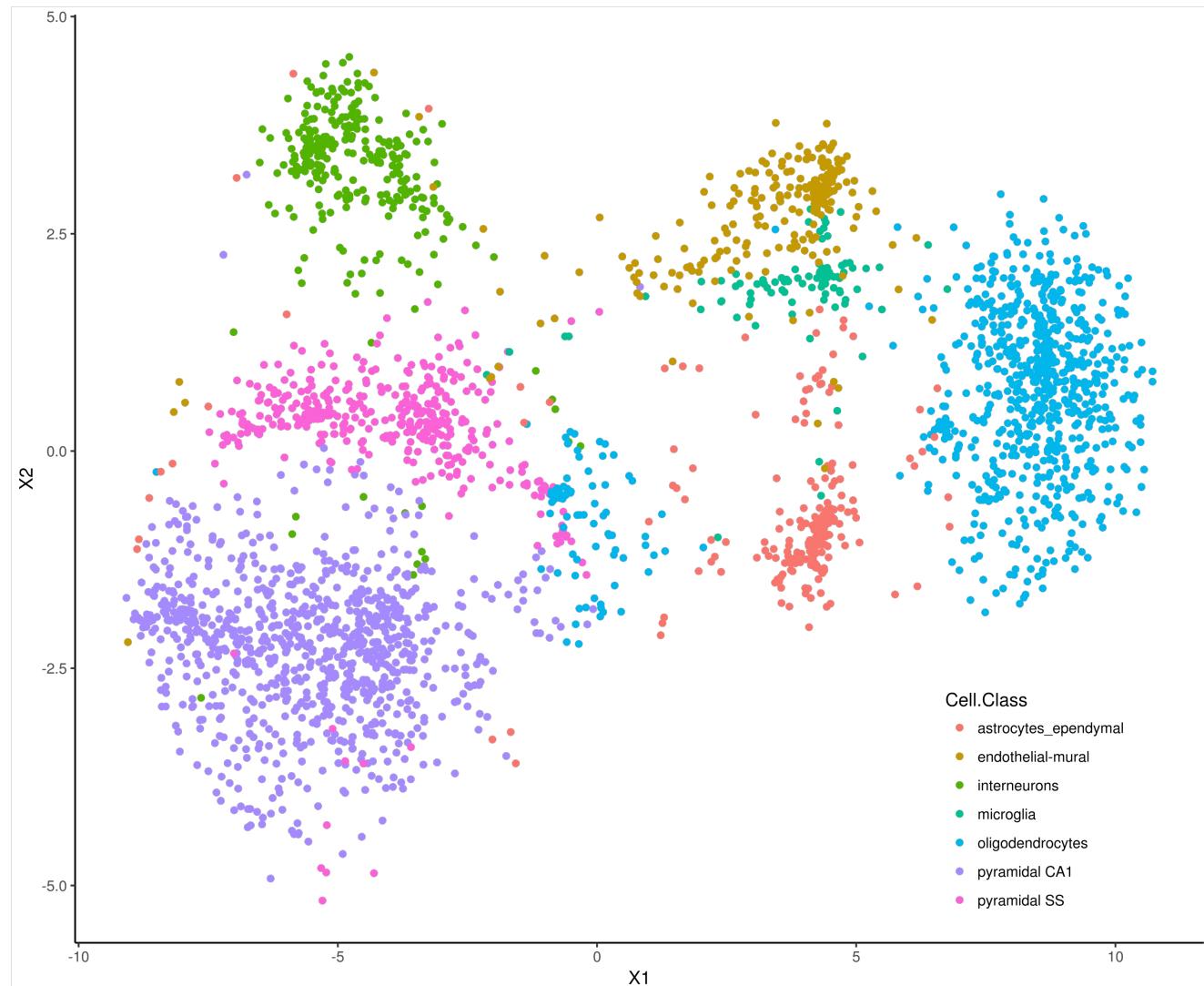


Pathway Activity in Breast Cancer



Critical Gaps in scRNA-seq data

- Complex data
 - New analysis challenges
- Interactive, Simple Analysis?
 - Inexperienced Users
 - Optimizing Parameters
 - Filtering Failed Samples
 - Filtering Low expression genes
- Some packages already exist
 - QC
 - Clustering
 - Full Analysis Portal?

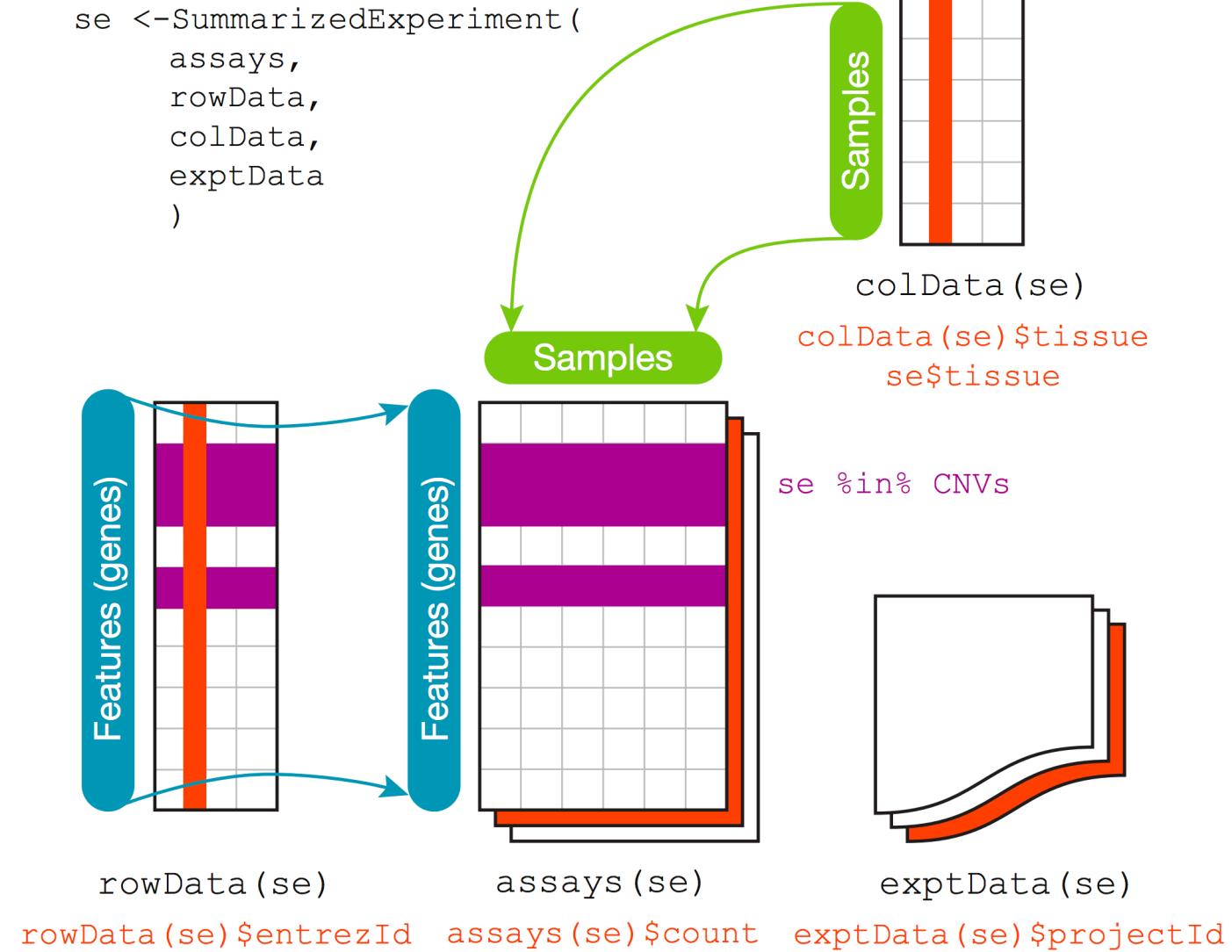


SingleCellTK

Single Cell Toolkit

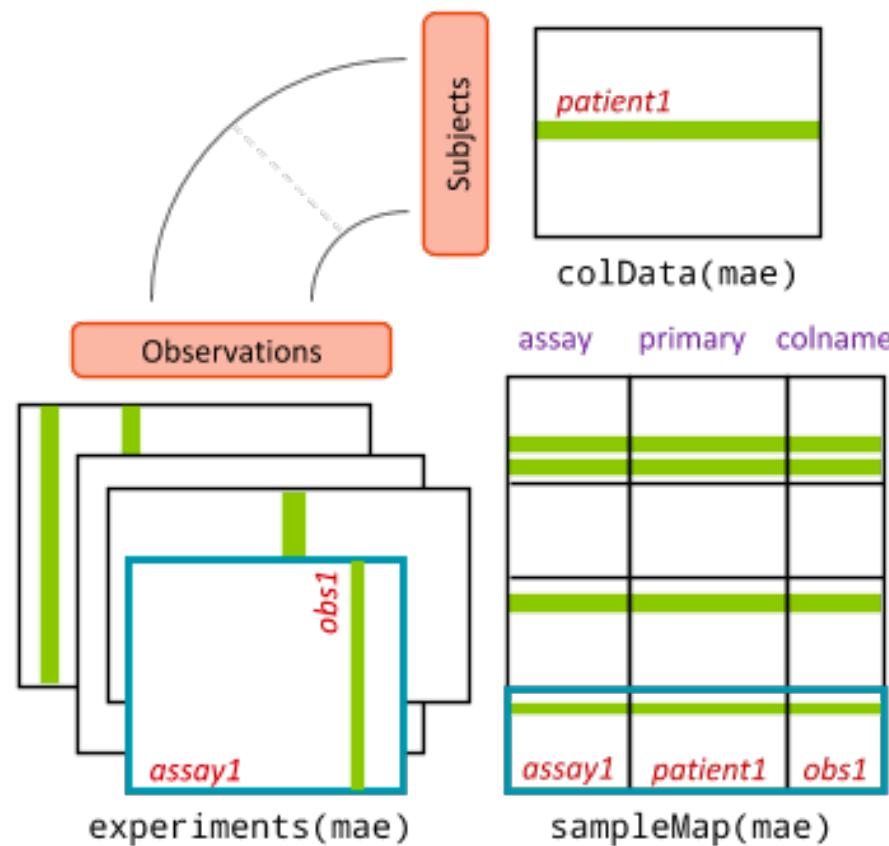
- Standard R package *with* a Shiny toolkit
- R functions on top of a SingleCellExperiment (SCE) object
 - Can download analysis performed in the toolkit and continue on the command line
- SCE object can be brought in/out of Shiny at any stage
 - Great for common tasks:
 - Interactive clustering/visualization
 - Differential gene/pathway analysis
- (Also works for bulk RNA-seq analysis)

SummarizedExperiment



MultiAssayExperiment

- Combine multiple data types in similar objects with functionality to query across features and cells.

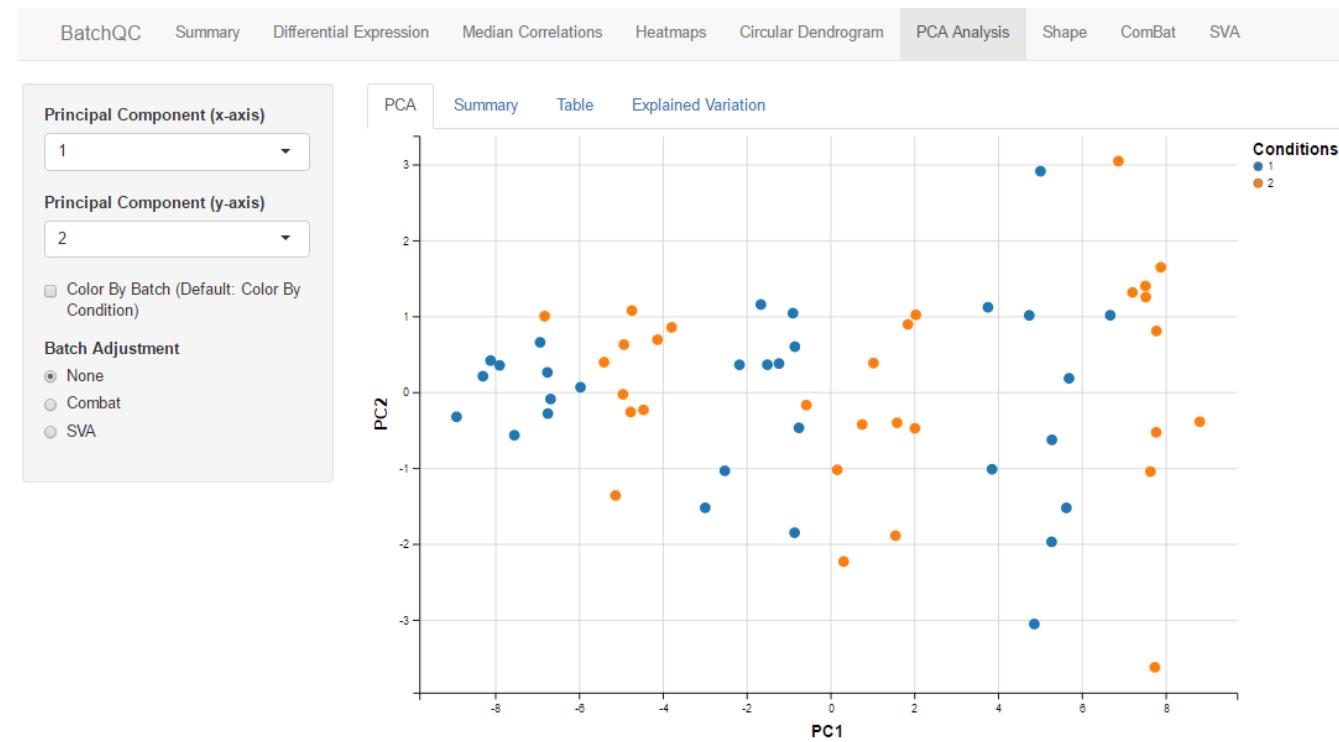
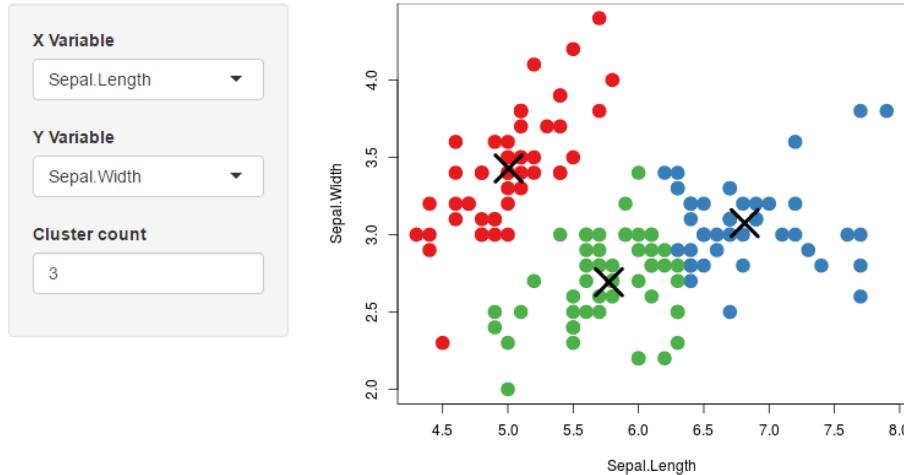


SingleCellExperiment

- <https://bioconductor.org/packages/devel/bioc/html/SingleCellExperiment.html>
- reducedDims
 - PCA, TSNE, any other sample size matrix of dimensionality reduction data
- isSpike – indicate which probes are spike ins
- sizeFactors – add scaling factors
- Named assays: counts, normcounts, logcounts, cpm, tpm
 - Just a convention
- DelayedArray – run common array functions without loading it into memory

Shiny

- Simple Web Apps Written in R
- Interactive and Reactive
- Iris k-means clustering



SingleCellITK

Data Upload

QC and Filtering

Visualization and
Clustering

ComBat Batch
Correction

Differential
Expression

MAST

Pathway Activity
Analysis

Experimental
Design

Data Export

[Home](#) » [Bioconductor 3.9](#) » [Software Packages](#) » [singleCellTK](#)

singleCellTK

platforms all

rank unknown

posts 0 in Bioc 1 year

build ok

updated before release

DOI: [10.18129/B9.bioc.singleCellTK](https://doi.org/10.18129/B9.bioc.singleCellTK)



Interactive Analysis of Single Cell RNA-Seq Data

Bioconductor version: Release (3.9)

Run common single cell analysis directly through your browser including differential expression, downsampling analysis, and clustering.

Author: David Jenkins

Maintainer: David Jenkins <dfj at bu.edu>

Citation (from within R, enter `citation("singleCellTK")`):

Jenkins D, Faits T, Khan MM, Briars E, Carrasco Pro S, Johnson WE (2019). *singleCellTK: Interactive Analysis of Single Cell RNA-Seq Data*. R package version 1.4.0, https://combiomed.github.io/sctk_docs/.

Installation

Documentation »

Bioconductor

- Package [vignettes](#) and manuals.
- [Workflows](#) for learning and use.
- [Course and conference](#) material.
- [Videos](#).
- Community [resources](#) and [tutorials](#).

R / [CRAN packages](#) and [documentation](#)

Support »

Please read the [posting guide](#). Post questions about Bioconductor to one of the following locations:

- [Support site](#) - for questions about Bioconductor packages
- [Bioc-devel](#) mailing list - for package developers



Single Cell Toolkit

Filter, cluster, and analyze single cell RNA-Seq data

Need help? [Read the docs.](#)

Upload

[\(help\)](#)

Choose data source:

- Upload files
- Upload SCtkExperiment RDS File
- Use example data

Upload data in tab separated text format:

Example count file:

Gene	Cell1	Cell2	...	CellN
------	-------	-------	-----	-------

Example sample annotation file:

Cell	Annot1	...
------	--------	-----

Example feature file:

Gene	Annot2	...
------	--------	-----

Single Cell Toolkit

Filter, cluster, and analyze single cell RNA-Seq data

Need help? [Read the docs.](#)

Upload

[\(help\)](#)

✓ Successfully Uploaded!



Choose data source:

- Upload files
- Upload SCtkExperiment RDS File
- Use example data

Choose Example Dataset:



Single Cell Toolkit

Filter, cluster, and analyze single cell RNA-Seq data

Need help? [Read the docs.](#)

Upload

[\(help\)](#)

✓ Successfully Uploaded!



Choose data source:

- Upload files
- Upload SCtkExperiment RDS File
- Use example data

Choose Example Dataset:

mouseBrainSubset



SCTK features - *Differential Expression analysis-MAST*

Single Cell Toolkit v0.3.9 Upload Data Summary and Filtering DR & Clustering Barcode Collection Differential Expression ▾ Pathway Activity Analysis Sample Size

MAST

Select Assay:
logcounts

Adaptive Thresholding:

Hurdle Model:
 Use Adaptive Thresholds

Select fold change threshold

Select expression threshold

Select Condition for Hurdle Model
condition

p-value (FDR) cutoff:

Adaptive thresholding Results Table Violin Plot Linear Model Heatmap

Plot	Threshold Range	N	Bandwidth
(0, 1)	(0.0144, 0.163)	19216	0.02479
(1, 2)	(0.163, 0.334)	68296	0.04462
(2, 3)	(0.334, 0.53)	74475	0.06479
(3, 4)	(0.53, 0.755)	74934	0.0838
(4, 5)	(0.755, 1.01)	62929	0.1036
(5, 6)	(1.01, 1.31)	57195	0.1237
(6, 7)	(1.31, 1.65)	47982	0.1538
(7, 8)	(1.65, 2.04)	37354	0.1831
(8, 9)	(2.04, 2.48)	25554	0.2486
(9, 10)	(2.48, 2.99)	17546	0.3426
(10, 11)	(2.99, 3.58)	17010	0.3729
(11, 12)	(3.58, 4.25)	11939	0.4304
(12, 13)	(4.25, 5.02)	10134	0.4633
(13, 14)	(5.02, 5.91)	7839	0.4981
(14, 15)	(5.91, 6.92)	6650	0.5434
(15, 16)	(6.92, 8.08)	5543	0.5622
(16, 17)	(8.08, 9.42)	5027	0.5832
(17, 18)	(9.42, 10.9)	4868	0.3081
(18, 19)	(10.9, 14.7)	3404	0.2805

MAST

Select Assay:

logcounts

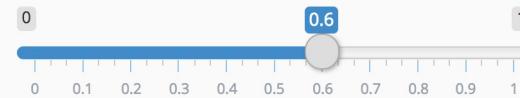
Adaptive Thresholding:

Run Thresholding

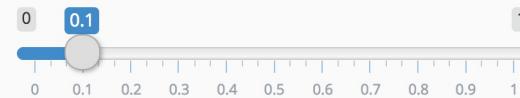
Hurdle Model:

 Use Adaptive Thresholds

Select fold change threshold



Select expression threshold



Select Condition for Hurdle Model

condition

p-value (FDR) cutoff:



Run DE Using Hurdle



SCTK features – *Pathway activity analysis - GSVA*



Enter Top 'N' Genes value:

Max genes: 19972

Apply p-value Cutoff

Apply logFC Cutoff

Scale Expression values?

Options

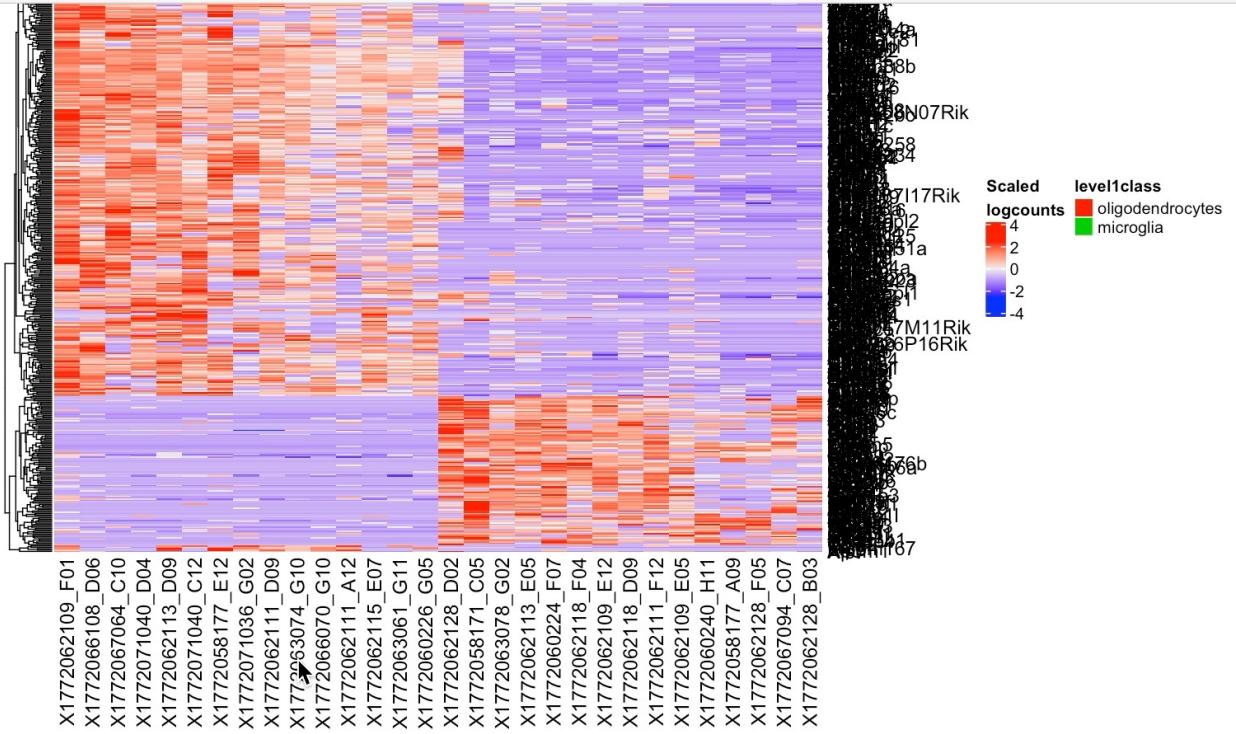
General Options

Row Labels Column Labels

Column Dendograms Row Dendograms

Cluster Rows Cluster Columns

Columns Title



Sequencing Depth

Number of cells

Snapshot

Minimum readcount to detect gene

10

Minimum number of cells with nonzero expression to detect gene

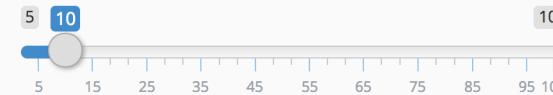
3

Number of bootstrap iterations.

10

Maximum log₁₀(number of simulated reads)

how many values to simulate



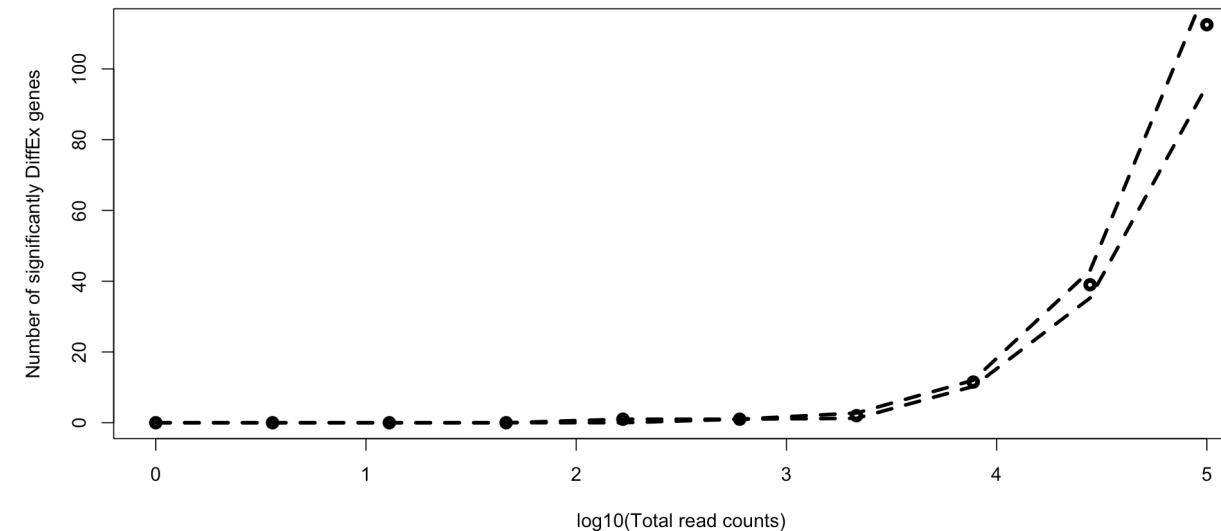
Condition for diffex

Run subsampler

Genes Detected

Minimum Detectable Effect Size

Number of Diffex Genes



Thanks

- Johnson Lab
 - David Jenkins
 - Mohammad Khan
 - Tyler Faits
 - Yuqing Zhang
 - Yue Zhao
- Additional Toolkit Contributors
 - Emma Briars
 - Sebastian Carrasco Pro
 - Steve Cunningham
 - Sean Corbett
- Bild Lab (University of Utah)
 - Andrea Bild
 - Mumtahena Rahman
 - Shelley Macneil
 - Sam Brady
- Single Cell Workgroup
 - Josh Campbell
 - Masanao Yajima

Common Questions

- How is this different from other tools/packages?
 - (Is this the right question?)
- Can we import data to/from other tools?
 - CellRanger object
 - To/from Seurat
- Is it scalable to larger datasets?

Decontamination of ambient RNA with DecontX

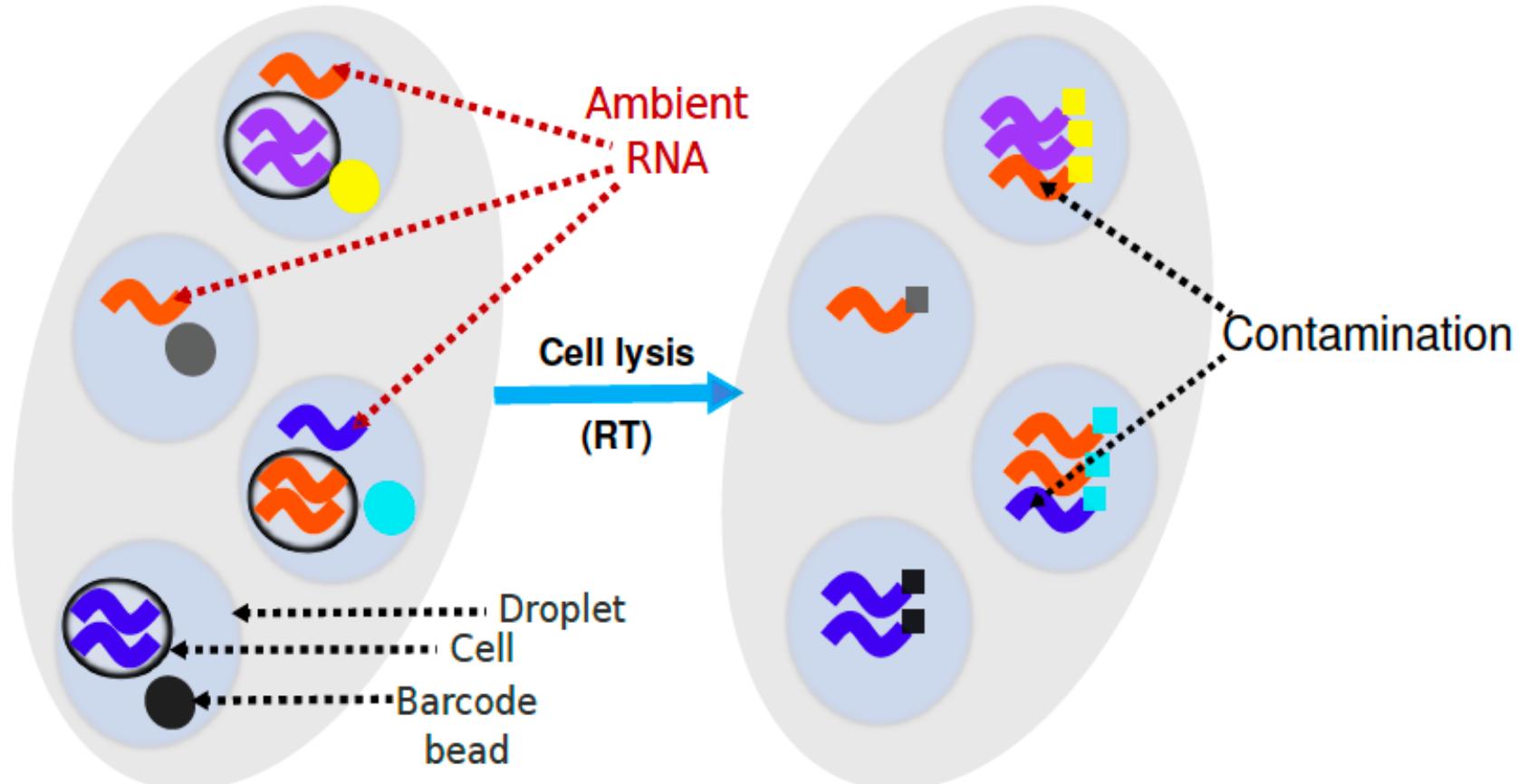
<https://www.biorxiv.org/content/10.1101/704015v1>



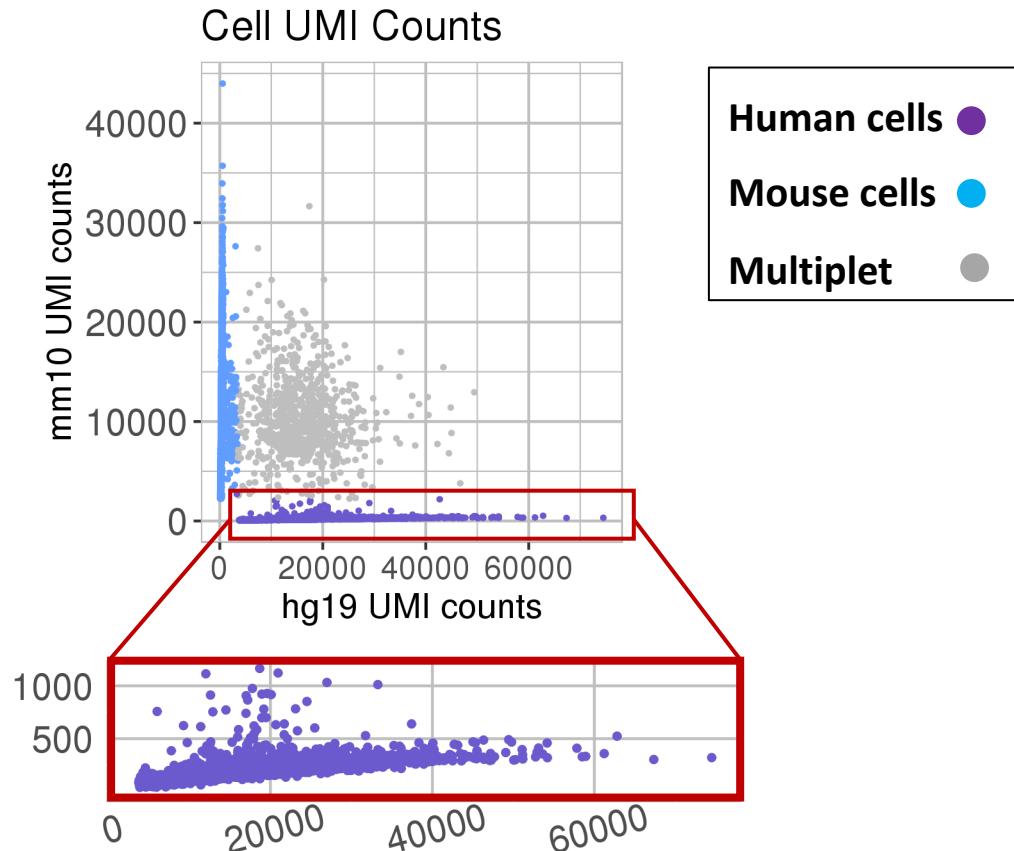
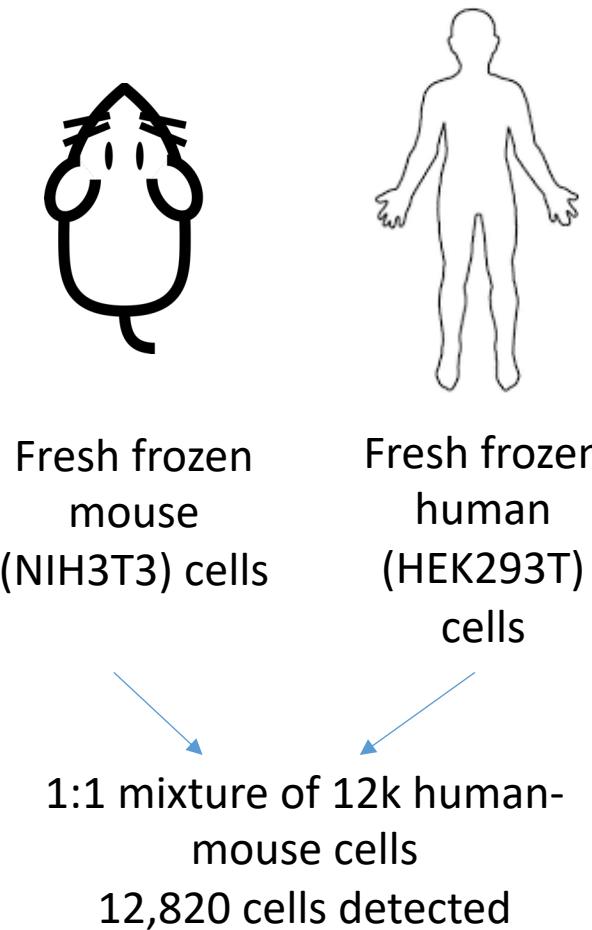
<https://github.com/campbio/celda>



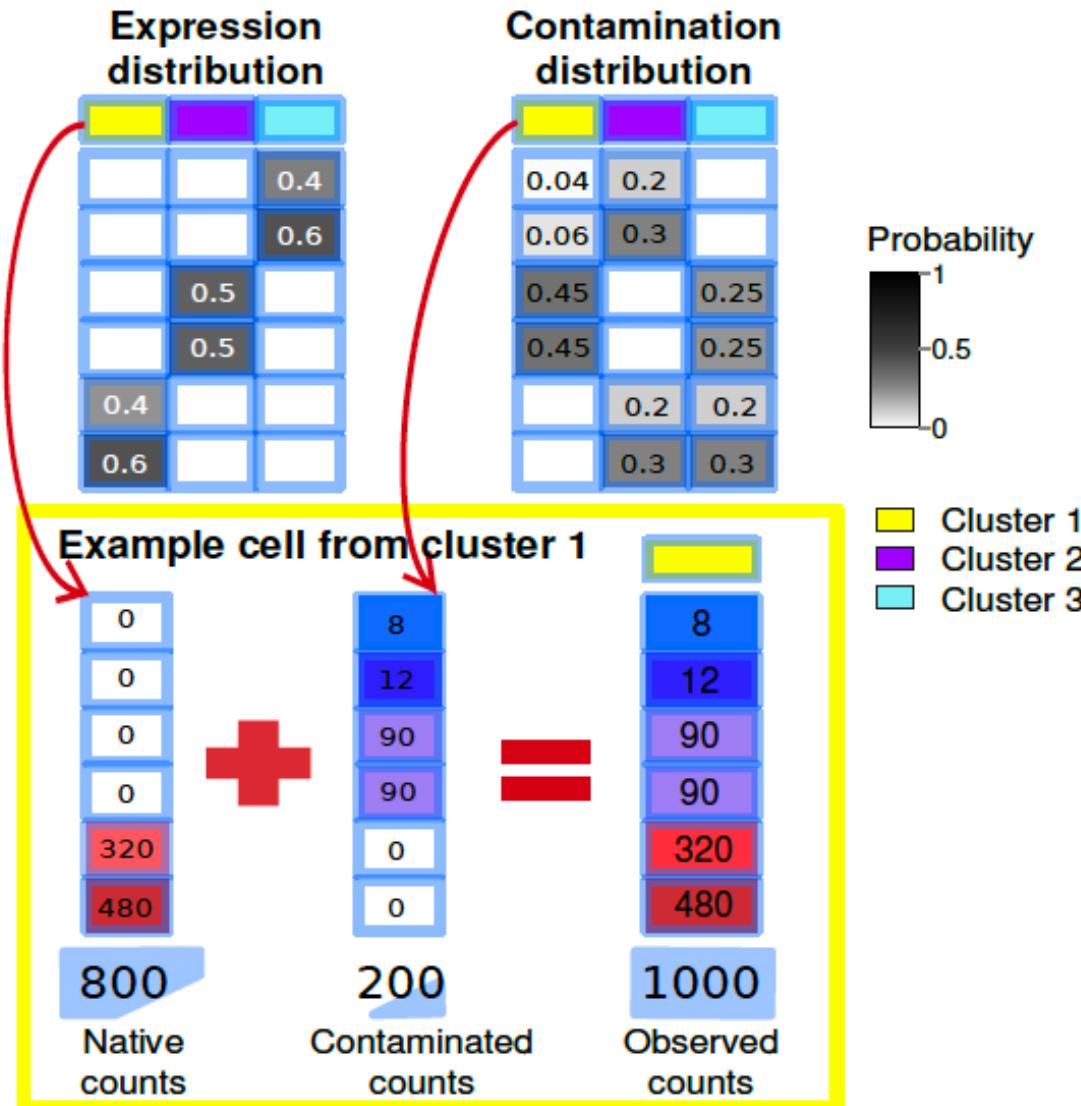
Contamination from ambient RNA in scRNA-seq data



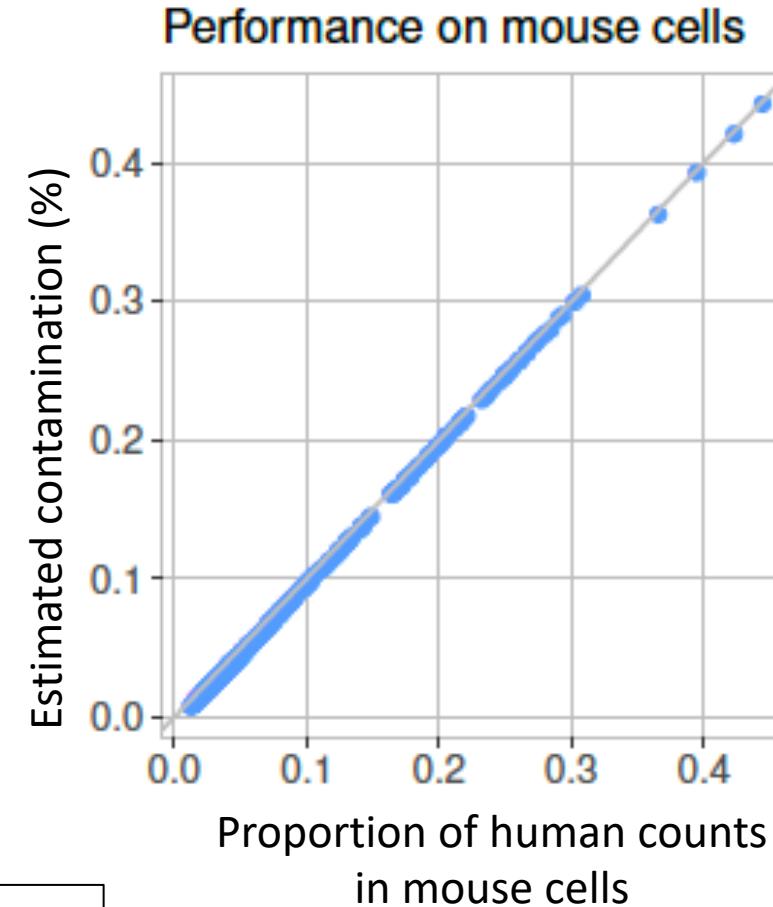
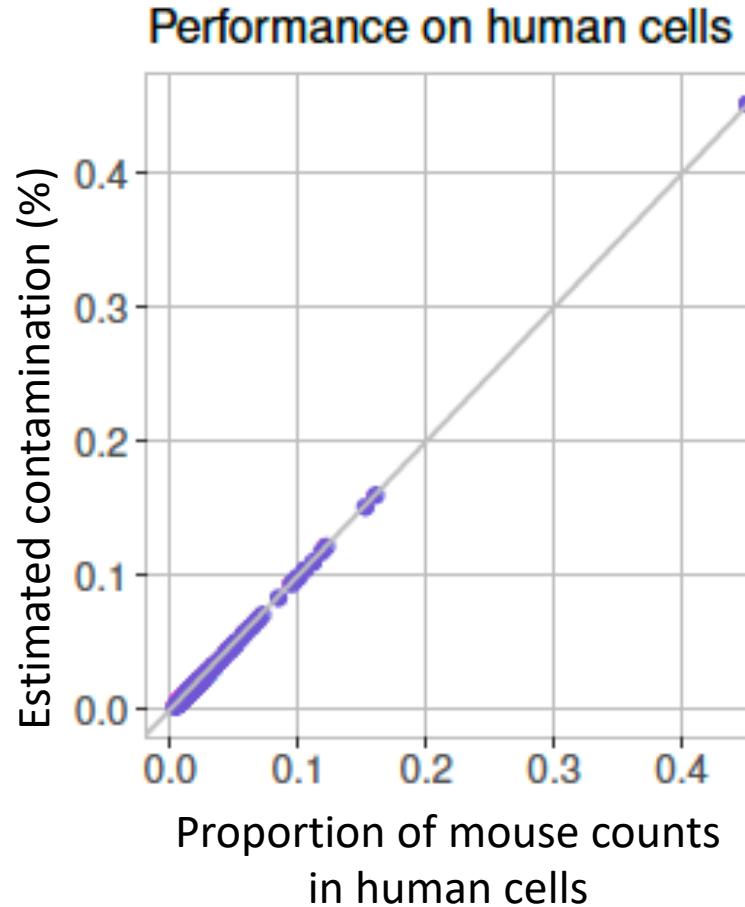
Example of contamination in a mouse-human cell mixture dataset from 10X.



Each cell is a mixture of transcripts from its native cell type and a weighted combination of all other cell types.

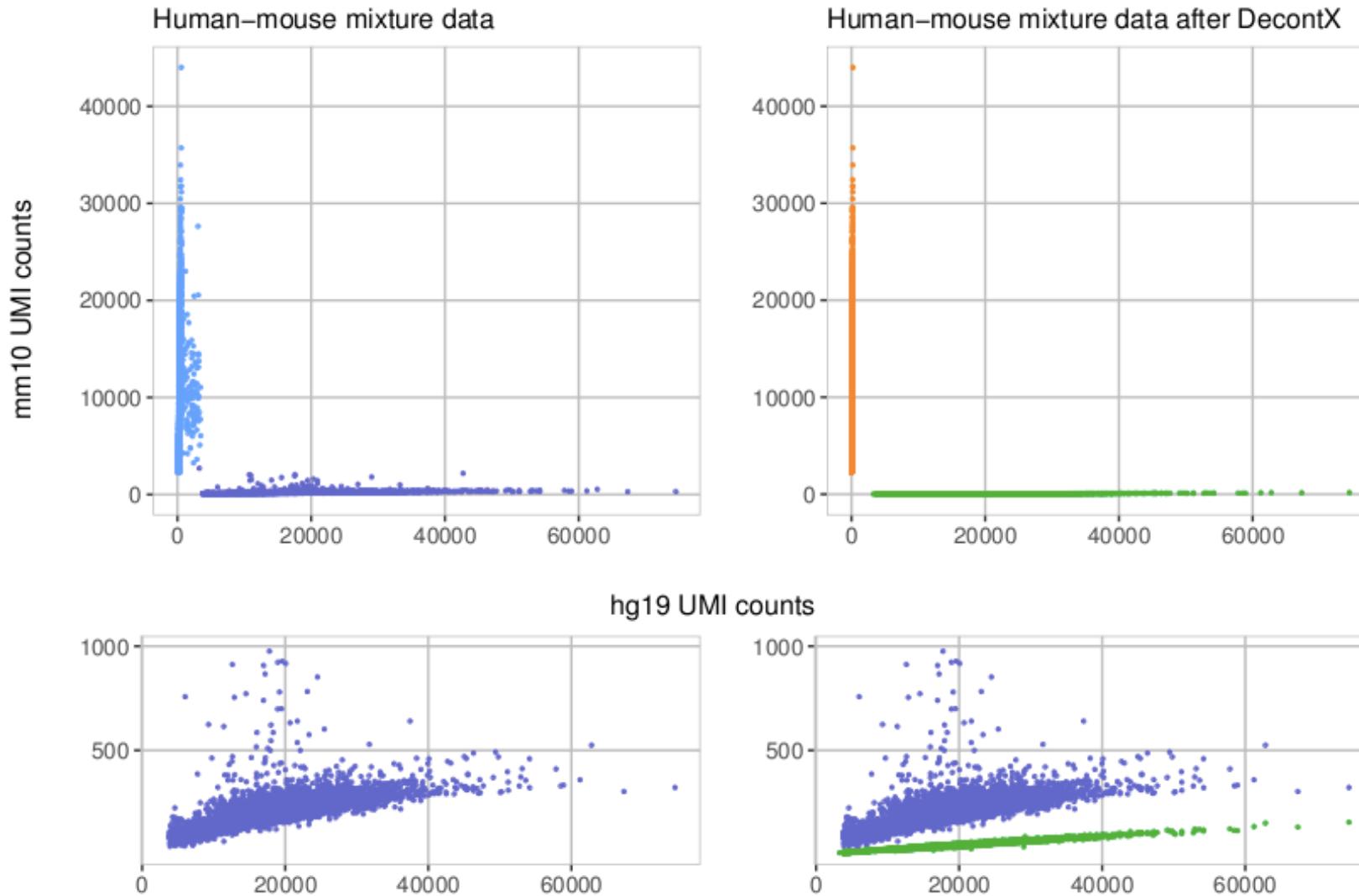


DecontX-estimated % contamination is highly correlated with the proportion of counts from the other organism.



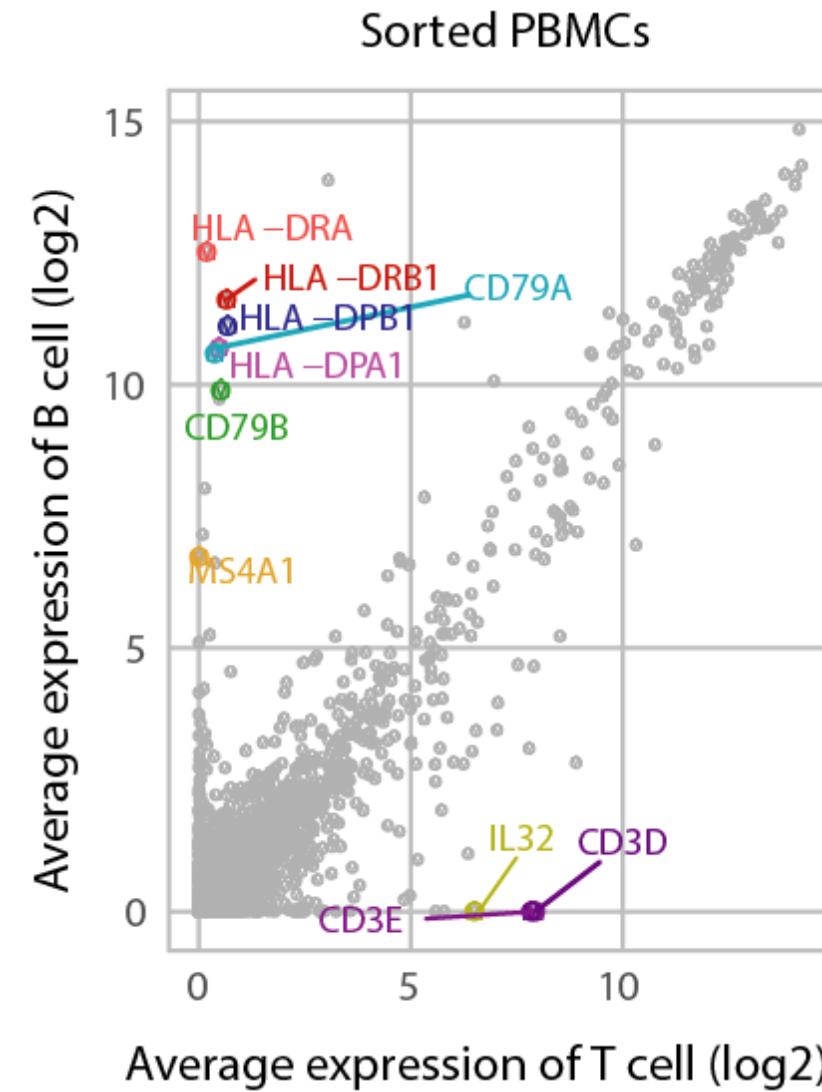
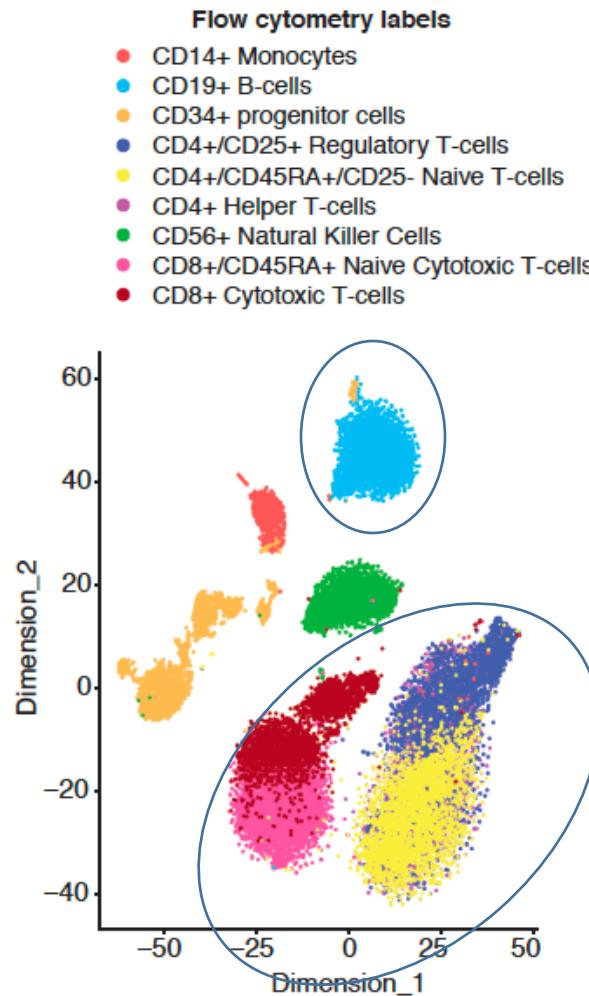
Human cells ●
Mouse cells ●

DecontX largely removed the contamination counts from ambient RNA of the other organism.

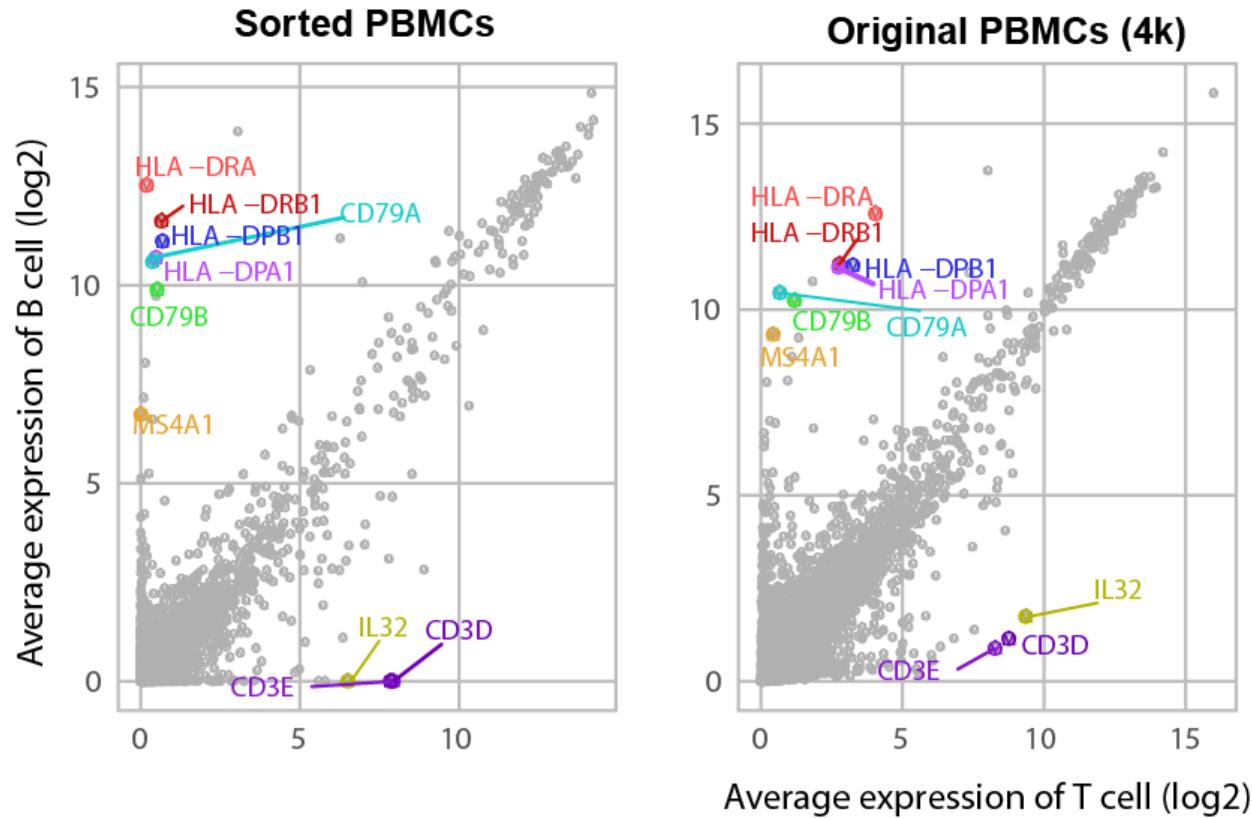


PBMCs sorted and profiled in separate channels of 10X showed little cross contamination.

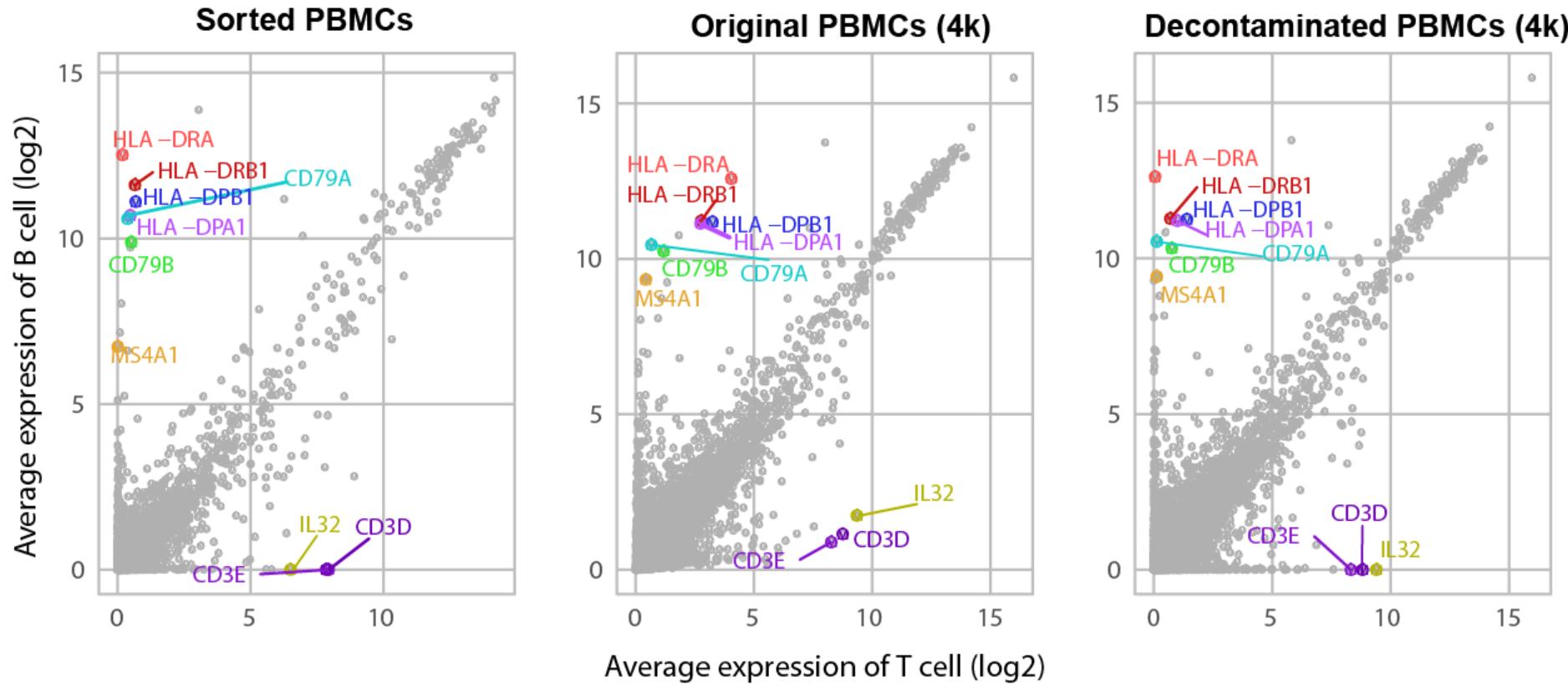
Number of cells: 84,431



PBMCs that were profiled in the same channel of 10X showed higher levels of marker gene cross contamination.



DecontX removed the majority of marker genes cross contamination from T- and B-cell populations.



When cell types are profiled in different channels, cell-type specific genes are lowly detectable in other populations.

B-cells

CD79A
CD79B
MS4A1

T-cells

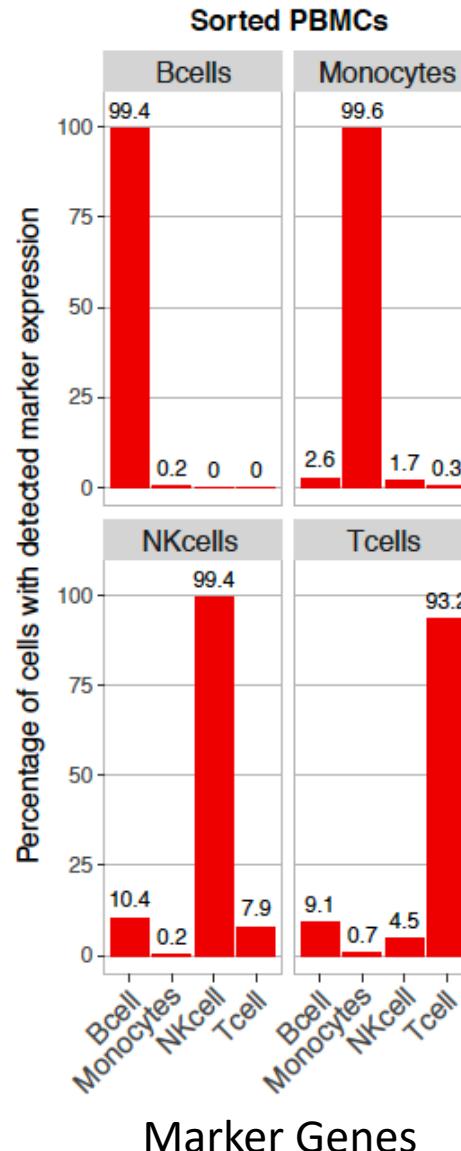
CD3E
CD3D

NK-cells

GNLY

Monocytes

S100A8
S100A9
VCAN



When cell types are profiled in the same channel, cell-type specific genes are detected at high levels in other populations.

B-cells

CD79A
CD79B
MS4A1

T-cells

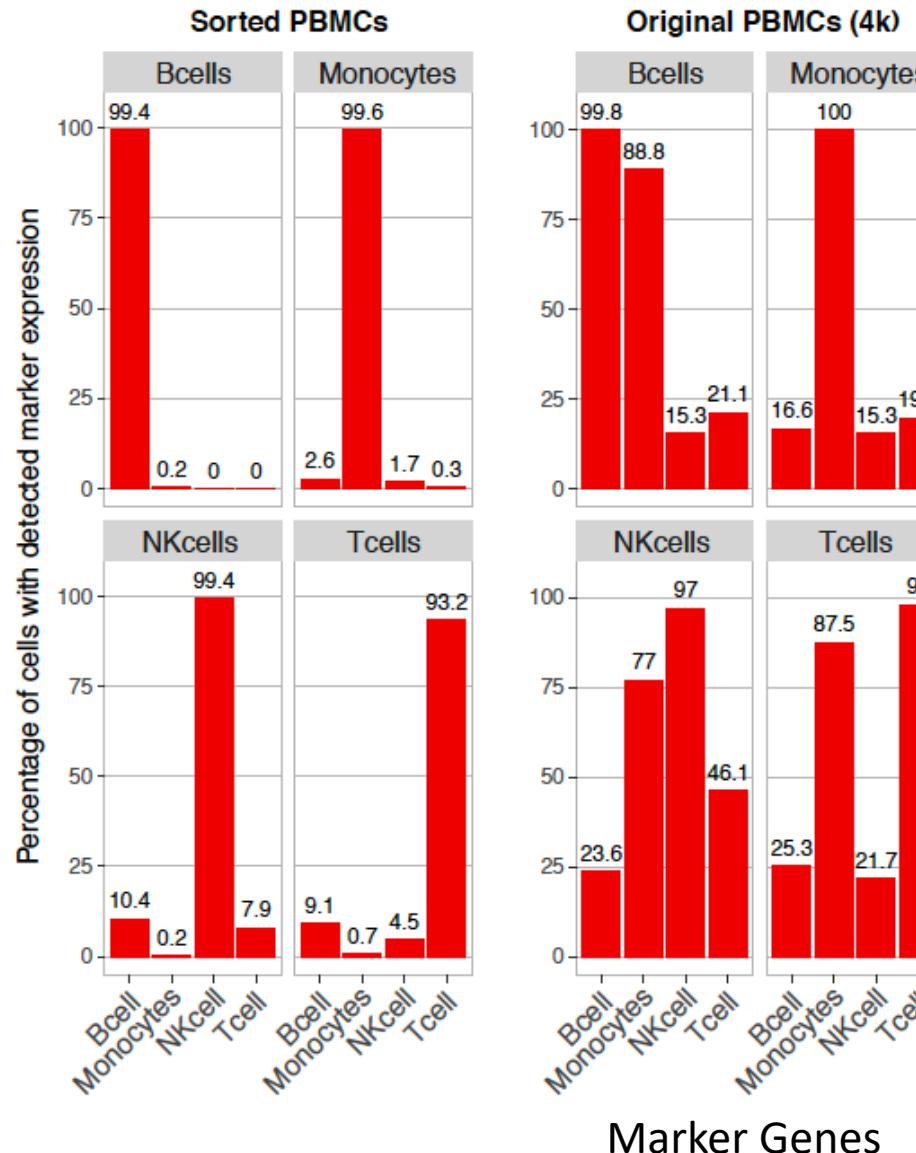
CD3E
CD3D

NK-cells

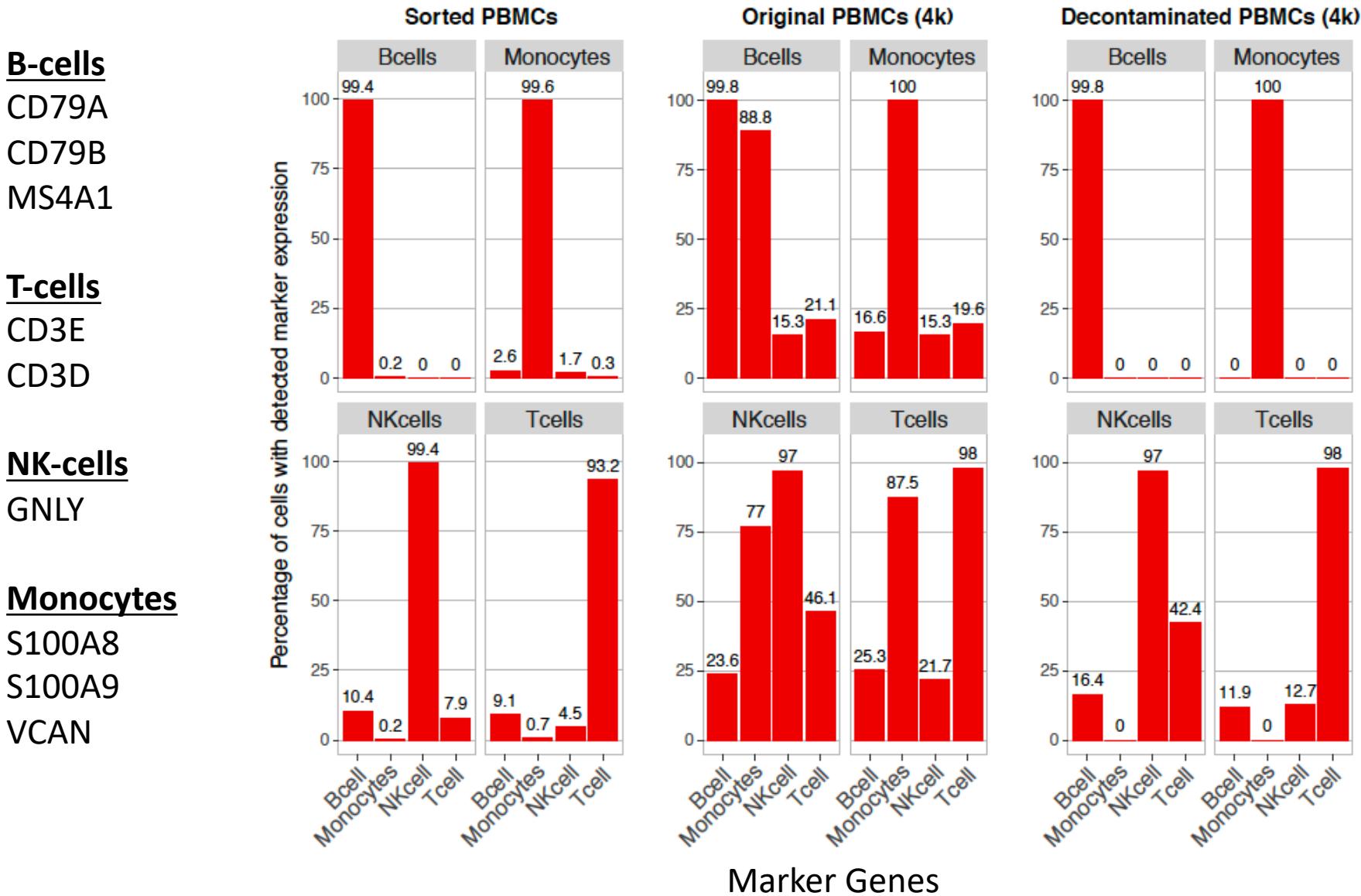
GNLY

Monocytes

S100A8
S100A9
VCAN

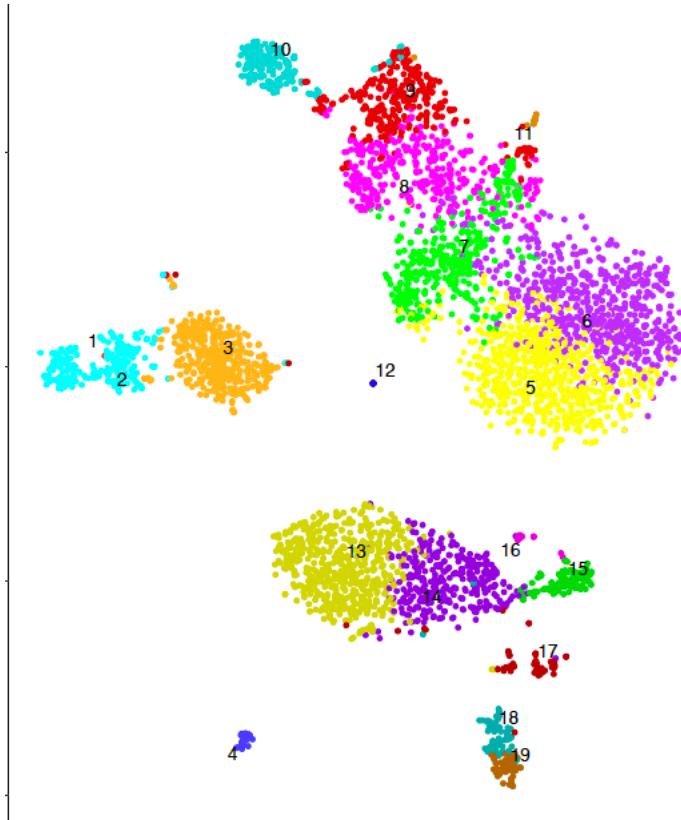


DecontX removed the majority of marker gene cross-contamination from all populations.

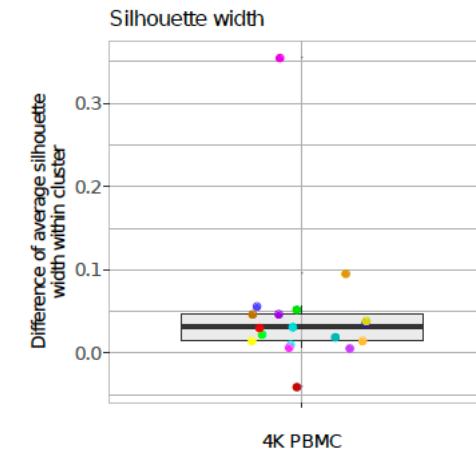
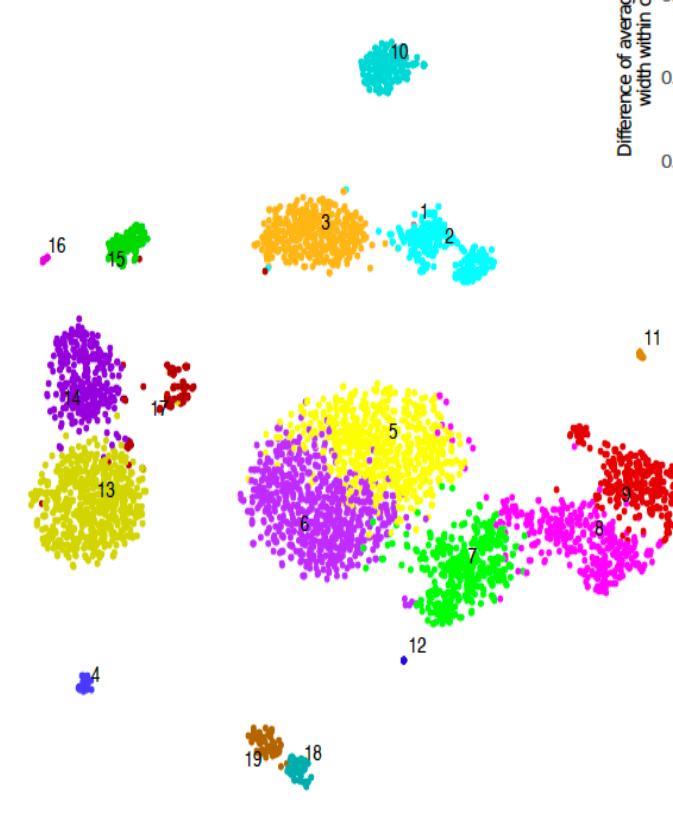


Decontamination improved overall cell cluster separation.

Original Counts
PBMC 4K



Decontaminated Counts
PBMC 4K



Average contamination level in PBMC 4K was 6.9%

Summary

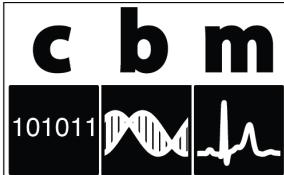
- Estimation of contamination from ambient RNA can be useful to quality control experimental pipelines.
- Decontaminated counts can improve downstream clustering and visualization.
- DecontX is available in the Celda R package on Bioconductor and Github:



<https://github.com/campbio/celda>



DecontX Acknowledgements



**Boston University School of Medicine
Computational Biomedicine**

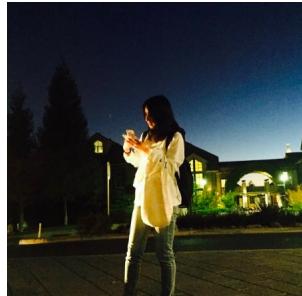
Masanao
Yajima



Evan
Johnson



Shiyi (Iris)
Yang



Zhe
Wang



Yusuke
Koga



Sean
Corbett



Eric Reed
Ahmed Youssef
Jing Zhang

Salam Alabdullatif
Ali Lashkaripour
Zhaorong Li
Timo Hu



<https://github.com/campbio/celda>



QC pipeline development

1. Raw Data

Single cell/nuc RNA-seq

10X
inDrops
CEL-seq2
Drop-seq
SMART-seq2

CITE-seq/Total-seq

10X
inDrops

TCR/BCR-seq
10X

2. Preprocessing

CellRanger
STARsolo
BUStools
dropEST
etc.

scruff
Import into R

SingleCellExperiment
MultiAssayExperiment

Export to Python object

3. Quality control

Standard metrics
Doublets
Ambient RNA
Seq. Depth/Saturation
Empty Drop Detection
Batch correction

Interactive
visualization and
analysis

SCTK

Summary of QC metrics

<u>Standard metrics</u>	<u>Doublet Detection</u>	<u>Ambient RNA estimation</u>
- Total UMI per cell	Scrublet	DecontX
- Total Genes per cell	doubletCells (Scran)	
- % Mitochondria per cell	DoubletFinder	
- % counts in top X features	DoubletDecon	
- 50, 100, 200, 500	DoubletDetection	

Additional items?

Interactive scRNA-Seq analysis with the Single Cell Toolkit (SCTK)

Evan Johnson and Josh Campbell
Division of Computational Biomedicine
Boston University School of Medicine

May 1, 2020

Single Cell Toolkit (SCTK)

- Standard R package w/Shiny toolkit
- Interacts and operates on an SCE (SingleCellExperiment) object
- SCE object in/out of Shiny any stage
- Great for common tasks:
 - Interactive clustering/visualization
 - Differential gene/pathway analysis
- (Also works for bulk RNA-seq analysis)

The screenshot shows the Bioconductor website with the following details:

- Header:** Bioconductor OPEN SOURCE SOFTWARE FOR BIOINFORMATICS, Home, Install, Help, Developers, About, Search.
- Breadcrumbs:** Home > Bioconductor 3.9 > Software Packages > singleCellTK
- Title:** singleCellTK
- Metrics:** platforms all, rank unknown, posts 0, in Bioc 1 year, build ok, updated before release.
- DOI:** [10.18129/B9.bioc.singleCellTK](https://doi.org/10.18129/B9.bioc.singleCellTK)
- Description:** Interactive Analysis of Single Cell RNA-Seq Data
- Author:** David Jenkins
- Maintainer:** David Jenkins <dfj@bu.edu>
- Citation:** (from within R, enter `citation("singleCellTK")`):
Jenkins D, Faits T, Khan MM, Briars E, Carrasco Pro S, Johnson WE (2019). *singleCellTK: Interactive Analysis of Single Cell RNA-Seq Data*. R package version 1.4.0, https://combiomed.github.io/sctk_docs/.
- Installation:** To install this package, start R (version "3.6") and enter:

```
if (!requireNamespace("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install("singleCellTK")
```
- Documentation:** To view documentation for the version of this package installed in your system, start R and enter:

```
knowc vignettes("singleCellTK")
```
- Documentation Sidebar:** Documentation, Bioconductor, R / CRAN packages and documentation.
- Support Sidebar:** Support, Posting guide, Support site - for questions about Bioconductor packages, Bioc-devel mailing list - for package developers.

Single Cell Toolkit v1.1.20 127.0.0.1

Upload Data Summary & Filtering Visualization & Clustering Batch Correction Differential Expression Enrichment Analysis Sample Size

Single Cell Toolkit

Filter, cluster, and analyze single cell RNA-Seq data

Need help? [Read the docs.](#)

Upload

(help)

Choose data source:

- Upload files
- Use example data

Choose Example Dataset:

mouseBrainSubset

Mouse Brain Subset: GSE60361

A subset of 30 samples from a single cell RNA-Seq experiment from Zeisel, et al. Science 2015. The data was produced from cells from the mouse somatosensory cortex (S1) and hippocampus (CA1). 15 of the cells were identified as oligodendrocytes and 15 of the cell were identified as microglia.

Upload

BOSTON UNIVERSITY c b m

Data Summary & Filtering

(help)

Data Summary Assay Details Annotation Data Visualize

Settings: Hide All Show All

Select Assay: counts

Delete Outliers
Filter Samples by Annotation
Filter Genes by Feature Annotation
Convert Gene Annotations
Delete an Annotation Column
Randomly Subset

Download SCKExperiment

Summary Contents:

Metric	Value
Number of Samples	30
Number of Genes	19972
Average number of reads per cell	8058
Average number of genes per cell	2395
Samples with <1700 detected genes	6
Genes with no expression across all samples	7923

Counts Histogram:

Single Cell Toolkit v0.3.9 Upload Data Summary and Filtering DR & Clustering Batch Correction Differential Expression Pathway Activity Analysis Sample Size

Pathway Activity Analysis

Select Assay: logcounts

Select Method: GSVA

Gene list source:

- Manual Input
- MSigDB c2 (Human, Entrez ID only)

Select Gene List(s): ALL

Number of top pathways: 25 (3,272)

Select Condition(s) of interest for plot: condition

Plot Type:

- Violin
- Heatmap

Run

Results Table

IZUKA, RECURRENT_LIVER_CANCER
OLSSON_EIF3_TARGETS_UP
KEGG_ETHER_LIPID_METABOLISM
JEON_SMAD5_TARGETS_DN
ZHAN_MULTIPLE_MYELOMA_MF_UP
KEGG_PEROXOSOME
KEGG_PRKR_SIGNALING_PATHWAY
WANG_SMARCF1_TARGETS_UP
KEGG_VEGF_SIGNALING_PATHWAY
KAAB_HEART_ATRIUM_VS_VENTRICLE_UP
KRIE_RESPONSE_TO_TOSEDOSTAT_6HR_UP
BROWNE_HCMV_INFECTION_12HR_DN
NAGASHIMA_NRG1_SIGNALING_DL
VERRECCHIA_EARLY_RESPONSE_TO_SGRB4A
ZHU_CMV_24_HR_DN
ZUCCHI_METASTASIS_DN
CHIBA_RESPONSE_TO_TSAL_DN
KEGG_BASE_EXCISION_REPAIR
KEGG_PROPANOATE_METABOLISM
KEGG_BLADDER_CANCER
ZHAN_MULTIPLE_MYELOMA_CD1_AND_CD2_D
ZHU_CMV_8_HR_UP
MOTHRA_POC
KEGG_ALZHEIMERS_DISEASE
KEGG_HUNTINGTONS_DISEASE

Single Cell Toolkit Upload Data Summary & Filtering Visualization & Clustering Batch Correction Differential Expression Enrichment Analysis Sample Size

Samplewise Visualization and Clustering

(help)

Run New Dimensional Reduction

Select: Assay: logtpm
Method: PCA
Run

DR Options: reducedDimName: logtpm
View More Options

Available Reduced Dims:

Reduced Dimension
TSNE_logtpm
PCA_logtpm

Remove a reducedDim: TSNE_logtpm Delete

Visualization Settings:

Select Reduced Dimension Data: PCA_logtpm
Axis Settings
X Axis: PC1
Y Axis: PC2
Style Settings
Color points by: condition
Shape points by: No Shape
Update Plot

Visualize

“inDrops” is custom microfluidics device that can process large numbers of cells with high capture rates.

Allon Klein

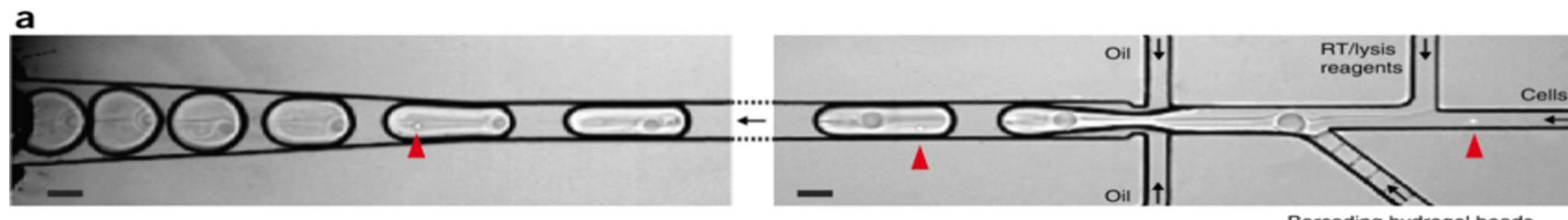


Harvard Medical School

PROTOCOL

Single-cell barcoding and sequencing using droplet microfluidics

Rapolas Zilionis^{1,2}, Juozas Nainys¹, Adrian Veres^{2–4}, Virginia Savova², David Zemmour⁵, Allon M Klein² & Linas Mazutis¹



Goals of parent R33:

- Aim 1. Optimize inDrops for microscopic samples and fixed cells.
- Aim 2. Optimize inDrops for low-cost, high-throughput, high sensitivity targeted transcriptomics.
- Aim 3. Integrate single cell genomics with histopathology.

Aims of ITCR-IMAT collaboration

IMAT (HMS - Klein)

Aim 1: Develop a Total-seq protocol for the inDrop system.

Aim 2: Develop single cell ATAC-seq for inDrop.

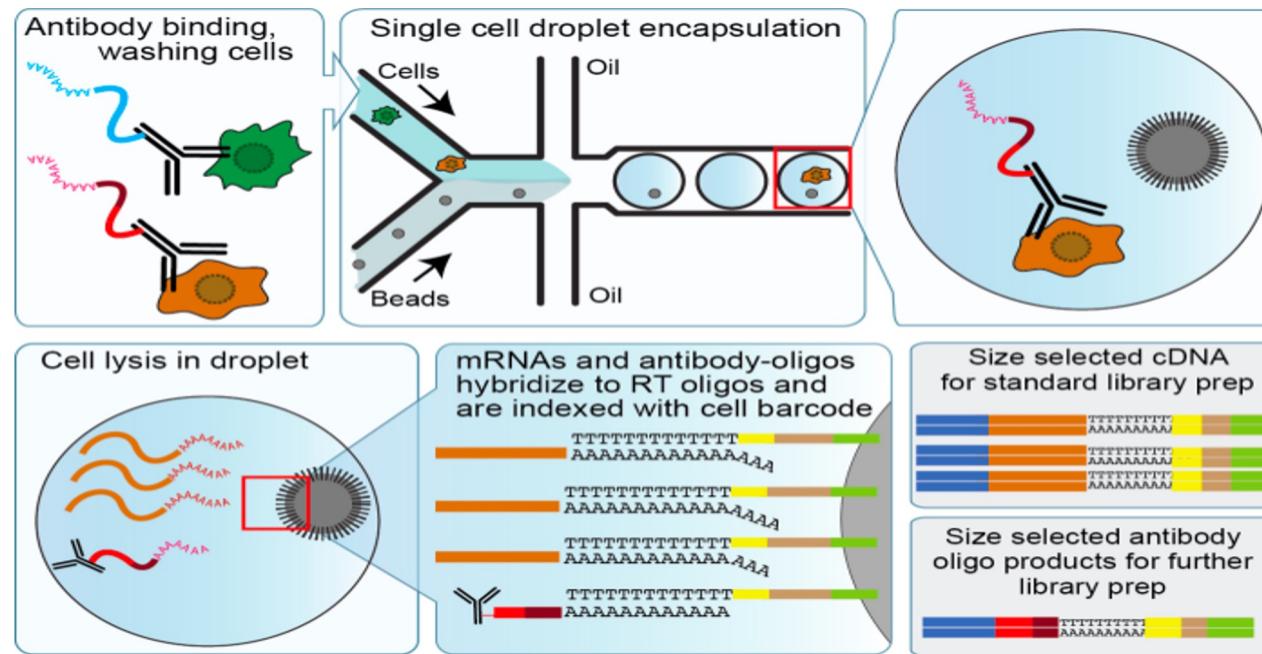
ITCR (BU – Johnson/Campbell)

Aim 1: Develop computational modules for analysis and display of single-cell Total/CITE-seq and ATAC-seq data generated from tumor specimens.

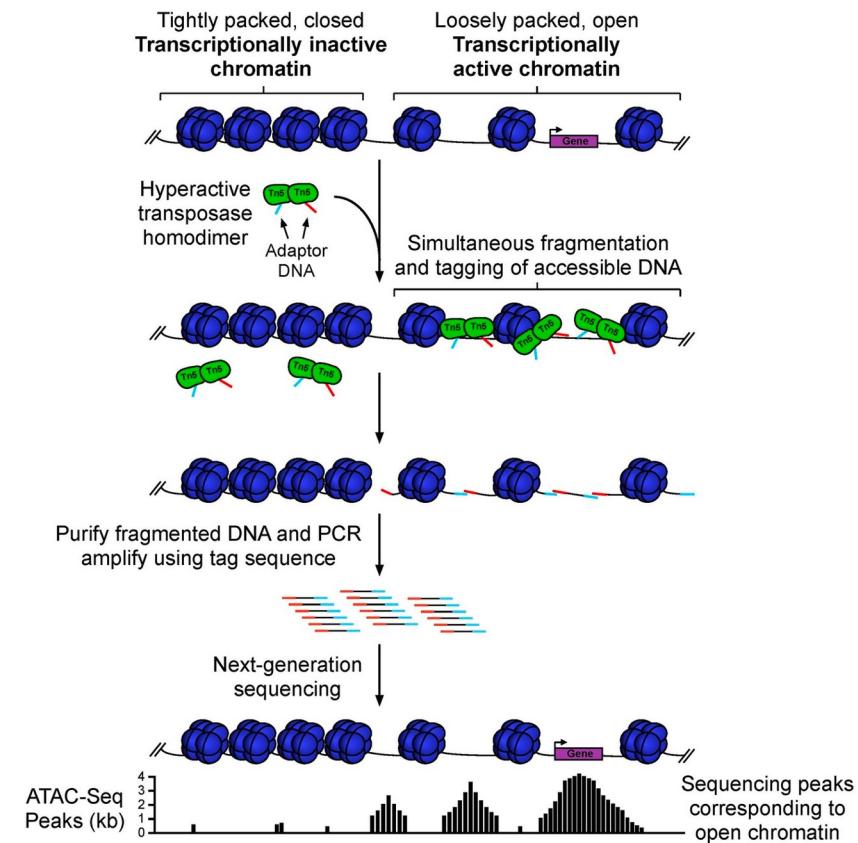
Aim 2: Develop computational modules for comprehensive assessment and correction of batch effects or sample-specific effects across tumor specimens.

Aim 1: Develop computational modules for analysis and display of single-cell CITE-seq and ATAC-seq data generated from tumors.

Total-Seq/CITE-seq antibody derived tags (ADTs) can be used to measure protein levels at single cell resolution



ATAC-seq (Assay for Transposase-Accessible Chromatin using sequencing) is a technique used in molecular biology to assess genome-wide chromatin accessibility.



Comprehensive scRNA-seq QC pipeline for the Human Tumor Atlas Network (HTAN)

1. Raw Data

Single cell/nuc RNA-seq



2. Preprocessing

CellRanger

STARsolo

BUSTools

dropEST

HCA Optimus

SEQC

SCTK

Import into R

SingleCellExperiment
MultiAssayExperiment

3. Quality control

Standard metrics

Doublets

Ambient RNA

Empty Drop Detection

Interactive
visualization and analysis

SCTK

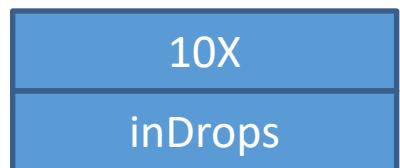
Expansion to include import functions for Total-seq/CITE-seq or scATAC-seq data in R data containers.

1. Raw Data

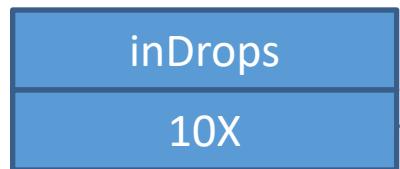
Single cell/nuc RNA-seq



CITE-seq/Total-seq



scATAC-Seq



2. Preprocessing

CellRanger

STARsolo

BUSTools

dropEST

HCA Optimus

SEQC

SCTK

Import into R

SingleCellExperiment
MultiAssayExperiment

3. Quality control

Standard metrics

Doublets

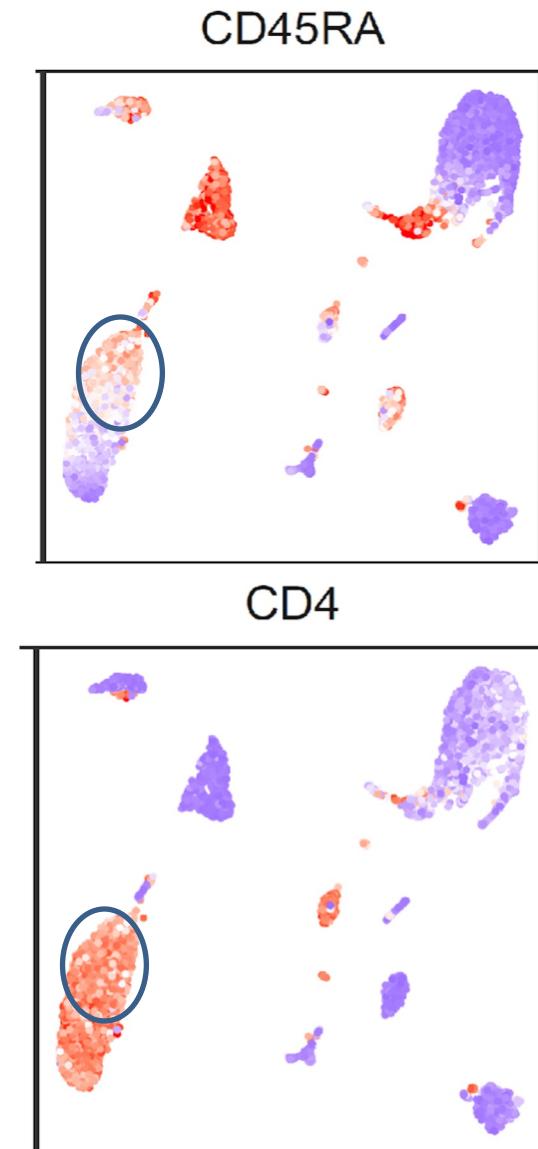
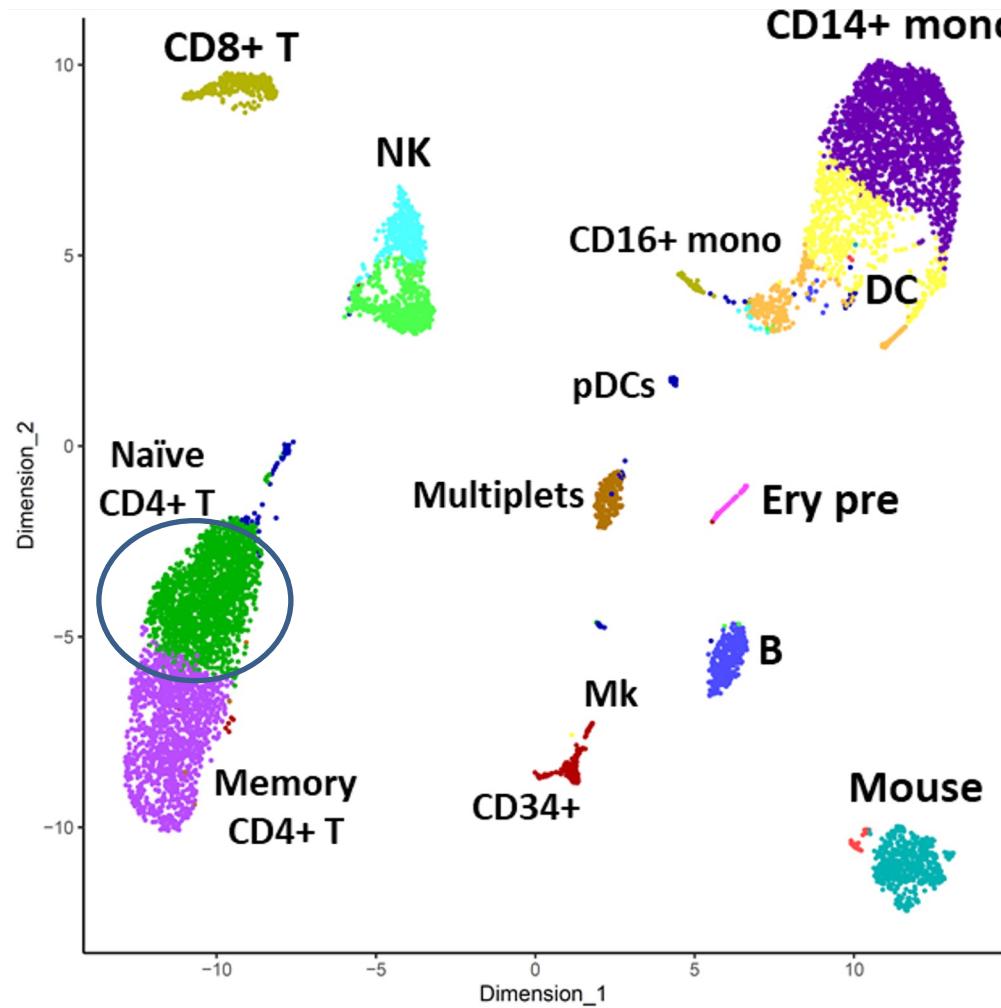
Ambient RNA

Empty Drop Detection

Interactive
visualization and analysis

SCTK

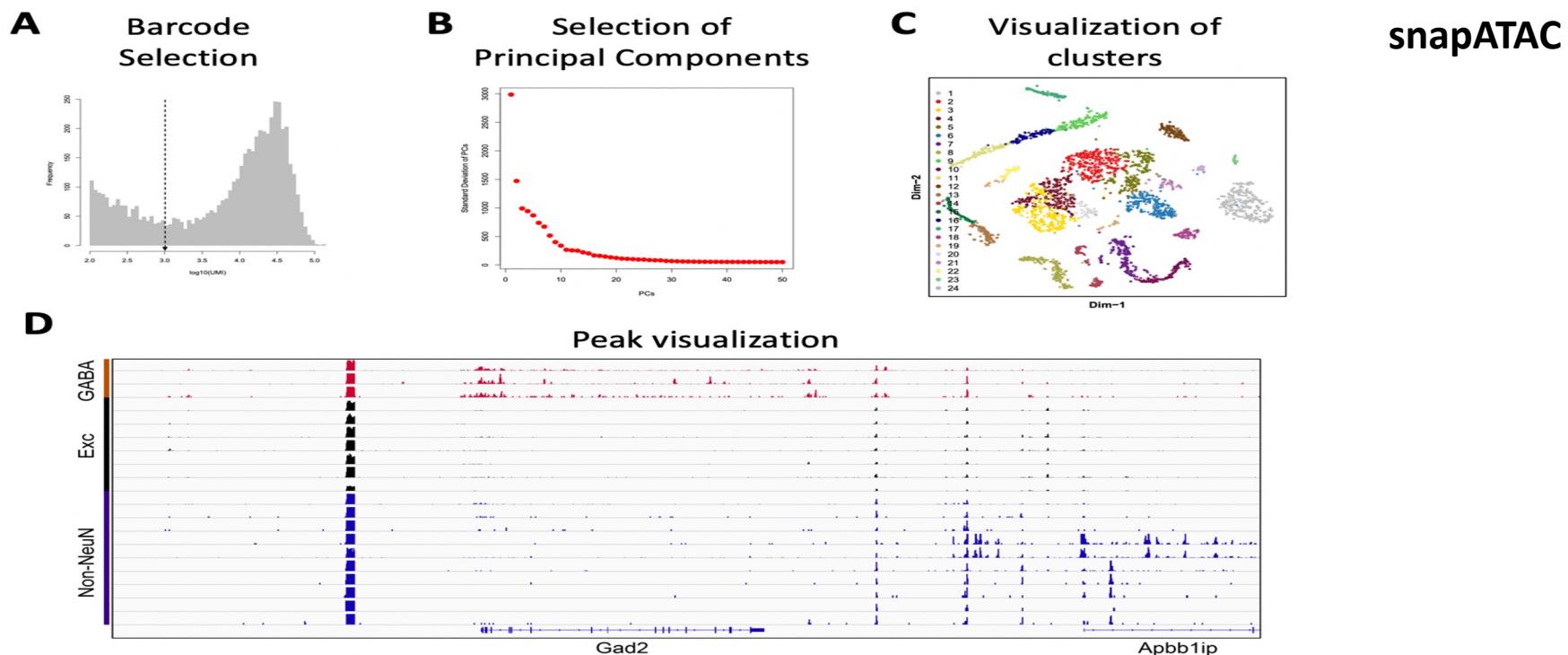
Developing novel statistical approaches for joint clustering of ADT (protein) and scRNA-seq data.



Develop computational modules for analysis and display of single-cell ATAC-seq data generated from tumors .

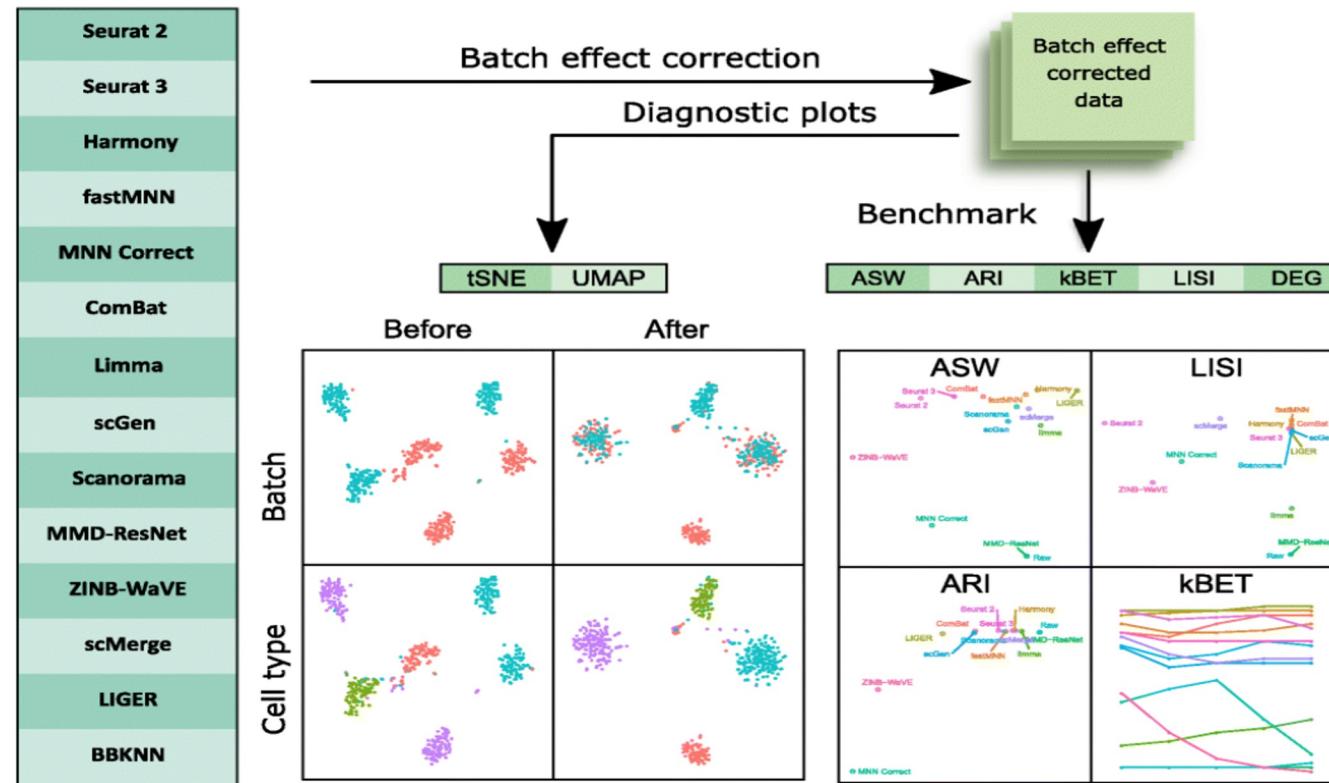
Fast and Accurate Clustering of Single Cell Epigenomes Reveals Cis-Regulatory Elements in Rare Cell Types

Rongxin Fang, Sebastian Preissl, Xiaomeng Hou, Jacinta Lucero, Xinxin Wang, Amir Motamed, Andrew K. Shiu, Eran A. Mukamel, Yanxiao Zhang, M. Margarita Behrens, Joseph Ecker, Bing Ren



Aim 2: Develop computational modules for comprehensive assessment and correction of batch effects or sample-specific effects across tumors.

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



Batch correction tools implemented in SCTK:

- 1) ComBat/ComBat-Seq
- 2) Seurat3 Integration
- 3) Harmony
- 4) scMerge
- 5) FastMNN
- 6) MNNcorrect
- 7) BBKNN
- 8) LIGER
- 9) scGEN
- 10) Scanorama
- 11) ZinB-wave

ComBat-Seq: Appropriate assumptions for count data

Mean-variance dependence in RNA-seq counts:

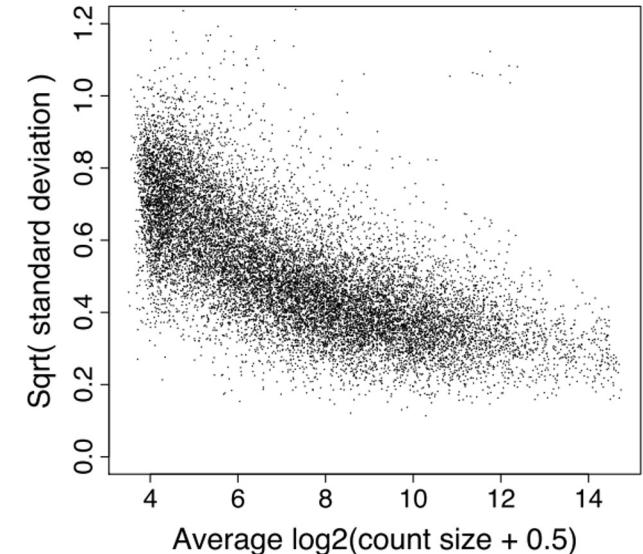
- Over-dispersion (variance > mean)
- Genes with smaller counts tend to have larger variance

Negative Binomial (NB):

$$y \sim NB(\mu, \phi)$$

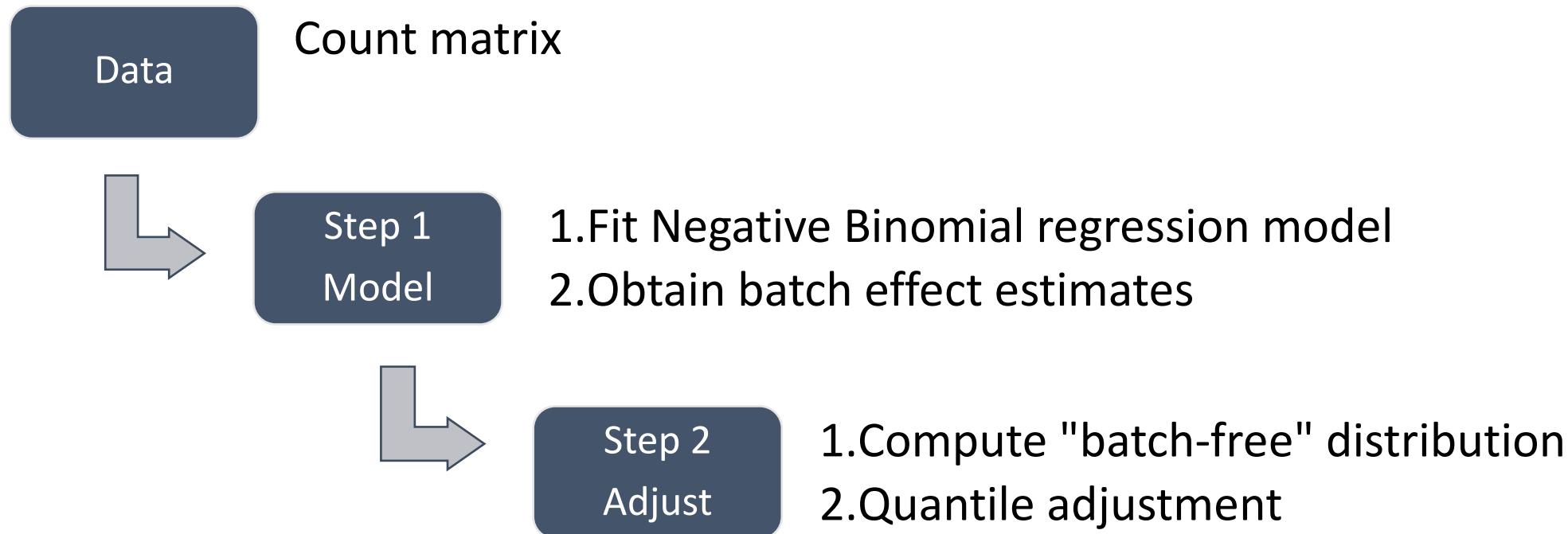
- Negative Binomial regression used in edgeR & DESeq2
- Variance is a function of mean

$$\text{var}(y) = \mu + \phi\mu^2$$



Law, Charity W., et al. "voom: Precision weights unlock linear model analysis tools for RNA-seq read counts." Genome biology 2014

ComBat-Seq algorithm



ComBat-Seq algorithm: Model

Negative Binomial regression

Gene-wise model: for a certain gene g , count in sample j from batch i $| y_{gij} \sim NB(\mu_{gij}, \phi_{gi})$

$$\log \mu_{gij} = \alpha_g + X_j \beta_g + \gamma_{gi} + \log N_j \quad Var(y_{gij}) = \mu_{gij} + \phi_{gi} \mu_{gij}^2$$

Decompose scaled counts
into 3 components

α_g	Average level for gene g (in “negative” samples)
$X_j \beta_g$	Biological condition of sample j
γ_{gi}	Mean batch effect

ϕ_{gi} | Variance batch effect

ComBat-Seq algorithm: Adjust

Adjust the data

- Calculate parameters for “batch-free” distribution: $y_{gj}^* \sim NB(\mu_{gj}^*, \phi_g^*)$

$$\log \mu_{gj}^* = \log \hat{\mu}_{gij} - \hat{\gamma}_{gi}$$

$$\phi_g^* = \frac{1}{N_{batch}} \sum_i \hat{\phi}_{gi}$$

- Map quantiles from empirical distribution to the batch-free distribution

ComBat-Seq algorithm: Adjust

Adjust the data

Original count matrix

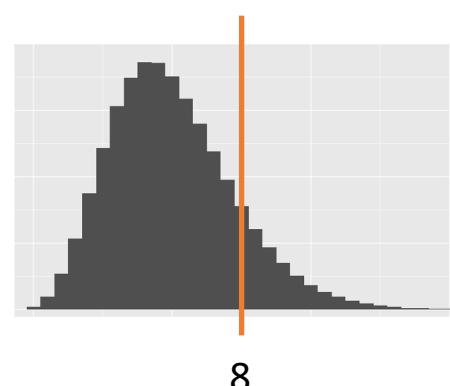
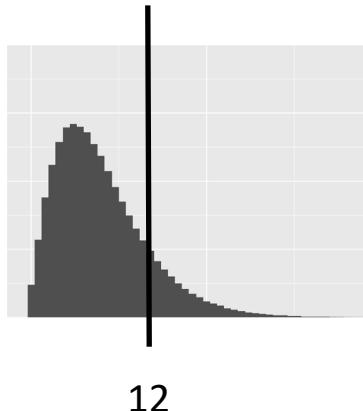
	S 1	S 2
G 1	14	12
G 2	0	0
G 3	112	11
...		

Adjusted count matrix

	S 1	S 2
G 1	9	8
G 2	0	0
G 3	60	47
...		

Empirical distribution
of original counts:

$$y_{gij} \sim NB(\hat{\mu}_{gij}, \hat{\phi}_{gi})$$



Batch-free distribution
for adjusted counts:

$$y_{gj}^* \sim NB(\mu_{gj}^*, \phi_g^*)$$

ComBat Methods for scRNA-seq Analysis

ComBat-Cell-Seq

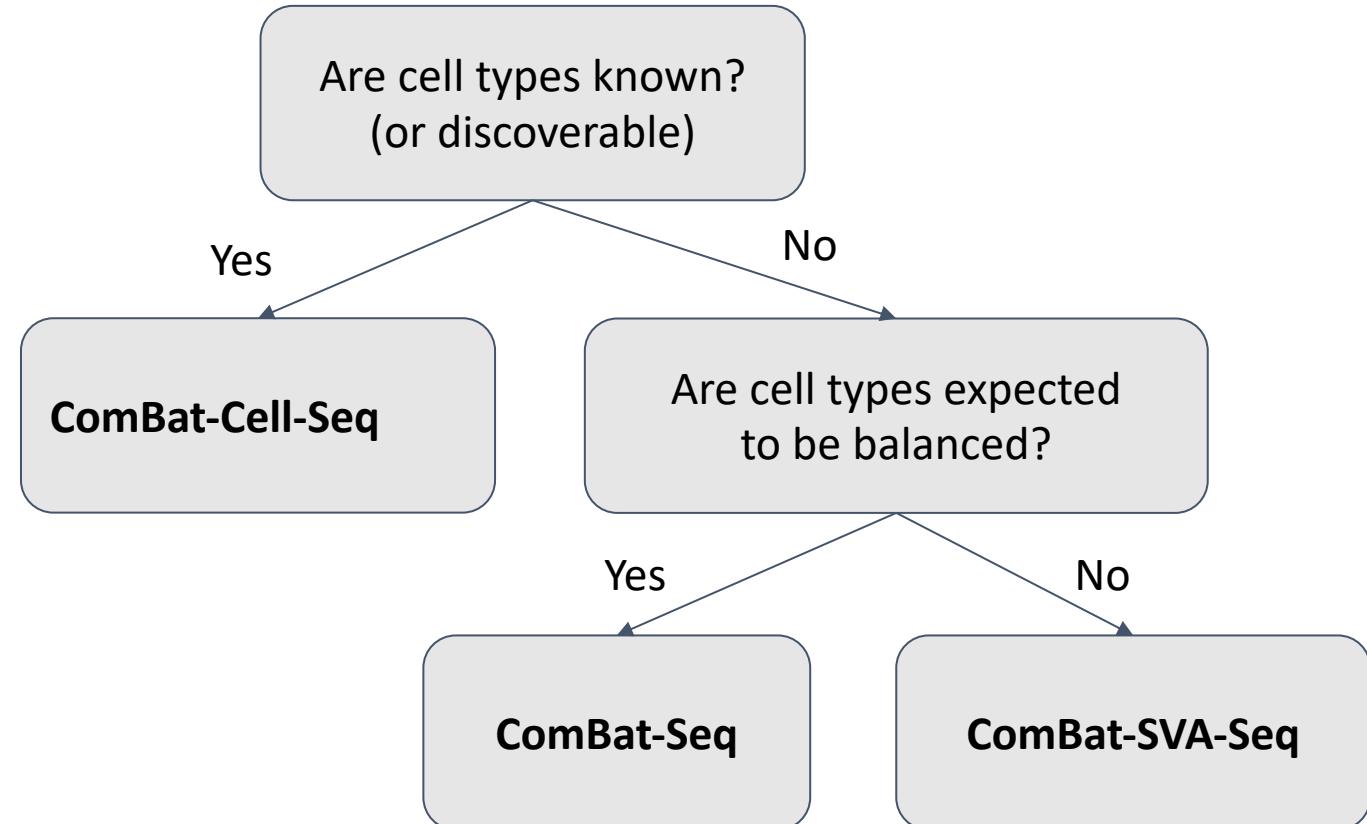
- Include cell types in ComBat design model

ComBat-Seq

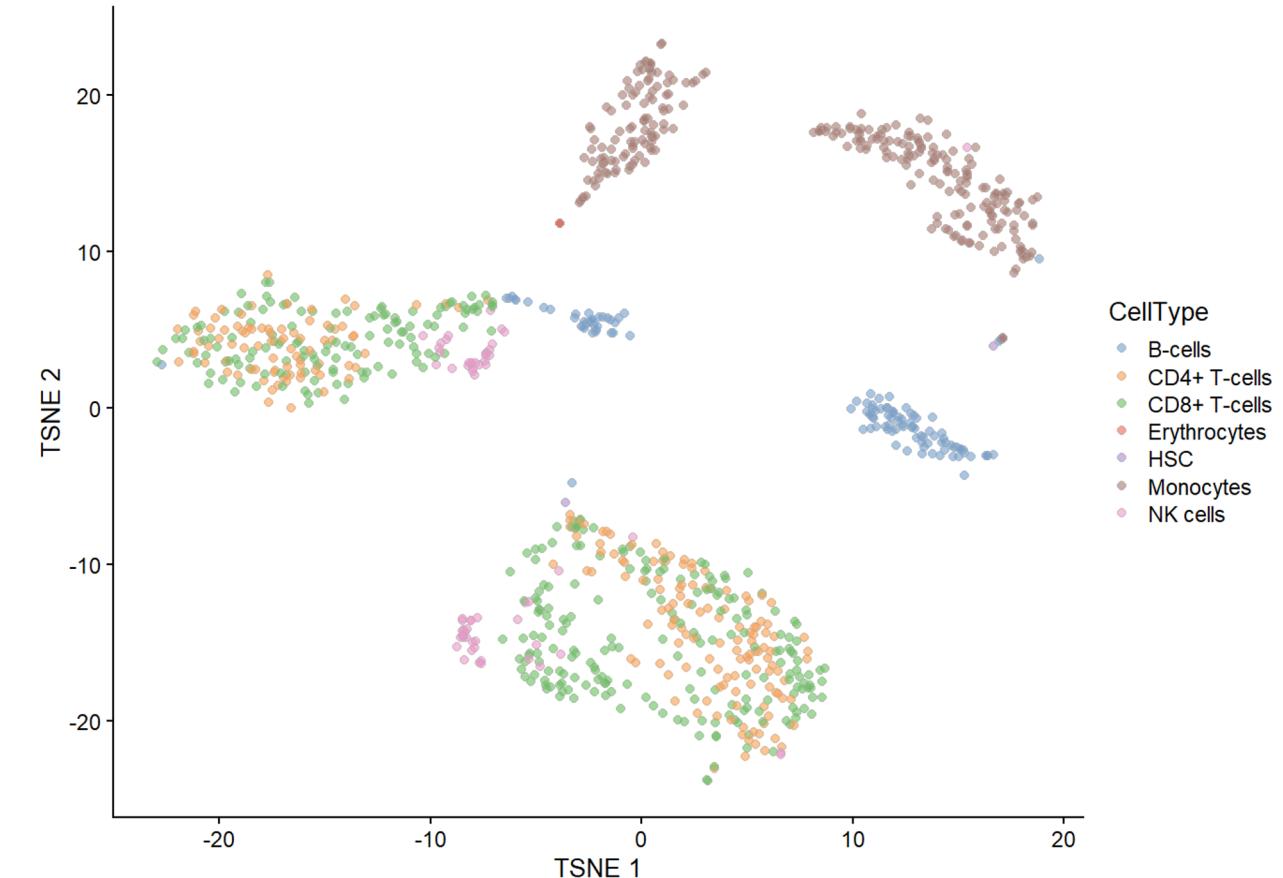
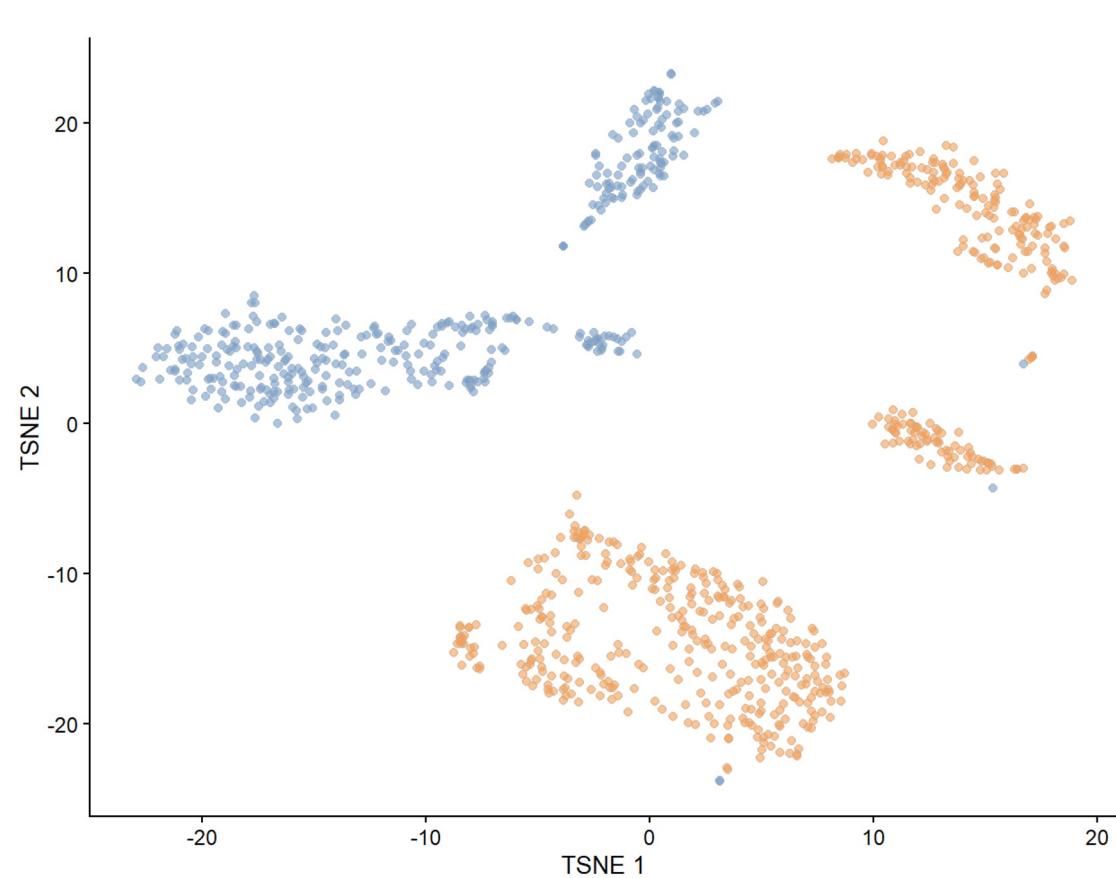
- Okay if combining batches with ‘balanced’ cell types

ComBat-SVA-Seq

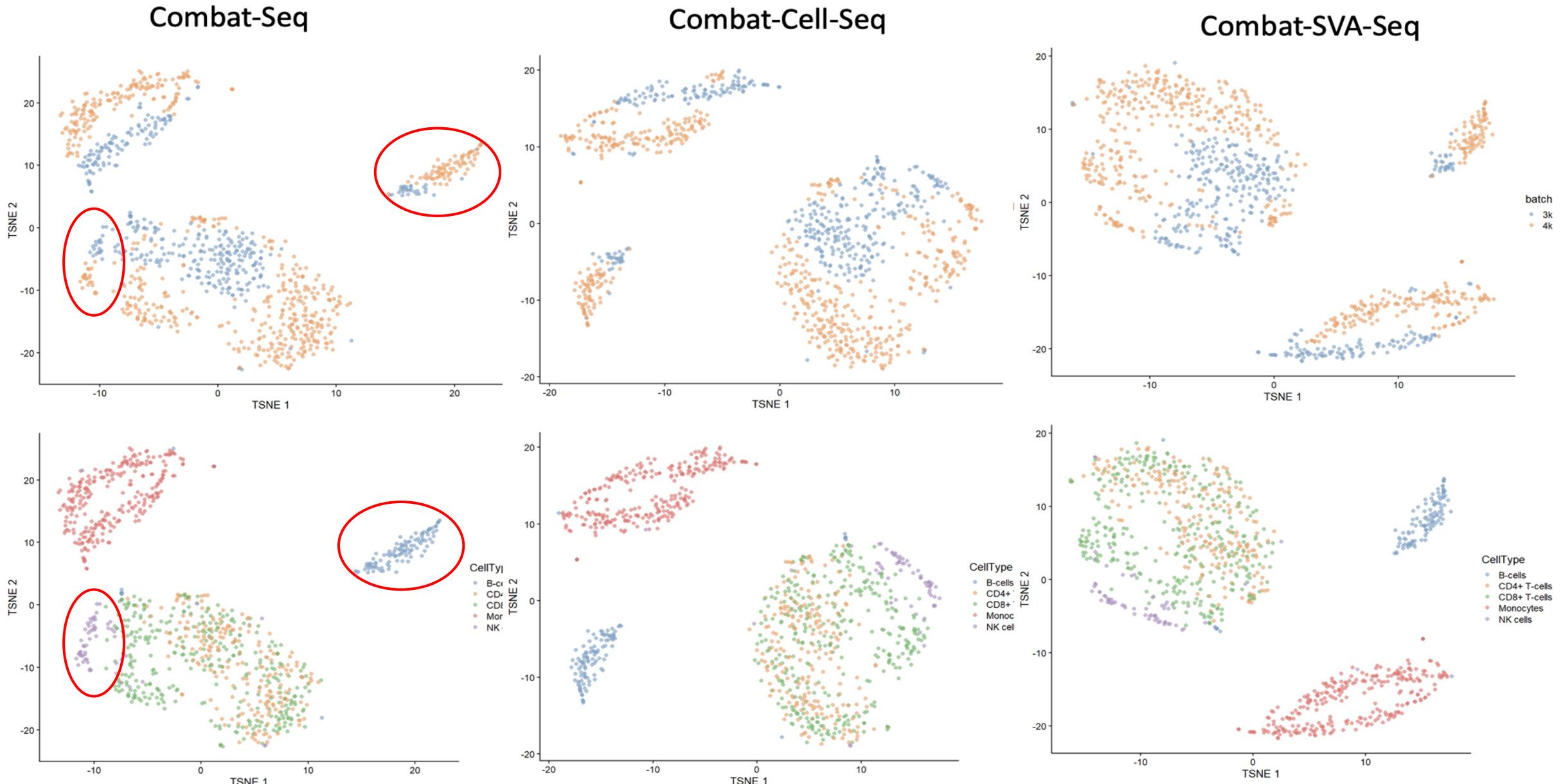
- Use SVA to identify surrogate cell-type variability



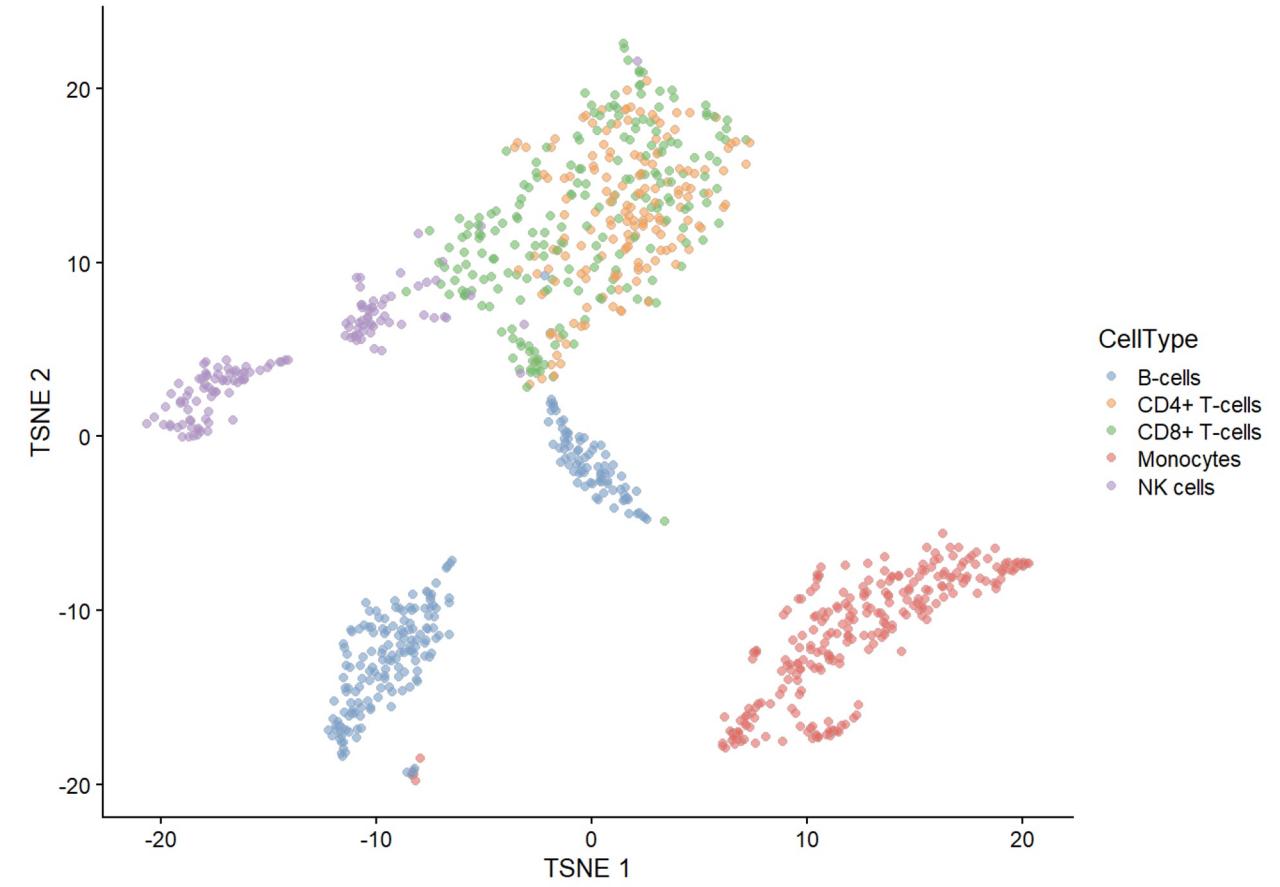
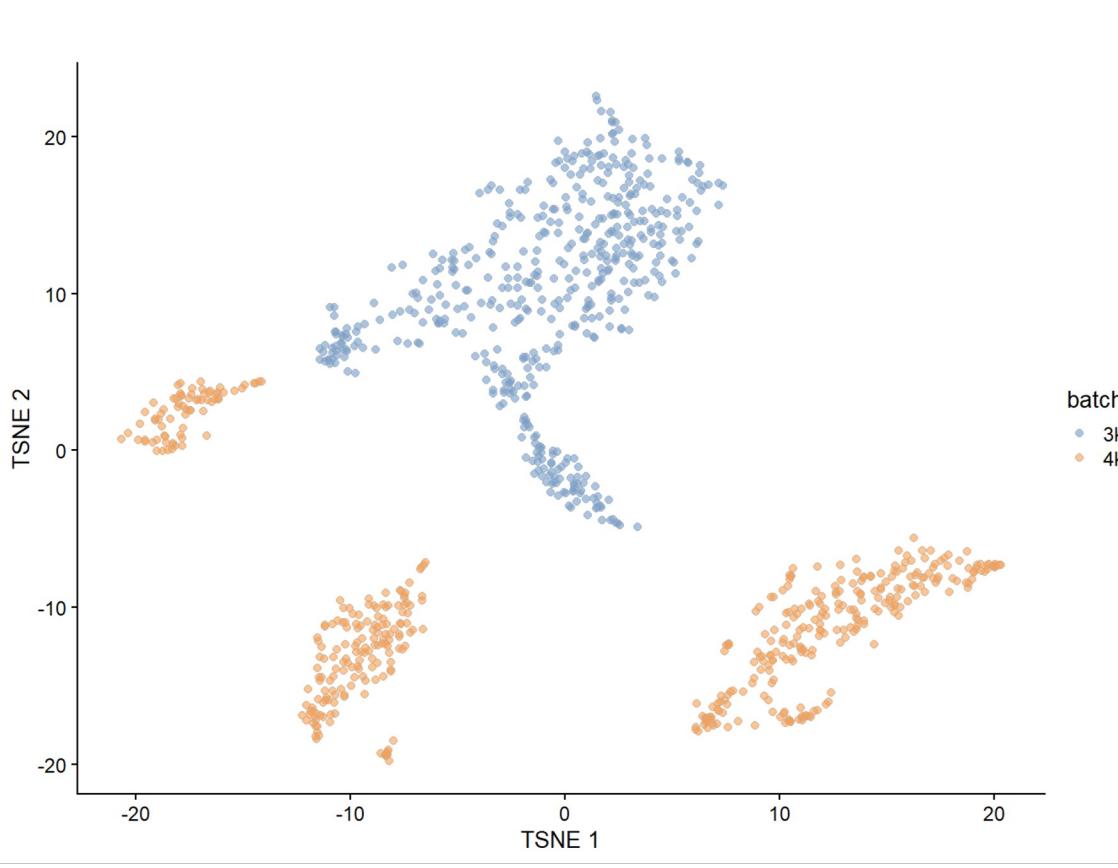
ComBat-Seq for Balanced Designs (ComBat-Cell-Seq, ComCombat-Seq)



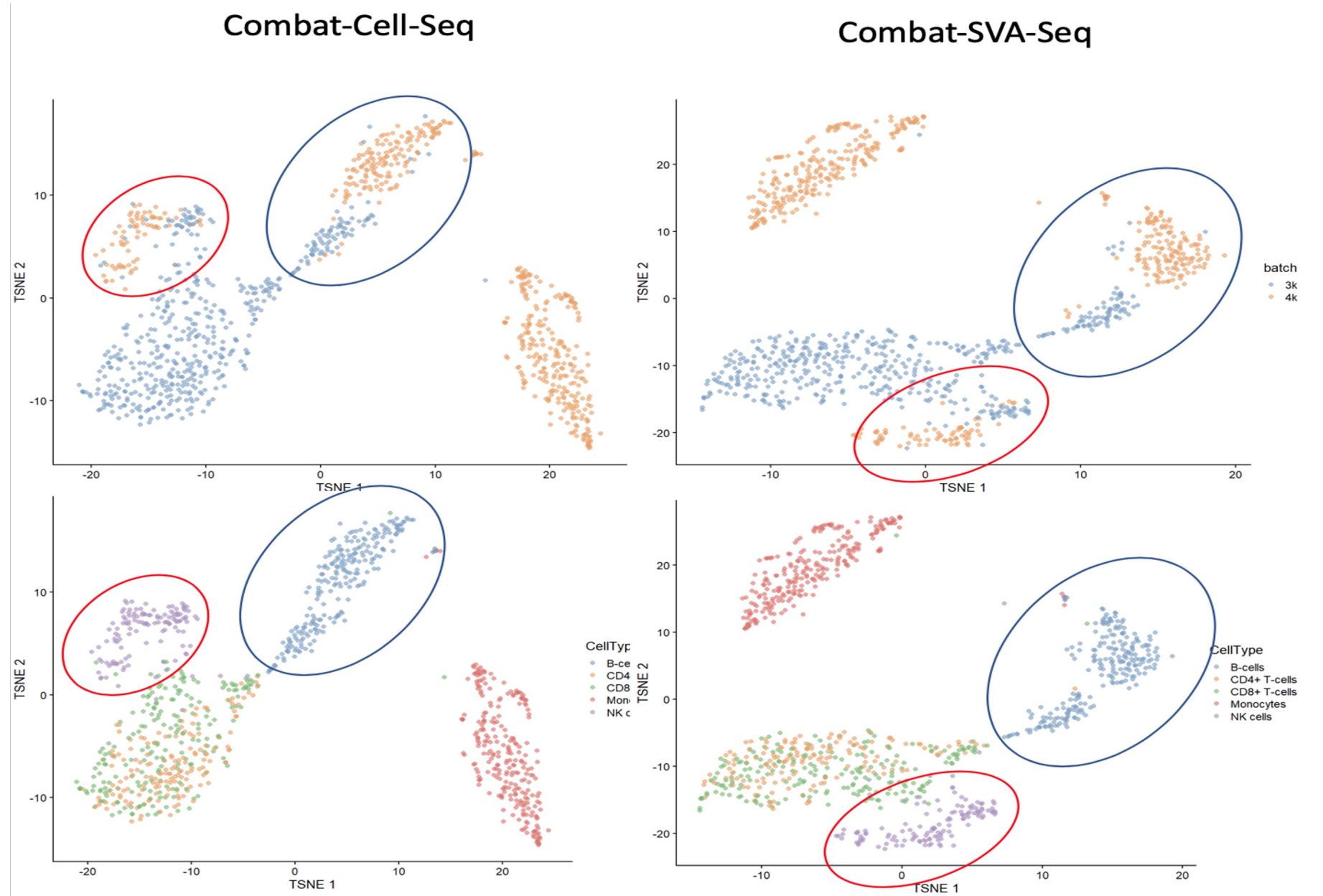
ComBat-Seq for Balanced Designs (ComBat-Cell-Seq, ComBat-Seq)



ComBat-Seq for Unbalanced Designs (ComBat-Cell-Seq, ComBat-SVA-Seq)



ComBat-Seq for Unbalanced Designs (ComBat-Cell-Seq, ComBat-SVA-Seq)



ComBat-seq Summary

For balanced designs:

- **ComBat-Seq, ComBat-Cell-Seq, ComBat-SVA-Seq** all work well!

For unbalanced designs:

- **ComBat-Seq**: May remove cell-type specific variation
- **Combat-Cell-Seq**: Performs extremely well
- **Combat-SVA-Seq**: Not quite as good as Combat-Cell-Seq, but performs well in both simulated and real-data examples