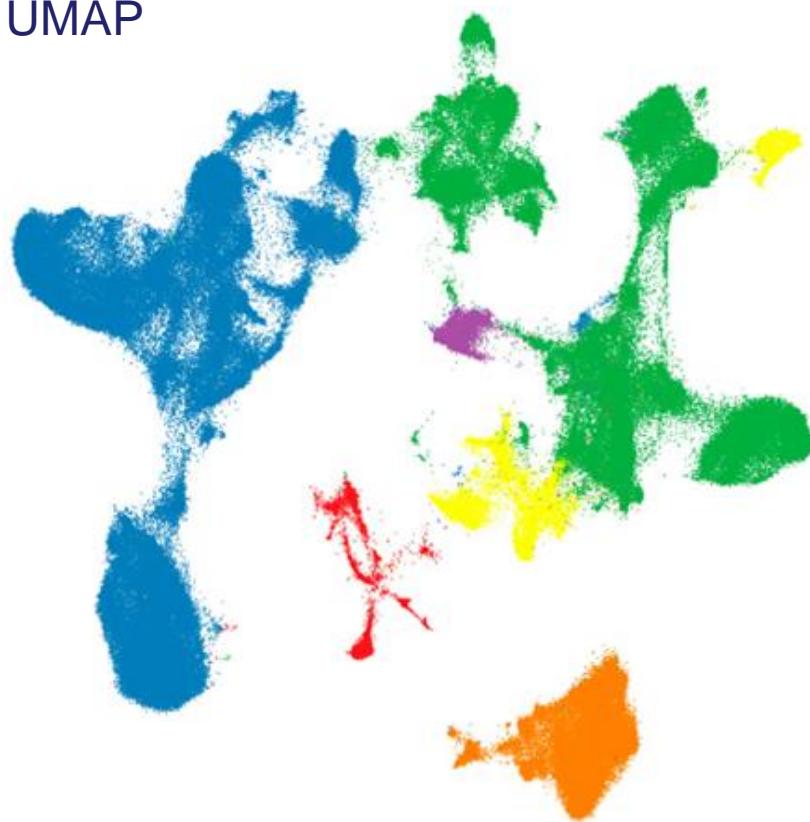
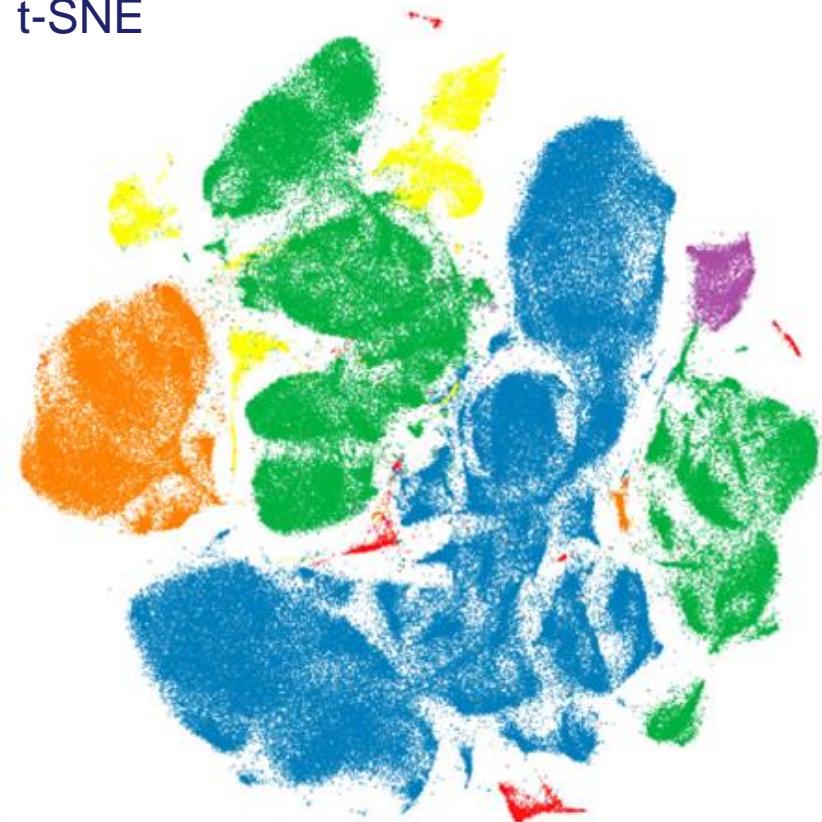


# Single Cell Biology & Data Science Methods

UMAP



t-SNE



Becht et al. 2018

*Jonathan Irish, Sierra Barone Lima, Cassidy Mayeda*

Associate Professor  
Cell & Developmental Biology  
Pathology, Microbiology & Immunology

Data Science  
Program Coordinator

Web Applications  
Research Assistant

## Code & Examples for Today:

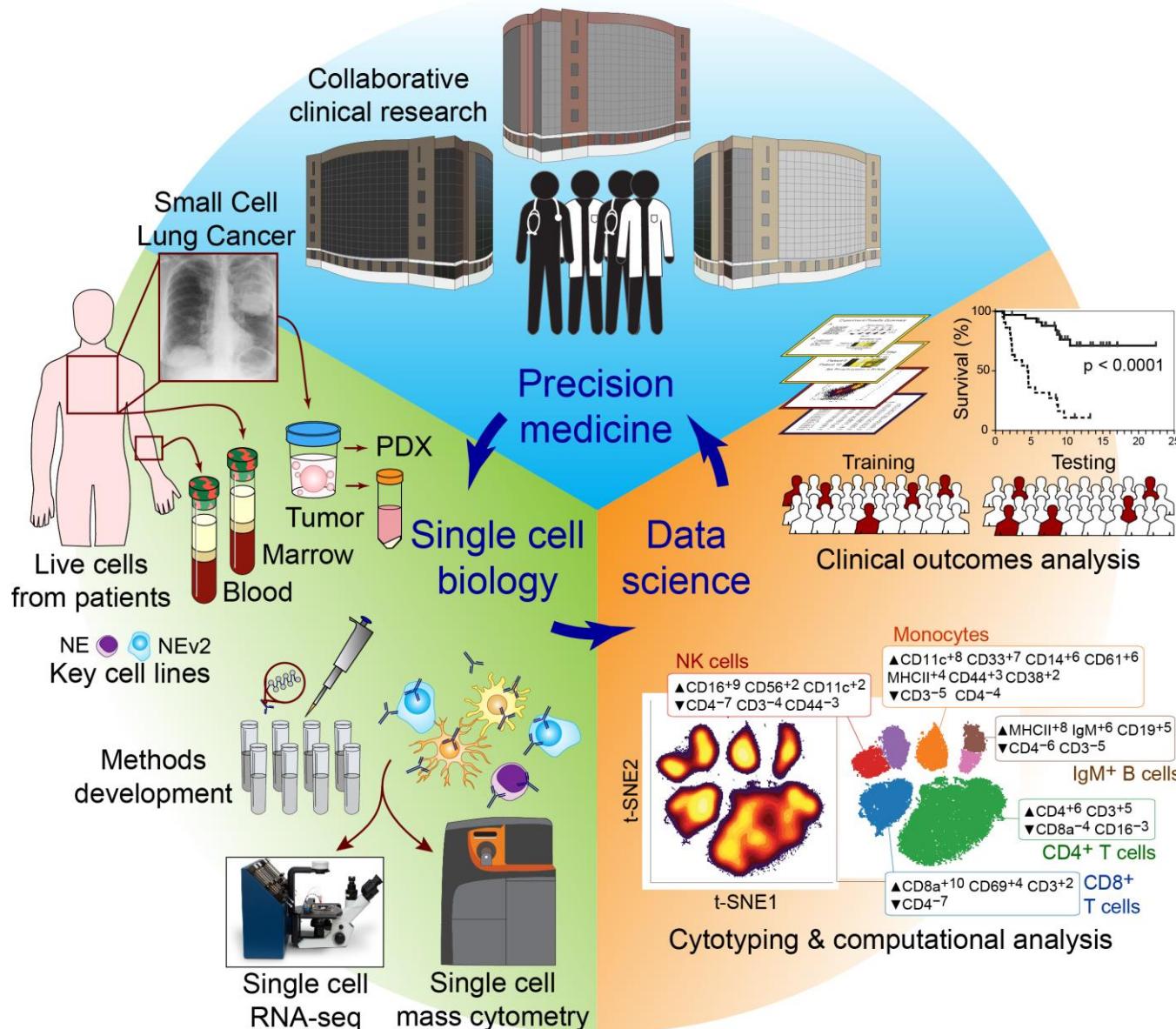
<https://cytolab.shinyapps.io/cGVHD/>

<https://github.com/cytolab/irish-data-science>

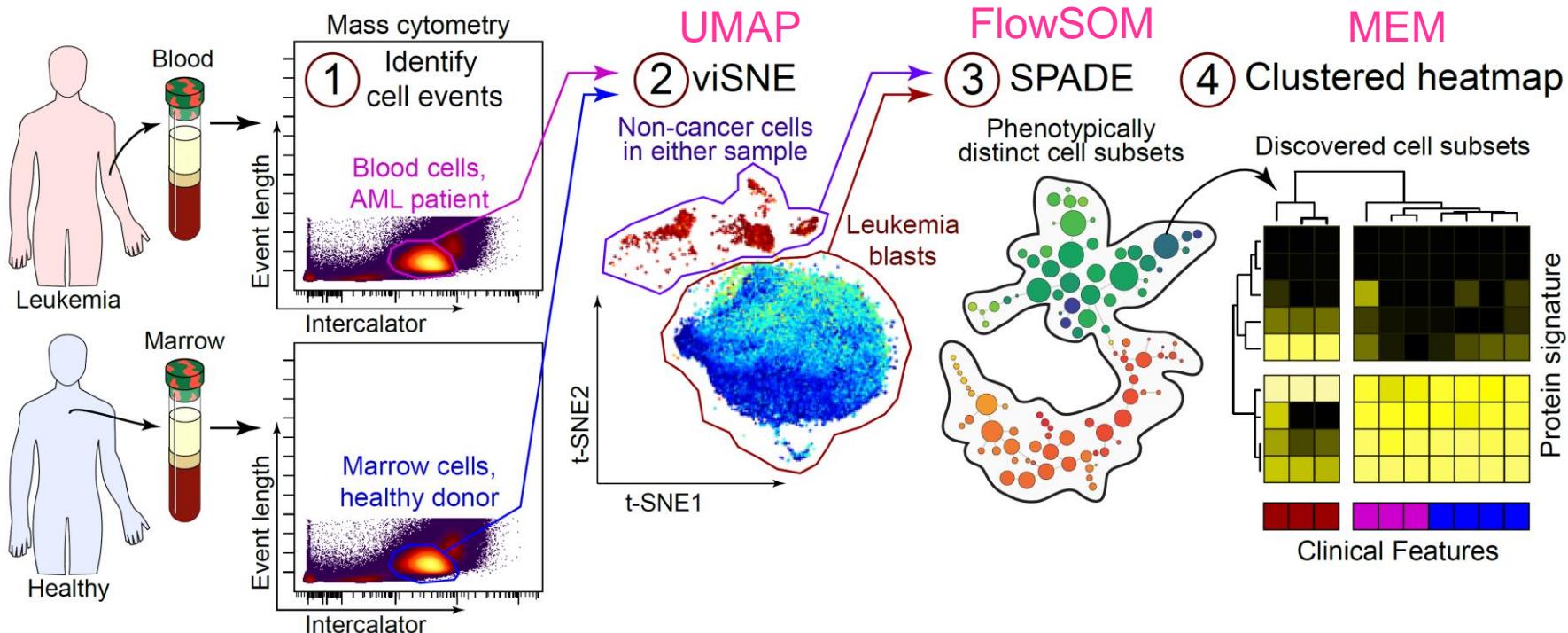
<https://github.com/cytolab/mem>

# Part 1: Introduction to Data Science

# Goal: Systematically Dissect Cellular Mechanisms Across Time, Treatments, Tissues, & Tumor Types



# Machine Learning Is a Key Skillset for Biologists (And the Tools Are Rapidly Evolving)



Protocol:  
Diggins et al., *Current Protocols in Cytometry* 2017  
Data files:  
<https://flowrepository.org/id/FR-FCM-ZZKZ>  
More great tools:  
Saeys et al., *Nature Reviews Immunology* 2016

Original research:  
Diggins et al., *Methods* 2015  
Diggins et al., *Nature Methods* 2017

Typical workflow and goal: learn & label cytotypes (cell identities),  
reveal and assess unexpected & abnormal cells

Need: human reference data (more examples) with annotations

Effective data analysis is critical in clinical research,  
& this now means working *with* computational tools  
that reveal and model patterns across data types

Tools from one area can be applied in others  
(economics, math, **patients**, **cells**, pixels, ...)

Data science workshop can be self-taught:

<https://github.com/cytolab/>

# Discussion Questions We Expect to Cover

- 1) What are key differences between tools (t-SNE/viSNE, SPADE, UMAP, FlowSOM, PCA, MEM, Citrus, etc.)? What is the difference between transforming, clustering, and modeling data? What type of modeling are we doing (if any)?
- 2) What do terms like linear or parametric analysis mean? Does the data's scale matter (arcsinh5, arcsinh15, linear)?
- 3) What do all the settings do (e.g., t-SNE iterations, perplexity, SPADE downsampling & node #)? When should they be changed?
- 4) How does one compare new samples with a prior analysis? How do we test tools with expert gating?
- 5) What are some “red flags” indicating problems? What does a good t-SNE, UMAP, FlowSOM, or other analysis run look like?

But first: what is data science?

## Irish lab view of data science:

Systematically varying analytical elements  
in order to test a hypothesis

(Varied analytical elements might be different data types, data sub-samples, different initial assumptions, contrasting analytical tools, input parameters, etc.)

It's relatively new that datasets are robust enough to enable mining & exploration.

# Rumsfeldian Data Science

Known knowns: What do you know about your system?

Known unknowns: What do you know remains to be learned?

Unknown unknowns: What don't you know you don't know?

Donald Rumsfeld (Feb 12, 2002): Reports that say that something hasn't happened are always interesting to me, because as we know, there are known knowns; there are things we know we know. We also know there are known unknowns; that is to say we know there are some things we do not know. But there are also **unknown unknowns – the ones we don't know we don't know**. And if one looks throughout the history of our country and other free countries, it is the latter category that **tend to be the difficult ones**.

# Socratic Data Science

Known knowns: What do you know about your system?

Known unknowns: What do you know remains to be learned?

Unknown unknowns: What don't you know you don't know?

Unknown knowns: What don't you know, but think you do?  
i.e. Which 'priors' are incorrect?

If you fear incorrect priors, unsupervised analysis may be able to help.

Socrates according to Plato's *Apology*: I am wiser than this man, for neither of us appears to know anything great and good; but he fancies he knows something, although he knows nothing; whereas I, as I do not know anything, do not fancy I do. In this trifling particular, then, I appear to be wiser than he, because I do not fancy I know what I do not know.

# Citrus & RAPID Connect Cell Clusters to Clinical Outcomes, RAPID is Designed for Unsupervised Analysis of Survival

---

## Citrus

Bruggner, Tibshirani, et al., PNAS 2014

Finding  
cell clusters

Unsupervised  
(hierarchical clustering,  
cells may be in 2+ clusters)

Determining number  
of cell clusters to seek

Unsupervised  
(must be >5% of sample)

Modeling  
cluster features

Supervised, multivariate  
(lasso regularized  
logistic regression,  
nearest shrunken centroid)

Splitting patients  
into groups

Supervised, happens at start  
(expert knows cut points,  
assigns patients to groups)

## RAPID

bioRxiv 2019

Unsupervised  
(various: FlowSOM, KNN,  
t-SNE + FlowSOM)

Unsupervised  
(seeks few clusters  
w/ low internal variation)

Unsupervised, univariate  
(median or MEM, simply a  
statistical description of cluster)

Unsupervised, happens at end  
(cluster abundance as cut point,  
Cox model of hazard)

# Defining Your System

## 1) Elements, the studied units of the system.

- ▶ Patients, cells, images, pixels, transcripts, genomes, peptides.
- ▶ We will envision elements as “rows” in a spreadsheet.

## 2) Features, the things measured for each element.

- ▶ Clinical outcomes, phospho-proteins, pixel density, nucleotides.
- ▶ We will envision features as “columns” in a spreadsheet.
- ▶ Feature selection may rely on hypotheses, rules, or prior knowledge.

## 3) Scales, the type & range of the measurements for each feature.

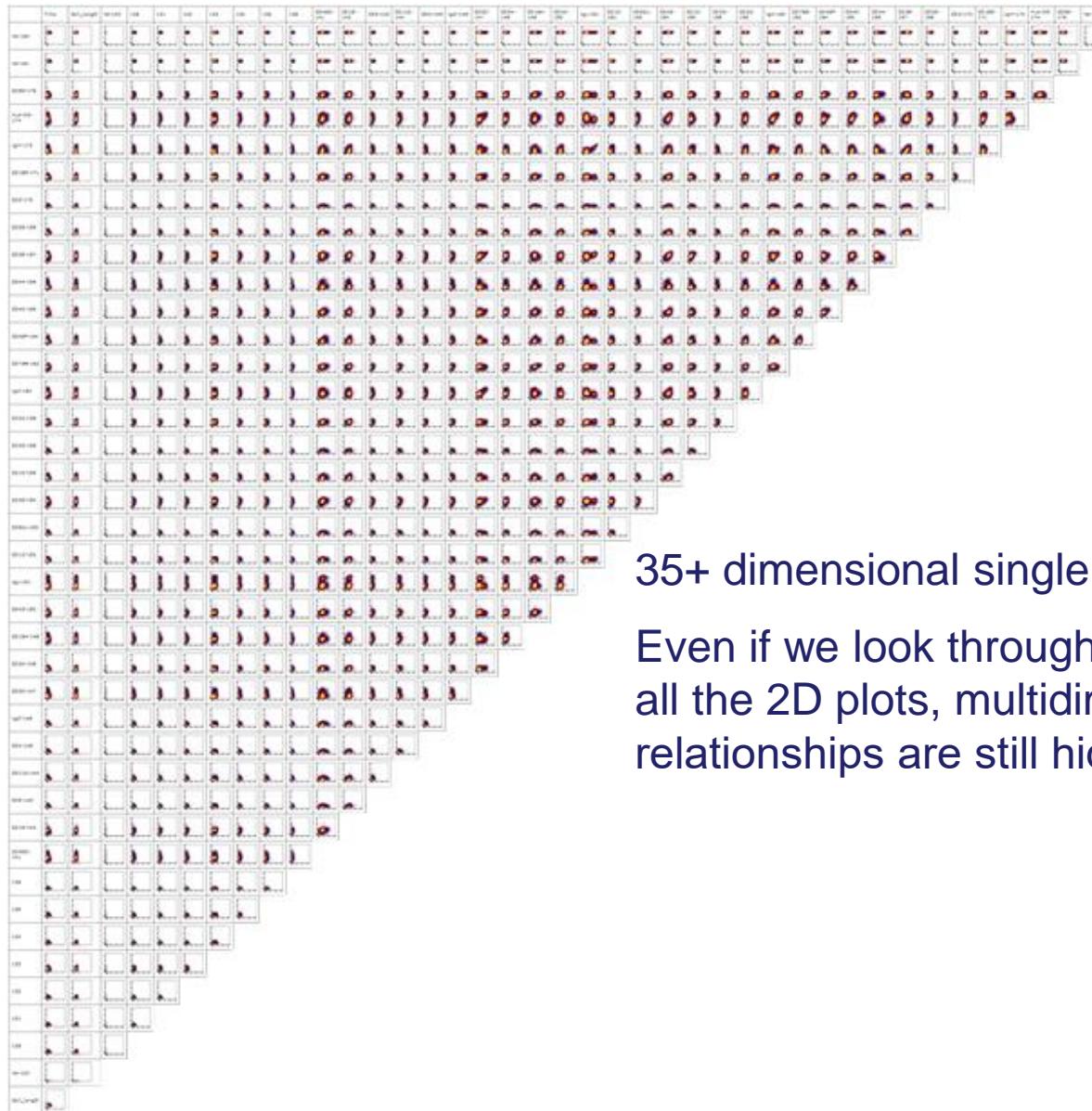
- ▶ Categorical, linear, log & base, arcsinh & cofactor.
- ▶ -150 to 262,144; 1 to 10,000; 0 to 50; 1 to 100; 0 to 1; NR, PR, CR.
- ▶ Will largely explore the data without units until we create reports.

## 4) Prior knowledge, the things assumed to be known for the system.

- ▶ Organization of elements (groups, order, etc.), feature relationships.
- ▶ Supervised analysis explicitly uses prior knowledge.
- ▶ Unsupervised analysis looks for patterns without prior knowledge.

## Part 2: Quantifying Cell Biology & Cytometry Tools

We Now Make Billions of Multi-D Single Cell Measurements  
=> Need for Machine Learning Tools & Human Readable Views



35+ dimensional single cell data:

Even if we look through  
all the 2D plots, multidimensional  
relationships are still hidden...

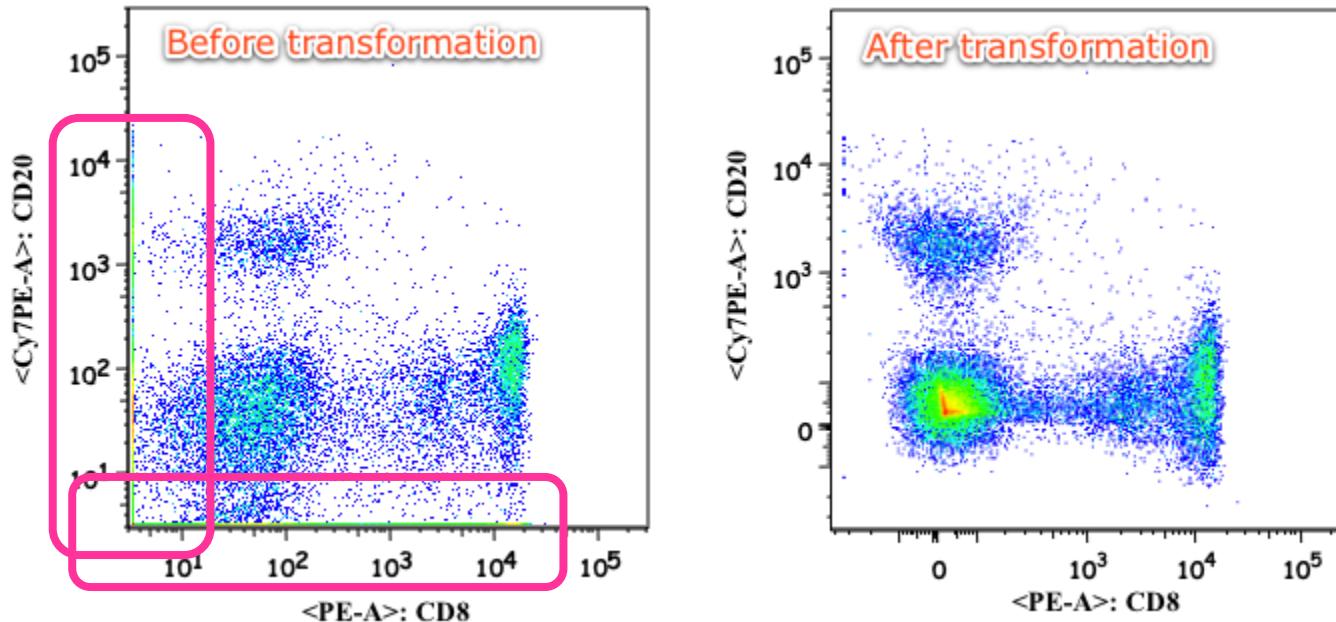
# Part 3: The Importance of Scales

Before we get to ‘desert’, some math ‘veggies’:

Scales matter: poorly or variably scaled data can destroy an analysis, most issues arise near zero

(pre-processing & normalization can also be critical)

Have you ever noticed two peaks within the cells that are biologically 100% negative for a marker?



<http://www.flowjo.com/v76/en/displaytransformwhy.html>

Results from bad scaling (poor transformation)  
and it can be an issue for computational analysis.

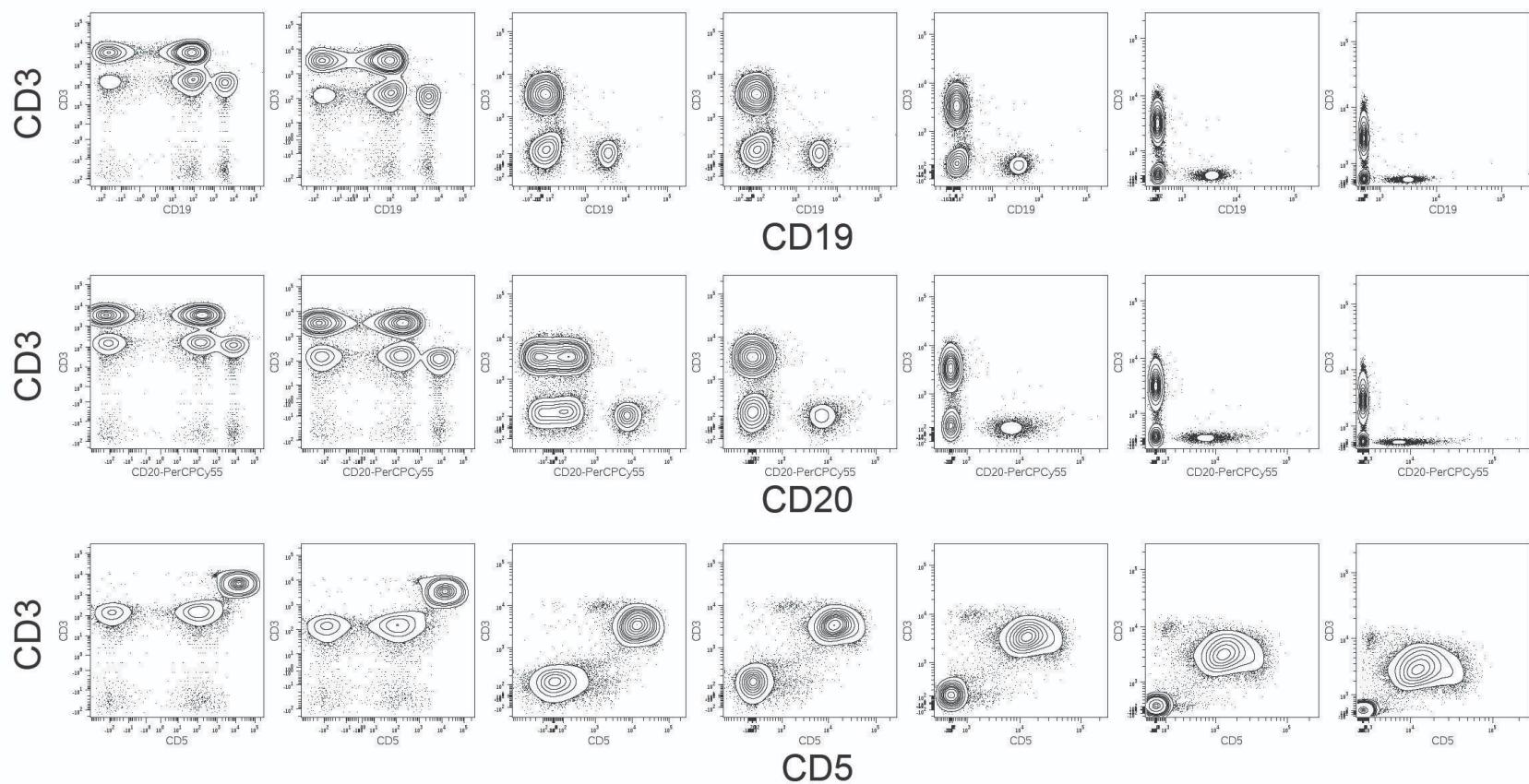
Scaling is important in both mass and fluorescence cytometry.

# Scaling Matters for Measuring Distance (Fluorescence Flow)

Healthy human PBMC, intact cells gate

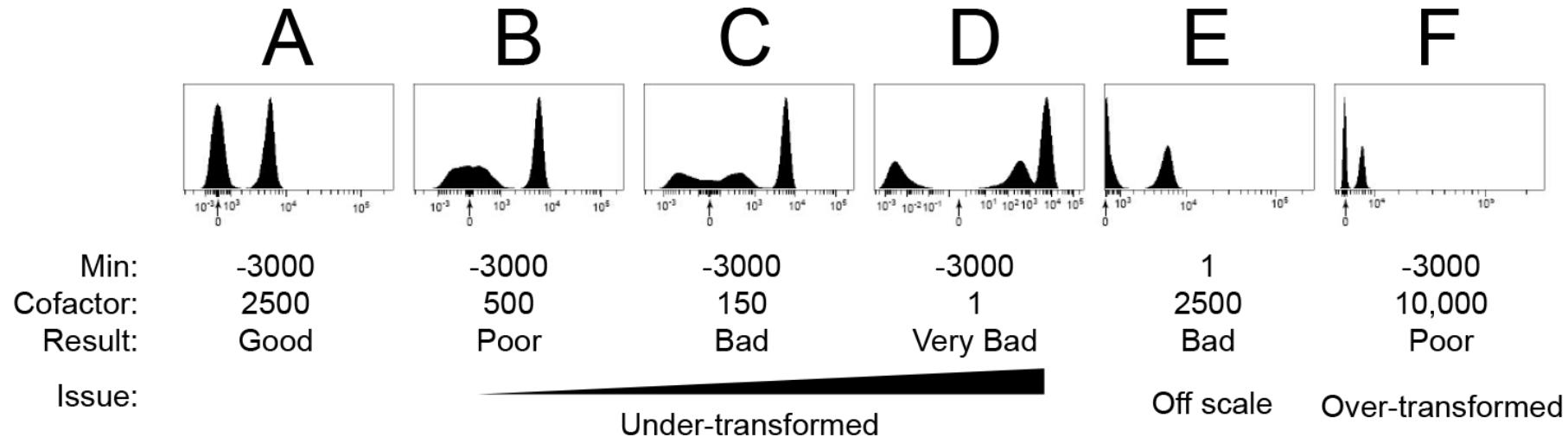
arcsinh(intensity / argument)

CD3 arg	1	5	150	150	300	750	1500
CD19 arg	1	5	150	150	300	750	1500
CD20 arg	1	5	150	500	1500	3000	6000
CD5 arg	1	5	150	500	1500	3000	6000



# Scaling Matters for Measuring Distance (Compensation Beads)

A 50:50 mix of + and - events stained only for PerCP-Cy5.5 is shown using different scales.



$$\text{arcsinh}(x) \text{ with cofactor } c = \ln\left(\frac{x}{c} + \sqrt{1 + \left(\frac{x}{c}\right)^2}\right)$$

For fluorescent flow cytometry data a biexponential or arcsinh transformation corrects the scale near zero.

Since computational analysis techniques compare distance similar to what a person does when looking at a plot, these techniques can identify artificial populations near zero (see C and D) if data are not appropriately transformed prior to analysis.

# Examples of Four Common Data Scales

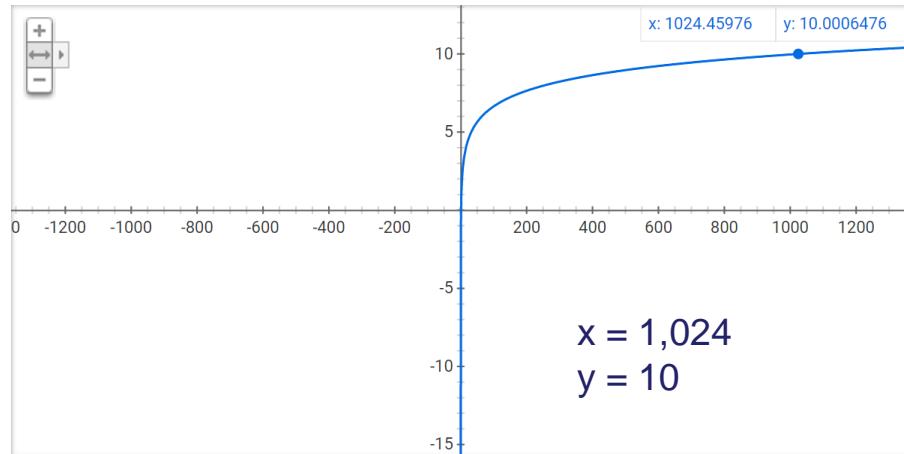
$$\log_2(x)$$

$$\log_{10}(x)$$

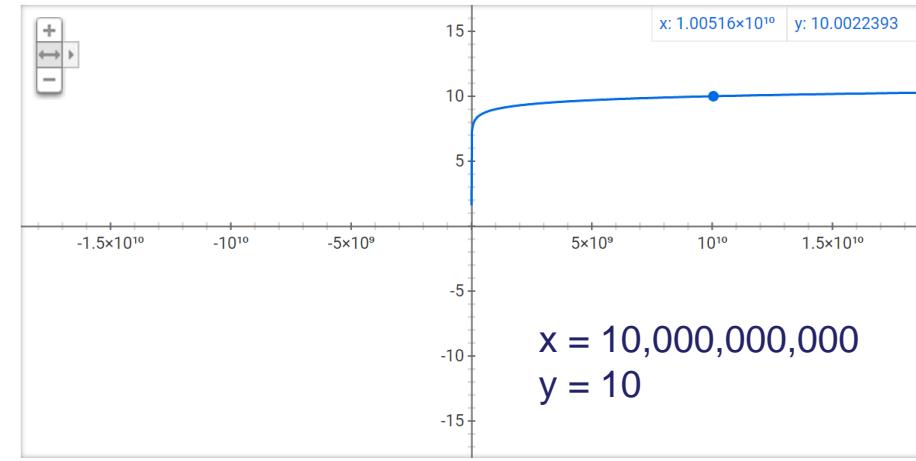
$$\operatorname{arcsinh}(x/5)$$

$$\operatorname{arcsinh}(x/150)$$

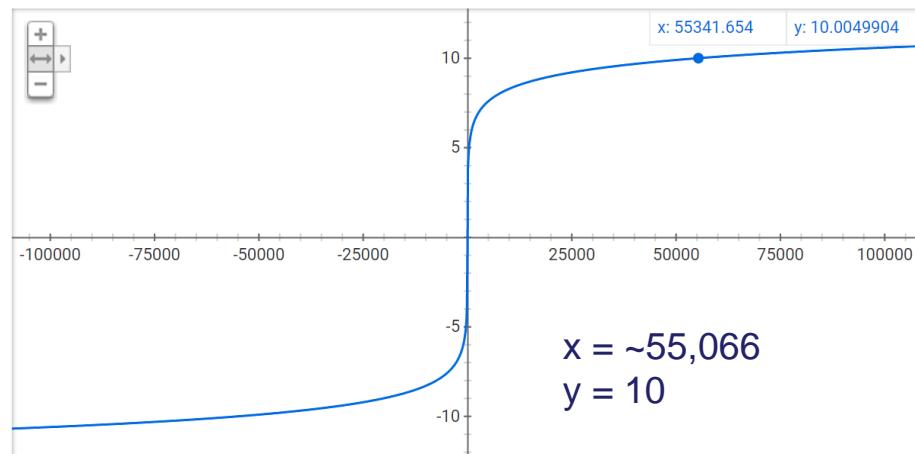
Graph for  $\ln(x)/\ln(2)$



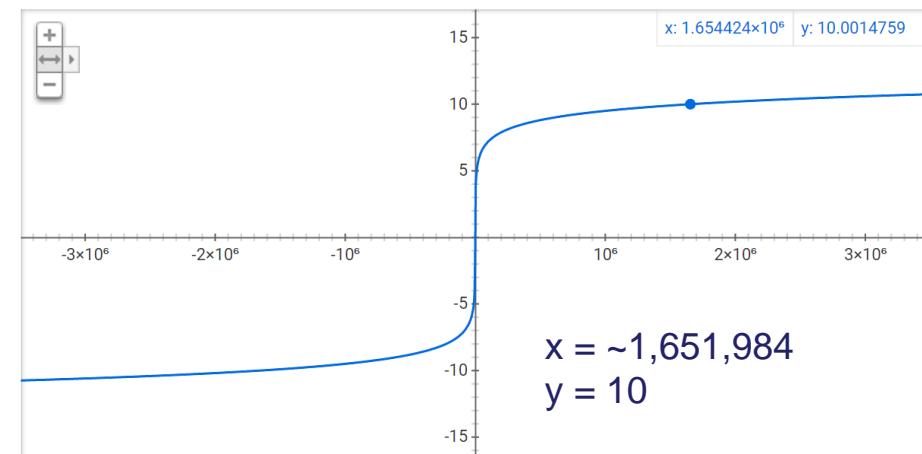
Graph for  $\ln(x)/\ln(10)$



Graph for  $\ln(x/5+\sqrt{1+(x/5)^2}) = \operatorname{arcsinh}(x/5)$

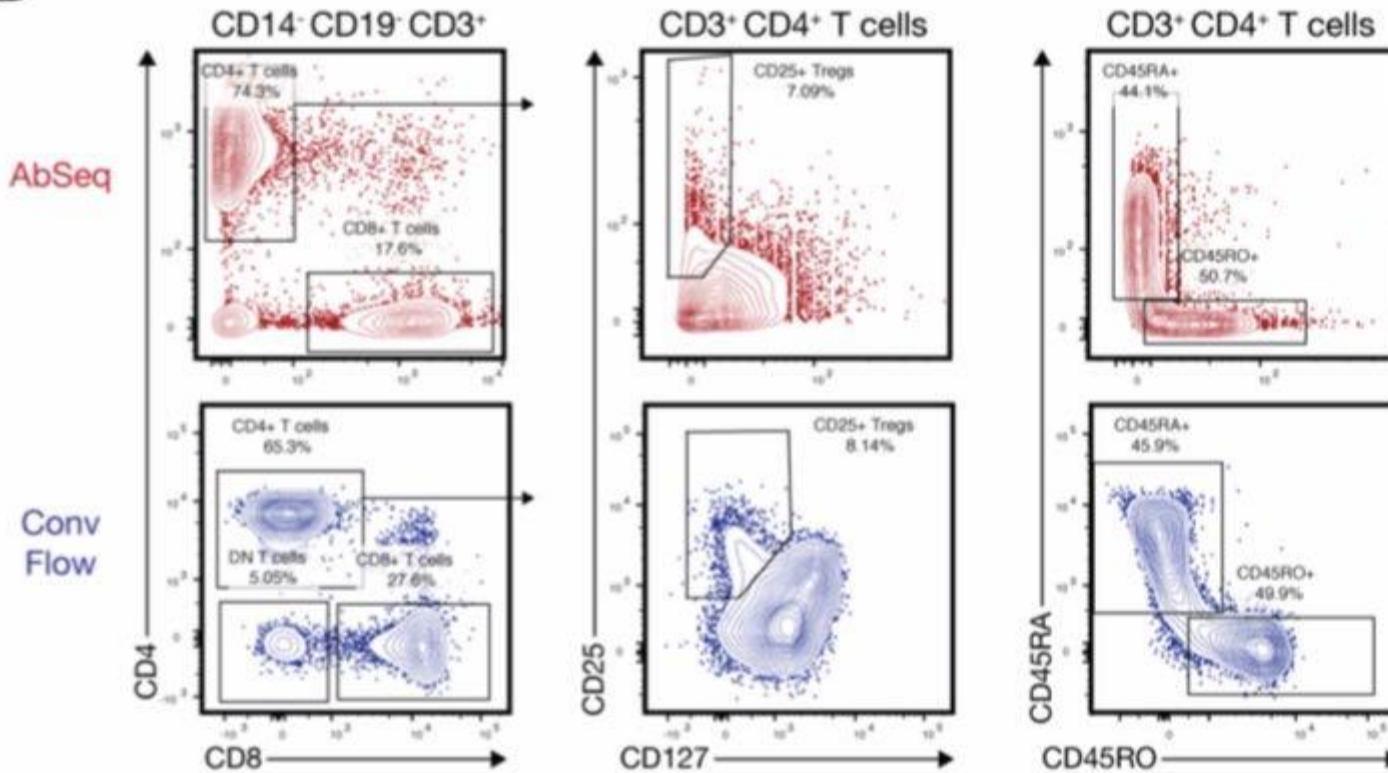


Graph for  $\ln(x/150+\sqrt{1+(x/150)^2}) = \operatorname{arcsinh}(x/150)$



# Scaling Matters for Measuring Distance (RNA seq vs. Fluorescence Flow)

D

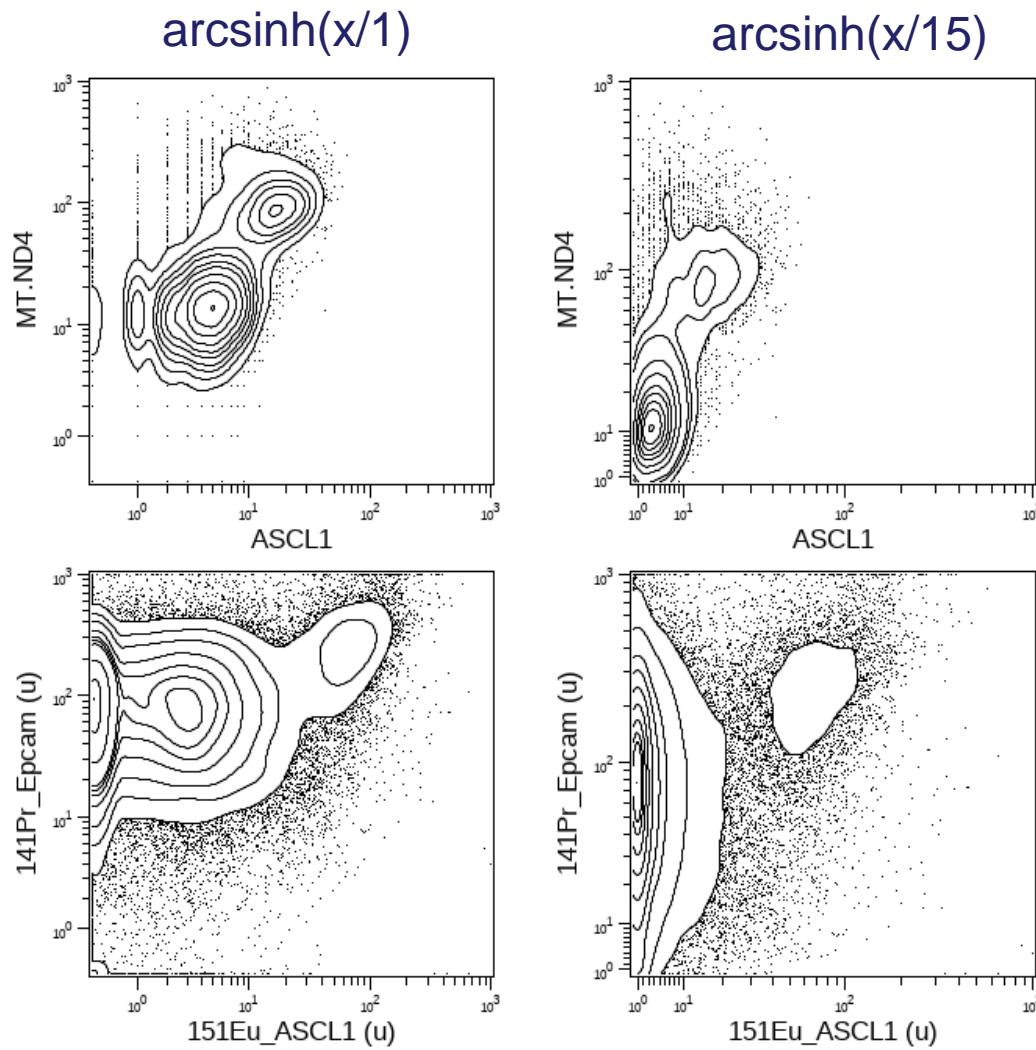


# Example Challenge, Scaling to Compare scRNA-seq vs. CyTOF (Here with one example cell line, DMS454, and ASCL1 on the x-axis)

DMS454  
SCLC cell line

scRNA-seq

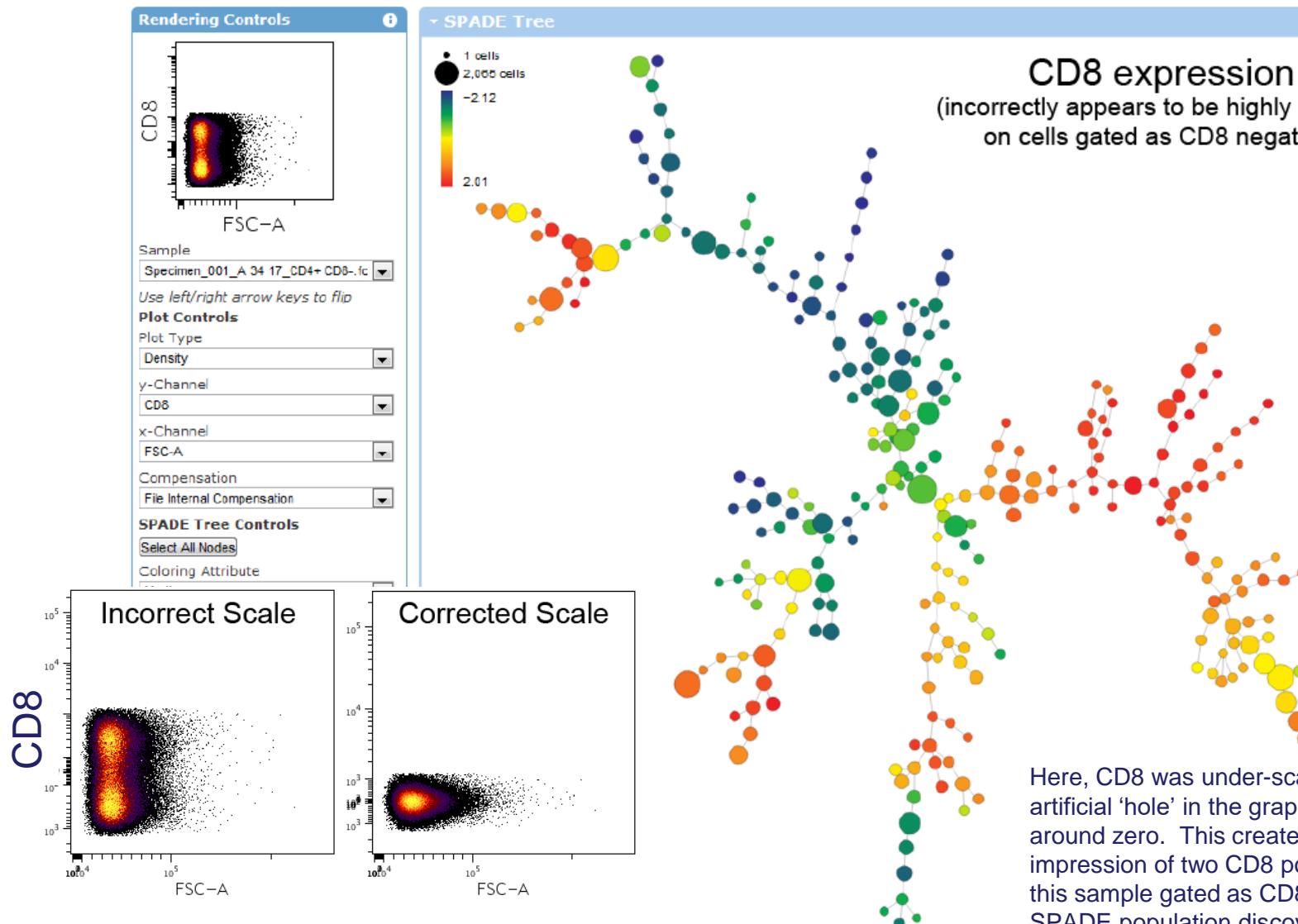
CyTOF



Qualitative: Is there a subset of cells that is ASCL1<sup>+</sup>?

Quantitative: How much ASCL1 is in each cell? How much does ASCL1 change after Tx?

# Inappropriate Scaling Can Lead to False Population Discovery

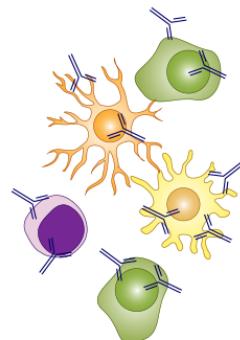
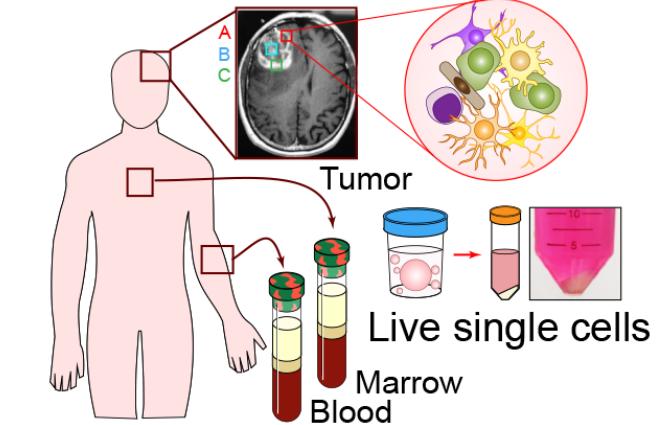


Here, CD8 was under-scaled so that an artificial 'hole' in the graph existed around zero. This created the false impression of two CD8 populations in this sample gated as CD8 negative. SPADE population discovery treated this as significant.

# Part 4: Dimensionality Reduction Tools: t-SNE & UMAP Change Everything

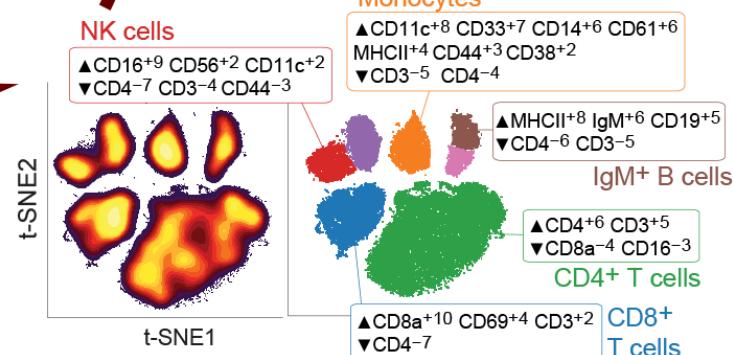
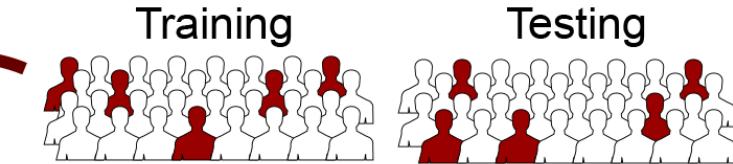
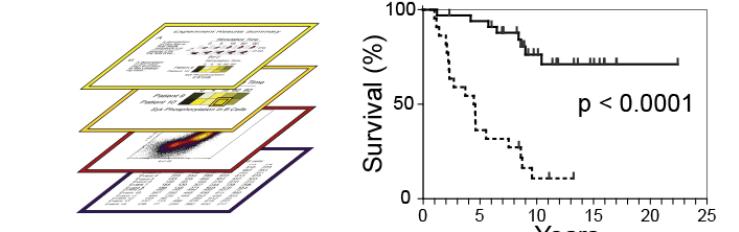
# Elements of Translational Data Science

Sample patients over time at key clinical decision points



Measure cell identity, signaling, biomarkers, and functional responses

Integrate systems biology data & clinical outcomes to guide treatment



Apply machine learning tools to reveal patterns & identify groups

# Key Topic Areas, Terms, and Cytometry Workflows

- 1: Field Changes: Data Science & Latest Tools
- 2: History: Non-linear, PCA, Trajectories, Supervised
- 3: Dimensionality Reduction: t-SNE & UMAP
- 4: Clustering: SPADE (k-means), KNN, FlowSOM, Citrus
- 5: Enriched Features: MEM, ΔMEM, RMSD
- 6: Cytometry:
  - 2004: Expert => Expert => Heatmap (Irish/Nolan)
  - 2011: Expert => SPADE => Heatmap (Bendall/Qiu)
  - 2013: t-SNE => Expert (viSNE/Pe'er, Van Der Maaten)
  - 2014: t-SNE => DensVM => Heatmap (Newell)
  - 2015: t-SNE => SPADE => Heatmap (Diggins/Irish)
  - 2015: “KNN” (=> t-SNE) => Heatmap (Phenograph)
  - 2015: FlowSOM (Van Gassen/Saeys)
  - 2017: t-SNE => SPADE => MEM (Diggins/Irish)
  - 2018: UMAP => Expert (Newell, McInnes)
  - 2019: UMAP => FlowSOM => MEM (Barone/Irish)

# Unsupervised Analysis: Not Using Prior Knowledge To Guide the Analysis

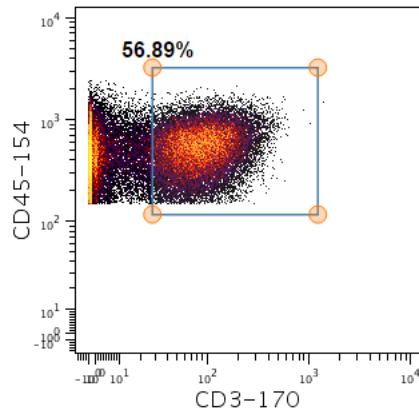
Prior knowledge examples: Stem cells express CD34, these samples were from patients that responded to drug

## Supervised Approaches

- Expert gating
- Citrus
- CellCNN (neural network)
- Wanderlust

## Unsupervised Approaches

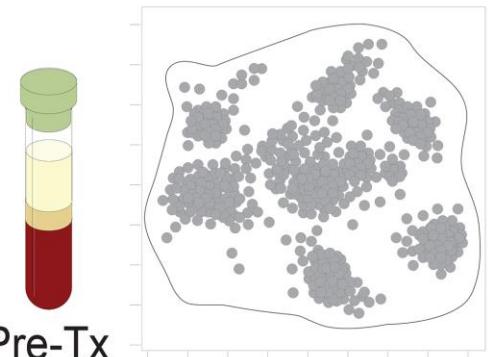
- Most heatmap clustering
- SPADE, FlowSOM
- t-SNE / viSNE, UMAP
- Phenograph



# Clinical Trial Monitoring: What Do We Need to Know? Automate Four Key Readouts vs. Clinical Outcomes

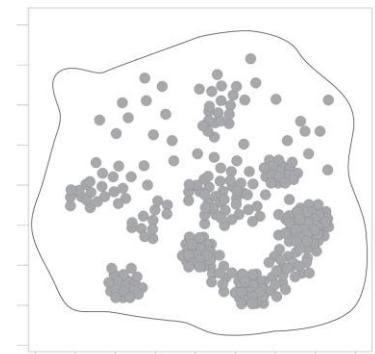
## Features of Dynamic Populations

### 1 Systems Plasticity



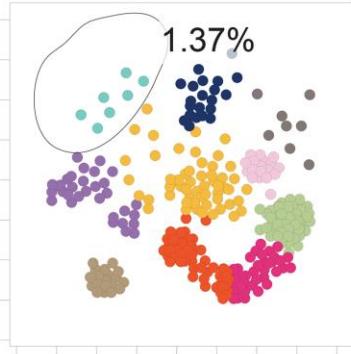
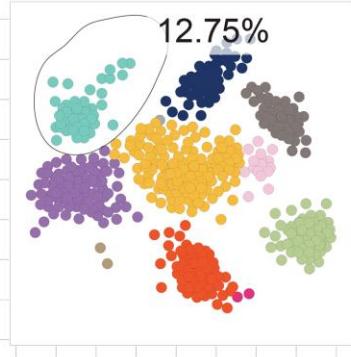
Pre-Tx

Time 1



Earth Mover's Distance  
on t-SNE or UMAP

### 2 Population abundance



Traditional gating  
or cluster frequency

### 3 Signature features

#### Pre-therapy

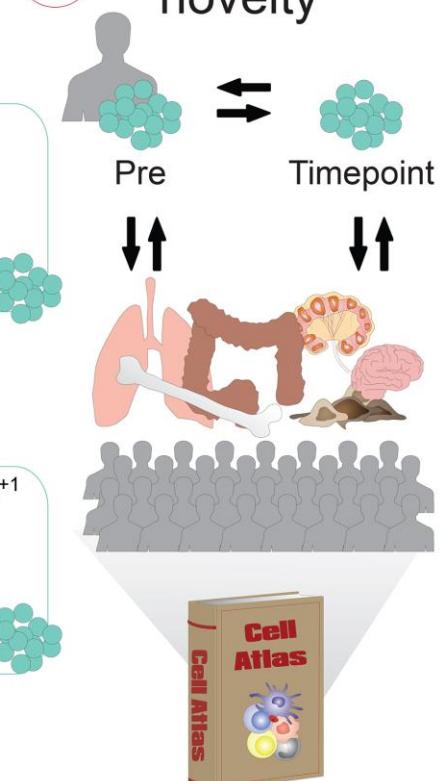
- ▲ HLA<sup>DR</sup><sup>+2</sup> CCR5<sup>+1</sup> CD38<sup>+1</sup>  
CD33<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD45RO<sup>-6</sup> CD3<sup>-6</sup>  
CD4<sup>-4</sup> CD45<sup>-2</sup> CCR4<sup>-1</sup>  
CCR7<sup>-1</sup> CD28<sup>-1</sup> CD27<sup>-1</sup>

#### Time point 1

- ▲ HLA<sup>DR</sup><sup>+2</sup> CD38<sup>+1</sup> CD45RA<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD4<sup>-6</sup> CD3<sup>-6</sup>  
CD45RO<sup>-5</sup> CCR5<sup>-2</sup> CD45<sup>-2</sup>  
CD28<sup>-2</sup> CD20<sup>-1</sup> CCR4<sup>-1</sup>  
CD27<sup>-1</sup>

Marker Enrichment  
Modeling (MEM)

### 4 Population novelty

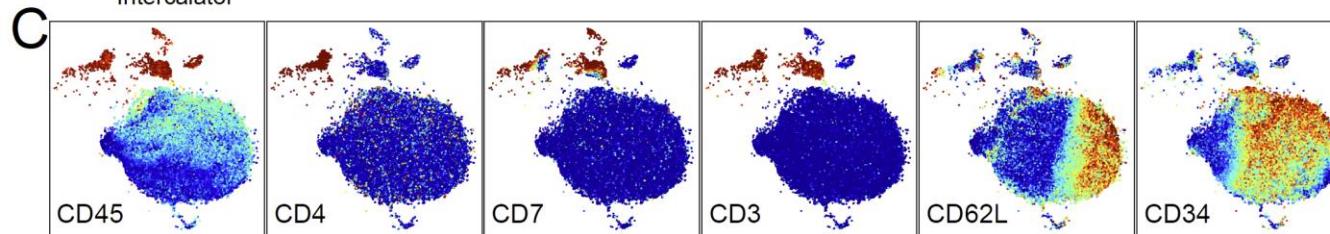
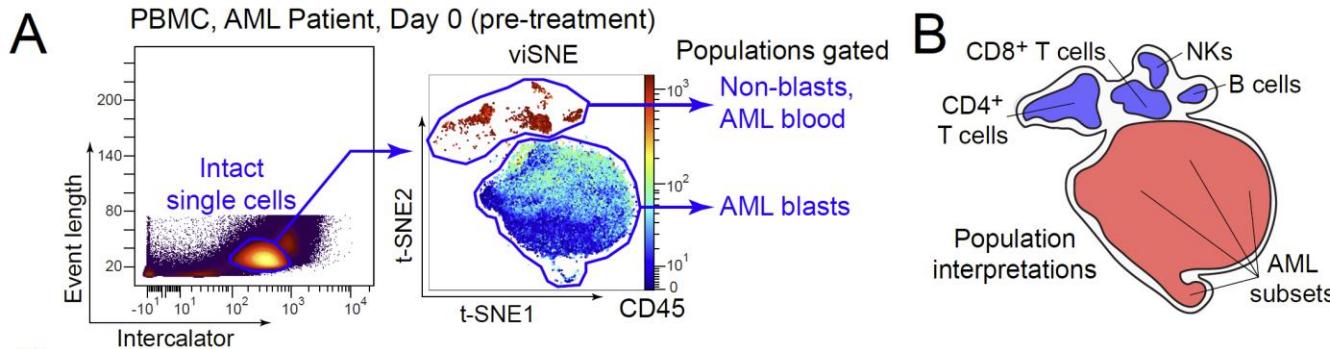


ΔMEM vs. Timepoint  
or Cell Atlas

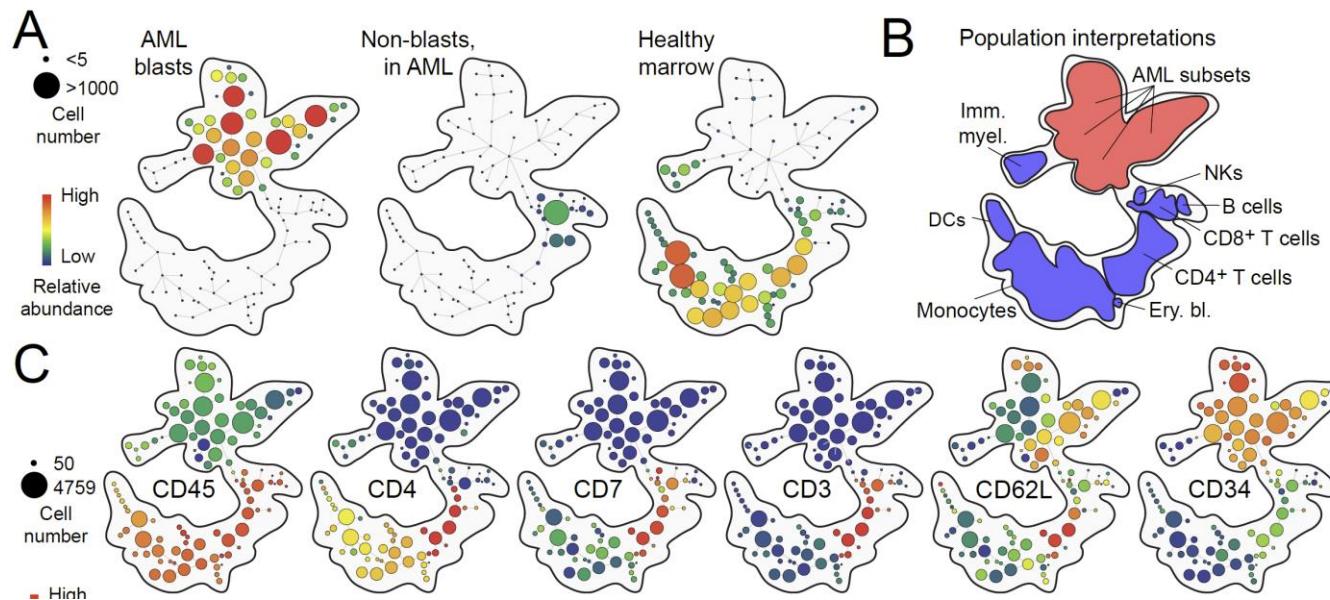
How we quantified

Greenplate et al., *Cancer Immunology Research* 2019

# Key Analysis Concepts: Dimensionality Reduction, Transformation, Clustering, Modeling, Visualization, & Integration



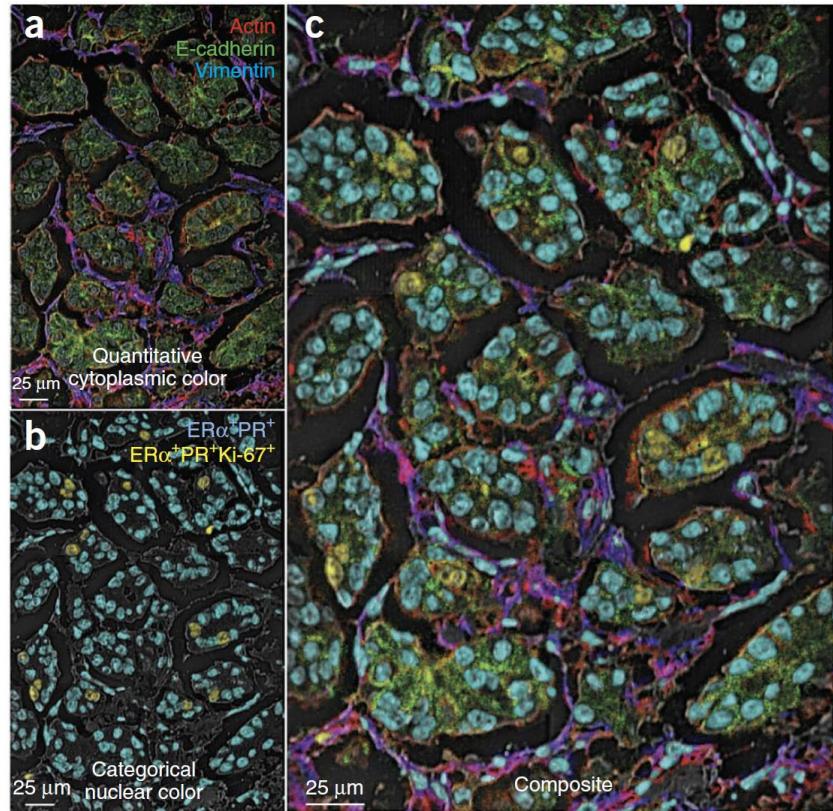
viSNE  
Amir et al.  
*Nature biotech* 2013



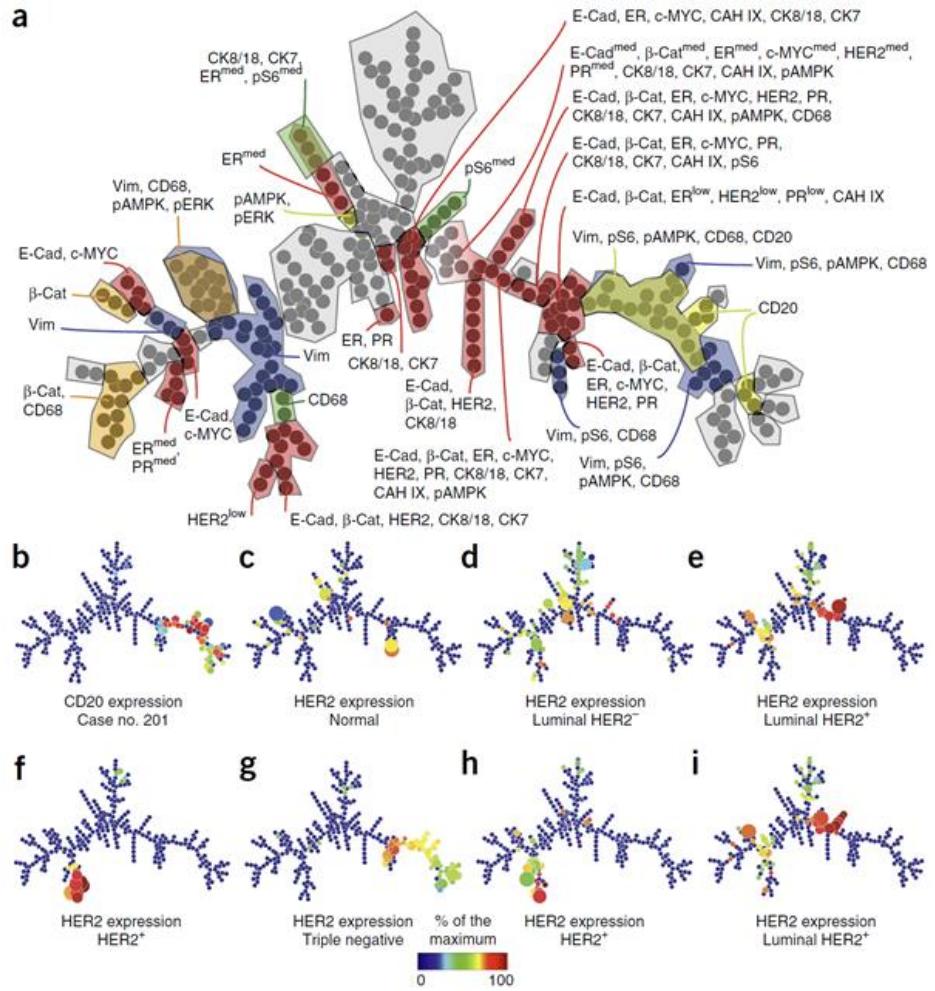
SPADE  
Qiu et al.  
*Nature biotech* 2011

Diggins et al., *Methods* 2015

# Cytomics (Cell Identity): Powered by Multiple Platforms: Imaging Mass Cytometry of Breast Cancer



Example MIBI breast cancer histology  
Angelo et al., *Nature Medicine* 2014



Analysis of IMC from 20+ breast cancer using SPADE  
Giesen et al., *Nature Methods* 2014

t-SNE was a game changer for single cell

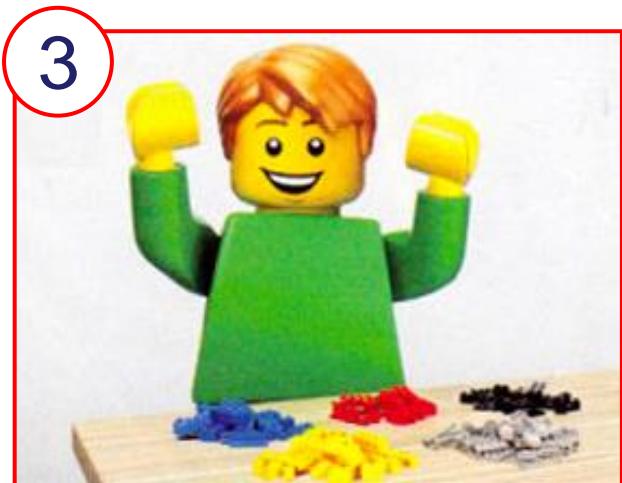
# Teaching Computers To Spot Useful Patterns : Grouping Cells by Selected Features (e.g. Protein Expression)



HD cytometry!!



Woah, that's a lot of data...

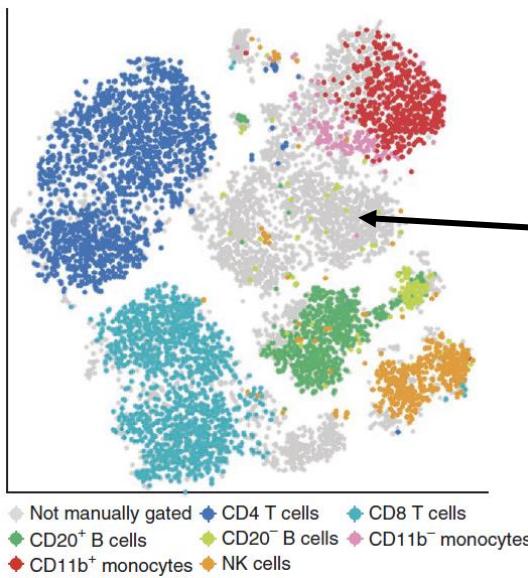


Computational tools



Biological knowledge

# Traditional Gating Overlooks Many Cells in Primary Samples

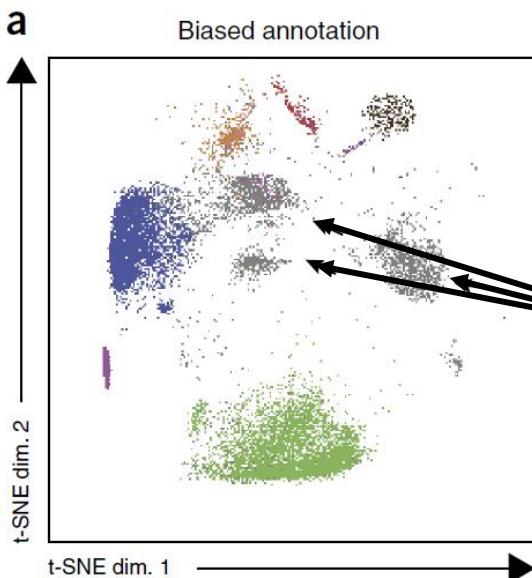


viSNE enables visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of leukemia

El-ad David Amir<sup>1</sup>, Kara L Davis<sup>2,3</sup>, Michelle D Tadmor<sup>1,3</sup>, Erin F Simonds<sup>2,3</sup>, Jacob H Levine<sup>1,3</sup>, Sean C Bendall<sup>2,3</sup>, Daniel K Shenfeld<sup>1,3</sup>, Smita Krishnaswamy<sup>1</sup>, Garry P Nolan<sup>2,4</sup> & Dana Pe'er<sup>1,4</sup>

nature biotechnology  
2013

In all cases, the viSNE gate included cells that were not classified by the expert manually gated biaxial plots; these cells are labeled in gray in the viSNE map. Examination of the marker expression of these cells reveals that they are typically just beyond the threshold of one marker, but the viSNE classification is strongly supported based on the expression of all other markers. For example, in **Figure 1d**, wherein cells are colored for CD11b marker expression, the cells in the gated region express the canonical monocyte marker CD33 (**Supplementary Fig. 1b**). However, only 47% of these cells were classified as monocytes by the manual gating (**Fig. 1b**).

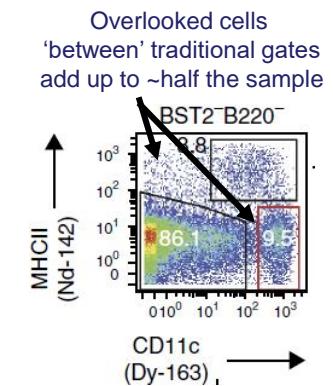


## High-dimensional analysis of the murine myeloid cell system

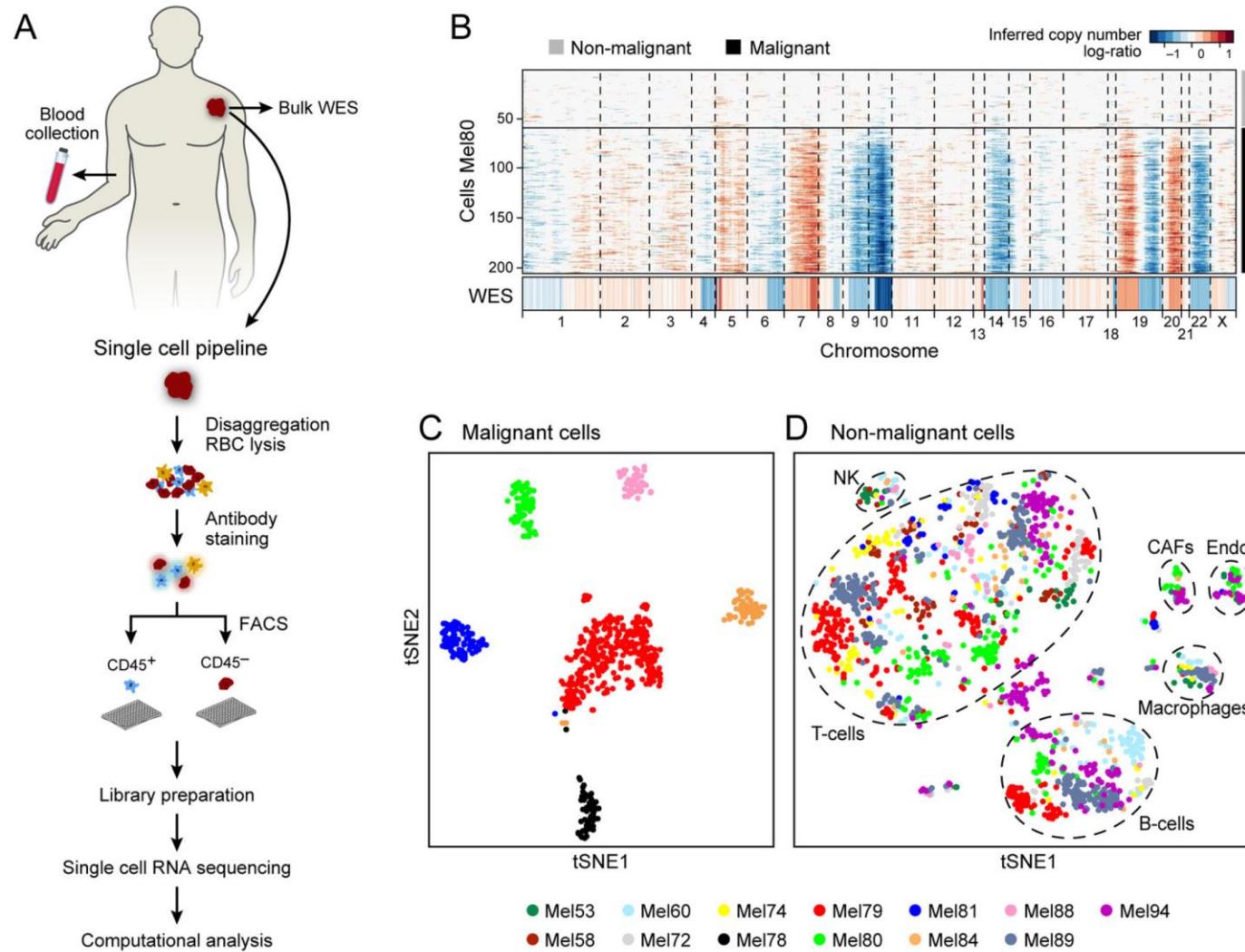
Burkhard Becher<sup>1,4,5</sup>, Andreas Schlitzer<sup>1,5</sup>, Jinmiao Chen<sup>1,5</sup>, Florian Mair<sup>2</sup>, Hermi R Sumatoh<sup>1</sup>, Karen Wei Weng Teng<sup>1</sup>, Donovan Low<sup>1</sup>, Christiane Ruedl<sup>3</sup>, Paola Riccardi-Castagnoli<sup>1</sup>, Michael Poidinger<sup>1</sup>, Melanie Greter<sup>2</sup>, Florent Ginhoux<sup>1</sup> & Evan W Newell<sup>1</sup>

nature immunology  
2014

Notably, whereas traditional biased gating strategies allowed for identification of only  $54.7 \pm 2.6\%$  (mean  $\pm$  s.e.m.,  $n = 3$  mice) of lung myeloid cells (different DC subsets, macrophages, monocytes, neutrophils), the automatic, computational approach identified nearly 100% of the cells ( $96.6 \pm 1.0\%$  (mean  $\pm$  s.e.m.,  $n = 3$  mice) accounted for by 14 predominant clusters).



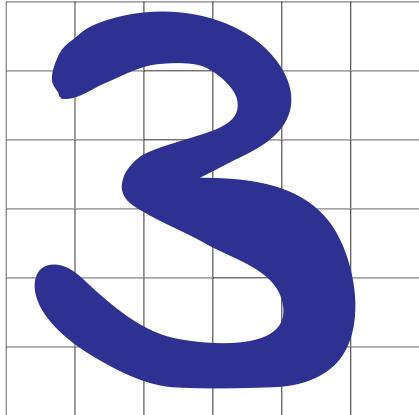
# Cytomics (Cell Identity): Powered by Multiple Platforms Melanoma Cell Diversity Based on scRNA-seq Data



# Stochastic Neighbor Embedding (SNE)

- SNE used for image recognition
- 60,000 handwritten greyscale images
- 28x28 pixels each

Example: 6x6 Pixel Image

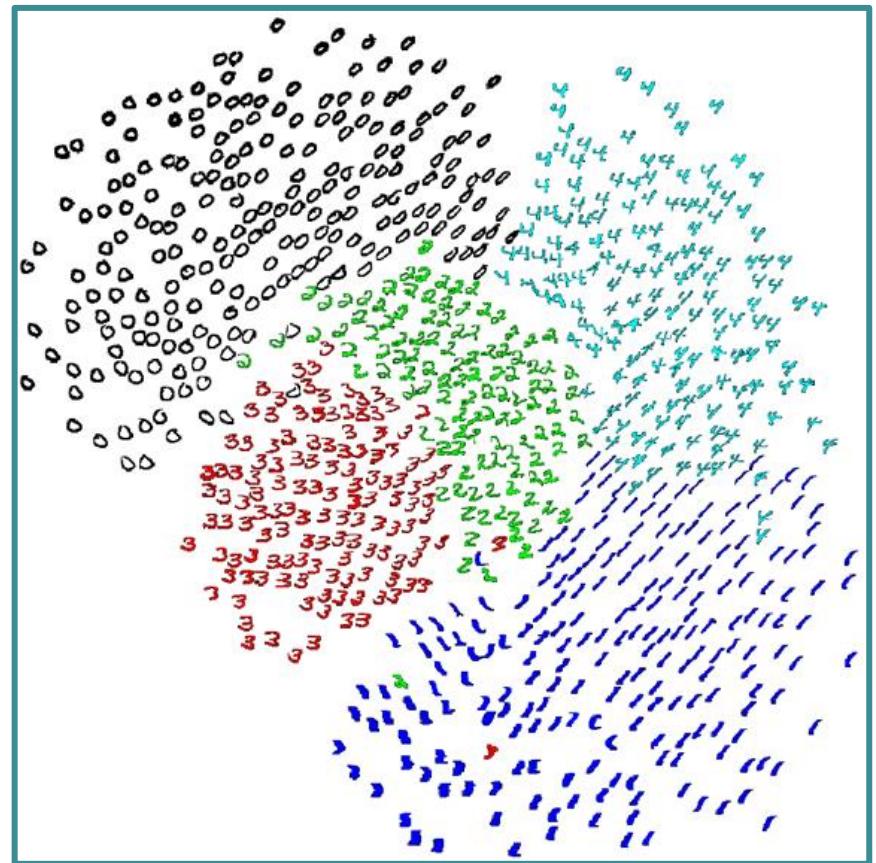


Vectorize (1x36)

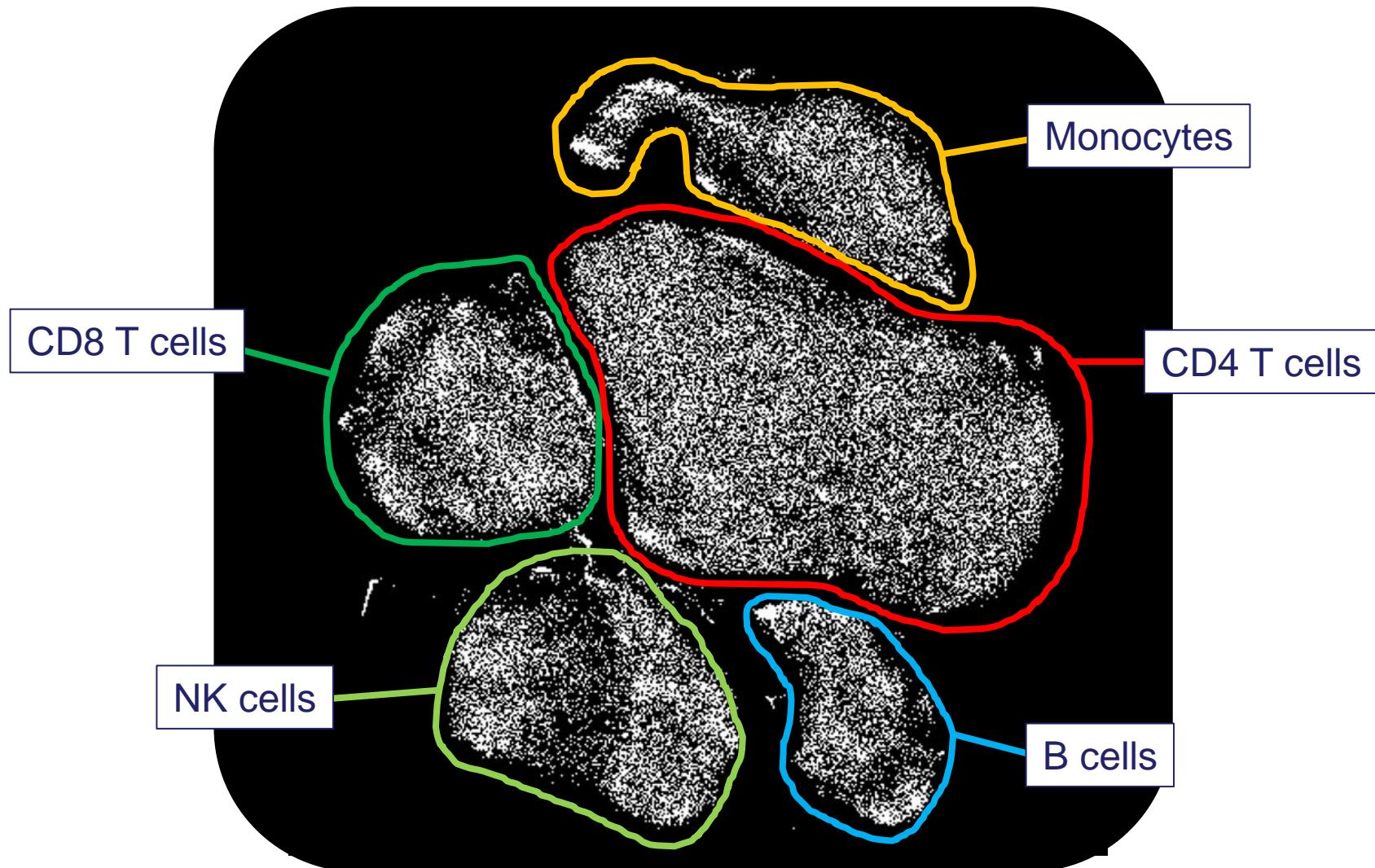


tSNE on all pixels

Hinton et al., "Advances in neural information processing systems." 2002.



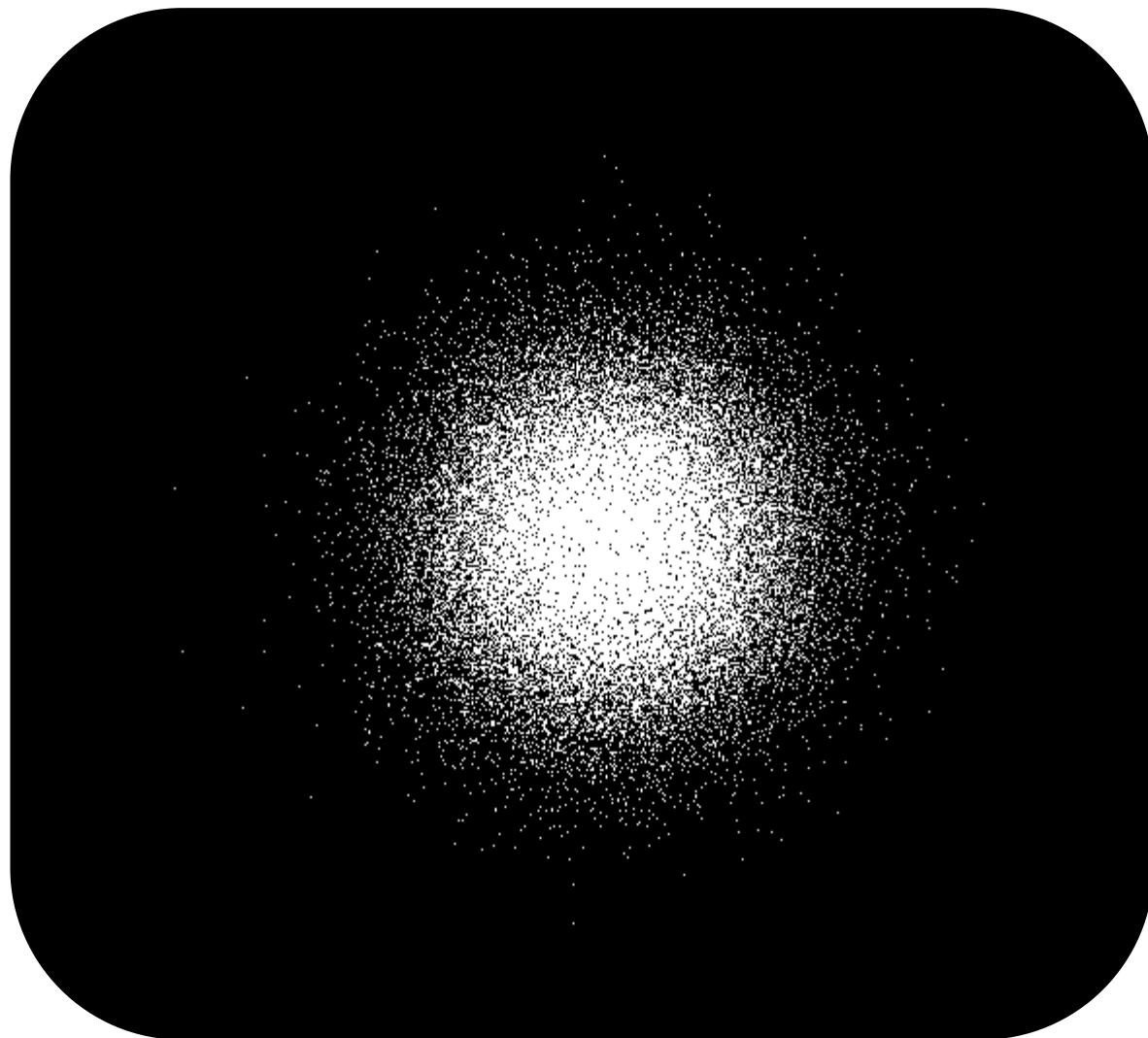
# viSNE / t-SNE Arranges Cells in 2D by Multi-D Similarity



Healthy human blood, mass cytometry,  
26 markers measured, viSNE analysis tool

Animation created by Cytobank team from iterations of viSNE / t-SNE using PBMC (26 features)

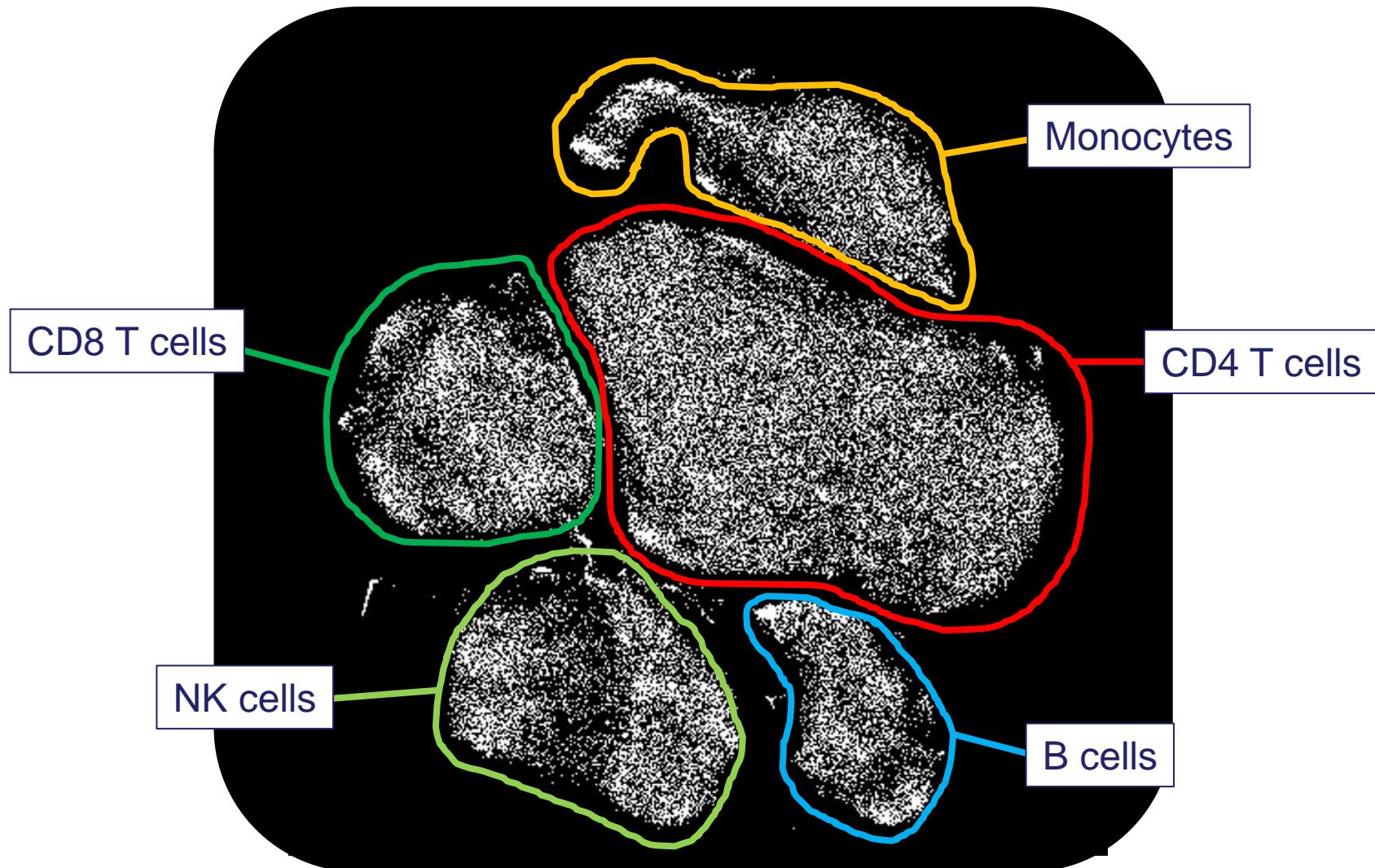
# viSNE / t-SNE Arranges Cells in 2D by Multi-D Similarity



Healthy human blood, mass cytometry,  
26 markers measured, viSNE analysis tool

Animation created by Cytobank team from iterations of viSNE / t-SNE using PBMC (26 features)

# viSNE / t-SNE Arranges Cells in 2D by Multi-D Similarity



Healthy human blood, mass cytometry,  
26 markers measured, viSNE analysis tool

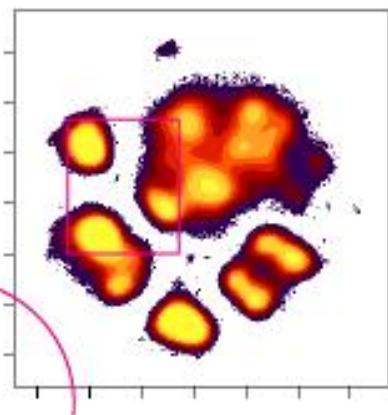
Animation created by Cytobank team from iterations of viSNE / t-SNE using PBMC (26 features)

# t-SNE Analysis Allows 2D Visualization of High Dimensional Single Cell Data

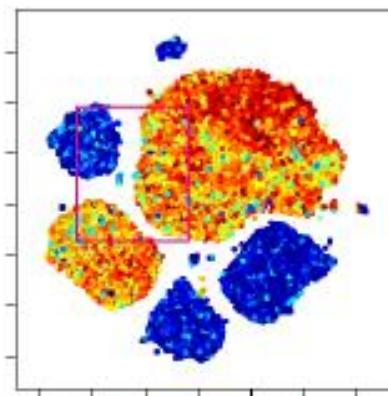
Same map, different information

Healthy Peripheral Blood Mononuclear Cells

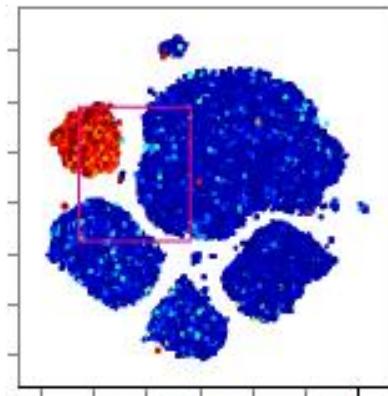
Cell Density



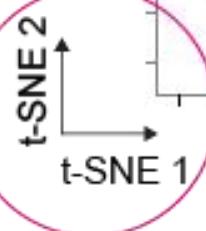
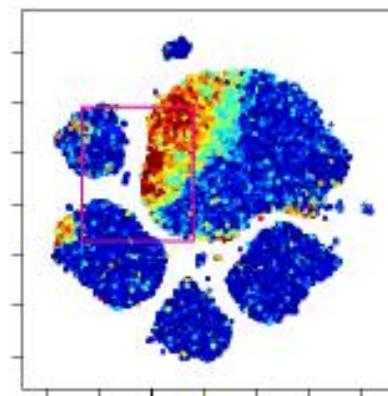
CD3



CD19



CD25

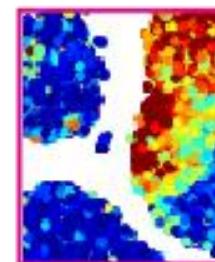
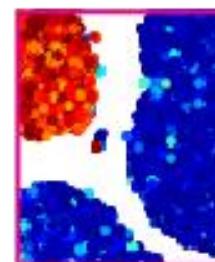
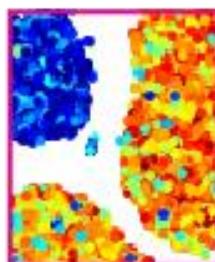


Cell density

Protein expression

min

max



New 2D axes that represent phenotypic similarities of single cells

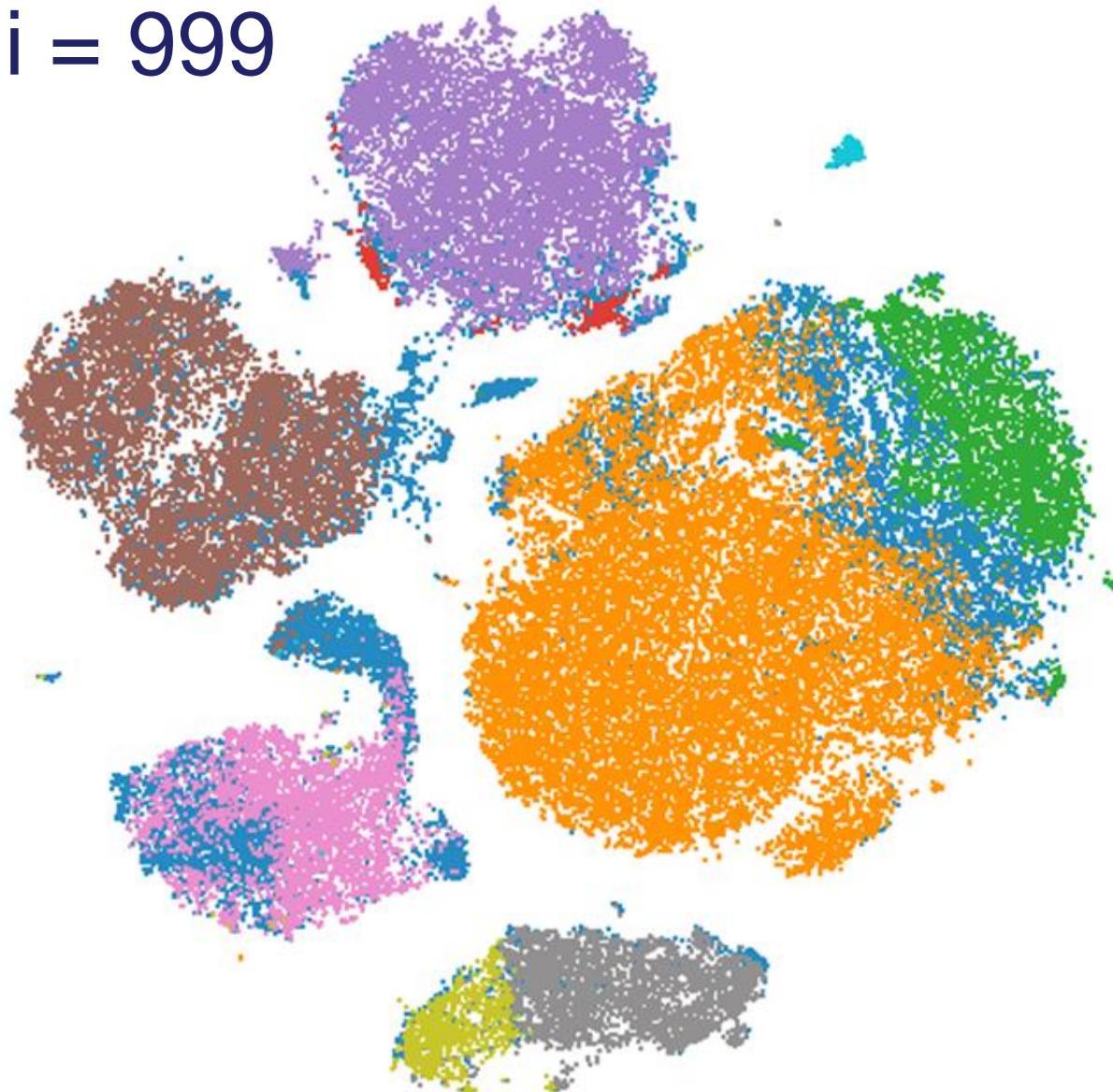
1 dot = 1 cell

# Viewing Expert Gates with viSNE Reveals Cyto Incognito

i = 999

Healthy human blood,  
mass cytometry,  
26D viSNE analysis

- Ungated
- CD45RA+ Naive CD4+
- CD45RA- CD4+ T cells
- CD45RA- Memory CD8+ T cells
- CD45RA+ Naive CD8+ T cells
- CD16+ NK cells
- CD14+ CD33+ Monocytes
- IgM+ B cells
- IgM- B cells
- CD123+ pDCs

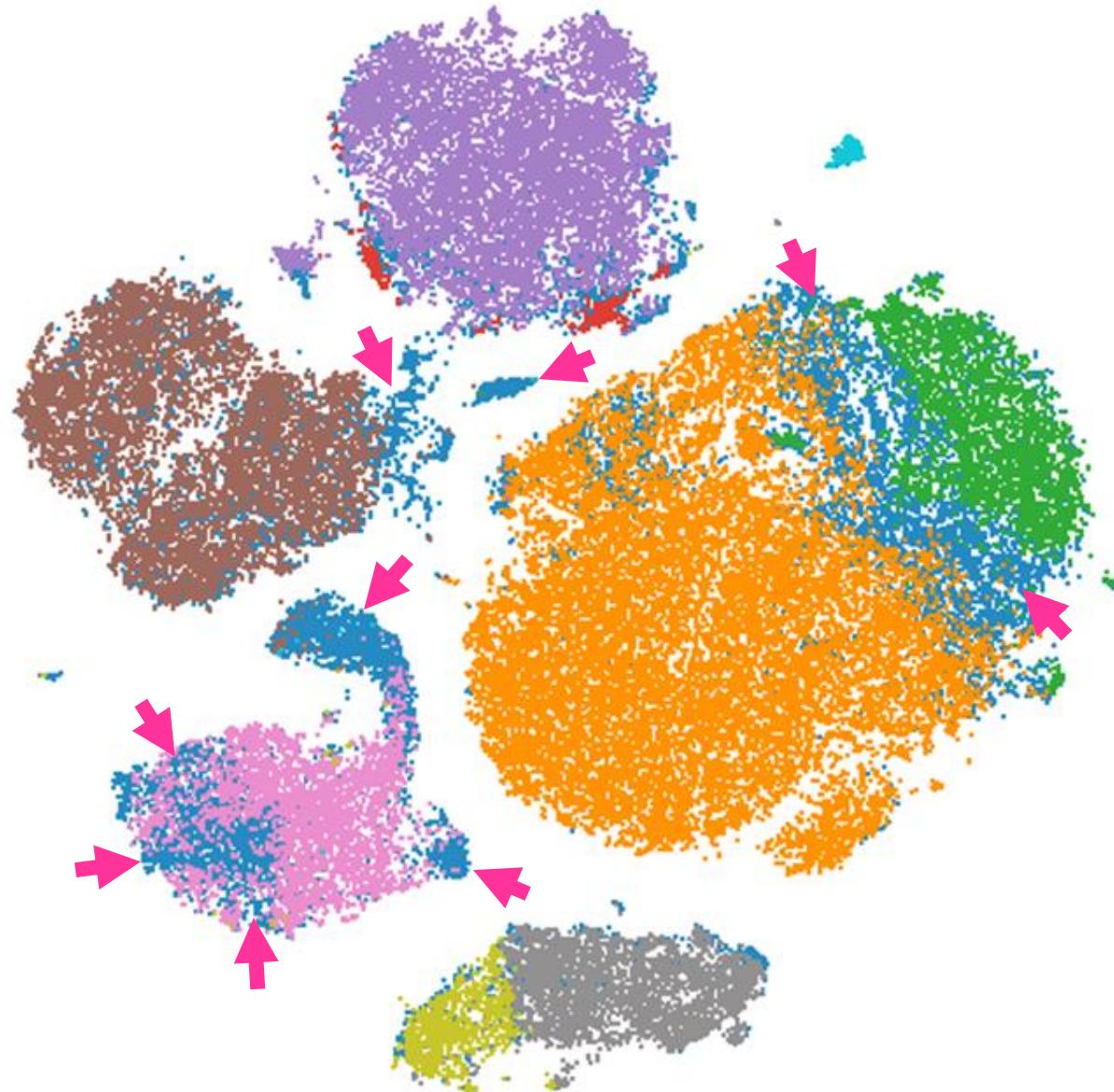


# Viewing Expert Gates with viSNE Reveals Cyto Incognito

Healthy human blood,  
mass cytometry,  
26D viSNE analysis

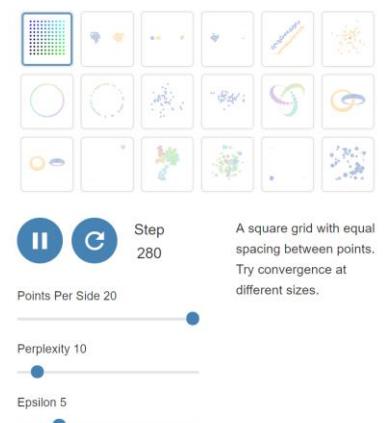
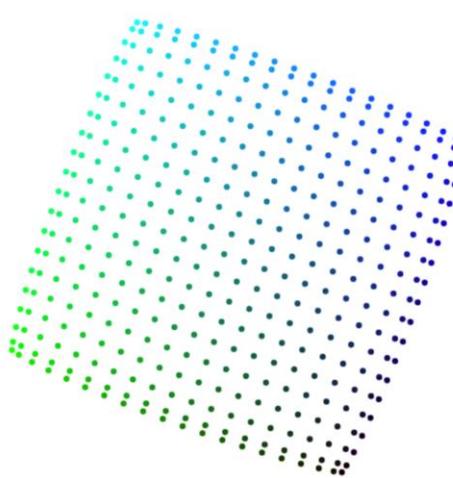
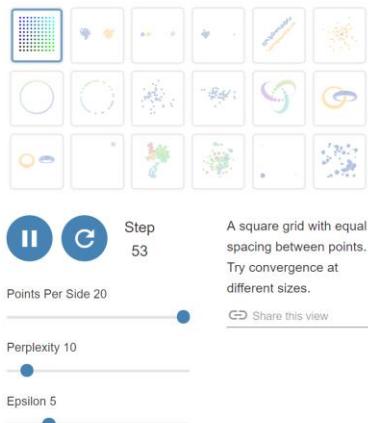
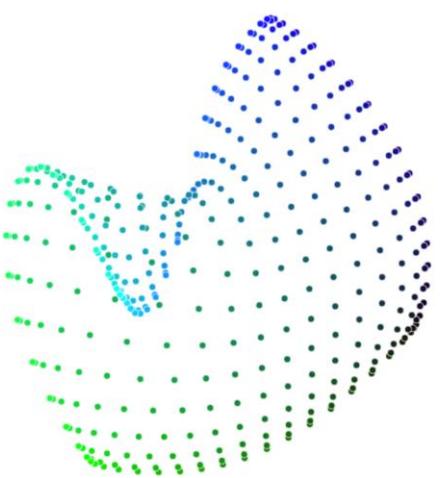
- Ungated
- CD45RA+ Naive CD4+
- CD45RA- CD4+ T cells
- CD45RA- Memory CD8+ T cells
- CD45RA+ Naive CD8+ T cells
- CD16+ NK cells
- CD14+ CD33+ Monocytes
- IgM+ B cells
- IgM- B cells
- CD123+ pDCs

➡ Cyto incognito  
(Cells overlooked or  
hidden in expert gating)



# t-SNE 2D Examples with Animations and Settings

<http://distill.pub/2016/misread-tsne/>



# opt-SNE Provides Automated Optimization of t-SNE Parameters and PCA Initialization (Fast, Reliable)

nature communications

[View all journals](#)

Search  [Login !\[\]\(68a1910447ad5088ad0b81e5d2b246b6\_img.jpg\)](#)

[Explore content](#)  [About the journal](#)  [Publish with us](#) 

[nature](#) > [nature communications](#) > [articles](#) > [article](#)

Article | [Open Access](#) | Published: 28 November 2019

## Automated optimized parameters for T-distributed stochastic neighbor embedding improve visualization and analysis of large datasets

Anna C. Belkina , Christopher O. Ciccolella, Rina Anno, Richard Halpert, Josef Spidlen & Jennifer E. Snyder-Cappione

*Nature Communications* **10**, Article number: 5415 (2019) | [Cite this article](#)

**12k** Accesses | **53** Citations | **67** Altmetric | [Metrics](#)

### Abstract

Accurate and comprehensive extraction of information from high-dimensional single cell datasets necessitates faithful visualizations to assess biological populations. A state-of-the-art algorithm for non-linear dimension reduction, t-SNE, requires multiple heuristics and fails to produce clear representations of datasets when millions of cells are projected. We develop opt-SNE, an automated toolkit for t-SNE parameter selection that utilizes Kullback-Leibler divergence evaluation in real time to tailor the early exaggeration and overall number of gradient descent iterations in a dataset-specific manner. The precise calibration of early exaggeration together with opt-SNE adjustment of gradient descent learning rate dramatically improves computation time and enables high-quality visualization of large cytometry and transcriptomics datasets, overcoming limitations of analysis tools with hard-coded parameters that often produce poorly resolved or misleading maps of fluorescent and mass cytometry data. In summary, opt-SNE enables superior data resolution in t-SNE space and thereby more accurate data interpretation.

# Check out Dr. Anna Belkina's opt-SNE Webinar on ISAC's CYTO U Learning Portal



[Donate](#) [Contact Us](#) [Sign In](#) [Join](#)

[About](#) [Membership](#) [CYTO U](#) [Publications](#) [Events](#) [Programs](#) [Resources](#) [News](#)

[« Go to Upcoming Event List](#)

A promotional graphic for a CYTO U webinar. The graphic has a dark blue background. On the left side, there is white text: "ISAC" with its logo above it, "NEW WEBINAR" in large letters, "Don't Leave Home Without a Map: Powerful Dimensionality Reduction Methods for Cytometry Data Visualization" in a larger font below that, "Presented by: Anna Belkina, MD, PhD Boston University School of Medicine" in a smaller font, and the CYTO UNIVERSITY logo at the bottom. On the right side, there is a portrait photo of a woman with long brown hair, identified as Dr. Anna Belkina. At the bottom right, there is text: "Wednesday August 25th 12pm EDT".

**ISAC**  
INTERNATIONAL SOCIETY FOR  
ADVANCEMENT OF CYTOMETRY

## NEW WEBINAR

**Don't Leave Home Without a Map:  
Powerful Dimensionality Reduction  
Methods for Cytometry Data Visualization**

Presented by:  
**Anna Belkina, MD, PhD**  
Boston University School of Medicine

**CYTO**  
UNIVERSITY  
Your Cytometry Learning Portal

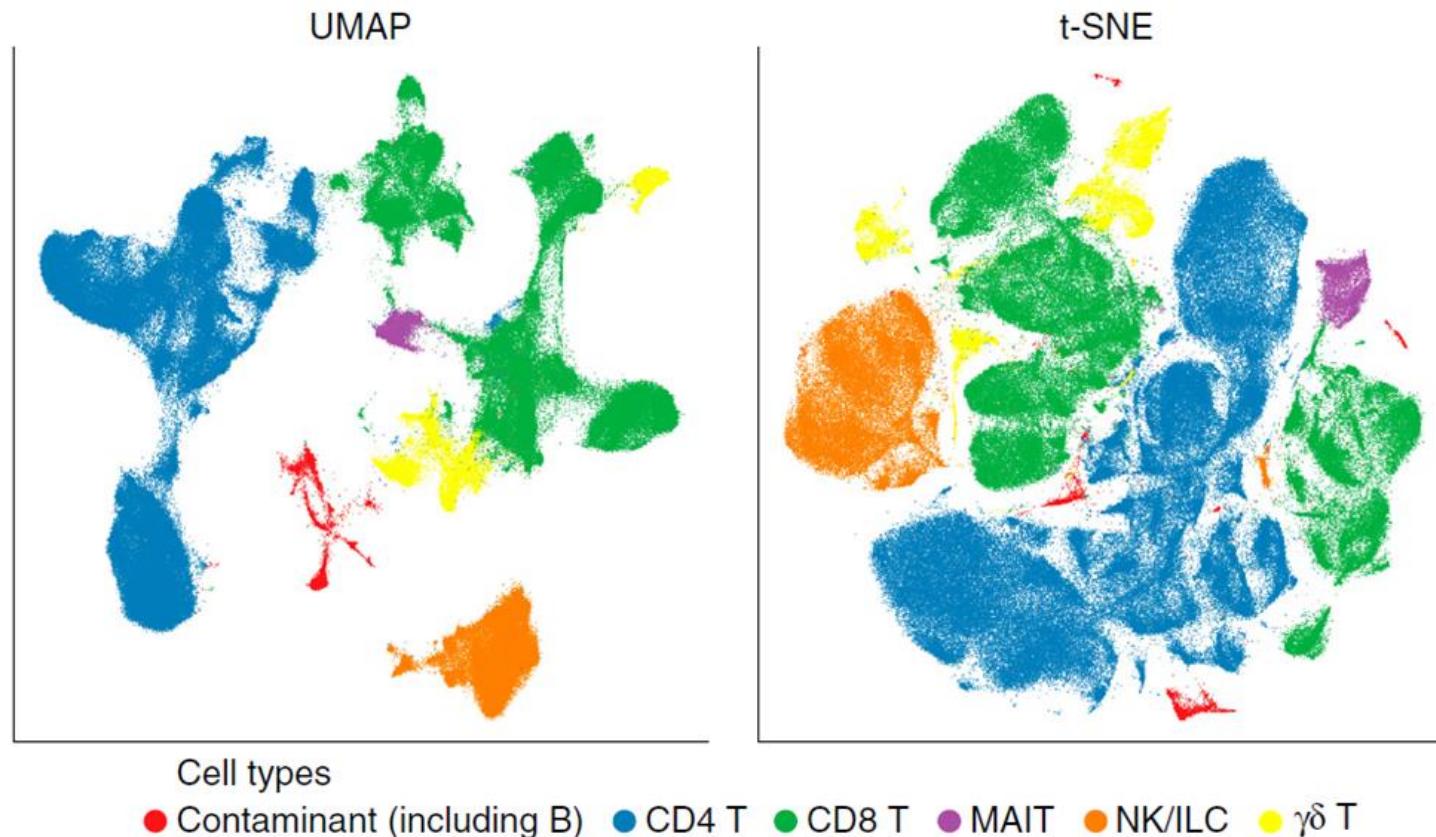
*Register today at <https://learning.isac-net.org>*

Wednesday  
August 25th  
12pm EDT

CYTO U Webinar: <https://learning.isac-net.org/products/dont-leave-home-without-a-map-powerful-dimensionality-reduction-methods-for-cytometry-data-visualization>

# Becht et al., UMAP Preserves Local and Global Structure (Analysis of Tissue T Cells; Color = Expert Knowledge / Source)

(a) UMAP better split CD8 T cells,  $\gamma\delta$  T cells, and contaminating cells

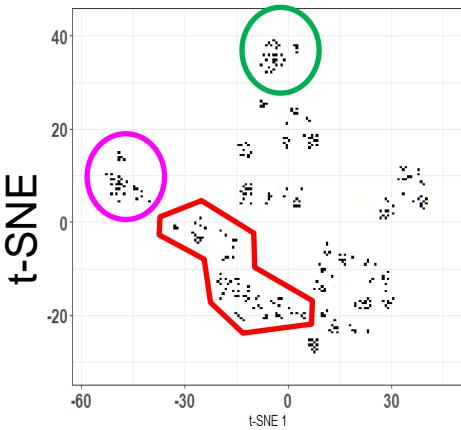


Dataset covering 35 samples originating from 8 distinct human tissues enriched for T and natural killer (NK) cells, of more than >300,000 cell events with 39 protein targets (Wong et al. dataset).

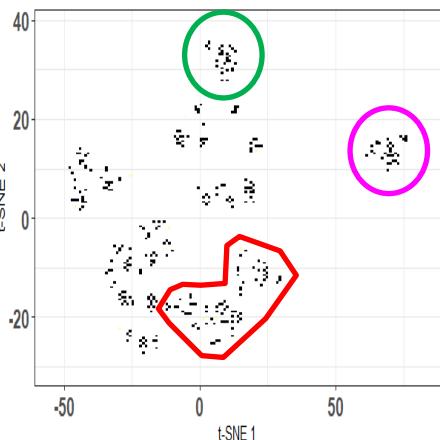
# Multiple Runs of t-SNE vs. UMAP on a Patient Dataset (n = 339)

Gandelman et al., cGVHD Patient Dataset

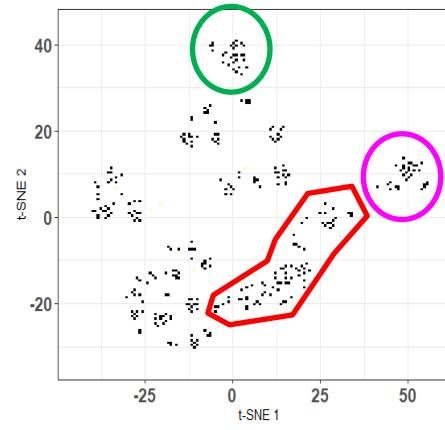
Run 1



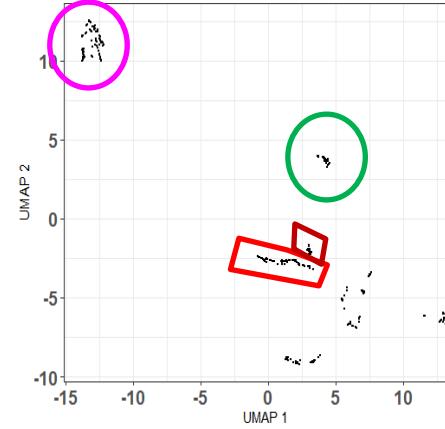
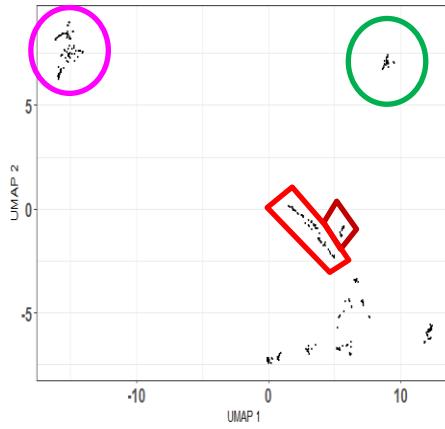
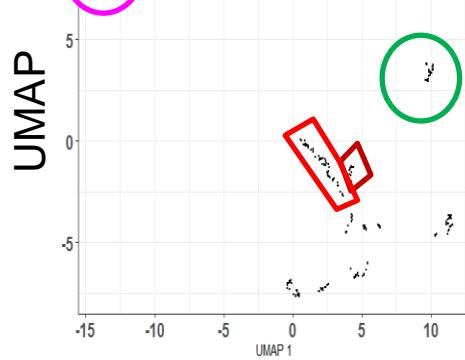
Run 2



Run 3



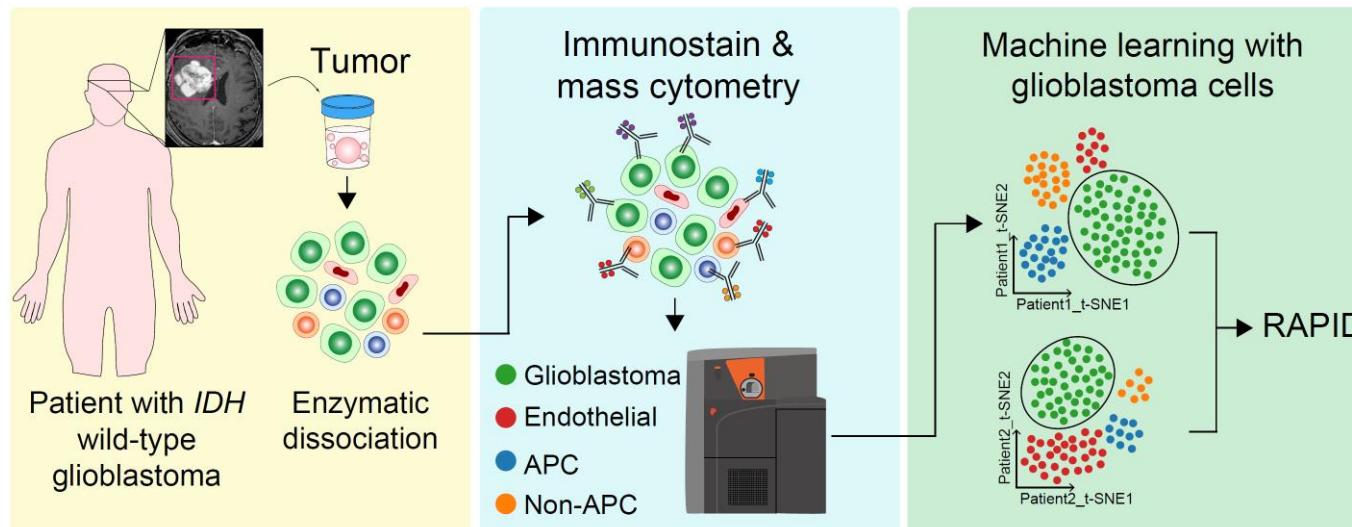
In the t-SNE plots, the relative relationship of the major islands (“global structure”) alters between runs; t-SNE focuses on local structure



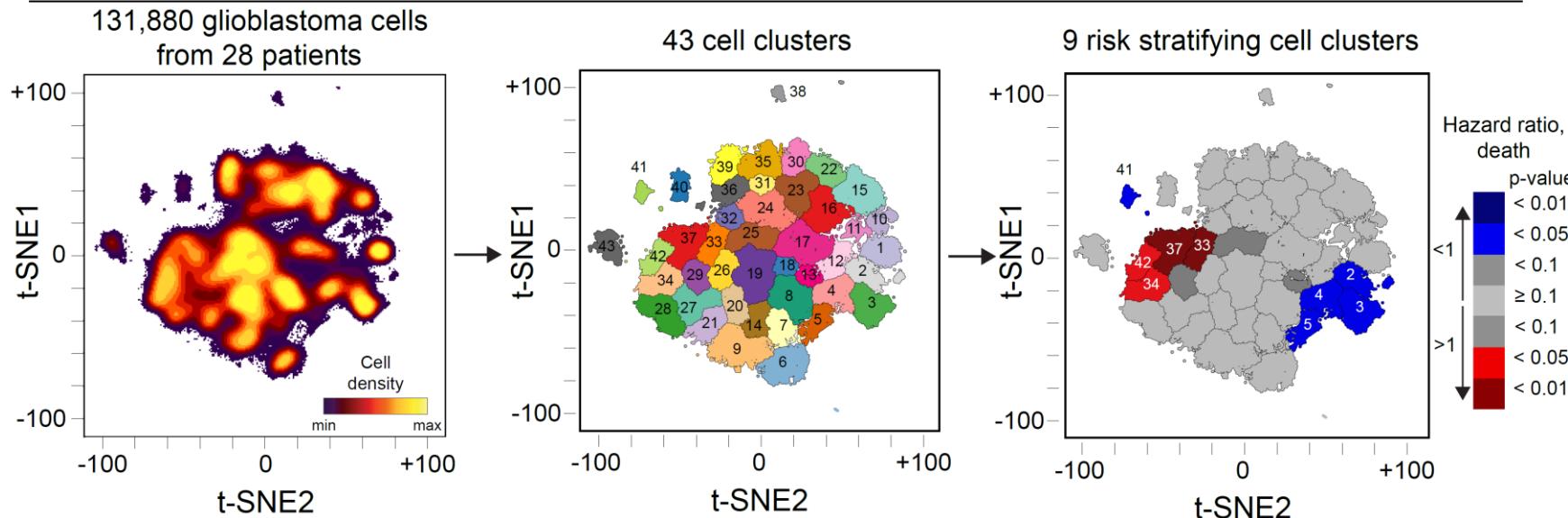
Relative island position (“global structure”) is more stable & reflects original measurements in UMAP

Principle Component Analysis (PCA) is linear and deterministic, meaning that it strictly preserves global structure (and can overlook significant local structures / paths / trajectories)

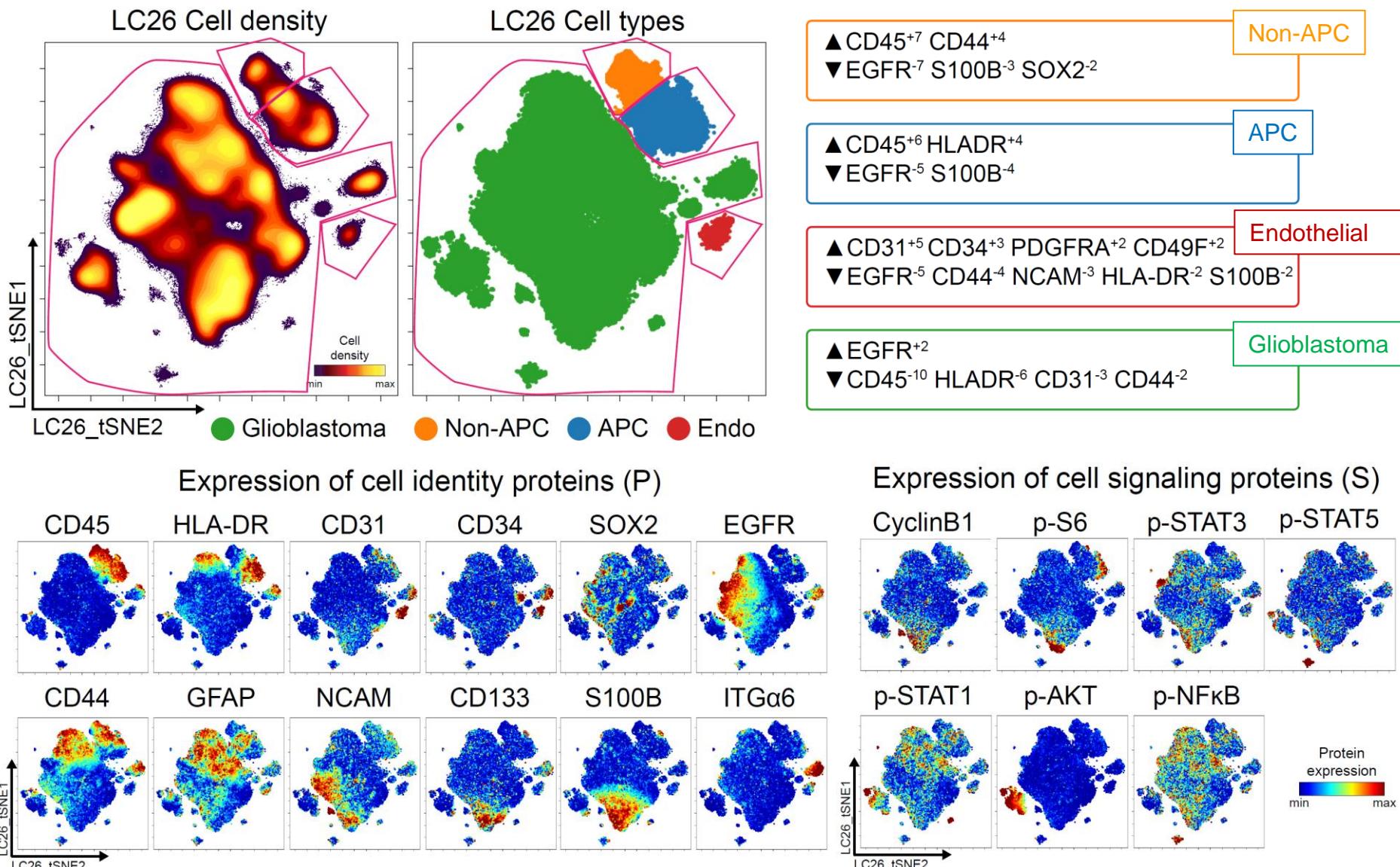
# RAPID Maps Clinical Outcomes Onto Clusters (in t-SNE, UMAP, 2D image, original features, PCA, etc.)



## Risk Assessment Population IDentification (RAPID) Maps Outcome onto t-SNE



# t-SNE Can Help Computationally Isolate Cancer Cells from Endothelial and Leukocyte Cell Subsets

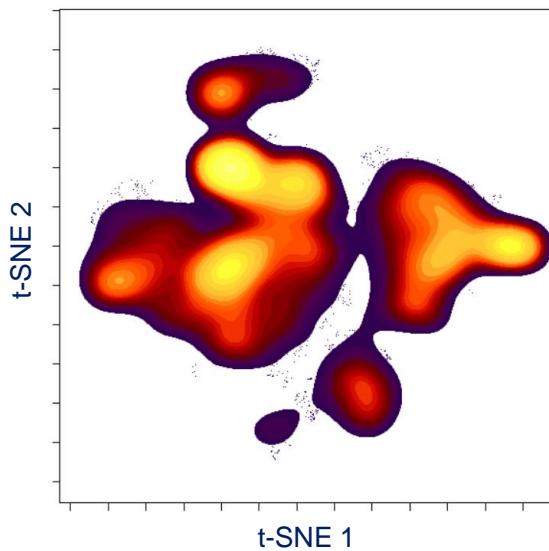


# UMAP (Uniform Manifold Approximation and Projection) is Another Dimensionality Reduction Tool

Superior run times

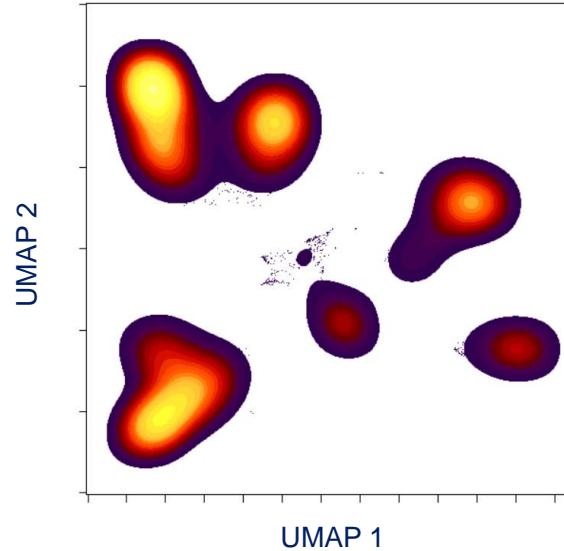
Emphasis on both global and local structure in the data

**t-SNE**

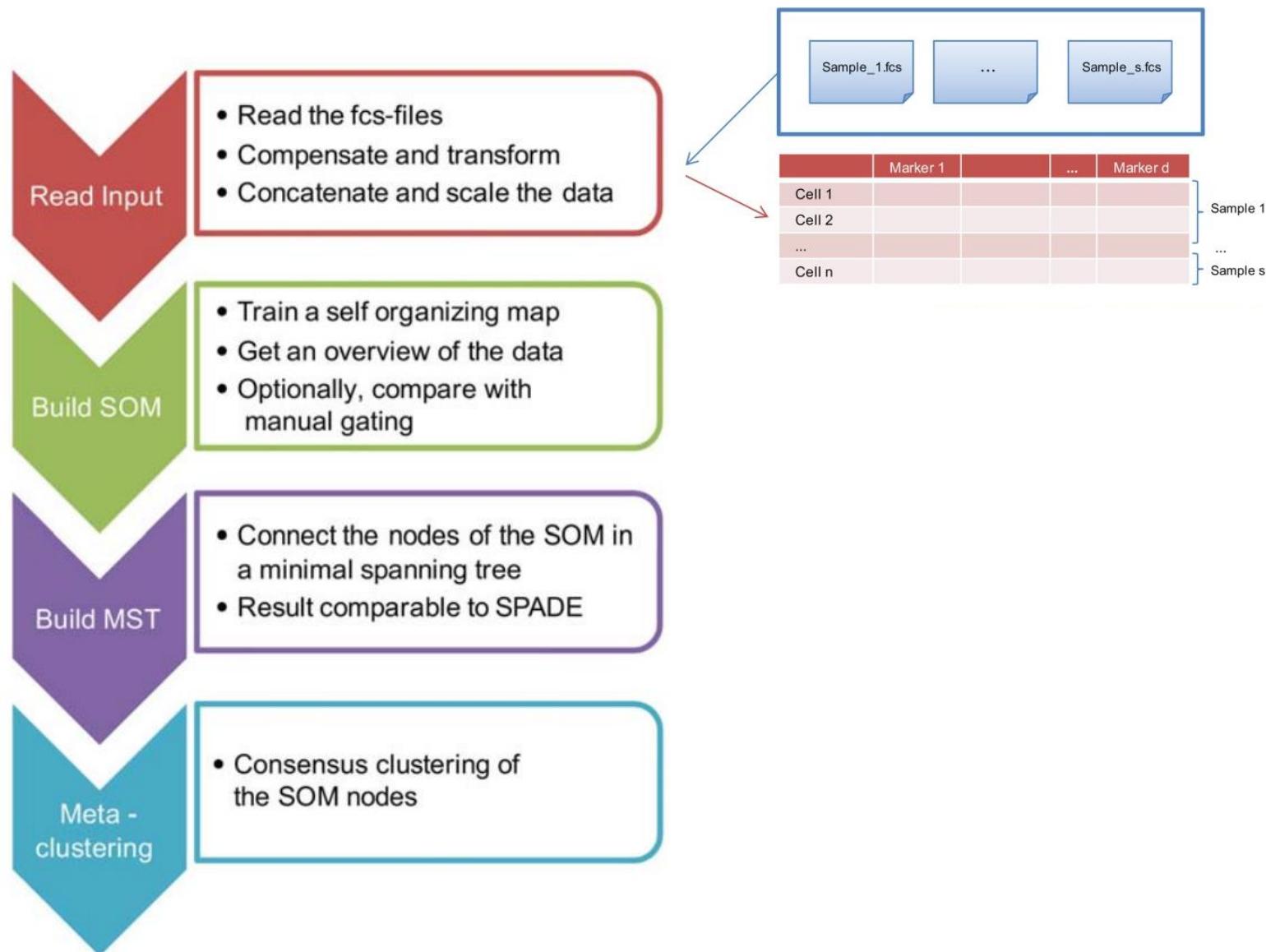


**vs.**

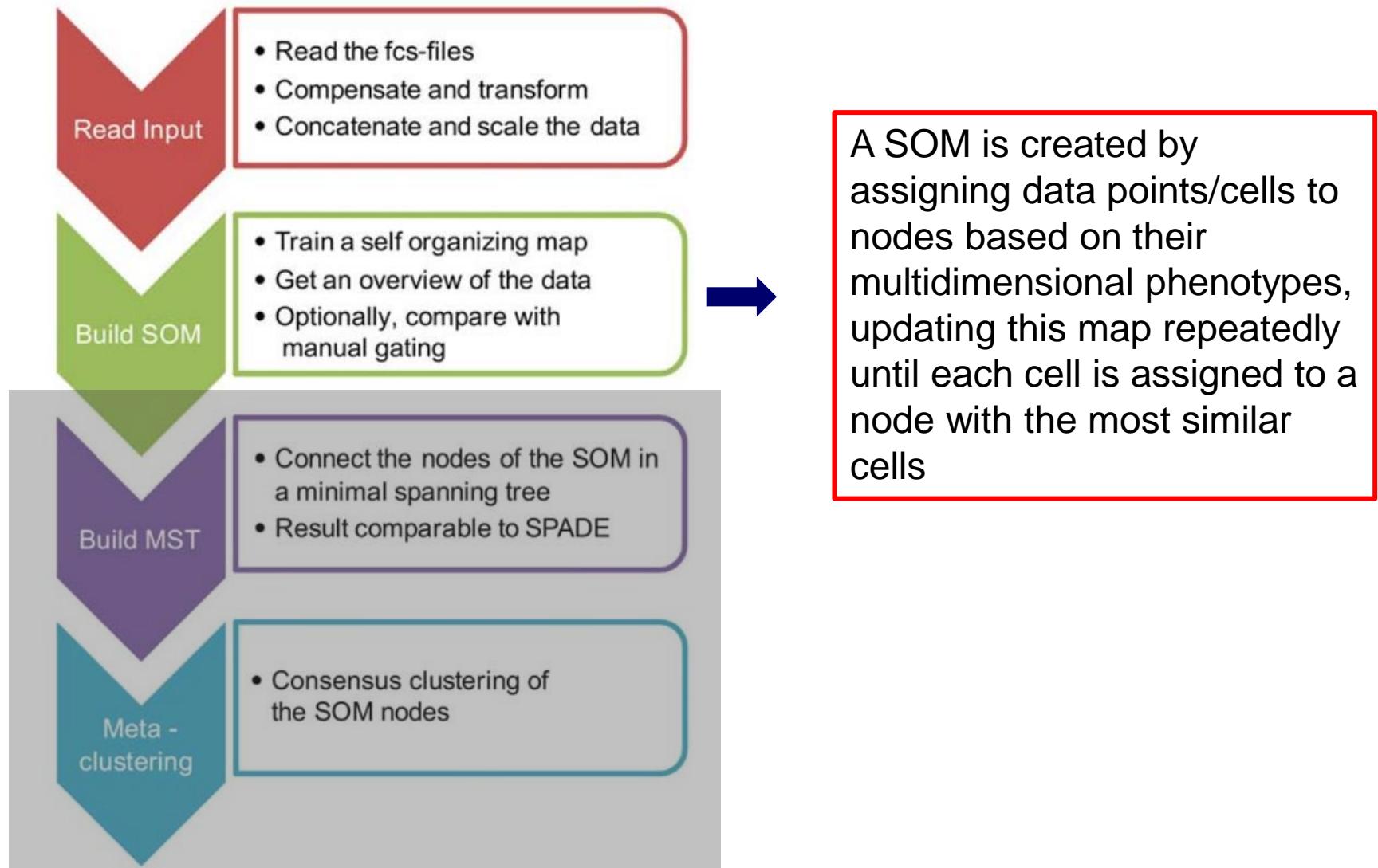
**UMAP**



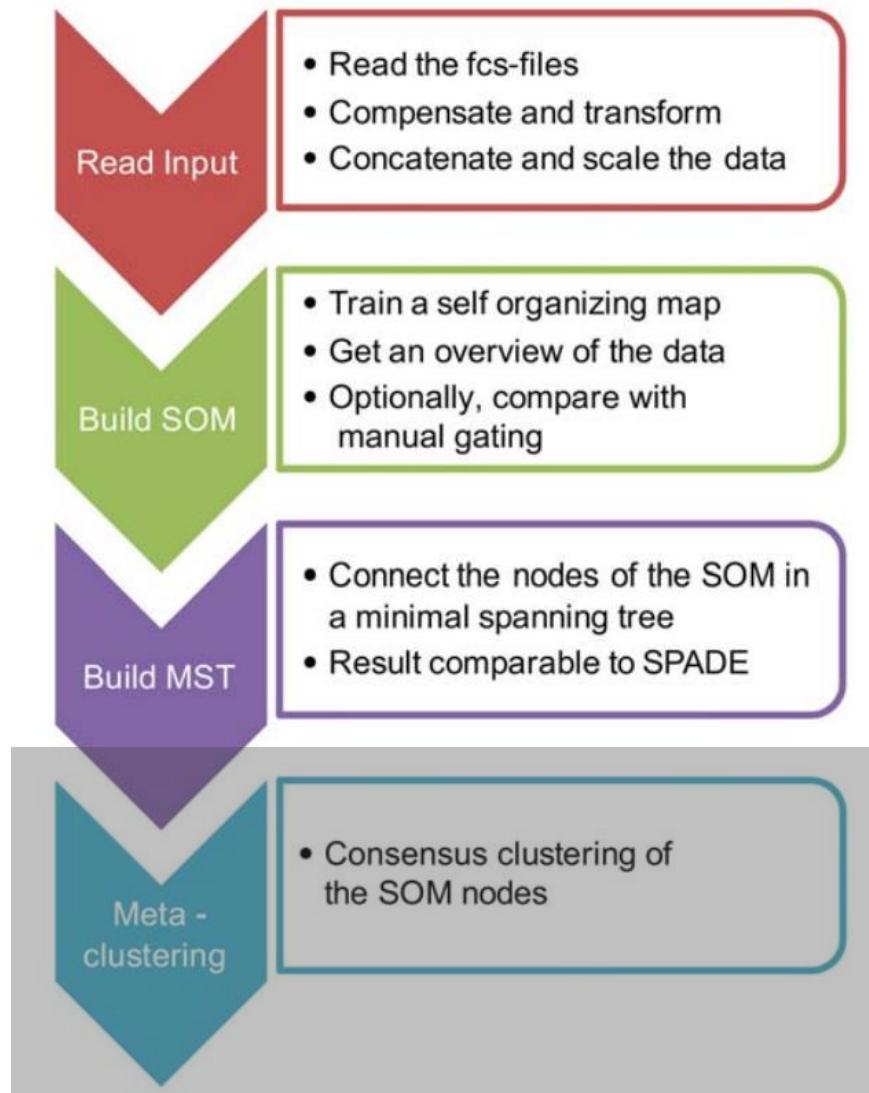
# Clustering with FlowSOM: Self-organizing Maps



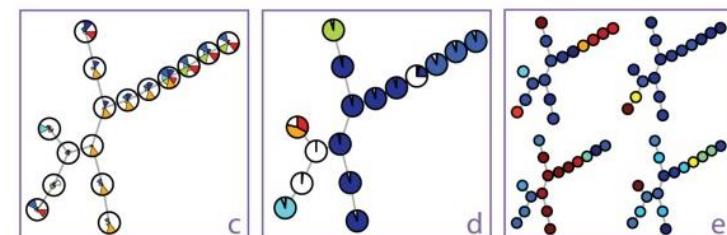
# Clustering with FlowSOM: Self-organizing Maps



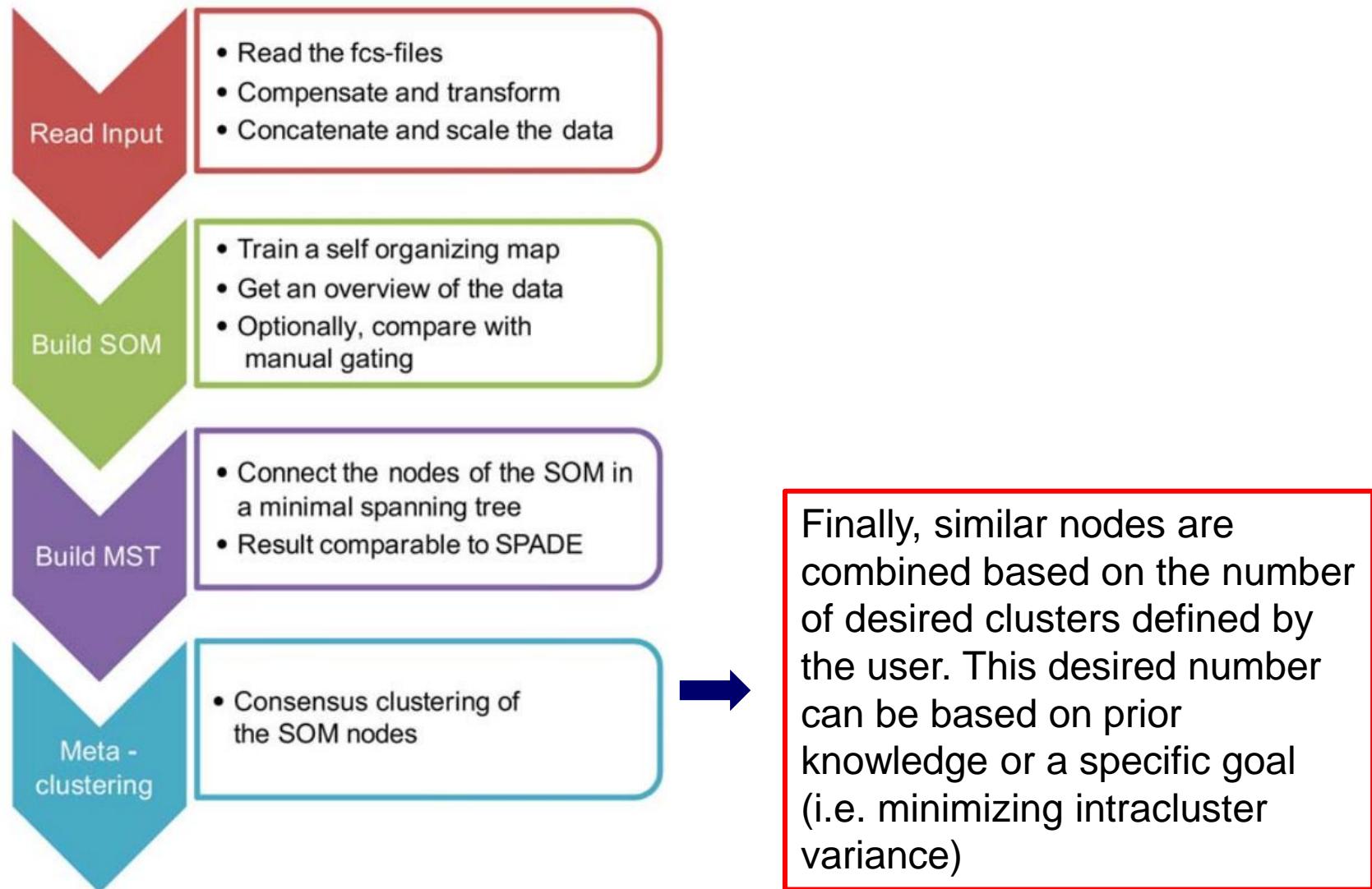
# Clustering with FlowSOM: Self-organizing Maps



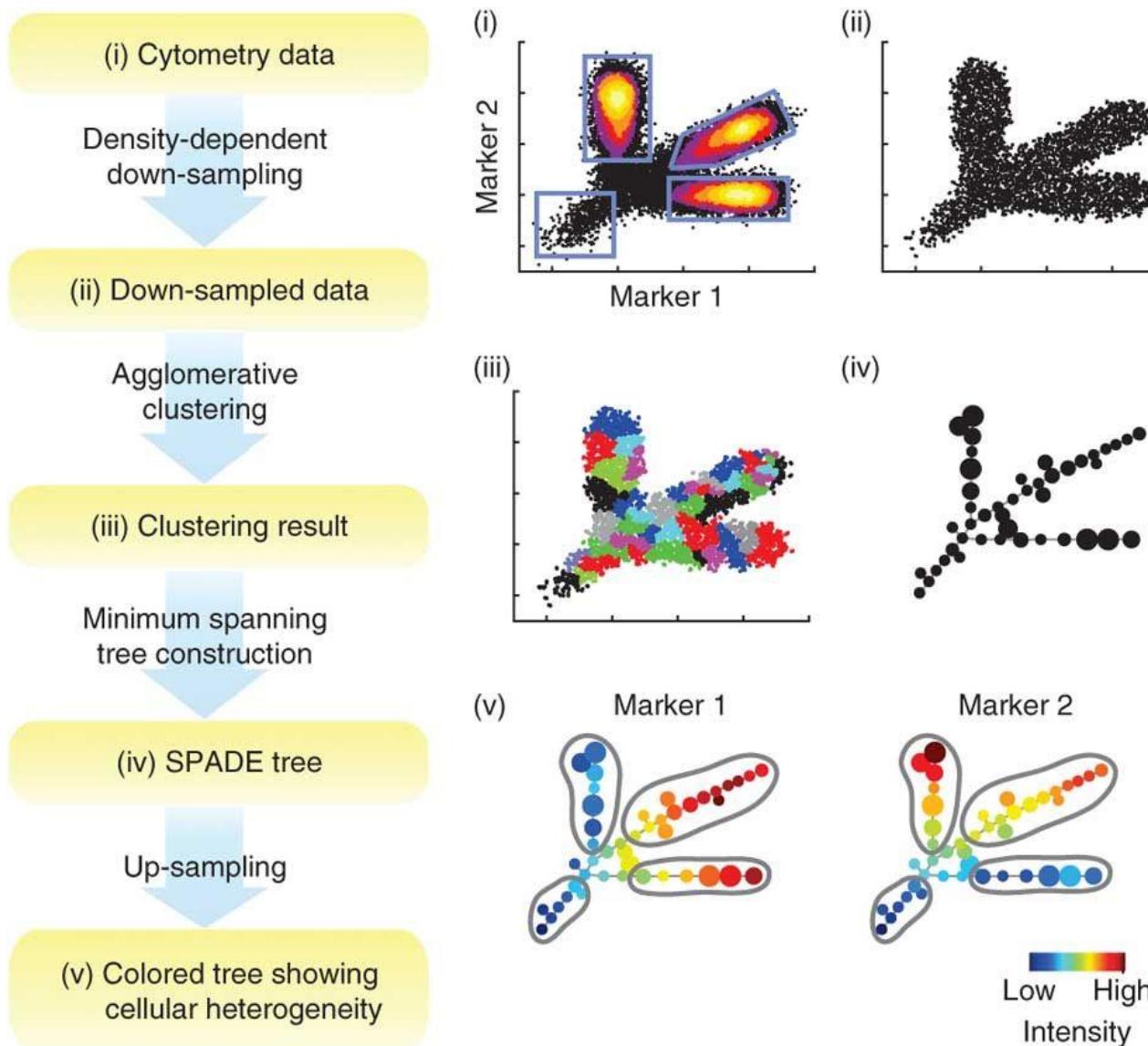
The next step is to arrange the nodes along a minimal spanning tree (MST), so that nodes that are most similar are closest on the tree  
\*not used in our visualization\*



# Clustering with FlowSOM: Self-organizing Maps

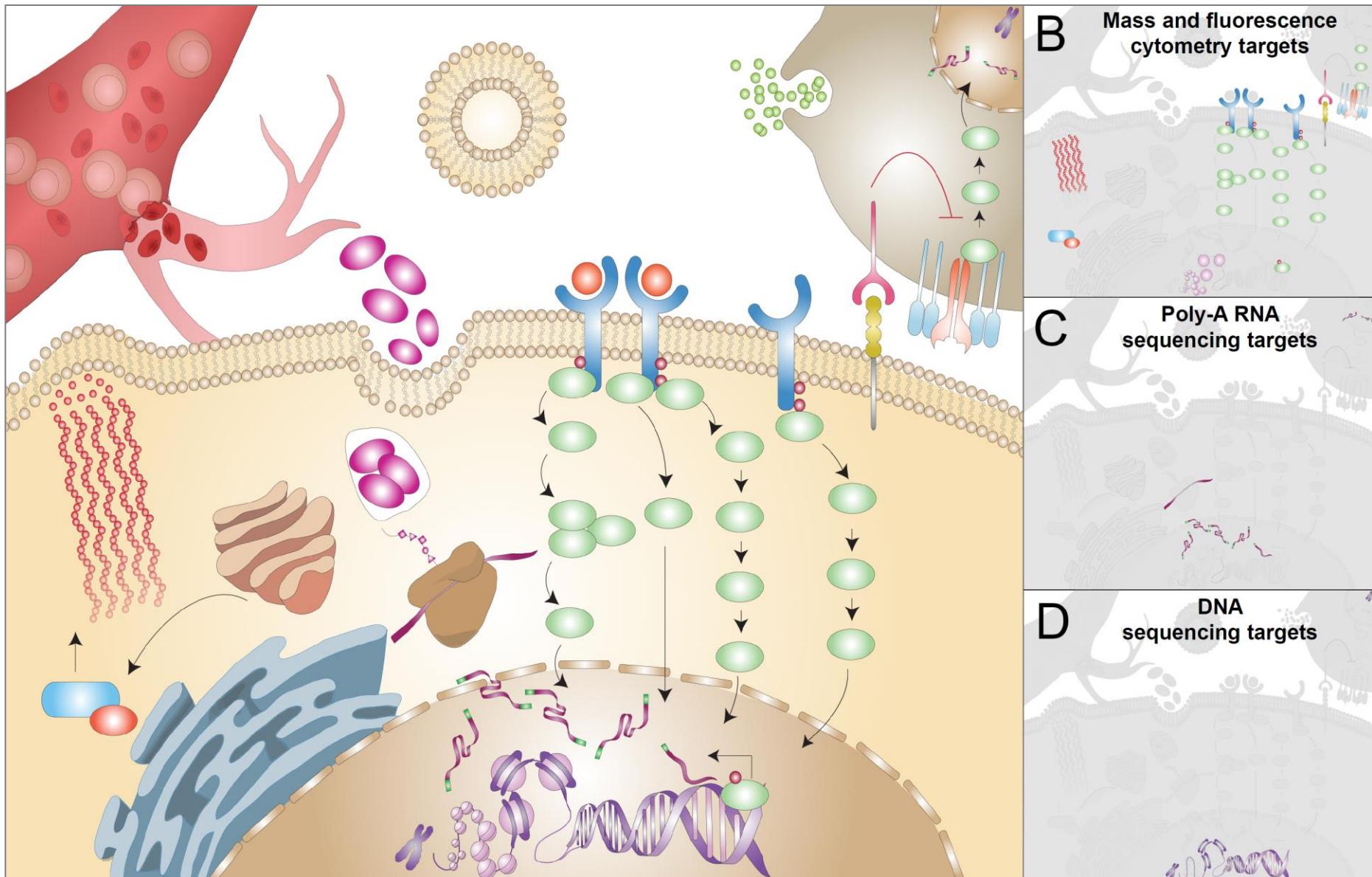


# Spanning-Tree Progression Analysis of Density-Normalized Events (SPADE) is an Alternative Clustering Tool



# Part 5: What Is This Cluster? MEM Labels & Identifies Group of Cells

# Challenge to the Field: Multiplex Across Cell Functions

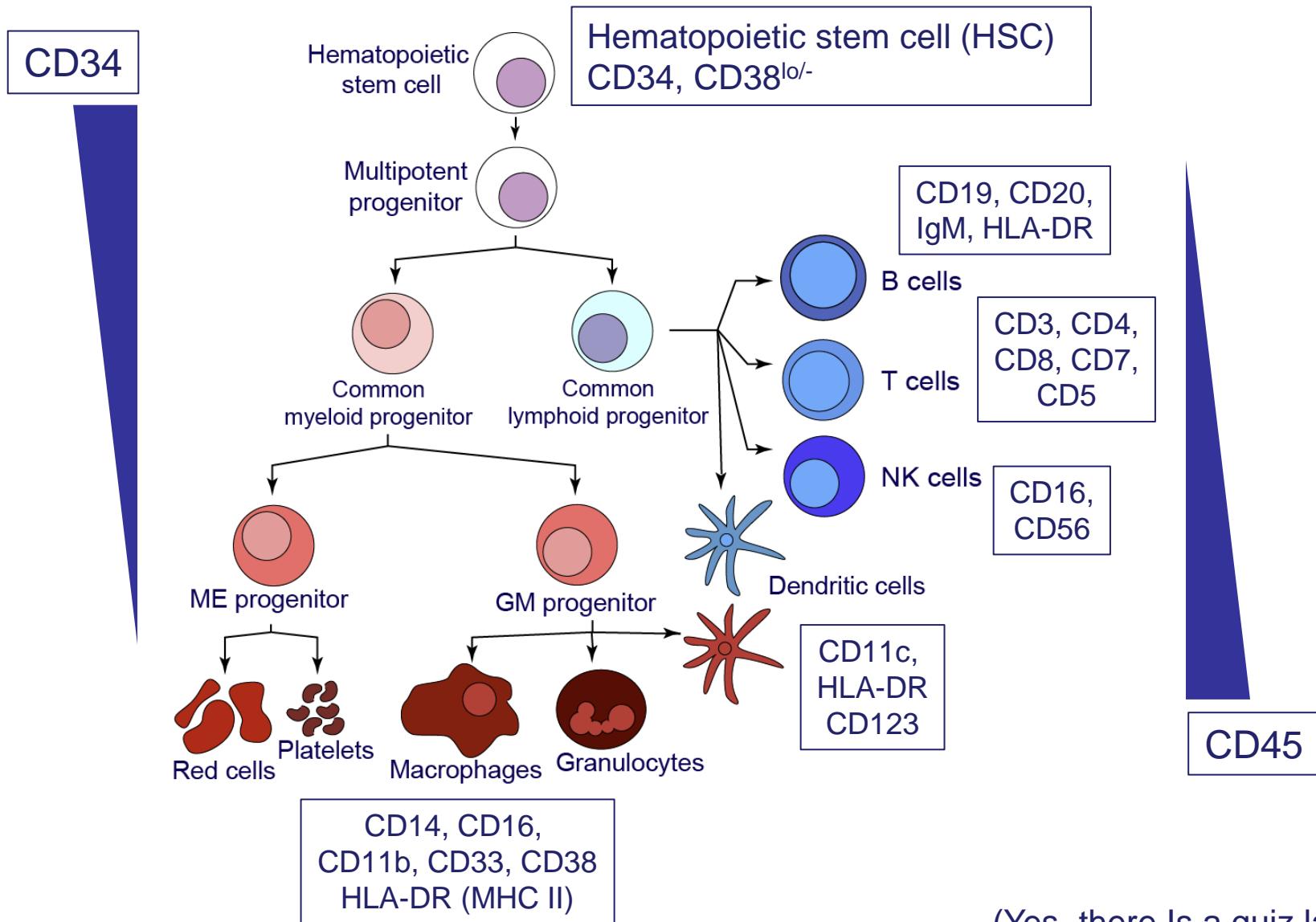


Beyond the Message: Advantages of Snapshot Proteomics with Single-Cell Mass Cytometry in Solid Tumors  
Mistry et al., *FEBS Journal* 2018

MEM summarizes a population's special features  
and is used in workflows "at the end"  
(in place of box and whisker plots or heatmaps)

[ So MEM complements tools from other steps, including  
t-SNE, SPADE, Citrus, FlowSOM, SCAFFOLD, Phenograph ]

# Human Bone Marrow Hematopoiesis & “Famous” Cell Identity Markers



Despite advances, no computational tools learn & label cell identity,  
a human must “stare and compare” using expert knowledge

Diggins et al., *Methods* 2015

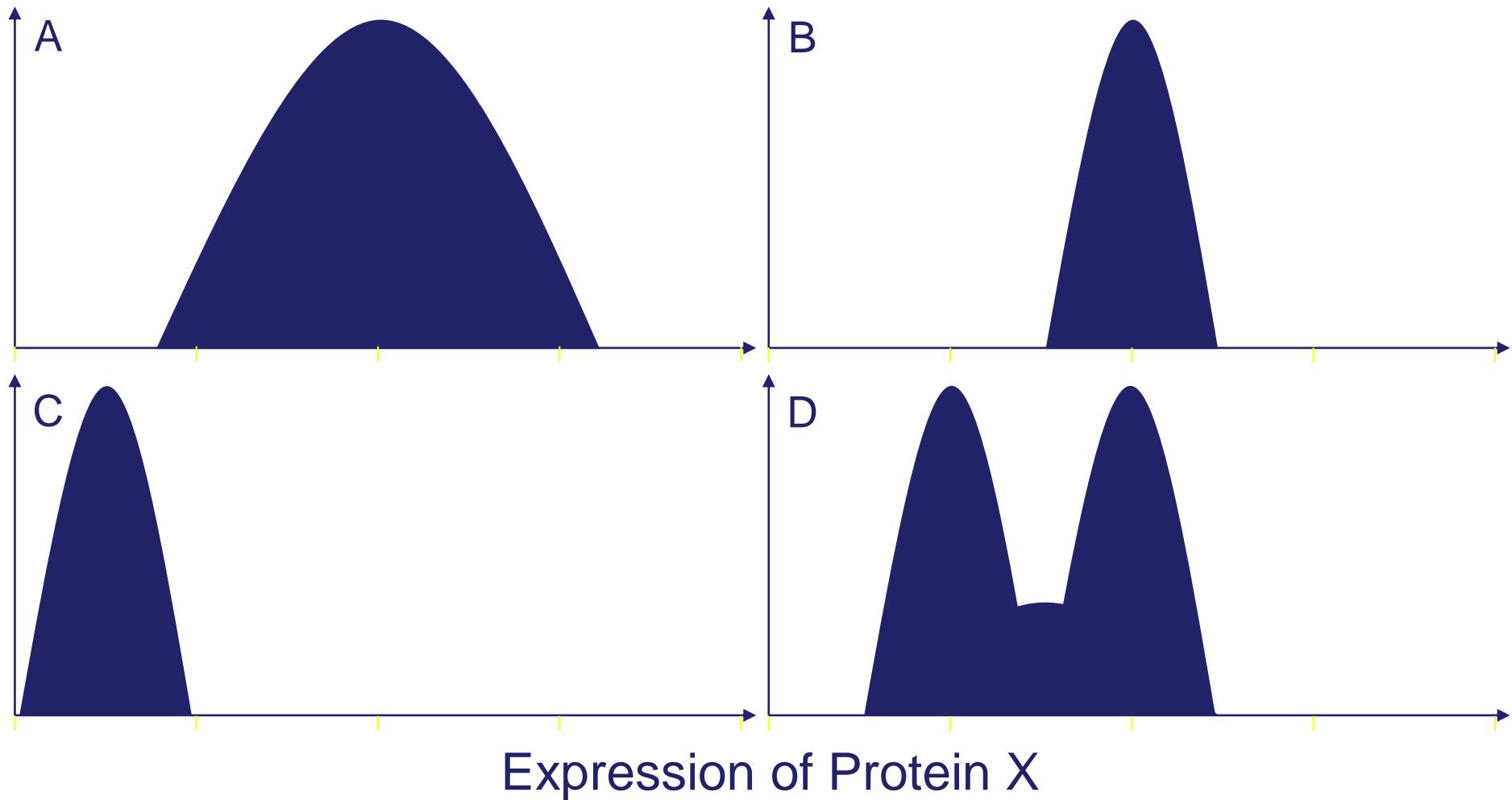
Populations are often labeled by metaphors of function  
 (“cancer stem cells”, “central memory T cells”)  
 or incomplete labels based on a few features (e.g. “PD-1+ CD8 T cells”).

We need an unbiased way to label & identify cells  
 (regardless of how they are found)

# Enrichment Tracks Feature Exclusivity In a Subset

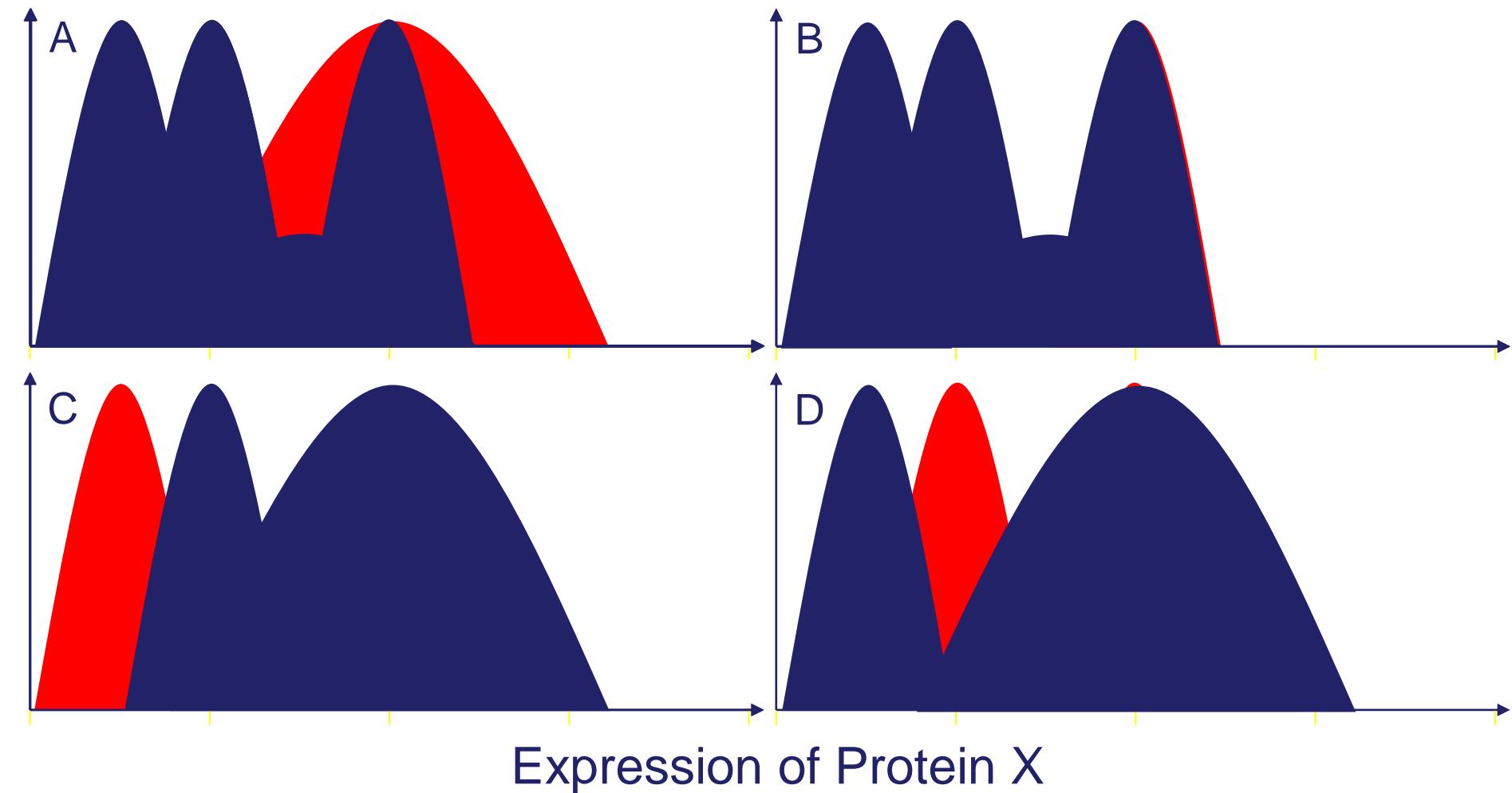
A, B, C, and D are 4 subsets where Protein X was measured.

In which subset is Protein X most distinct? (Which would be easiest to gate?)

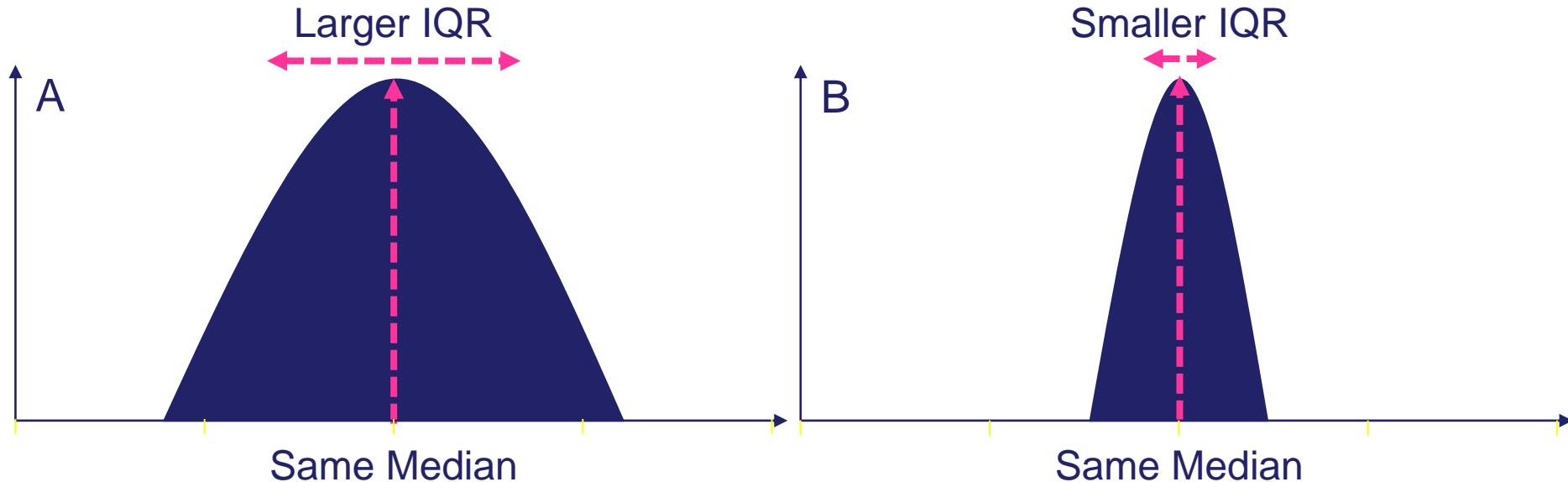


# Enrichment Tracks Feature Exclusivity In a Subset

A, B, C, and D are 4 subsets where Protein X was measured.  
In which subset is Protein X most distinct? (Which would be easiest to gate?)



# Median (50%) and Interquartile Range (25%-75%) Represent Key Features of Distributions



Core idea in MEM: given two protein distributions with equal medians, a smaller interquartile range (IQR) indicates greater enrichment

Not captured by median & IQR are other elements of shape  
(skewness, symmetry, # peaks, outliers, etc.)

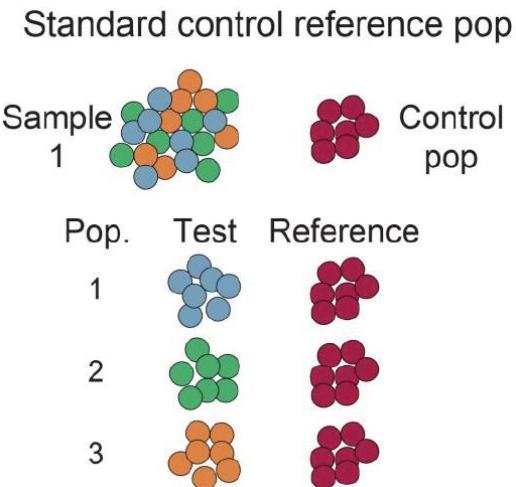
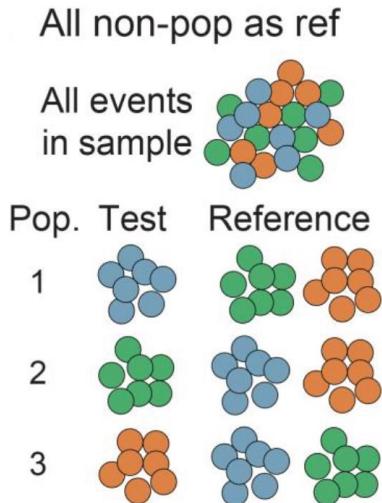
# MEM Quantifies Relative Enrichment By Combining Magnitude & Interquartile Range

$$MEM = |MAG_{test} - MAG_{ref}| + \frac{IQR_{ref}}{IQR_{test}} - 1$$



Linear transformation to -10 to +10  
(d20 scale, cause that's how we roll)

If  $MAG_{test} - MAG_{ref} < 0$ ,  $MEM = -MEM$



# Quiz Time: What Are These Cell Subsets & What Is This Tissue?

## Stem cells (HSCs)

▲ CD34<sup>+6</sup> CD33<sup>+4</sup> CD15<sup>+3</sup> CD38<sup>+3</sup>  
MHCII<sup>+3</sup> CXCR4<sup>+2</sup>  
▼ CD44<sup>-5</sup> CD45<sup>-5</sup> CD7<sup>-3</sup> 0.07%

## Natural killer cells

▲ CD16<sup>+9</sup> CD7<sup>+6</sup> CD38<sup>+5</sup> CD56<sup>+4</sup> CD161<sup>+4</sup>  
CD45RA<sup>+3</sup> CD8<sup>+2</sup> CD11b<sup>+2</sup> CD47<sup>+2</sup> 5.27%

## Progenitors

▲ MHCII<sup>+10</sup> CD33<sup>+7</sup> CD38<sup>+5</sup> CD123<sup>+3</sup>  
CD117<sup>+3</sup> CD19<sup>+2</sup> CD34<sup>+2</sup> CD13<sup>+2</sup>  
CD14<sup>+2</sup> CXCR4<sup>+2</sup>  
▼ CD45<sup>-3</sup> CD15<sup>-2</sup> 0.002%

## CD8<sup>+</sup> T cells

▲ CD8<sup>+8</sup> CD7<sup>+5</sup> CD3<sup>+3</sup>  
CD45RA<sup>+3</sup> CXCR4<sup>+2</sup> 9.25%

## Early myeloid cells

▲ MHCII<sup>+9</sup> CD33<sup>+8</sup> CD38<sup>+5</sup>  
CD4<sup>+3</sup> CD15<sup>+2</sup> CD14<sup>+2</sup>  
▼ CD45<sup>-2</sup> CD7<sup>-2</sup> 0.02%

## CD4<sup>+</sup> T cells

▲ CD4<sup>+7</sup> CD7<sup>+5</sup> CD3<sup>+5</sup>  
CD47<sup>+2</sup> CD45RA<sup>+2</sup> 8.12%

## Monocytes

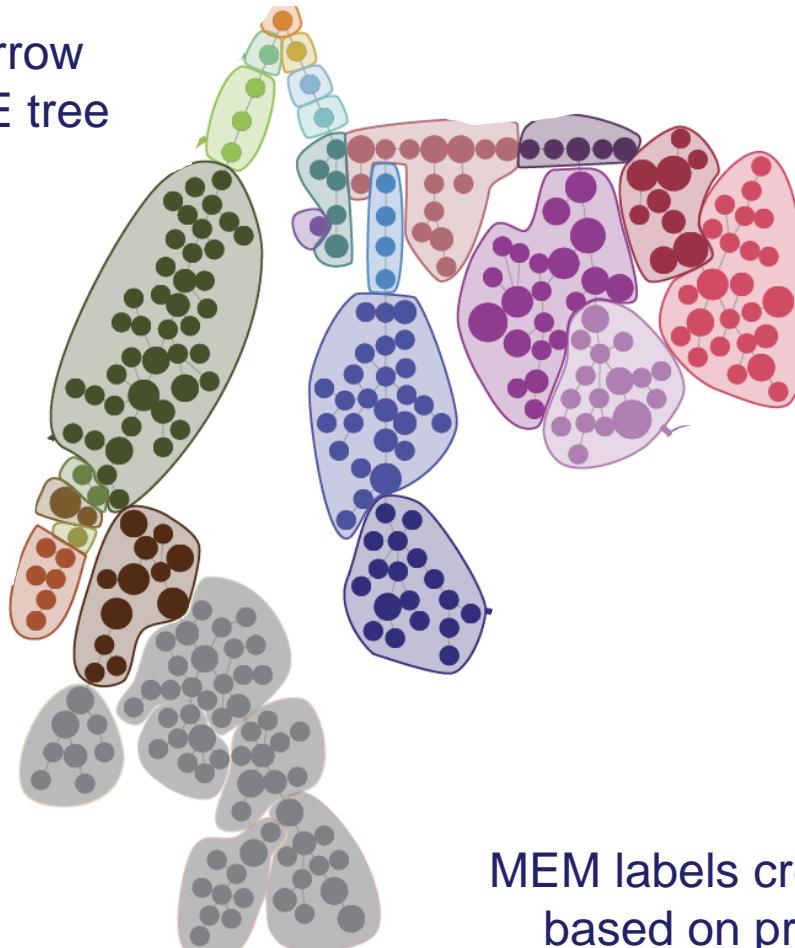
▲ CD33<sup>+10</sup> CD14<sup>+8</sup> CD11b<sup>+7</sup>  
MHCII<sup>+5</sup> CD4<sup>+4</sup> CD11c<sup>+4</sup>  
CD38<sup>+4</sup> CD13<sup>+3</sup>  
▼ CXCR4<sup>-2</sup> CD47<sup>-2</sup> 10.57%

## B cells

▲ MHCII<sup>+10</sup> CD20<sup>+9</sup> CD19<sup>+7</sup> IgM<sup>+5</sup> CD34<sup>+3</sup>  
CD45RA<sup>+3</sup> CXCR4<sup>+2</sup> CD47<sup>+2</sup> CD33<sup>+2</sup>  
▼ CD7<sup>-2</sup> 2.44%

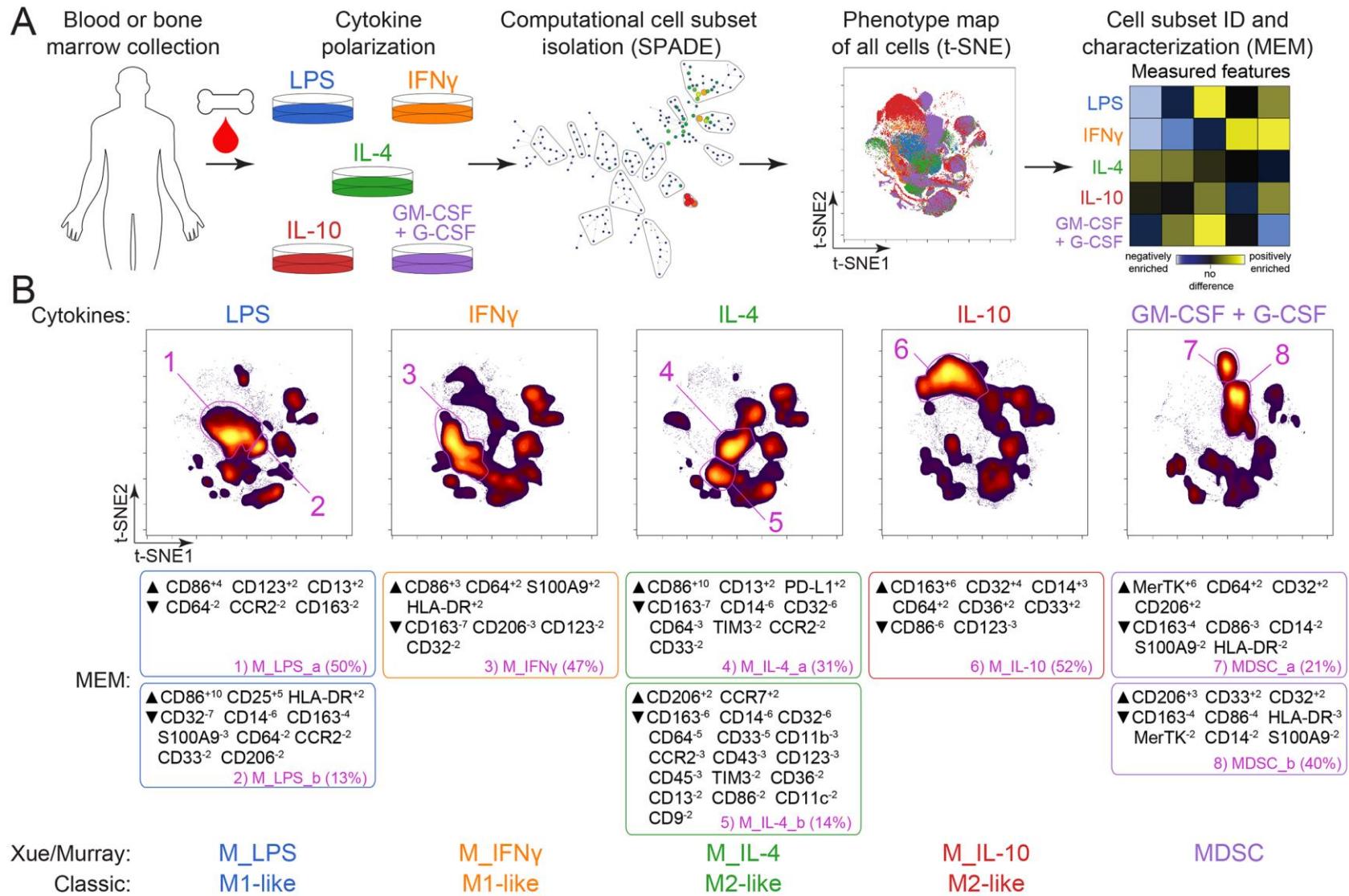
# Marker Enrichment Modeling Automatically Labels Cell Types in Human Bone Marrow Using -10 to +10 Enrichment Values

Cells from bone marrow grouped in a SPADE tree



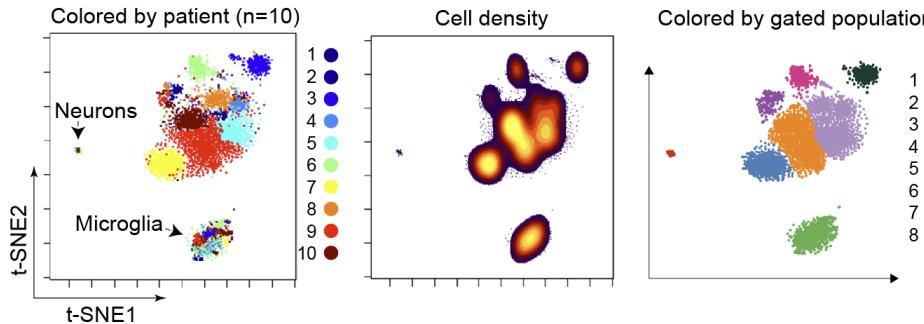
MEM labels created automatically based on protein enrichment

# Functional Profiling: Polarization Capacity of Monocytes



# MEM Classification of Cell Subsets in IDH-A Brain Tumors Based on scRNA-seq Transcript Expression

IDH-A brain tumor cells from 10 patients  
Gating for 8 cell types: t-SNE using 500 most variable transcripts, gating on local density



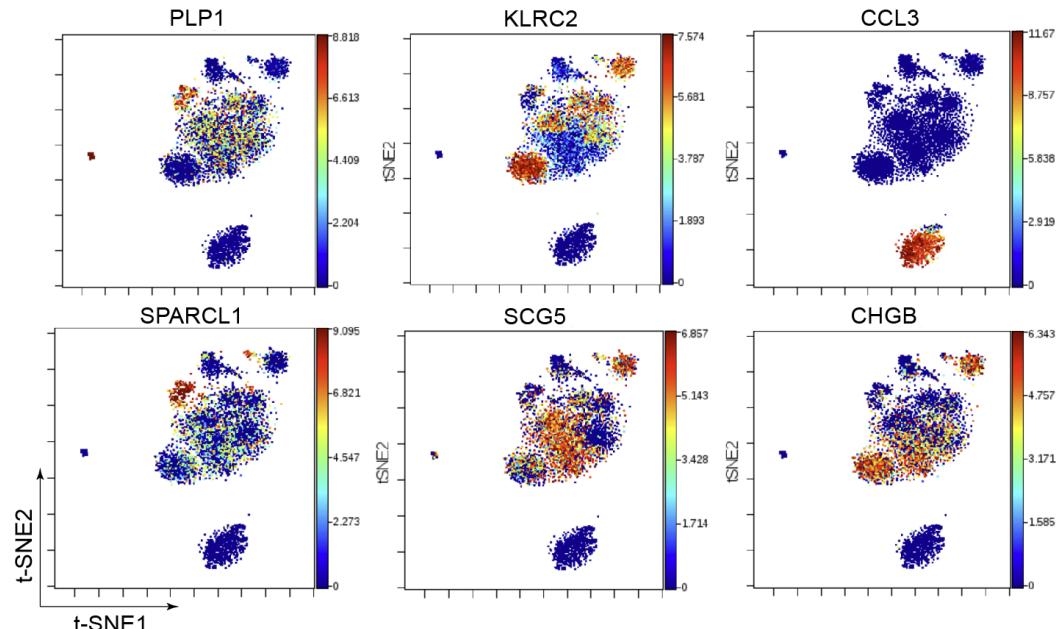
Per-cell transcript expression, top 6 most positively enriched in gated populations

MEM Label, Population 3:

▲ SAT1<sup>+6</sup> CCL3<sup>+5</sup> CCL4<sup>+5</sup>  
CD74<sup>+5</sup> HLA-DRA<sup>+5</sup>  
RGS1<sup>+5</sup> SPP1<sup>+5</sup>

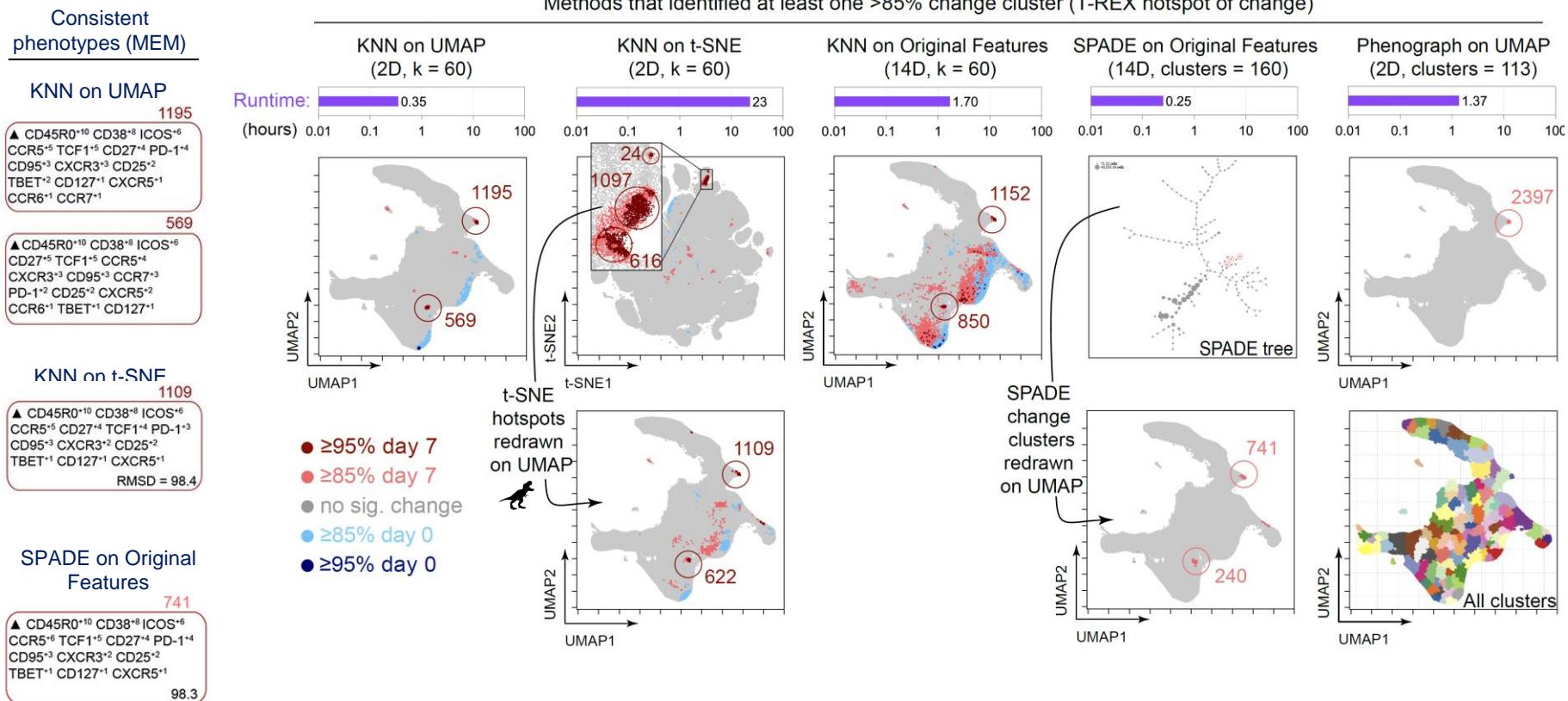
▼ GFAP<sup>-8</sup> OMG<sup>-8</sup> APOD<sup>-7</sup>  
KLRC2<sup>-7</sup> SCG5<sup>-7</sup> UCHL1<sup>-7</sup>

(ID'd as microglia)



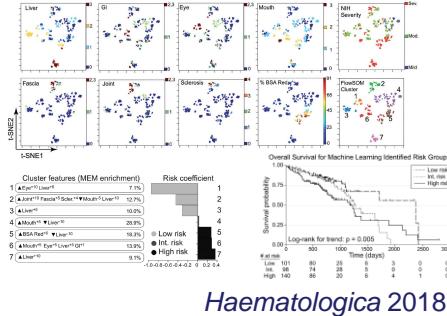
Data from leukemia patient blood, Venteicher et al., *Science* 2017  
MEM for scRNA-seq: Diggins et al., *in preparation*

# T-REX Worked with Other Algorithms to Identify Comparable Cells, But KNN on UMAP or t-SNE Outperformed KNN on Original Features



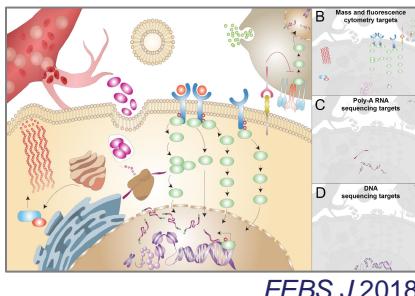
# Part 6: Example Data with cGVHD

# Adapting Advances in Machine Learning & Single Cell Biology for Immuno-Oncology



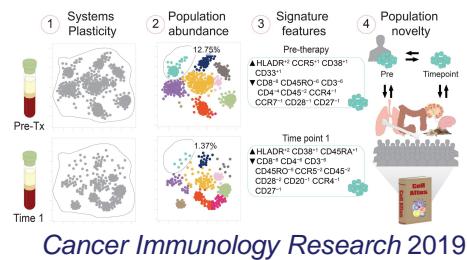
1. Gandelman et al., [Machine Learning Reveals Chronic Graft-Versus-Host Disease Phenotypes and Stratifies Survival After Stem Cell Transplant for Hematologic Malignancies](#). *Haematologica* 2018 PMC6312024.

Machine learning for patient medical phenotypes including t-SNE, FlowSOM, and MEM. Workflow inspired RAPID algorithm shown here for solid tumor cells.



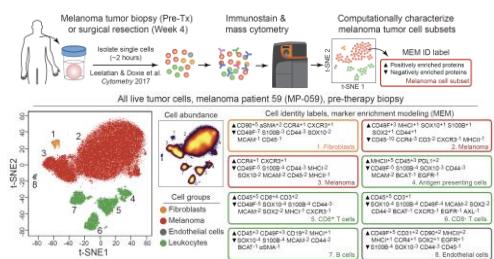
2. Mistry et al., [Beyond the Message: Advantages of Snapshot Proteomics with Single-Cell Mass Cytometry in Solid Tumors](#). *FEBS J* 2018 Dec 13. PMID: 30549207 DOI: 10.1111/febs.14730.

Reviews single cell mass cytometry with an emphasis on solid tumor and tissue research. Provides a detailed comparison with complementary single cell technologies, including single cell RNA-sequencing and quantitative multiplex imaging.



3. Greenplate et al., [Computational immune monitoring reveals abnormal double negative T cells present across human tumor types](#). *Cancer Immunology Research* 2019 Jan;7(1):86-99. PMC6318034.

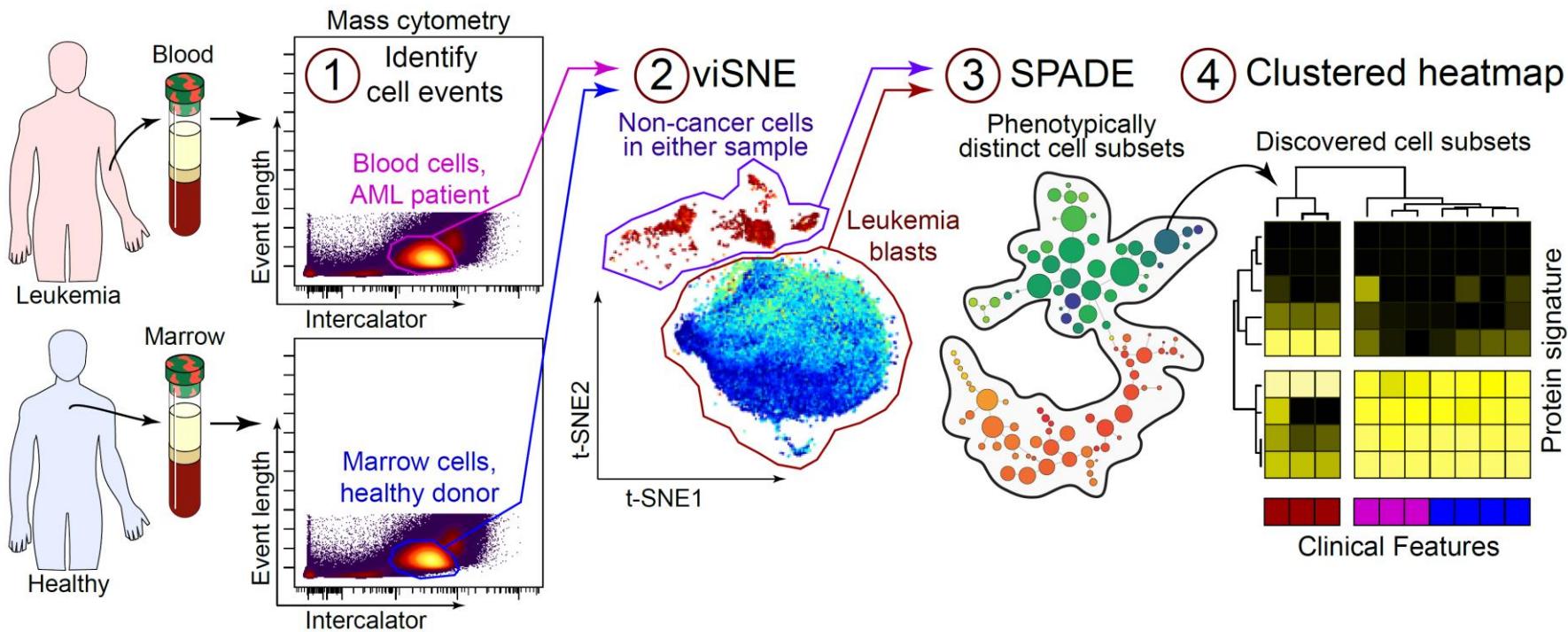
New computational workflow for longitudinal single cell tumor immunology. Revealed abnormal immune cells present in multiple tumor types. Will be applied to solid tumors to track changes in immune microenvironment.



4. Doxie et al., [BRAF and MEK inhibitor therapy eliminates Nestin-expressing melanoma cells in human tumors](#). *Pigment Cell Melanoma Research*. 2018 Jun 28. PubMed PMID: 29778085; PubMed Central PMCID: PMC6188784.

Compared biopsies over time of melanoma patients on treatment with mass cytometry to reveal *in vivo* cellular changes on treatment following targeted therapy.

# Machine Learning Tools Can Automate Sub-Population (Cluster / Subtype) Characterization



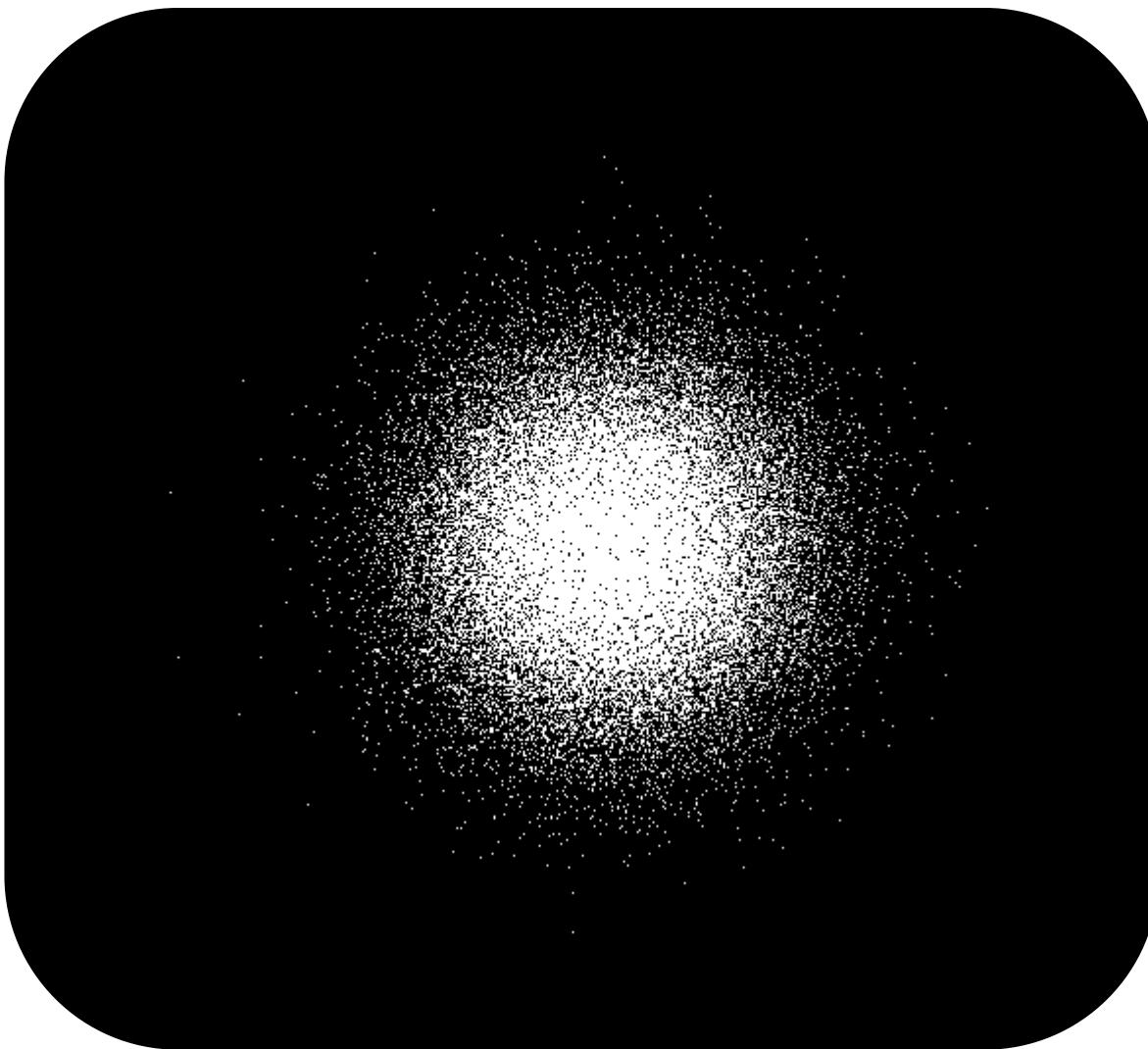
Protocol: Diggins et al., *Current Protocols in Cytometry* 2017  
Data files: <https://flowrepository.org/id/FR-FCM-ZZKZ>  
More great tools: Saeys et al., *Nature Reviews Immunology* 2016

Workflow: Diggins et al., *Methods* 2015  
viSNE/t-SNE: Amir et al., *Nat Biotech* 2013  
SPADE: Qiu et al., *Science* 2011

For sub-populations: quantify the known, reveal the unexpected, characterize the abnormal, & evaluate any associated risk

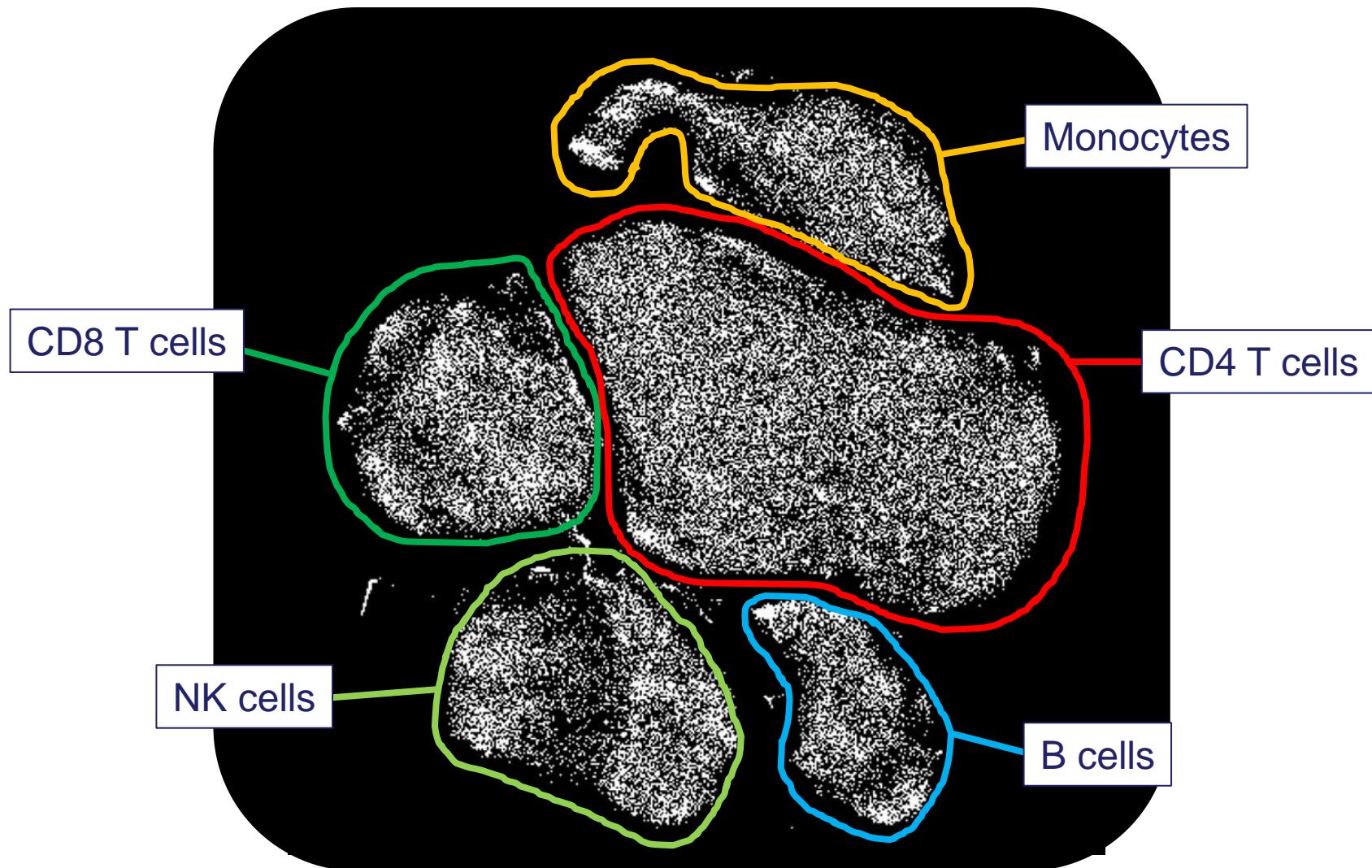
# Key Algorithm Concept: t-SNE Visualizes (in 2D) Close Phenotypic Relationships Using Many Markers

---



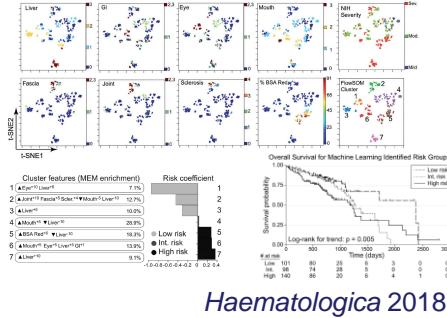
Healthy human blood, mass cytometry, 26 markers measured, viSNE / t-SNE analysis tool  
(Animation by Cytobank: each movie frame shows one iteration of t-SNE algorithm)

# Key Algorithm Concept: t-SNE Visualizes (in 2D) Close Phenotypic Relationships Using Many Markers



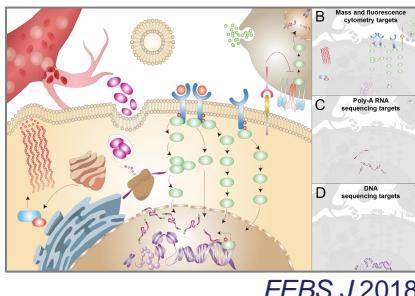
Healthy human blood, mass cytometry, 26 markers measured, viSNE / t-SNE analysis tool  
(Animation by Cytobank: each movie frame shows one iteration of t-SNE algorithm)

# Adapting Advances in Machine Learning & Single Cell Biology for Immuno-Oncology



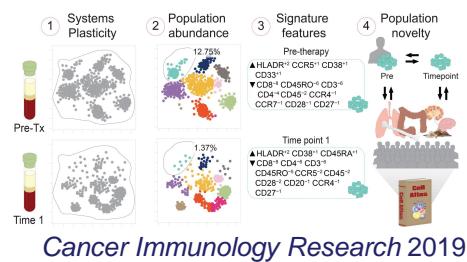
1. Gandelman et al., **Machine Learning Reveals Chronic Graft-Versus-Host Disease Phenotypes and Stratifies Survival After Stem Cell Transplant for Hematologic Malignancies.** *Haematologica* 2018 PMC6312024.

Machine learning for patient medical phenotypes including t-SNE, FlowSOM, and MEM. Workflow inspired RAPID algorithm shown here for solid tumor cells.



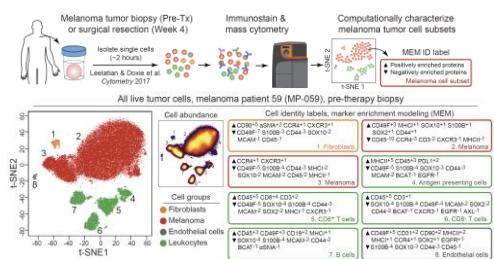
2. Mistry et al., **Beyond the Message: Advantages of Snapshot Proteomics with Single-Cell Mass Cytometry in Solid Tumors.** *FEBS J* 2018 Dec 13. PMID: 30549207 DOI: 10.1111/febs.14730.

Reviews single cell mass cytometry with an emphasis on solid tumor and tissue research. Provides a detailed comparison with complementary single cell technologies, including single cell RNA-sequencing and quantitative multiplex imaging.



3. Greenplate et al., **Computational immune monitoring reveals abnormal double negative T cells present across human tumor types.** *Cancer Immunology Research* 2019 Jan;7(1):86-99. PMC6318034.

New computational workflow for longitudinal single cell tumor immunology. Revealed abnormal immune cells present in multiple tumor types. Will be applied to solid tumors to track changes in immune microenvironment.



4. Doxie et al., **BRAF and MEK inhibitor therapy eliminates Nestin-expressing melanoma cells in human tumors.** *Pigment Cell Melanoma Research*. 2018 Jun 28. PubMed PMID: 29778085; PubMed Central PMCID: PMC6188784.

Compared biopsies over time of melanoma patients on treatment with mass cytometry to reveal *in vivo* cellular changes on treatment following targeted therapy.

# cGVHD Arises from Immune Dysregulation, Is Clinically Heterogeneous, and Is Scored “Multidimensionally”



## Areas of Potential Organ Involvement (“Dimensions”):

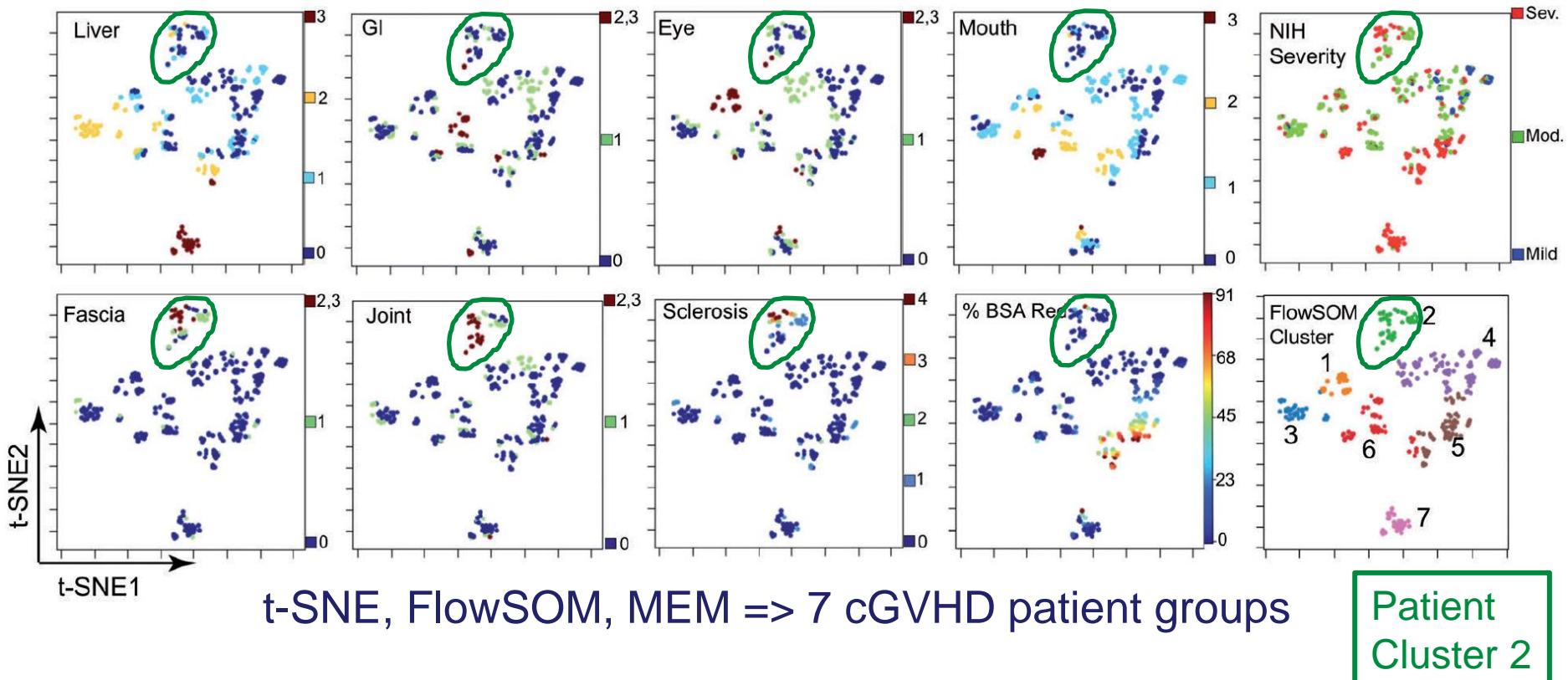
1. Liver
2. Gastrointestinal (GI)
3. Eye
4. Mouth
5. Fascia
6. Joint
7. Skin sclerosis
8. Skin redness
9. Lung

Each is scored (e.g. 0 to 3)

Sum of scores = NIH severity

# Data Science Strategies from Single Cell Analysis Reveal cGVHD Patient Groups Using Medical Phenotypes

Dots = 339 cGVHD patients  
t-SNE = 8 dimensional medical phenotype  
FlowSOM = identifies 7 cGVHD patient groups



Chronic graft-versus-host disease (cGVHD)  
t-SNE on 8 organ domain scores (e.g., liver involvement from 0 to 3)

Gandelman et al., *Haematologica* 2018

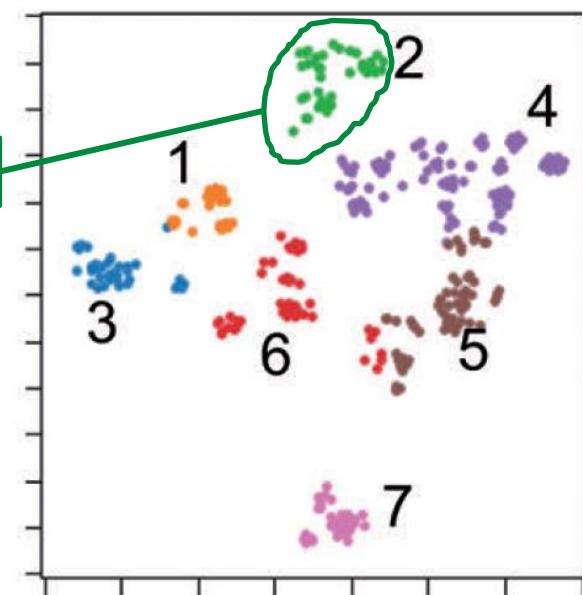
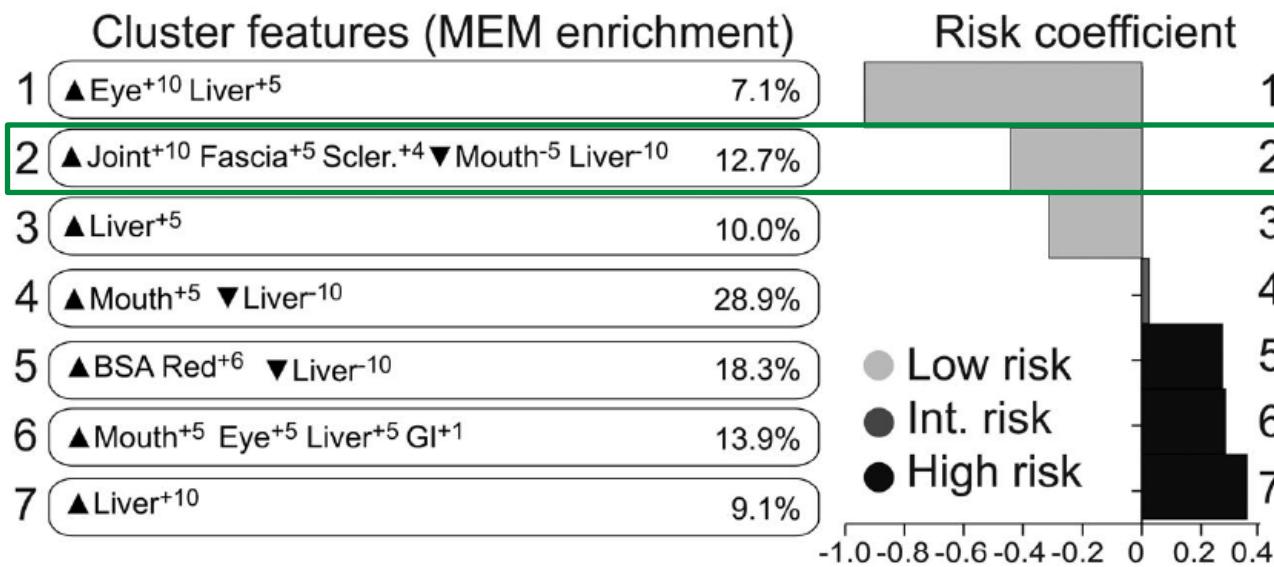
# Data Science Strategies from Single Cell Analysis Reveal cGVHD Patient Groups Using Medical Phenotypes

Dots = 339 cGVHD patients

t-SNE = 8 dimensional medical phenotype

FlowSOM = identifies 7 cGVHD patient groups

Machine learned cGVHD patient groups (clusters)



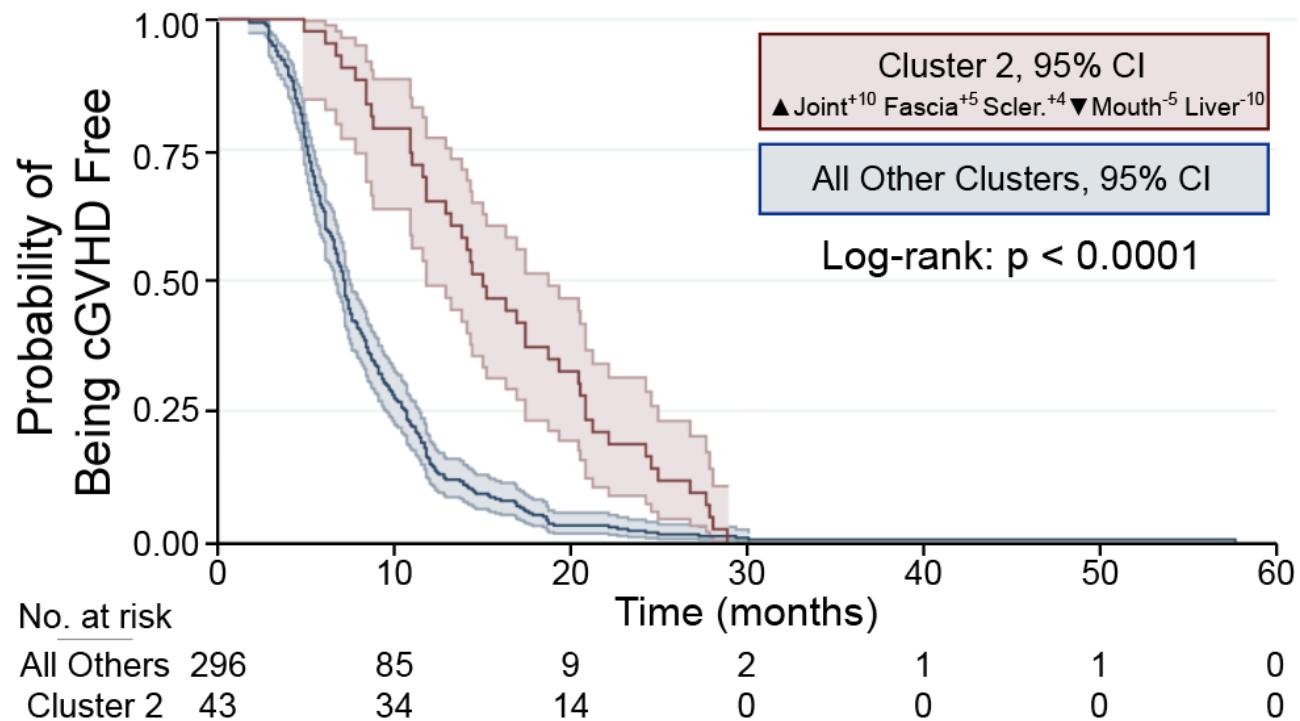
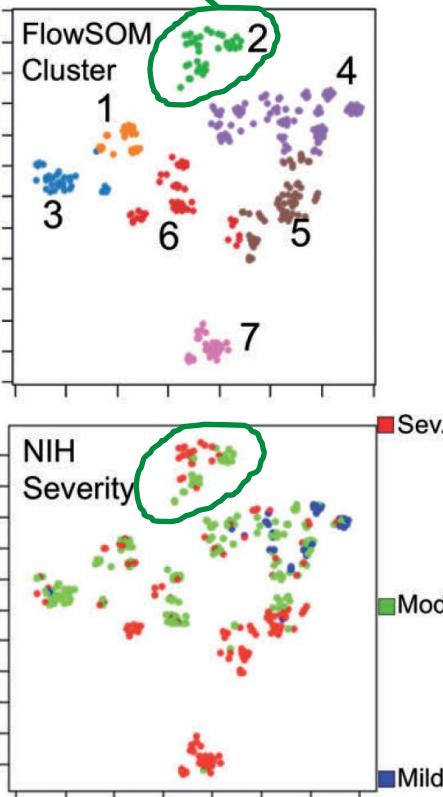
Chronic graft-versus-host disease (cGVHD)

t-SNE on 8 organ domain scores (e.g., liver involvement from 0 to 3)

Gandelman et al., *Haematologica* 2018

# Cluster 2 Was a Distinct Subtype of cGVHD Patients with Longer Time from Stem Cell Transplant to cGVHD

Patient Cluster 2



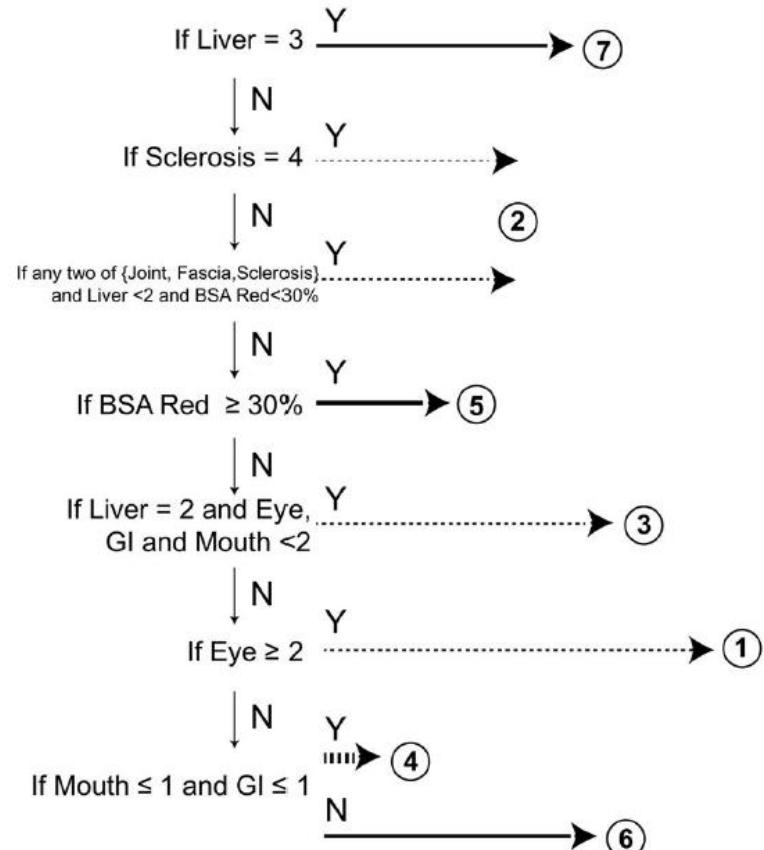
Phenotype of cGVHD Patient Cluster 2 (“MEM label”: +10 / -10)

2 ▲Joint<sup>+10</sup> Fascia<sup>+5</sup> Scler.<sup>+4</sup> ▼Mouth<sup>-5</sup> Liver<sup>-10</sup> 12.7%

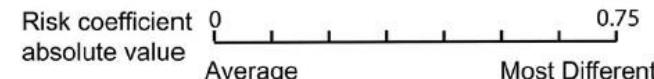
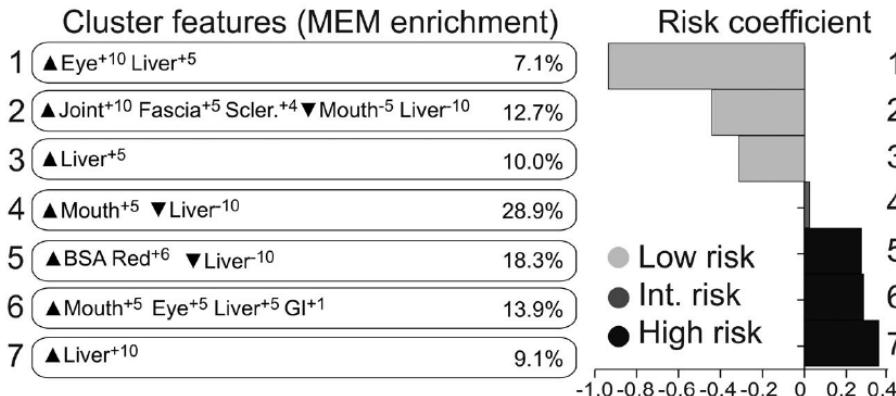
# Once Groups Are Revealed by Machine Learning, Clinically Accessible Strategies May Detect Them

7 Comparable Patient Groups  
Identified by Simple Yes/No Question Series

## DECISION TREE

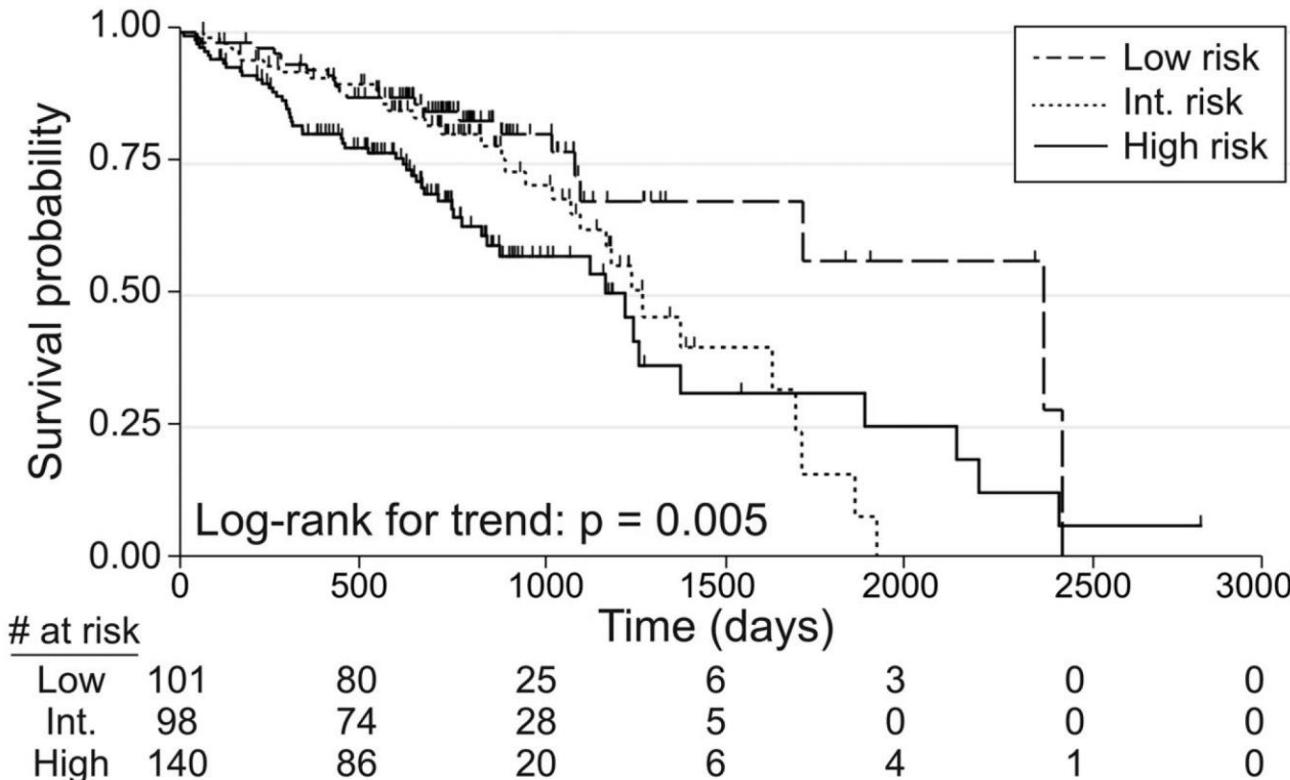


## 7 Machine-Identified cGVHD patient groups



# Decision Tree cGVHD Patient Groups Are Also Risk Stratified for Overall Survival

Overall Survival for Machine Learning Identified Risk Groups

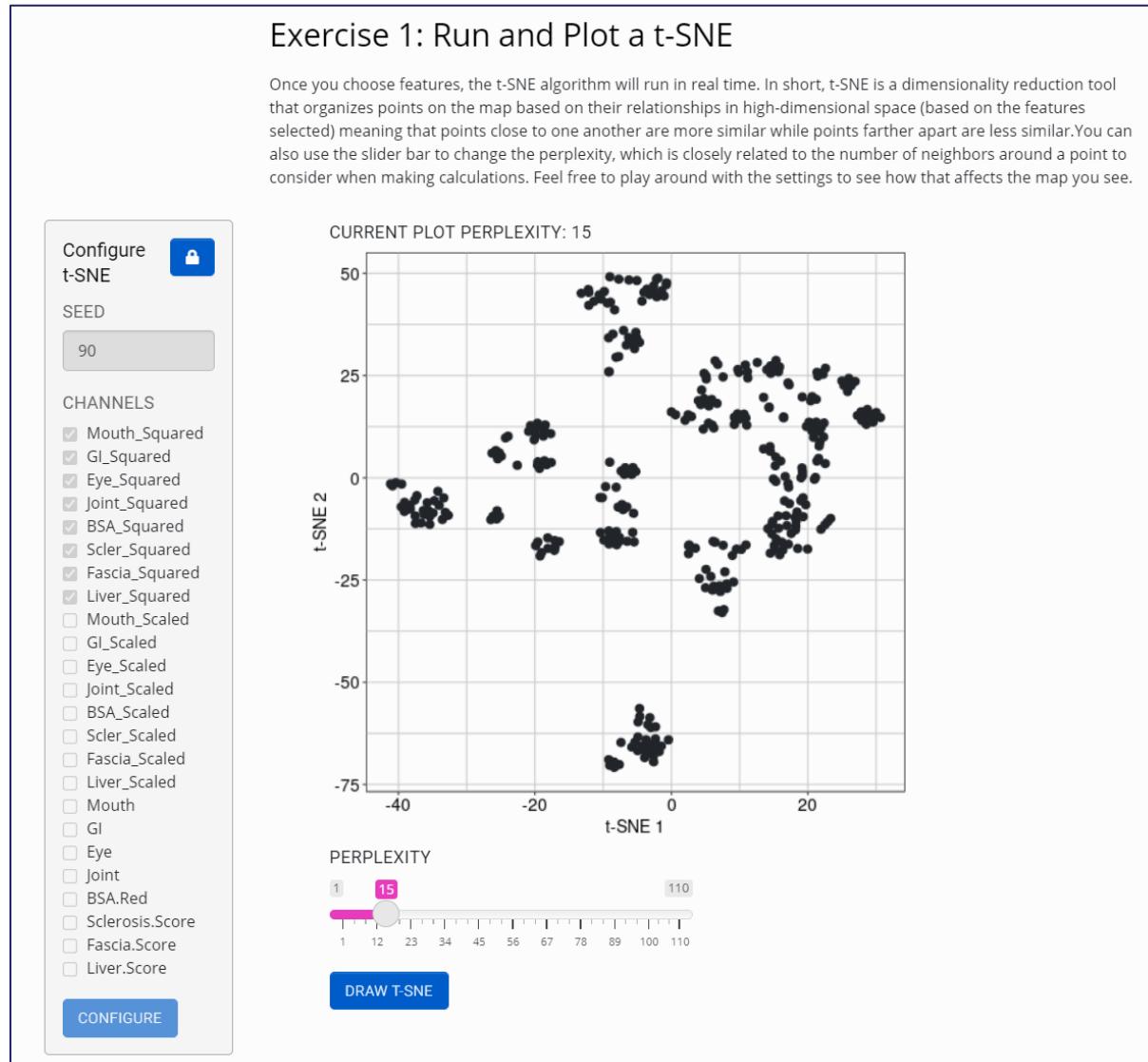


Low risk: HR = 2.79,  
95% CI: 1.58-4.91;  $p < 0.001$

High risk: HR = 2.65,  
95% CI: 1.42-4.94;  $p < 0.0001$

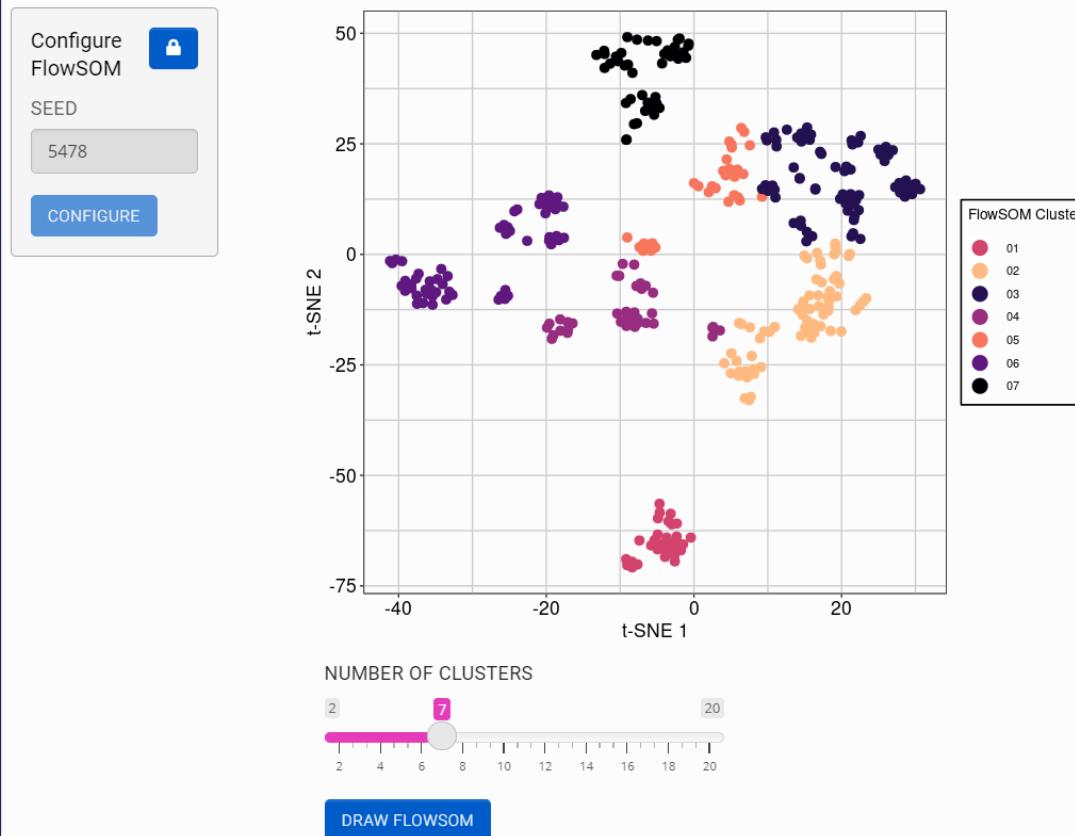
Decision tree cGVHD patient groups (clusters)

	Risk	Freq.	
1	-0.76	8.6%	Low
2	-0.34	10.0%	
3	-0.53	8.6%	
4	-0.06	33.6%	Int.
5	+0.21	17.4%	
6	+0.49	11.8%	
7	+0.46	10.0%	High



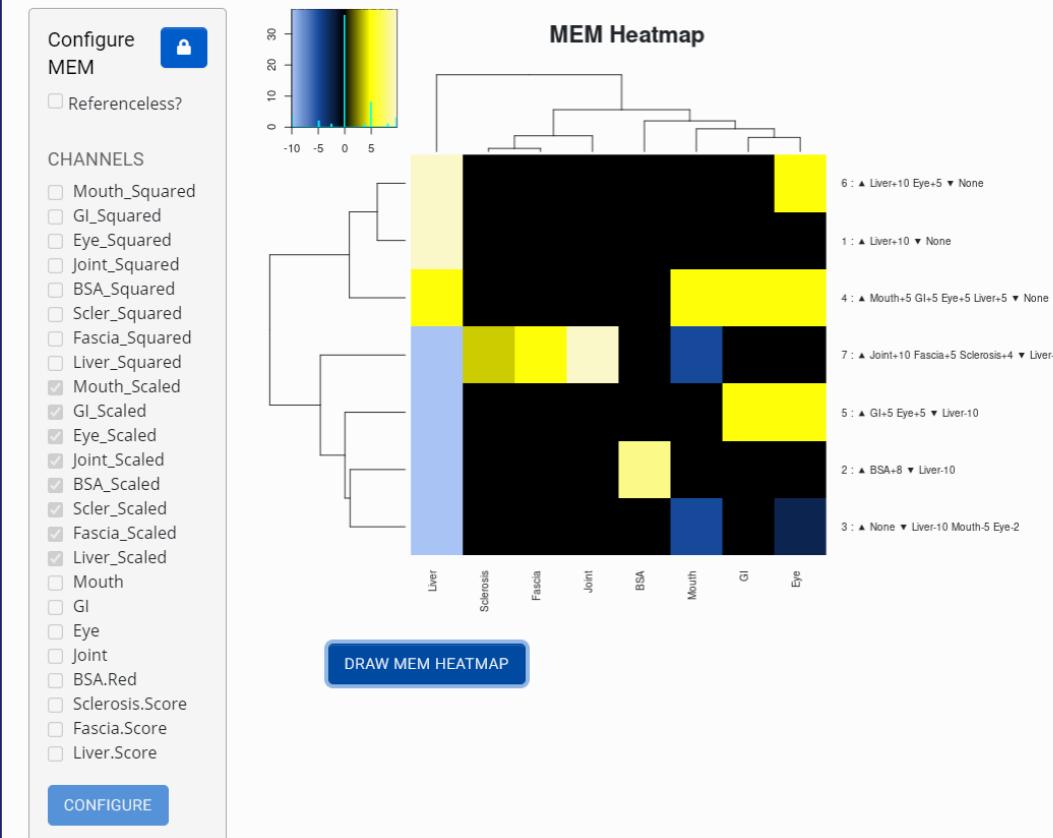
## Exercise 2: Run flowSOM on the t-SNE axes

Looking at the t-SNE plot, you may be able to visualize islands of dots (or patients in this case) and therefore identify some populations or groups yourself based on the location on the plot. However, that method is rather subjective and can come with certain biases, so we're going to go ahead and cluster groups together in an automatic, unsupervised way. The algorithm we will use is called flowSOM, a tool that uses self organizing maps to cluster and detect subsets. You as the user can choose the target number of clusters for flowSOM by changing the value using the slider.

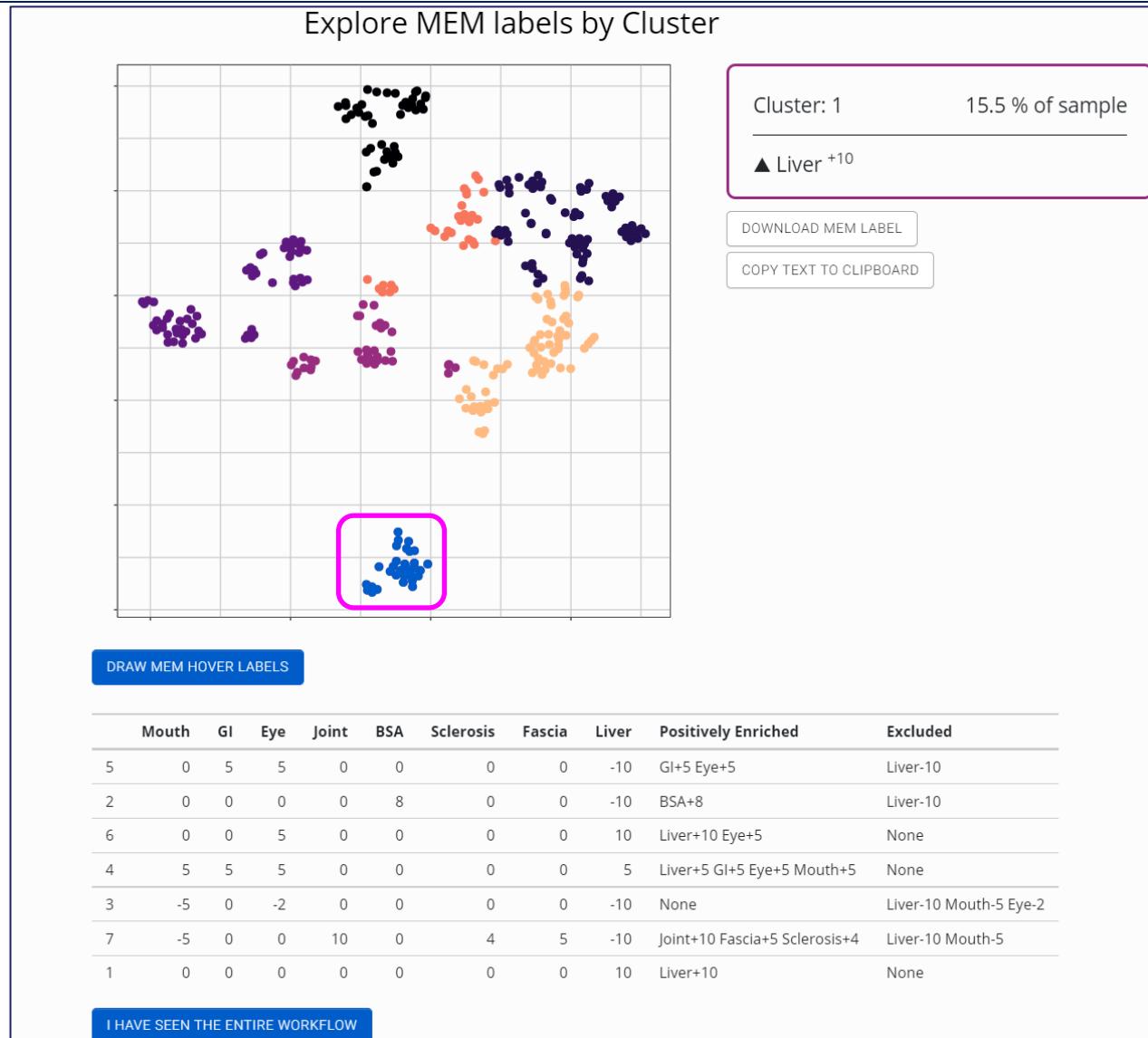


## Exercise 3: Run MEM on the flowSOM clusters

Next we need to characterize these flowSOM clusters. How are patients being grouped together? What clinical features are separating patients? There are a few ways to answer these questions. We could use median expression, but median expression does not necessarily reveal other differences samples might have, namely variance in feature expression. We can investigate differences like this using marker enrichment modeling (MEM) which takes into account the variance of a given feature and its median expression. Scores are linearly transformed on a scale from 0 to 10. A positive score (ex. +5) identifies a feature that is specifically expressed on a population while a negative score (ex. -3) identifies a feature that is specifically lacking from a population.

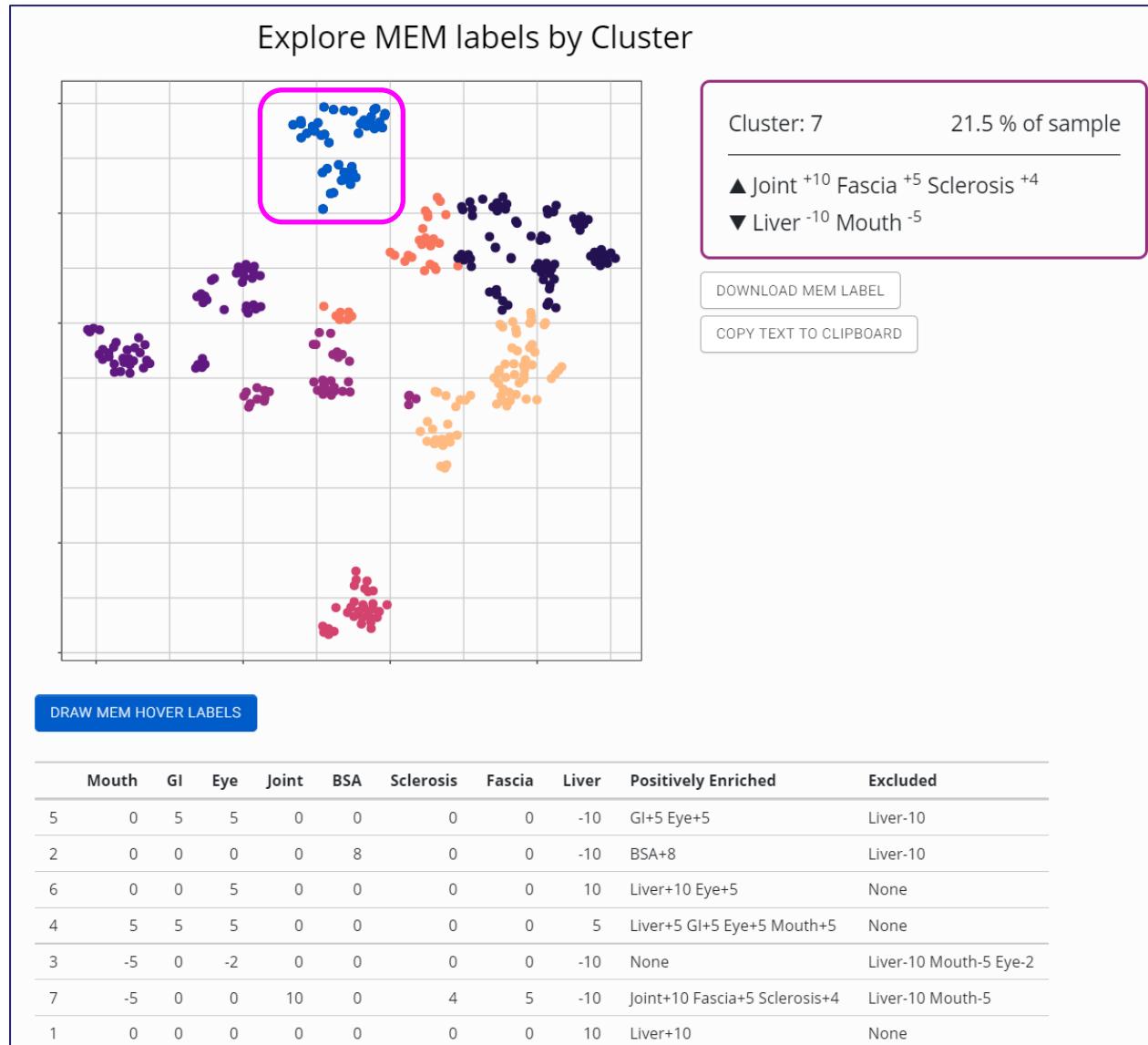


<https://cytolab.shinyapps.io/cGVHD/>



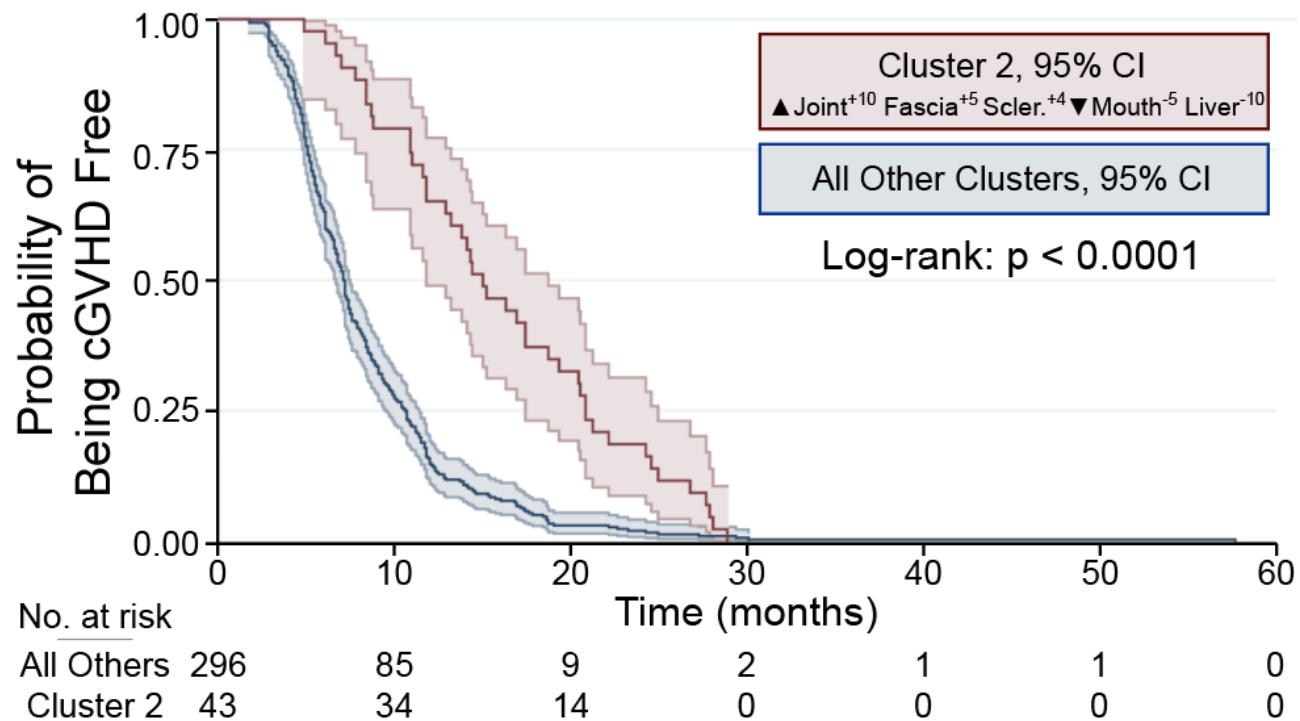
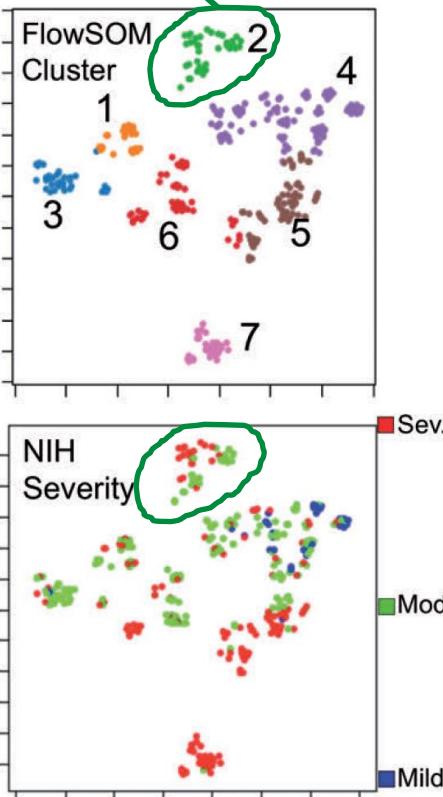
cGVHD web app by Mayeda et al.  
based on Gandelman et al., *Haematologica* 2018

<https://cytolab.shinyapps.io/cGVHD/>



# Cluster 2 Was a Distinct Subtype of cGVHD Patients with Longer Time from Stem Cell Transplant to cGVHD

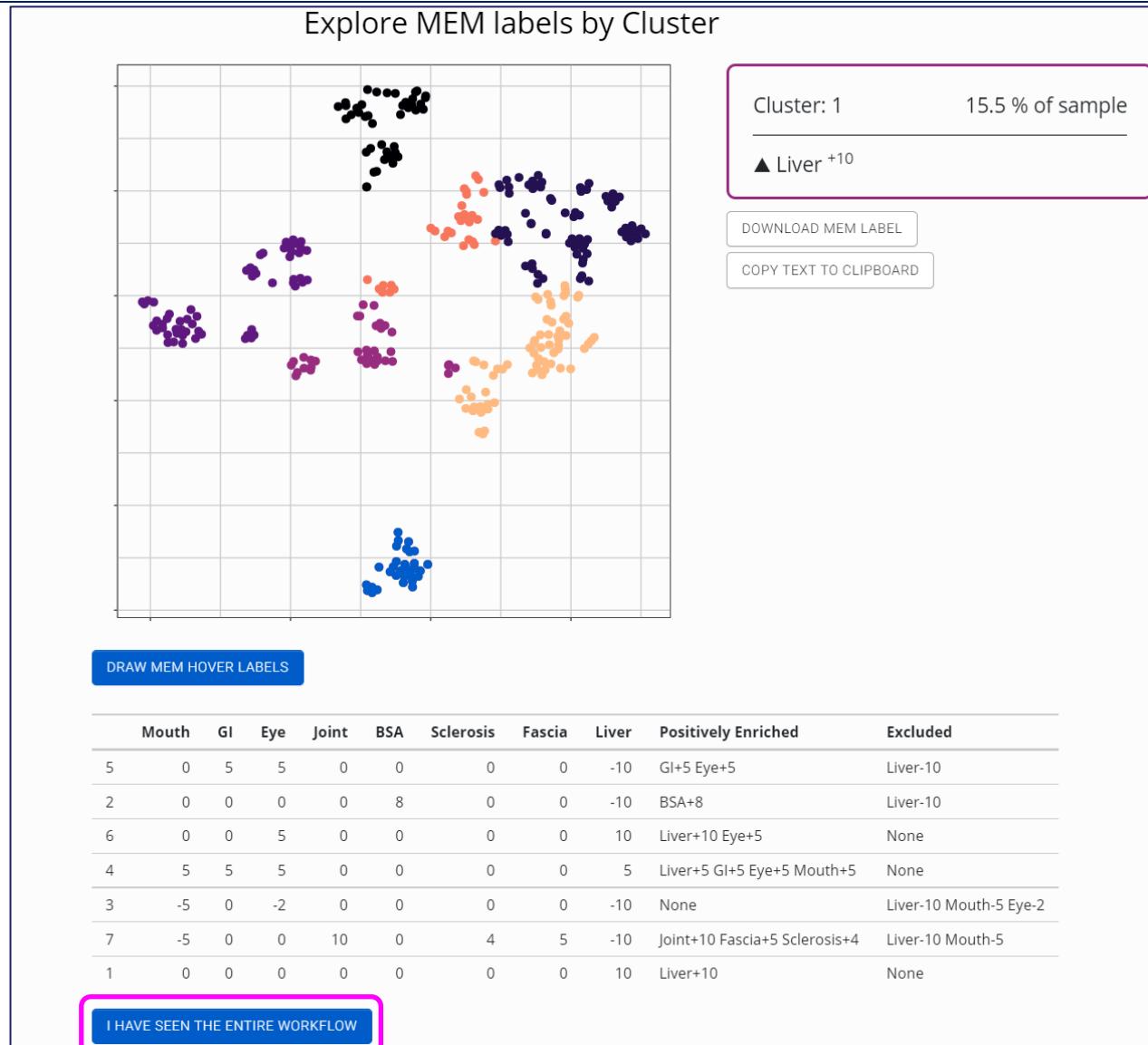
Patient Cluster 2

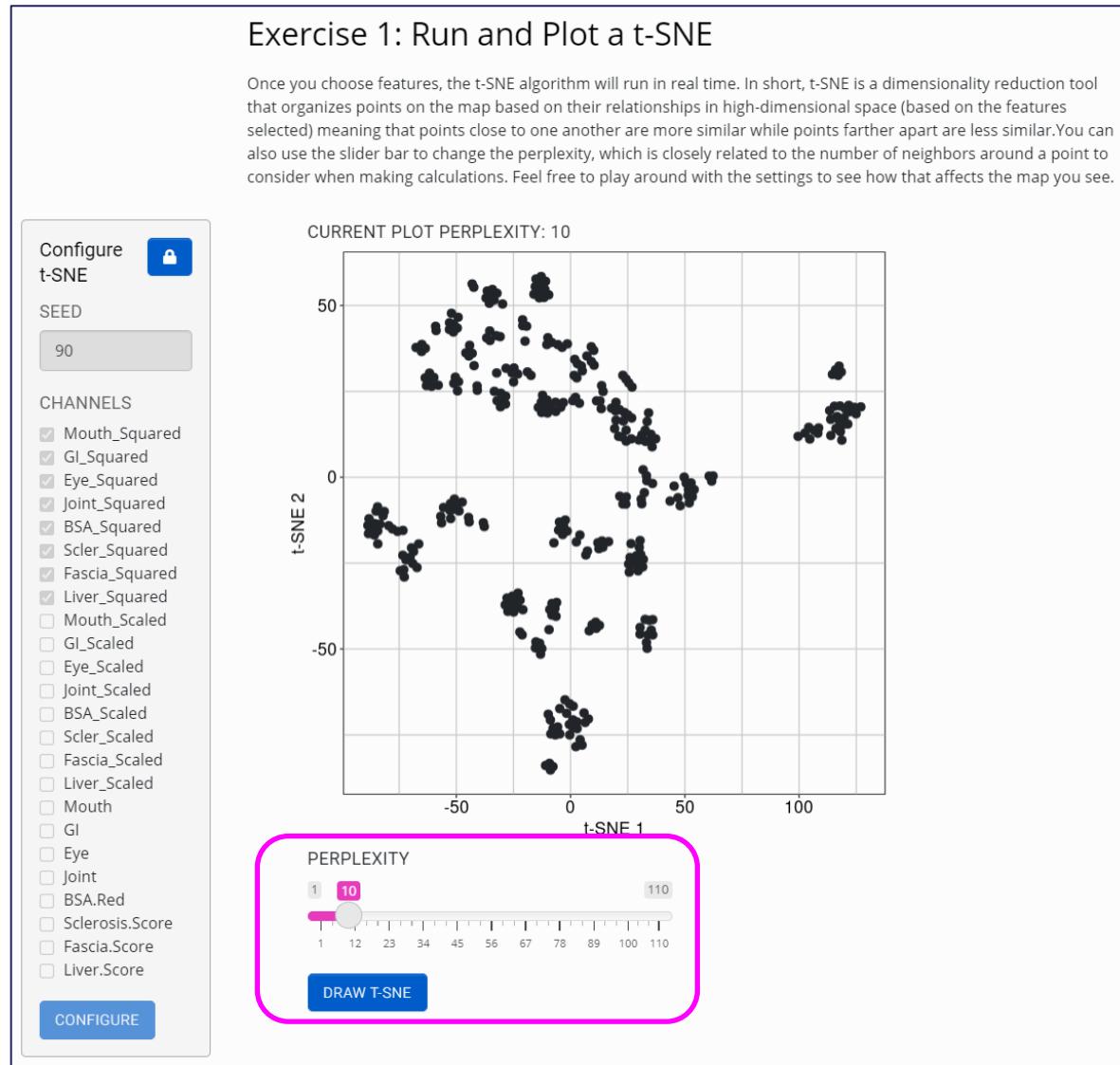


Phenotype of cGVHD Patient Cluster 2 (“MEM label”: +10 / -10)

2 ▲Joint<sup>+10</sup> Fascia<sup>+5</sup> Scler.<sup>+4</sup> ▼Mouth<sup>-5</sup> Liver<sup>-10</sup> 12.7%

<https://cytolab.shinyapps.io/cGVHD/>





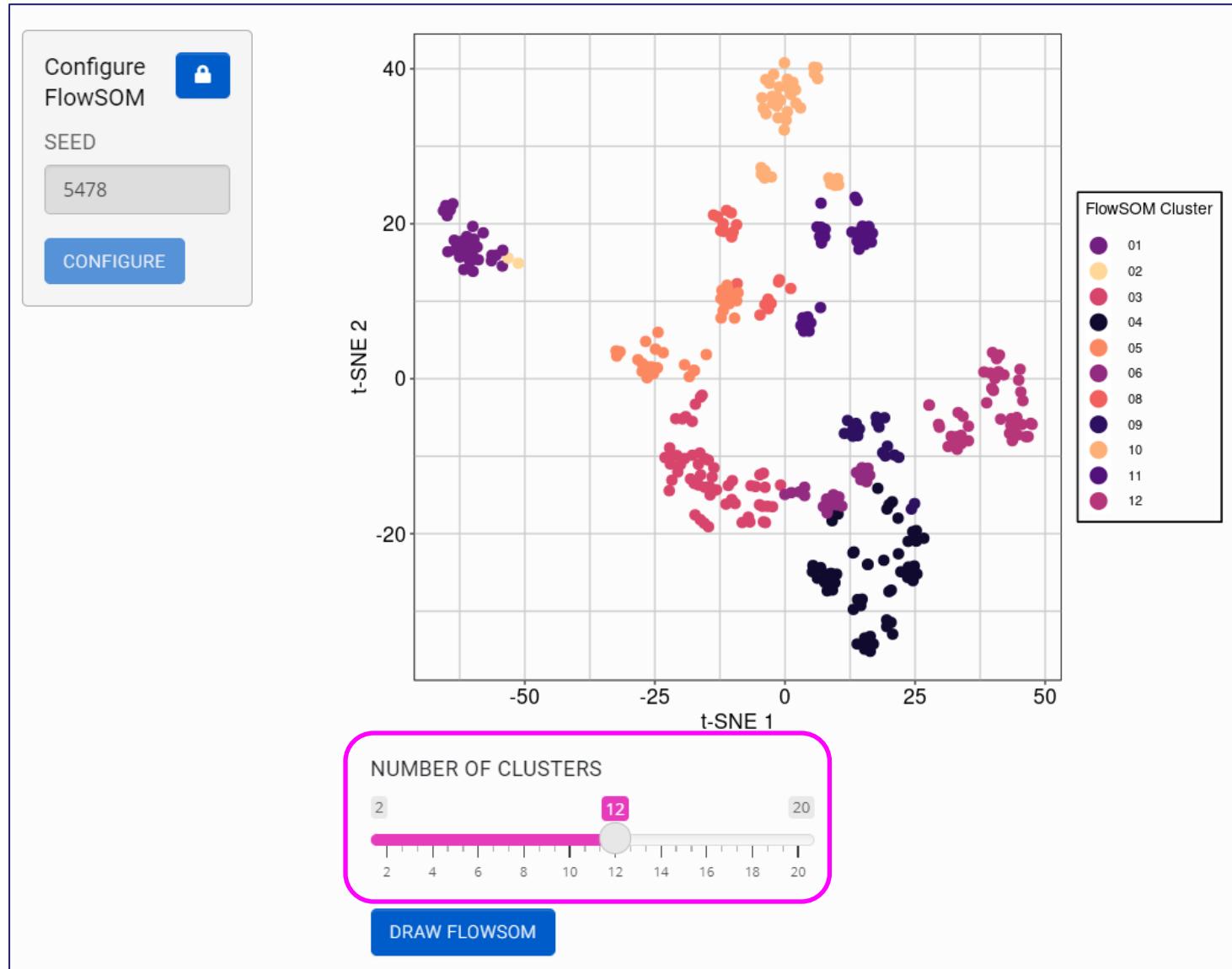
<https://cytolab.shinyapps.io/cGVHD/>



<https://cytolab.shinyapps.io/cGVHD/>

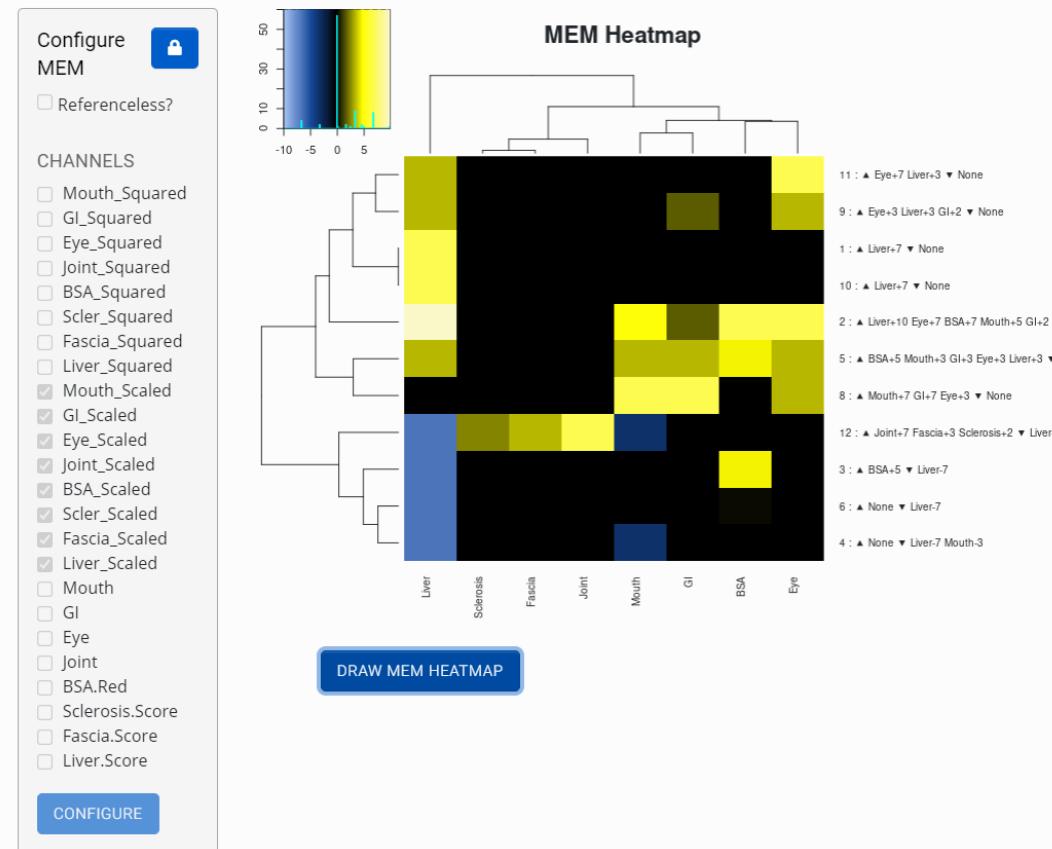


<https://cytolab.shinyapps.io/cGVHD/>

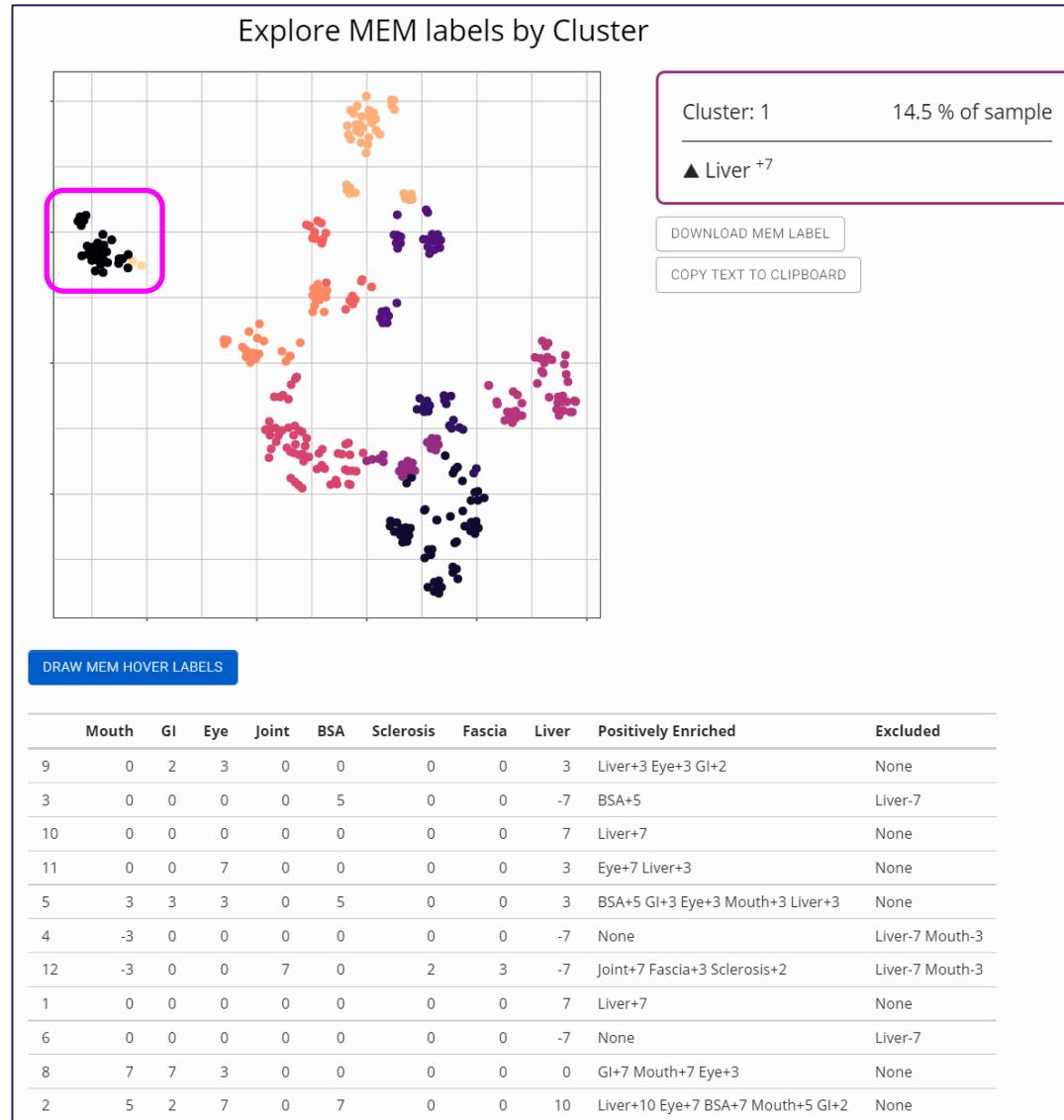


## Exercise 3: Run MEM on the flowSOM clusters

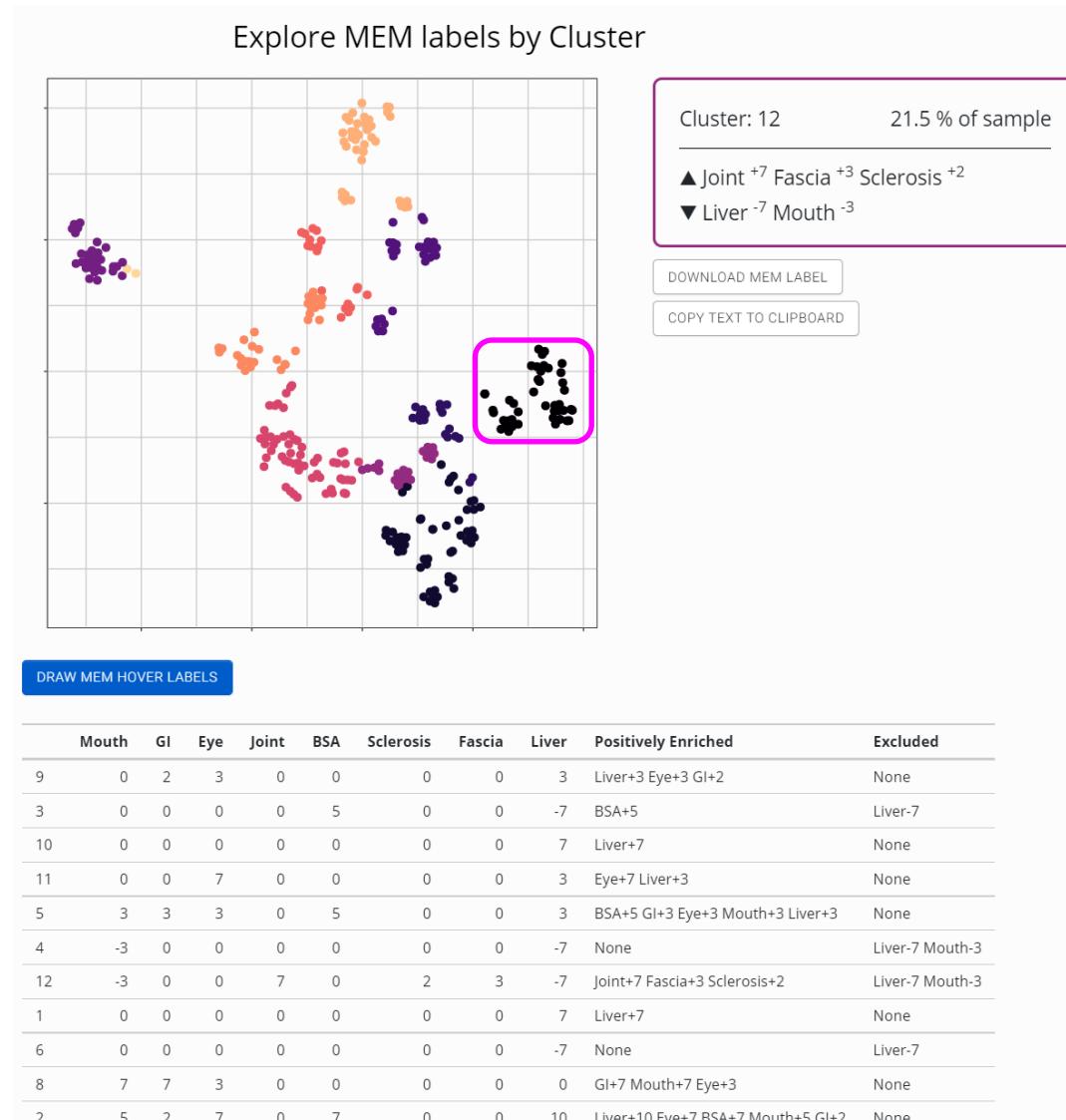
Next we need to characterize these flowSOM clusters. How are patients being grouped together? What clinical features are separating patients? There are a few ways to answer these questions. We could use median expression, but median expression does not necessarily reveal other differences samples might have, namely variance in feature expression. We can investigate differences like this using marker enrichment modeling (MEM) which takes into account the variance of a given feature and its median expression. Scores are linearly transformed on a scale from 0 to 10. A positive score (ex. +5) identifies a feature that is specifically expressed on a population while a negative score (ex. -3) identifies a feature that is specifically lacking from a population.



<https://cytolab.shinyapps.io/cGVHD/>



<https://cytolab.shinyapps.io/cGVHD/>



# Part A: Installing Tools in R

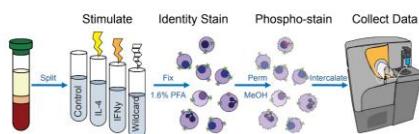
<https://github.com/cytolab>

Irish lab examples, scripts, & walkthroughs  
(MEM, RAPID, T-REX, & single cell courses)

# Single Cell Biology Workflow

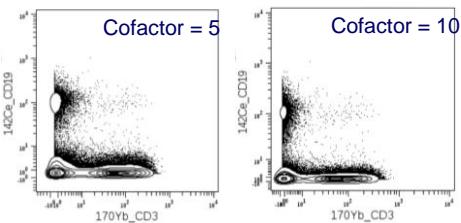
## Data collection

- 1) Panel design
- 2) Data collection



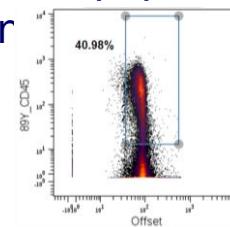
## Data processing

- 3) Cell event parsing
- 4) Scale transformation



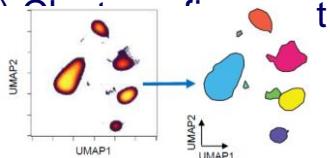
## Distinguishing initial populations

- 5) Live single cell gating
- 6) Focal population gating



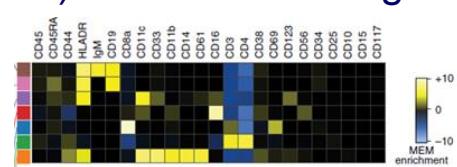
## Revealing cell subsets

- 7) Feature selection
- 8) Dimensionality reduction
- 9) Identify cell clusters
- 10) Cluster analysis



## Characterizing cell subsets

- 11) Feature comparison
- 12) Model populations
- 13) Learn cell identity
- 14) Statistical testing



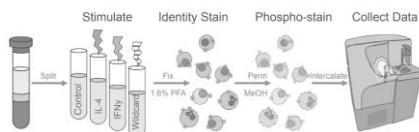
How much  
can be  
automated?

How do we  
select tools

# Single Cell Biology Workflow

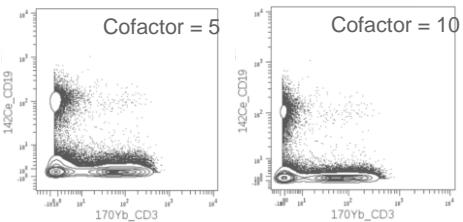
## Data collection

- 1) Panel design
- 2) Data collection



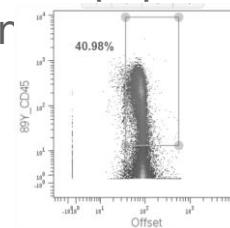
## Data processing

- 3) Cell event parsing
- 4) Scale transformation



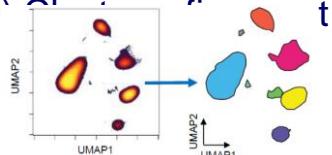
## Distinguishing initial populations

- 5) Live single cell gating
- 6) Focal population gating



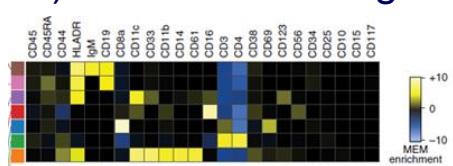
## Revealing cell subsets

- 7) Feature selection
- 8) Dimensionality reduction
- 9) Identify cell clusters
- 10) Cluster analysis



## Characterizing cell subsets

- 11) Feature comparison
- 12) Model populations
- 13) Learn cell identity
- 14) Statistical testing



How much  
can be  
automated?

How do we  
select tools



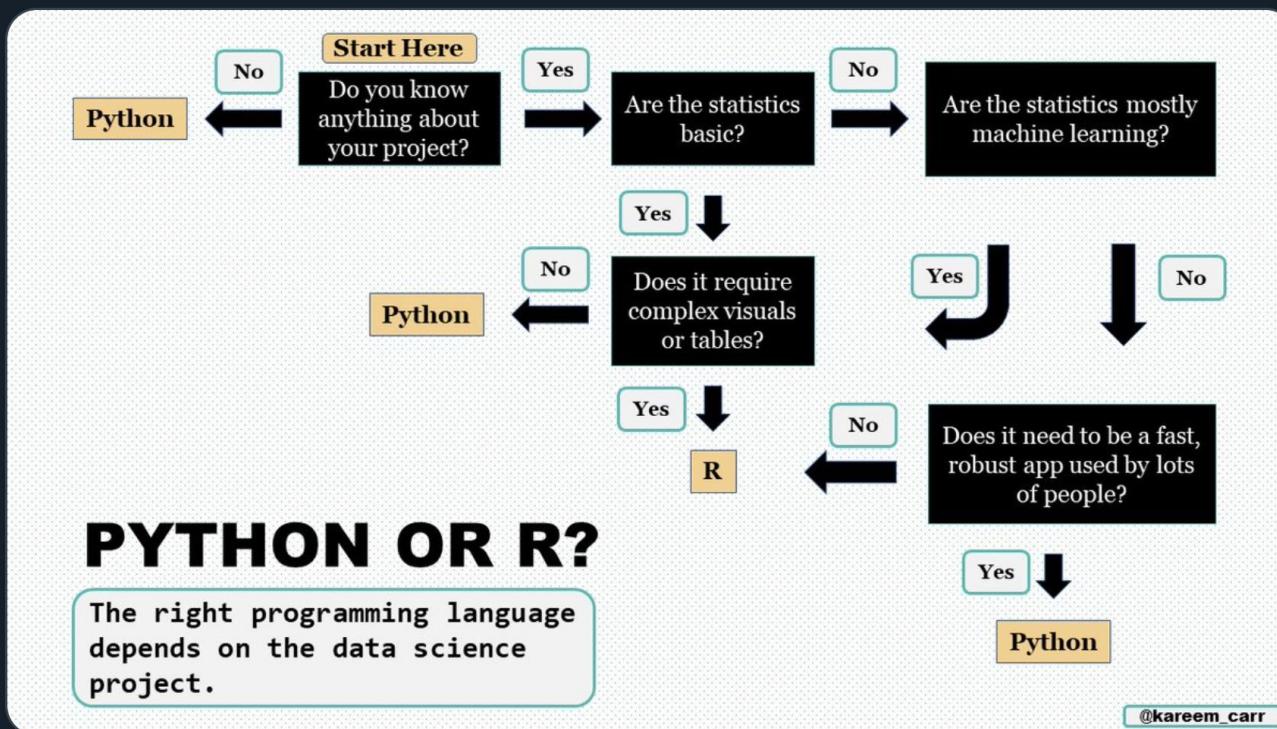
Tweet

You Retweeted



Kareem🔥data science thirst trap🔥Carr  
@kareem\_carr

People often ask me whether to learn R or Python for data science so I made a flowchart of my thought process.



# Installation and Intro to Working in R

# Download scripts, data, and R packages from GitHub

1

Go to link below and download repository:

[https://github.com/JonathanIrish/datasci\\_london](https://github.com/JonathanIrish/datasci_london)

No description, website, or topics provided.

28 commits 1 branch 0 releases 1 contributor View license

Branch: master New pull request Create new file Upload files Find File Clone or download

sierrabarone condensed code and updated plots

R initial commit 10 days ago  
data initial commit 10 days ago  
datafiles removed output files folder 9 days ago  
figures reworked examples 10 days ago  
man initial commit 10 days ago

2

Clone with HTTPS ?  
Use Git or checkout with SVN using the web URL.  
<https://github.com/cytolab/irish-data-science>

Open in Desktop

Download ZIP

3

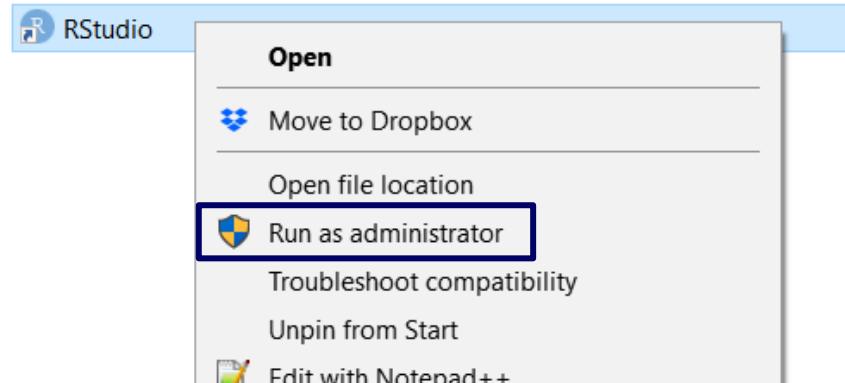
# Github repository contents

- 0) installation script (R markdown, .rmd)
- 1) example analysis scripts (.rmd)
- 2) Data files (.fcs)
- 3) MEM package (.r, .rproj, etc.)
- 4) Documentation files (.rd, .md)
- 5) Other misc. files (.txt, .pdf, .rdata, etc.)

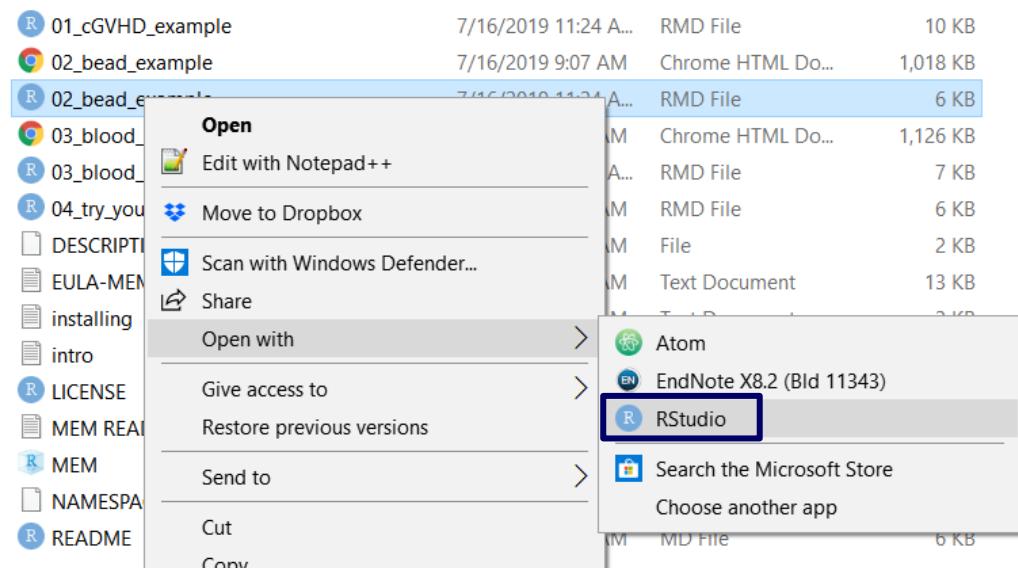
\*make sure you **unzip** downloaded folder

# We Will Open RStudio (Not R) For Interacting with Scripts

For PC users, run  
RStudio as  
administrator



For all, open .Rmd files  
with RStudio



\*make sure you **unzip** downloaded folder

## Scripts:

00\_install\_tools.rmd (run only once)

01\_cGVHD\_example.Rmd

02\_bead\_example.Rmd

03\_blood\_cell\_PBMC\_example.html

04\_try\_your\_data.Rmd

05\_scRNA-seq\_example.Rmd

06\_duncan\_tsne\_analysis.Rmd

\*make sure **console** is open

# Working Script and Code

# Environment

The screenshot shows the RStudio interface with the following components:

- Script Editor:** Displays the R script `01_PBMC_workflow_example.Rmd`. The code includes sections for setup, loading libraries, reading FCS files, and combining data.
- Environment View:** Shows the global environment with objects like `combined.data`, `data`, and `transformed.chos...`.
- Packages View:** Shows the user library with various Bioconductor packages installed, such as `Biobase`, `BiocGenerics`, and `BiocManager`.
- Console:** Displays the R session history, including commands for reading files, combining data, and transforming it.

```
24 ````{r setup, include=FALSE}
25 # Time <10 sec
26
27
28 # Load all libraries
29 # If you get an error message, you will need to try re-installing packages by
30 # going back to the 00_install_tools.RMD script
31 library(FlowsOM)
32 library(flowCore)
33 library(Biobase)
34 library(ggplot2)
35 library(hexbin)
36 library(MEM)
37 library(tidyverse)
38 library(Rtsne)
39 library(uwot)
40 library(viridis)
41 library(ggExtra)
42
43 ````{r data_preparation, warning=FALSE}
44 # Time <10 sec
45
46
47 # read files into R by setting working directory and directing R to the fcs files
48 setwd(paste(getwd(), "/datafiles/PBMC", sep = ""))
49 files <- dir(pattern = "*.fcs")
50
51 # convert and combine data for use in downstream analysis
52 data <- lapply(lapply(files, read.FCS), exprs)
53 combined.data = as.data.frame(do.call(rbind, data))
54
```

```
1:1 Data Analysis Workflow Example on PBMC Data (t-SNE, UMAP, FlowSOM, MEM) R Markdown
Console Terminal Jobs
C:/Users/Sierra/Desktop/irish-data-science/ 
> files <- dir(pattern = "*.fcs")
>
> # convert and combine data for use in downstream analysis
> data <- lapply(lapply(files, read.FCS), exprs)
> combined.data = as.data.frame(do.call(rbind, data))
>
> # choose channels with markers to use for downstream analysis and apply arcsinh
> # transformation with a cofactor of 15
> transformed.chosen.markers <- combined.data %>%
+   select(contains("-"), !contains("Ir")) %>%
+   mutate_all(function(x)
+     asinh(x / 15)) # cofactor here is 15; this can be changed
>
> # set seed for reproducible results (43 is chosen below)
> overall_seed = 43
>
```

Console

Plots, Files, Help, etc.

# Working Script and Code

# Environment

In this window, you can see the prepared script. Any text following # is a comment that is not part of the code, but can help explain what different lines of code are doing. The rest of the text is the actual code.

The screenshot shows the RStudio interface with three main panes:

- Script Pane:** Displays the R script `01_PBMC_workflow_example.Rmd`. The code includes comments explaining the purpose of various library imports and data processing steps.
- Environment Pane:** Shows a list of installed packages and their versions. A red box highlights the package `Biobase`, which is checked.
- Console Pane:** Displays the R command-line history, showing the execution of the script and its resulting data manipulation commands.

Package	Description	Version
ape	Analyses of Phylogenetics and Evolution	5.3
askpass	Safe Password Entry for R, Git, and SSH	1.1
assertthat	Easy Pre and Post Assertions	0.2.1
backports	Reimplementations of Functions Introduced Since R-3.0.0	1.1.4
base64enc	Tools for base64 encoding	0.1-3
BH	Boost C++ Header Files	1.69.0-1
bibtex	Bibtex Parser	0.4.2
<b>Biobase</b>	Biobase: Base functions for Bioconductor	2.44.0
BiocGenerics	S4 generic functions used in Bioconductor	0.30.0
BiocInstaller	Install/Update Bioconductor, CRAN, and github Packages	1.30.0
BiocManager	Access the Bioconductor Project Package Repository	1.30.4
BiocParallel	Bioconductor facilities for parallel evaluation	1.18.0
BiocVersion	Set the appropriate version of Bioconductor packages	3.9.0
biocViews	Categorized views of R package repositories	1.52.0
bit	A Class for Vectors of 1-Bit Booleans	1.1-14
bit64	A S3 Class for Vectors of 64bit Integers	0.9-7
bitops	Bitwise Operations	1.0-6
bmp	Read Windows Bitmap (BMP) Images	0.3

Console

Plots, Files, Help, etc.

# Working Script and Code

# Environment

The screenshot shows the RStudio interface. On the left, the 'Code' tab displays an R Markdown script titled '01\_PBMC\_workflow\_example.Rmd'. The code includes sections for setup, loading libraries (FlowsOM, flowCore, Biobase, ggrepel, hexbin, MEM, tidyverse, Rtsne, uwot, viridis, ggExtra), data preparation, and reading FCS files. The 'Console' tab at the bottom shows the execution of the script, including the creation of 'combined.data' and 'transformed.chosen.markers' data frames. To the right, the 'Environment' tab lists variables like 'combined.data' and 'transformed.chosen.markers' with their respective types and values. The 'Packages' tab shows a list of installed packages such as acepack, ape, askpass, assertthat, backports, base64enc, and BH.

In this window, you can see the code running. Errors and warnings will display here. You can type in the console without changing the base code above.

Console

Plots, Files, Help, etc.

\*make sure console is open

# Working Script and Code

# Environment

The environment contains the data you have loaded into R as well as any variables you have defined.

Code in the script pane:

```
24 25 ````{r setup, include=FALSE}
26 # Time <10 sec
27
28 # Load all libraries
29 # If you get an error message, you will need to try re-installing packages by
30 # going back to the 00_install_tools.RMD script
31 library(FlowsOM)
32 library(flowCore)
33 library(Biobase)
34 library(ggplot2)
35 library(hexbin)
36 library(MEM)
37 library(tidyverse)
38 library(Rtsne)
39 library(uwot)
40 library(viridis)
41 library(ggExtra)
42
43 ````{r data_preparation, warning=FALSE}
44 # Time <10 sec
45
46 # read files into R by setting working direc
47 setwd(paste(getwd(), "/datafiles/PBMC", sep
48 files <- dir(pattern = "*.fcs")
49
50
51 # convert and combine data for use in downstream analysis
52 data <- lapply(lapply(files, read.FCS), exprs)
53 combined.data = as.data.frame(do.call(rbind, data))
54
```

Console output:

```
> files <- dir(pattern = "*.fcs")
>
> # convert and combine data for use in downstream analysis
> data <- lapply(lapply(files, read.FCS), exprs)
> combined.data = as.data.frame(do.call(rbind, data))
>
> # choose channels with markers to use for downstream analysis and apply arcsinh
> # transformation with a cofactor of 15
> transformed.chosen.markers <- combined.data %>%
+   select(contains("-"), !contains("Ir")) %>%
+   mutate_all(function(x)
+     asinh(x / 15)) # cofactor here is 15; this can be changed
>
> # set seed for reproducible results (43 is chosen below)
> overall_seed = 43
>
```

Environment pane:

combined.data	49651 obs. of 46 variables
data	List of 7
transformed.chosen.markers	49651 obs. of 25 variables
values	
files	chr [1:7] "CD4Tcells_PBMC.fcs" "CD8Tcells_PBMC.fcs" "Monocytes_PBMC.fcs" "Neutrophils_PBMC.fcs" "Plasmacytoid_Dendritic_Cells_PBMC.fcs" "Regulatory_T_Cells_PBMC.fcs" "T_H17_Cells_PBMC.fcs"
overall_seed	43

Packages pane:

Description	Version
ACE and AVAS for Selecting Multiple Regression Transformations	1.4.1
Analyses of Phylogenetics and Evolution	5.3
askpass	Safe Password Entry for R, Git, and SSH
assertthat	Easy Pre and Post Assertions
backports	Reimplementations of Functions Introduced Since R-3.0.0
base64enc	Tools for base64 encoding
BH	Boost C++ Header Files
bibtex	Bibtex Parser
Biobase	Biobase: Base functions for Bioconductor
BiocGenerics	S4 generic functions used in Bioconductor
BiocInstaller	Install/Update Bioconductor, CRAN, and github Packages
BiocManager	Access the Bioconductor Project Package Repository
BiocParallel	Bioconductor facilities for parallel evaluation
BiocVersion	Set the appropriate version of Bioconductor packages
biocViews	Categorized views of R package repositories
bit	A Class for Vectors of 1-Bit Booleans
bit64	A S3 Class for Vectors of 64bit Integers
bitops	Bitwise Operations
bmp	Read Windows Bitmap (BMP) Images

Console

Plots, Files, Help, etc.

# Working Script and Code

# Environment

The screenshot shows the RStudio interface. On the left, the script editor displays a script named '01\_PBMC\_workflow\_example.Rmd' with code for setting up packages and reading FCS files. Below it, the console window shows the execution of this code. On the right, the environment browser shows variables like 'combined.data' and 'transformed.chos...', and the package manager shows a list of installed packages such as 'pack', 'ace', 'avas', 'base', etc.

This window will display files in your working directory, plots you have created, as well as packages you have installed and loaded. You can also access help pages for each package in this window.

Name	Description	Version
pack	ACE and AVAS for Selecting Multiple Regression Transformations	1.4.1
ape	Analyses of Phylogenetics and Evolution	5.3
base	Safe Password Entry for R, Git, and SSH	1.1
assertthat	Easy Pre and Post Assertions	0.2.1
BiocGenerics	Reimplementations of Functions Introduced Since R-3.0.0	1.1.4
bit	Tools for base64 encoding	0.1-3
bit64	Boost C++ Header Files	1.69.0-1
bitops	Bibtex Parser	0.4.2
base	Biobase: Base functions for Bioconductor	2.44.0
BiocManager	S4 generic functions used in Bioconductor	0.30.0
BiocInstaller	Install/Update Bioconductor, CRAN, and github Packages	1.30.0
BiocParallel	Access the Bioconductor Project Package Repository	1.30.4
BiocVersion	Bioconductor facilities for parallel evaluation	1.18.0
biocViews	Set the appropriate version of Bioconductor packages	3.9.0
bit	Categorized views of R package repositories	1.52.0
bit64	A Class for Vectors of 1-Bit Booleans	1.1-14
bitops	A S3 Class for Vectors of 64bit Integers	0.9-7
bmp	Bitwise Operations	1.0-6
Read Windows Bitmap (BMP) Images	0.3	

Console

Plots, Files, Help, etc.

# Open 00\_install\_tools.rmd

1

```
1 ---  
2 title: "Check Paths and Install Packages"  
3 author: "Copyright (c) 2016-2019 by Kirsten Diggins, Sierra Barone, and  
4 Jonathan Irish, All Rights Reserved; see EULA-MEM.text for MEM license  
5 information"  
6 date: "July 2019"  
7 output: html_document  
8 ---  
9   
10 cat("This section checks to see if files and paths are working correctly. You  
11 should see lists of files below. If it outputs character(0), something is  
12 wrong.\n\n")  
13 # Check the MEM code path  
14 cat("\n\nThe /MEM folder contains the MEM source code for install and related  
15 files:\n")  
16 list.files(getwd())  
17 # Check for datasets  
18 cat("\n\nCourse FCS format files are in subdirecties of the /datafiles  
19 folder:\n")  
20 list.files(paste(getwd(), "/datafiles", sep=""))  
21   
22   
23 # This only works for PC users
```

Header

2

Code

3

# Open 00\_install\_tools.rmd and begin installing required packages

## Code Section Title

```
`{r check_paths echo=FALSE, results = "markup"}  
# Check to make sure FCS files, documentation, and MEM code are available  
cat("This section checks to see if files and paths are working correctly. You  
should see lists of files below. If it outputs character(0), something is  
wrong.\n\n")  
  
# Check the MEM code path  
cat("\n\nThe /MEM folder contains the MEM source code for install and related  
files:\n")  
list.files(getwd())  
  
# Check for datasets  
cat("\n\nCourse FCS format files are in subdirecties of the /datafiles  
folder:\n")  
list.files(paste(getwd(), "/datafiles", sep=""))  
  
```{r installation_notes, echo=FALSE, results = "markdown"}  
# Print the contents a help file that explains installing packages  
writeLines(readLines(paste(getwd(), "installing.txt", sep="/")))  
```
```

CNTL-ENTER or  
COMMAND-  
RETURN to run a  
single line of code

OR

Press play to run  
entire section of  
code

This section checks  
that the files we will  
need are accessible  
in our working  
directory

This section prints  
installation text

# Open 00\_install\_tools.rmd and begin installing required packages

```
```{r install_bioconductor_packages, echo=FALSE, results = "hide"}  
# install bioconductor and flow cytometry tools for R  
cat("If this works, you should see 4 sets of messages about downloading files  
that end in a message saying something like package 'BiocManager' successfully  
unpacked and MD5 sums checked. You should see this for BiocManager, Biobase,  
flowCore, and FlowSOM.\n\n")  
install.packages("BiocManager", repos = "http://cran.us.r-project.org")  
  
if (!requireNamespace("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("flowCore")  
BiocManager::install("FlowSOM")  
```
```

This section downloads Bioconductor and flow cytometry tools we will need

```
```{r test_flow_installs, echo=FALSE, results = "markdown"}  
# Load and test whether bioconductor and flow packages are installed  
cat("If this works, you may see Attaching Package messages or no message at  
all; that's good. If you get a warning, go back to the last CHUNK.\n\n")  
library(FlowSOM)  
library(flowCore)  
library(Biobase)  
```
```

This section tests to make sure Bioconductor and flow cytometry tools are installed

```
```{r install_ggplots, echo=FALSE, results = "markup"}  
# install plotting packages  
cat("If this works, you will see text about packages being downloaded.\n\n")  
install.packages("gplots", repos = "http://cran.us.r-project.org")  
install.packages("ggplot2", repos = "http://cran.us.r-project.org")  
install.packages("hexbin", repos = "http://cran.us.r-project.org")  
install.packages("viridis", repos = "http://cran.us.r-project.org")  
install.packages("ggExtra", repos = "http://cran.us.r-project.org")  
```
```

```
```{r load_gplots, echo=FALSE, results = "markup"}  
# Load and test whether gplots and ggplot2 packages are installed  
cat("If this works, you may see Attaching Package messages or no message at  
all; that's good. If you get a warning, go back to the last CHUNK.\n\n")  
library(gplots)  
library(ggplot2)  
library(hexbin)  
library(viridis)  
library(ggExtra)  
```
```

The next sections install and load the tools to make plots

# Open 00\_install\_tools.rmd and begin installing required packages

```
```{r install_MEM, echo=FALSE, results = "markup"}  
# install MEM, load it, and test if it is all set  
cat("If this works, you should see several lines about installing files, then  
DONE (MEM) near the end. The MEM help page will also open in the Help menu in  
RStudio.\n\n")  
  
# If you have previously installed MEM, you may get an error message. If this  
is the case, try restarting your RStudio session  
install.packages(getwd(), type="source", repos=NULL)  
library(MEM)  
?MEM  
  
# OR  
# install.packages("devtools", repos = "http://cran.us.r-project.org")  
# devtools::install_github("cytolab/mem")  
...```

```

This section installs and loads the marker enrichment modeling tool

```
```{r install_last_packages, echo=FALSE, results = "markup"}  
# install the last packages for UMAP, t-SNE and other tools  
print("You may see a bunch of messages, this is OK as long as they are not  
errors.\n\n")  
install.packages("tidyverse", repos = "http://cran.us.r-project.org")  
install.packages("Rtsne", repos = "http://cran.us.r-project.org")  
install.packages("uwot", repos = "http://cran.us.r-project.org")  
install.packages("RColorBrewer", repos = "http://cran.us.r-project.org")  
...```

```

```
```{r load_last_packages, echo=FALSE, results = "markup"}  
# Load and test the last libraries  
library(tidyverse)  
library(Rtsne)  
library(uwot)  
library(RColorBrewer)  
...```

```

These sections install and load the other tools we will use for analysis

## Part B: Example with Healthy Blood

Open **03\_blood\_cell\_PBMC\_example.html** and  
work through the example

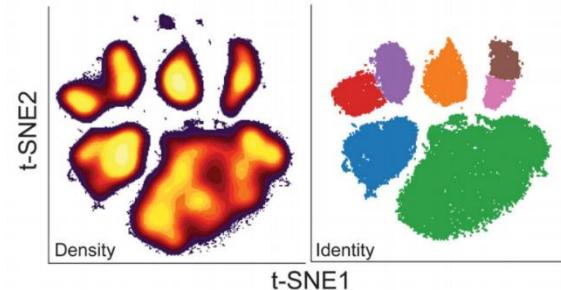
# 03\_blood\_cell\_PBMC\_example.html

```
01_PBMC_workflow_example.Rmd x
ABC Knit Insert Run
1 ---  
2 title: "Data Analysis Workflow Example on PBMC Data (t-SNE, UMAP, FlowsOM,  
MEM)"  
3 author: "Copyright (c) 2016-2019 by Kirsten Diggins, Sierra Barone, and  
Jonathan Irish, All Rights Reserved; see EULA-MEM.text for MEM license  
information"  
4 date: "July 2019"  
5 output:  
6   pdf_document:  
7     latex_engine: xelatex  
8   html_document:  
9     df_print: paged  
10  editor_options:  
11    chunk_output_type: inline  
12 ---  
13  
14 This data set contains 7 FCS (flow cytometry standard) files. Each FCS file  
15 contains single cell data for one cell subset that is a well-established,  
16 phenotypically distinct population. This is mass cytometry data for healthy  
17 human PBMC (peripheral blood mononuclear cells). The populations were expert  
18 gated following a t-SNE analysis. The first section of the code will run two  
19 dimensionality reduction tools, UMAP and t-SNE, on the data set. Next, you  
20 will run FlowsOM on the UMAP axes to cluster the islands of cell populations.  
21 Finally, you will run MEM to see enrichment scores for each of the FlowsOM  
22 clusters. The goal of this exercise is to run several computational tools on  
23 a single cell data set to get a feel for the workflow used in the Irish lab.  
24  
25 ```{r setup, include=FALSE}  
26 # Time <10 sec  
27  
28 # Load all libraries  
29 # If you get an error message, you will need to try re-installing packages by  
30 # going back to the 00_install_tools.RMD script  
31 library(FlowsOM)  
32 library(flowCore)  
33 library(BioBase)  
34 library(ggplot2)  
35 library(hexbin)  
36 library(MEM)  
37 library(tidyverse)  
38 library(Rtsne)  
39 library(uwot)  
40 library(viridis)  
41 library(ggExtra)  
42 ```

This code block is a R Markdown document titled '01_PBMC_workflow_example.Rmd'. It includes a header with metadata like title, author, date, and output type. The main content is a descriptive text about the PBMC dataset followed by R code. The R code section starts with ````{r setup, include=FALSE}```` and ends with ````{r}````. The R code itself loads various libraries including FlowsOM, flowCore, BioBase, ggplot2, hexbin, MEM, tidyverse, Rtsne, uwot, viridis, and ggExtra.
```

A description of the code and its purpose

a Identification of 7 canonical cell types in healthy human blood, 25D mass cytometry



This section loads the necessary libraries

# 03\_blood\_cell\_PBMC\_example.html

## Data Preparation

```
43
44 ~`{r data_preparation, warning=FALSE}
45 # Time <10 sec
46
47 # read files into R by setting working directory and directing R to
# files
48 setwd(paste(getwd(), "/datafiles/PBMC", sep = ""))
49 files <- dir(pattern = "*.fcs")
50
51 # convert and combine data for use in downstream analysis
52 data <- lapply(lapply(files, read.FCS), exprs)
53 combined.data = as.data.frame(do.call(rbind, data))
54
55 # choose channels with markers to use for downstream analysis and apply arcsinh
# transformation with a cofactor of 15
56 transformed.chosen.markers <- combined.data %>%
57   select(contains("-"), -contains("Ir")) %>%
58   mutate_all(function(x)
59     asinh(x / 15))      # cofactor here is 15; this can be changed
60
61 # set seed for reproducible results (43 is chosen below)
62 overall_seed = 43
63 ...
64
65
```

Read the data files into R and format for analysis

Select channels and scale the data

Choose parameters

# 03\_blood\_cell\_PBMC\_example.html

## Run t-SNE

```
66 ~ ````{r run_t-SNE}
67 # Time ~5 min
68
69 set.seed(overall_seed)
70
71 # the line below will run t-SNE on the scaled surface markers (to see help p...
72 # for t-SNE, type "?Rtsne -- enter" in console)
73 # you can view t-SNE progress by opening up the console below
74 mytsNE = Rtsne(
75   transformed.chosen.markers,                      # input scaled data
76   dims = 2,   # number of final
77   dimensions   # number of initial
78   initial_dims = length(transformed.chosen.markers), # number of initial
79   dimensions   # number of final
80   perplexity = 30,                                  # perplexity (similar to # of nearest neighbors,
81   # will scale with data sets, cannot be greater than
82   # the number of events minus 1 divided by 3)
83   check_duplicates = FALSE,
84   max_iter = 1000,                                   # number of iterations
85   verbose = TRUE
86 )
87 tsne.data = as.data.frame(mytsNE$Y)
88 ````{r plot_t-SNE}
89 # Time <10 sec
90
91 # setting aspect ratio for plots
92 range <- apply(apply(tsne.data, 2, range), 2, diff)
93 graphical.ratio <- (range[1] / range[2])
94
95 # t-SNE flat dot plot and density dot plot (1 dot = 1 cell)
96 tsne.plot <- data.frame(x = tsne.data[, 1], y = tsne.data[, 2])
```

This section will run a t-SNE analysis on the PBMC data with set parameters

You can choose the resulting numbers of dimensions, the perplexity, and the iterations

This section will plot the t-SNE results. Two plots will appear, a “flat” dot plot and a density plot

Go to t-SNE talk while this runs

# 03\_blood\_cell\_PBMC\_example.html

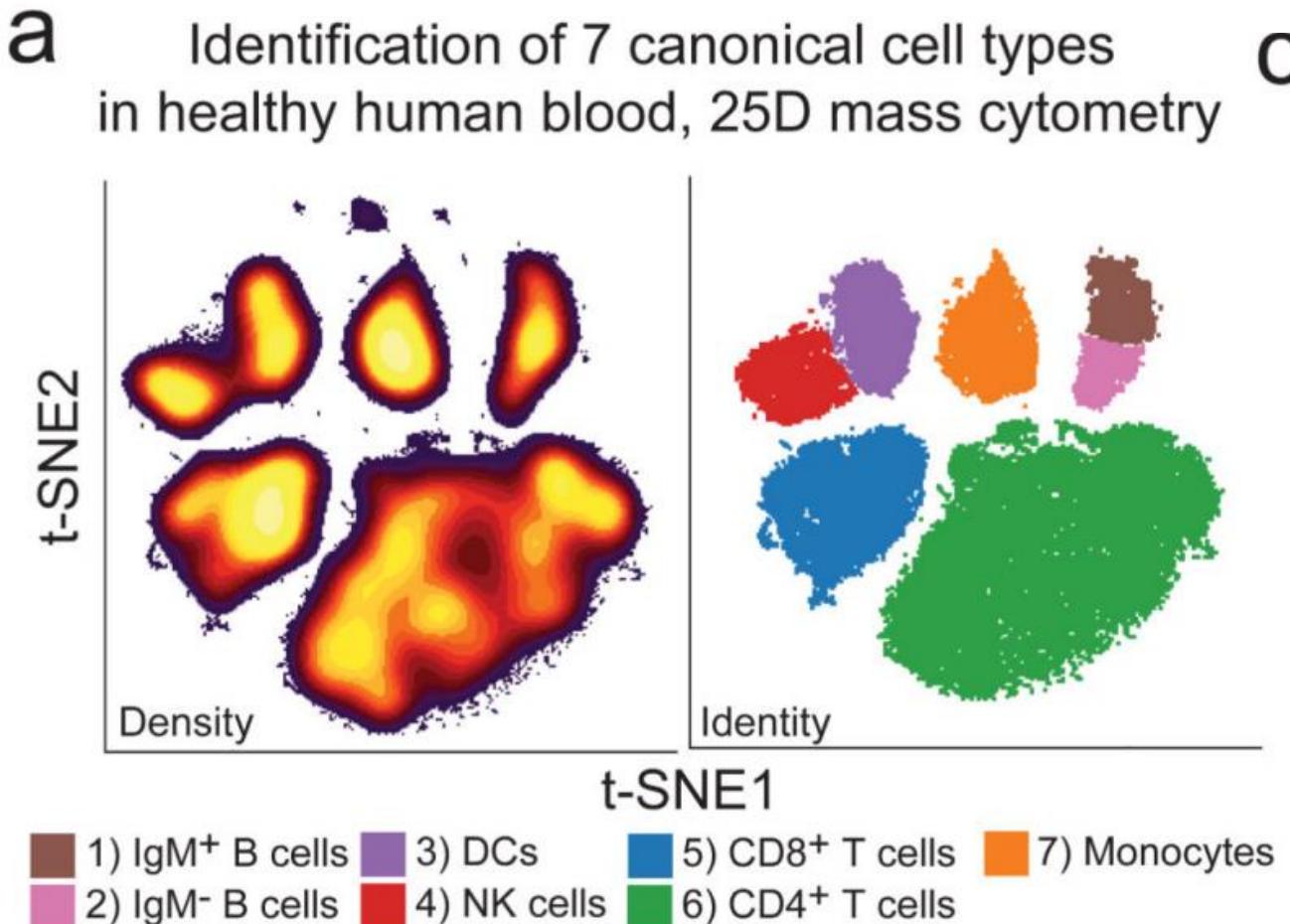
## Run UMAP

```
119 ````{r run_UMAP}
120 # Time ~1 min
121
122 # Run UMAP on all scaled surface markers
123 set.seed(overall_seed)
124
125 # the line below will run UMAP on the data set (to see help page for UMAP
# "UMAP -- enter" in console)
126 # you can view UMAP progress by opening up the console below
127 myumap <-
128   umap(transformed.chosen.markers, # input scaled data
129         n_neighbors = 15,           # number of nearest neighbors to look at,
scales with data set|
130         n_threads = 1,            # this argument makes UMAP reproducible
131         verbose = TRUE)
132 umap.data = as.data.frame(myumap)
133
134
135 ````{r plot_UMAP}
136 # Time <10 sec
137
138 # setting aspect ratio for plots
139 range <- apply(apply(umap.data, 2, range), 2, diff)
140 graphical.ratio <- (range[1] / range[2])
141
142 # UMAP flat dot plot and density dot plot (1 dot = 1 cell)
143 UMAP.plot <- data.frame(x = umap.data[, 1], y = umap.data[, 2])
144
145 # dot plot
146 ggplot(UMAP.plot) + coord_fixed(ratio = graphical.ratio) +
147   geom_point(aes(x = x, y = y), cex = 0.5) + labs(x = "UMAP 1", y = "UMAP 2",
title = "PBMC Data on UMAP
Axe
148
149 theme_bw() +
```

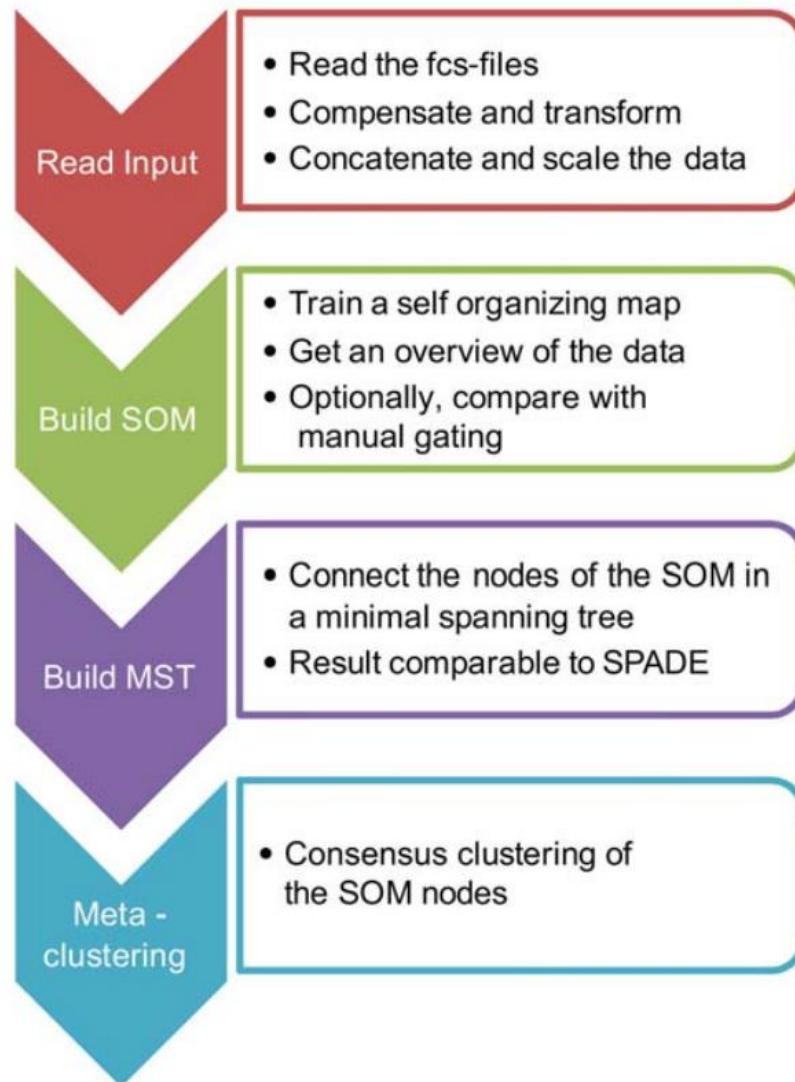
This section will run a UMAP analysis on the PBMC data using set parameters

This section will plot the UMAP results. Two plots will appear, a “flat” dot plot and a density plot

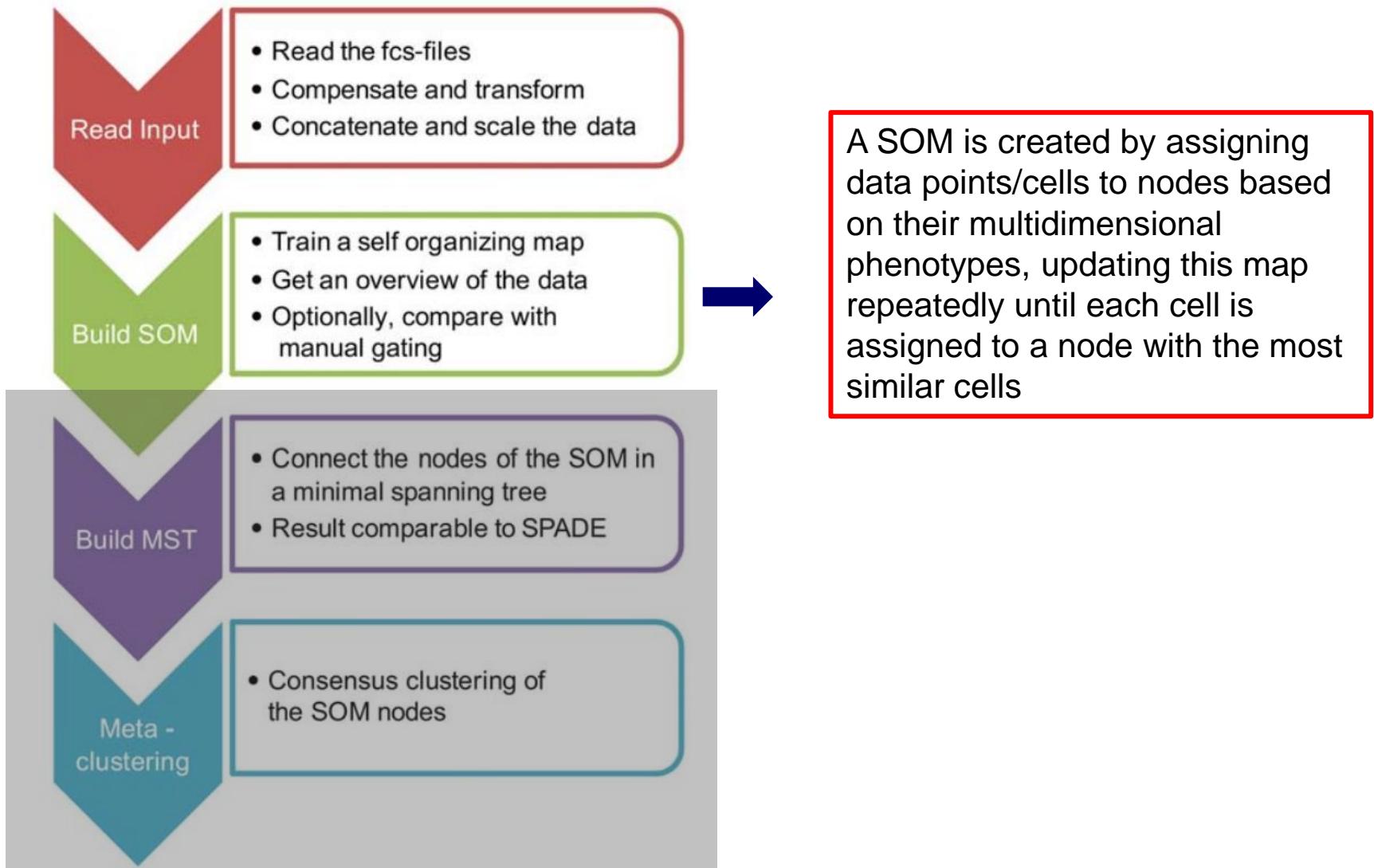
# Clusters can be Identified Based on Dimensionality Reduction Results



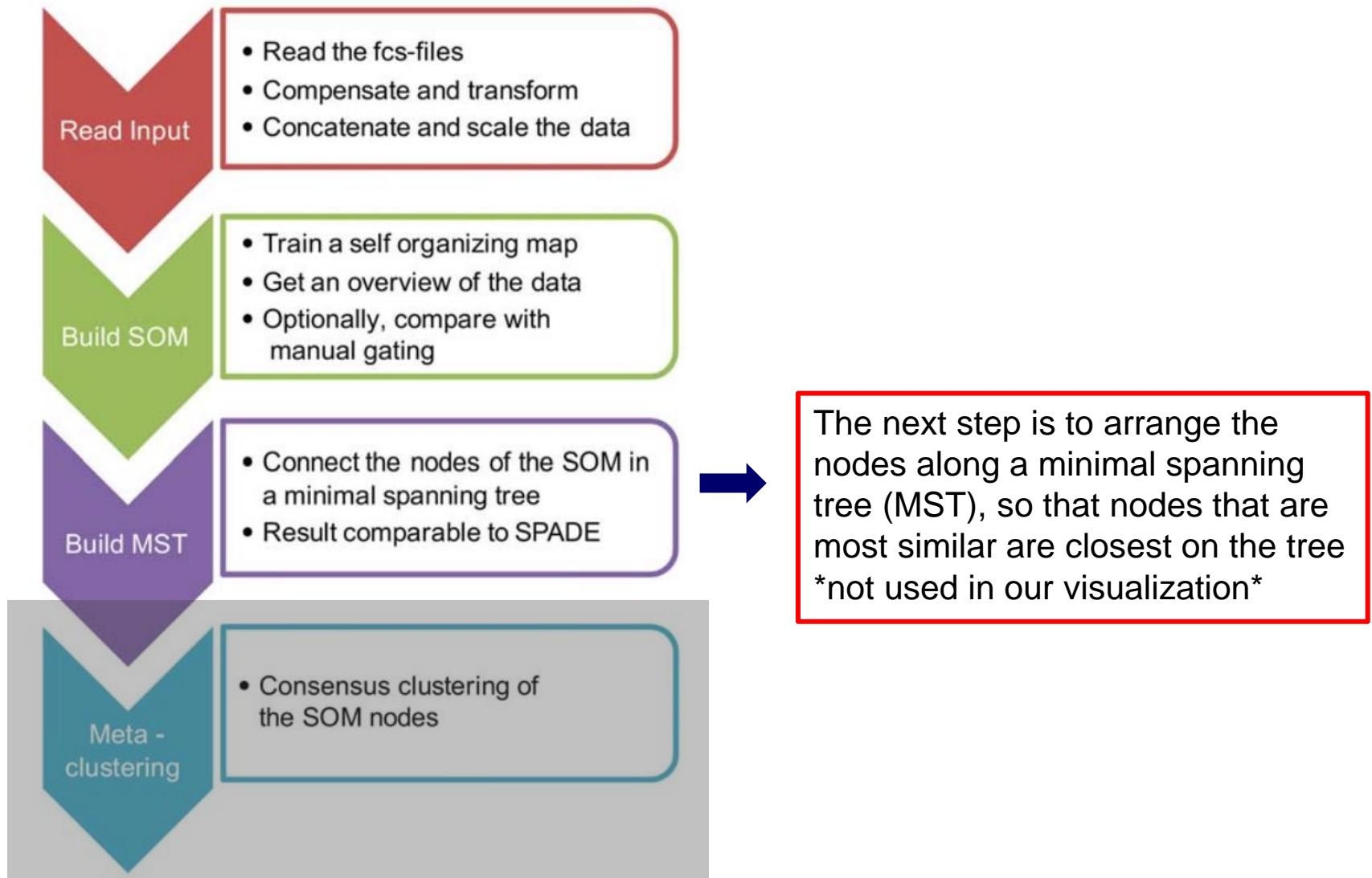
# Clustering with FlowSOM: Self-organizing Maps



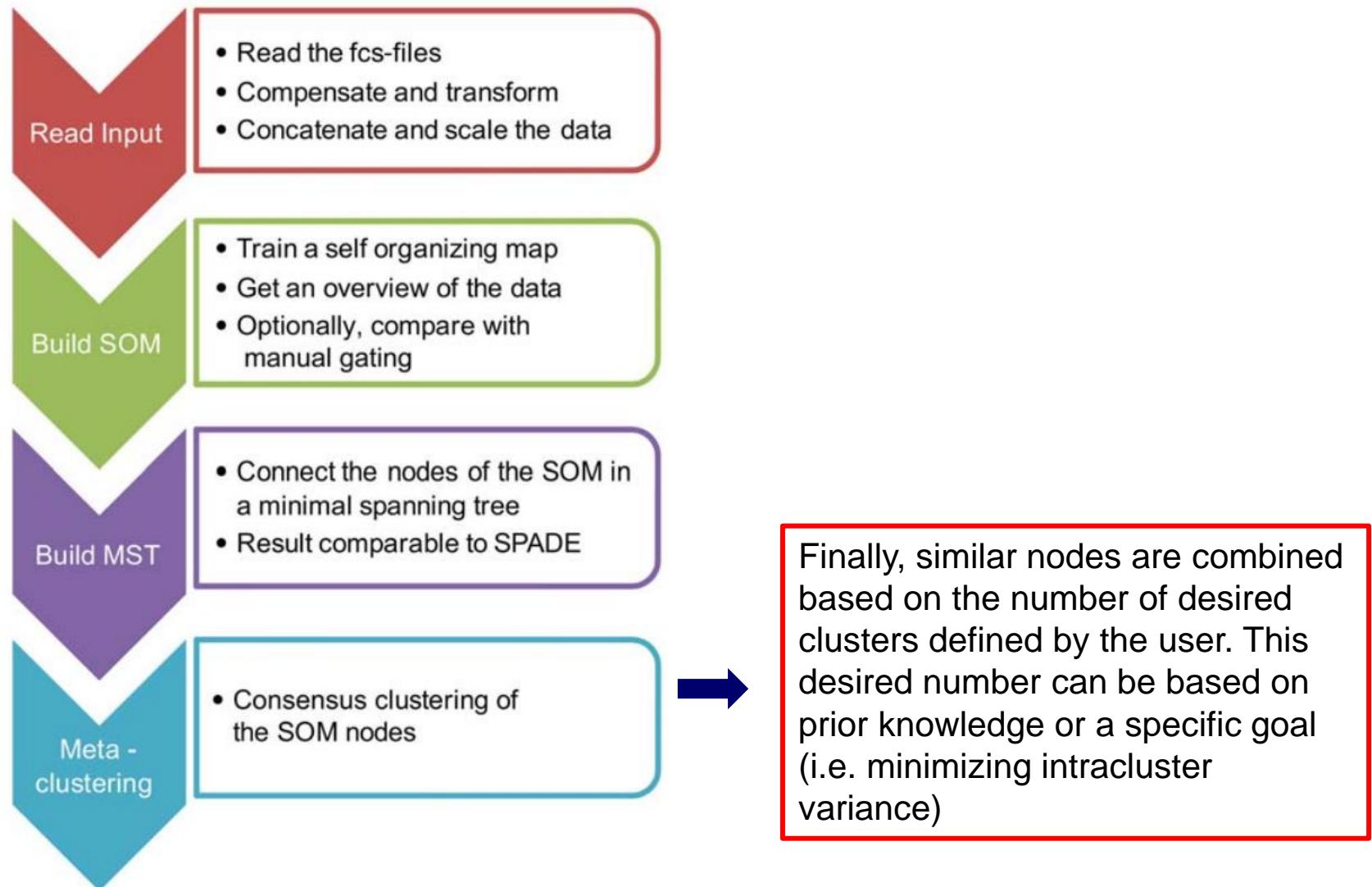
# Clustering with FlowSOM: Self-organizing Maps



# Clustering with FlowSOM: Self-organizing Maps

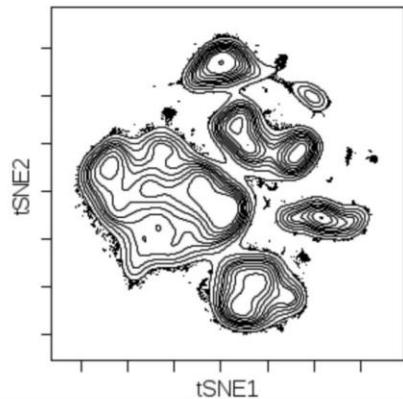


# Clustering with FlowSOM: Self-organizing Maps

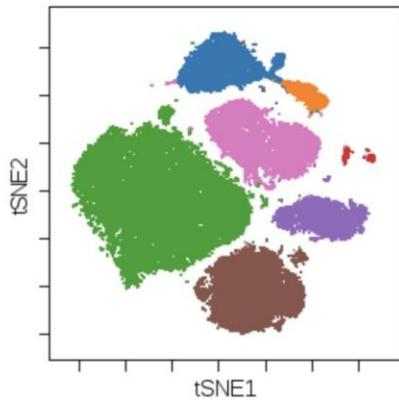


# Clustering with FlowSOM: Self-organizing Maps

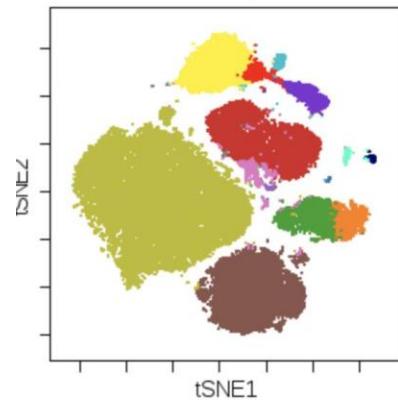
Contour plot of viSNE map



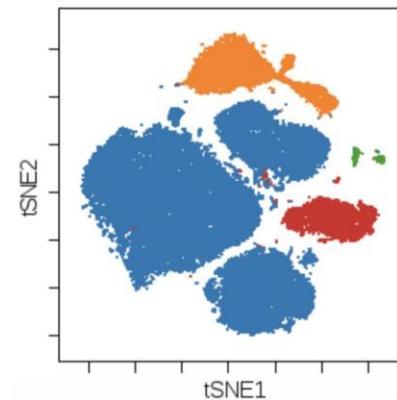
FlowSOM metaclusters overlaid on viSNE map



# metaclusters = 7



# metaclusters = 15



# metaclusters = 4

# 03\_blood\_cell\_PBMC\_example.html

## Run FlowSOM

```
166 ````{r run_FlowSOM}
167 # Time <10 sec
168
169 # create flowFrame for FlowsOM input (using umap axes as input)
170 matrix <- as.matrix(umap.data)
171 metadata <-
172   data.frame(name = dimnames(matrix)[[2]],
173             desc = paste('UMAP', dimnames(matrix)[[2]]))
174 metadata$range <- apply(apply(matrix, 2, range), 2, diff)
175 metadata$minRange <- apply(matrix, 2, min)
176 metadata$maxRange <- apply(matrix, 2, max)
177 flowframe <- new("flowFrame",
178   exprs = matrix,
179   parameters = AnnotatedDataFrame(metadata))
180
181 # implement the FlowsOM on the data by running the line below (to see help
182 # for FlowsOM, type "FlowsOM -- enter" in console)
183 fsom <-
184   FlowsOM(
185     flowframe,      # input flowframe
186     colstouse = c(1:2), # columns to use
187     nclus = 10,       # target number of clusters (this can be changed)
188     seed = overall_seed # set seed for reproducibility
189   )
190 FlowsOM.clusters <-
191   as.matrix(fsom[[2]][fsom[[1]]$map$mapping[, 1]])
192 ...
193 ````{r plot_clusters}
194 # Time <10 sec
195
196 # plot FlowsOM clusters on UMAP axes
197 ggplot(UMAP.plot) + coord_fixed(ratio=graphical.ratio) +
198   geom_point(aes(x=x, y=y, color=FlowsOM.clusters), cex = 0.5) +
199   labs(x = "UMAP 1", y = "UMAP 2", title = "FlowsOM clustering on UMAP Axes",
200         color = "Cluster") + theme_bw() +
201   guides(colour = guide_legend(override.aes = list(size=4)))+
202   labs(caption = "Data from Diggins et al., Nat Methods 2017, 14: 275-278
203 \nFlow Repository: FR-FCM-ZY63") +
204   theme(panel.grid.major = element_blank(),
205         panel.grid.minor = element_blank())
206 ...
207
```

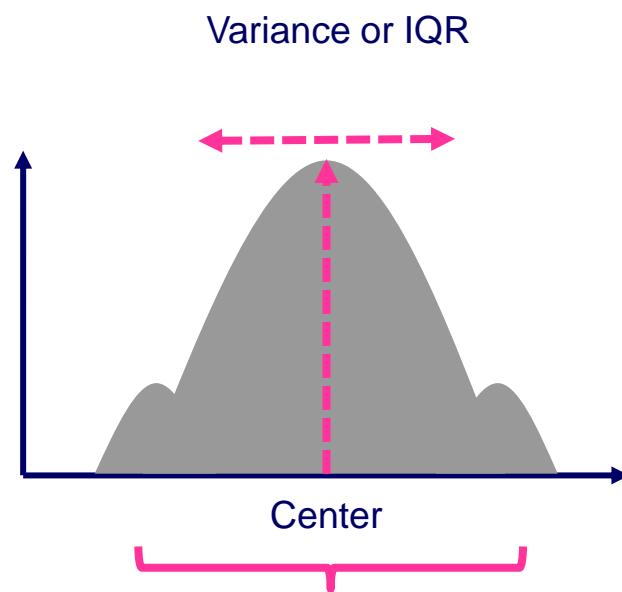
This section performs FlowSOM clustering on the UMAP results

You can choose the parameters the clustering is performed on (UMAP axes vs. measured markers) as well as a seed and desired number of clusters

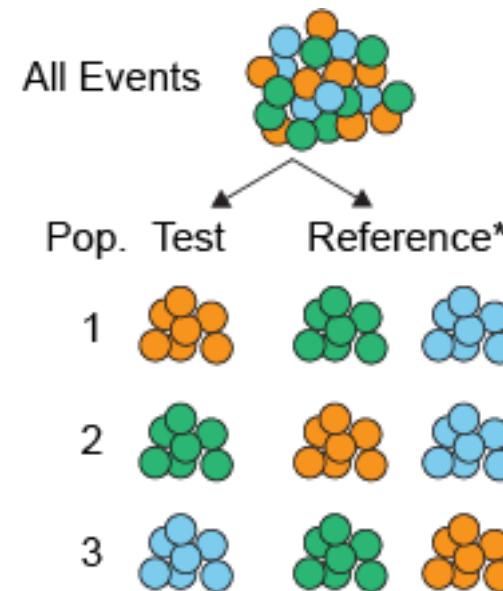
This section plots the identified clusters back onto the UMAP axes and generates a plot (a colored version of the UMAP plot from before)

# Marker Enrichment Modeling Analysis Identifies Markers that are Specifically Expressed or Lacking on Populations

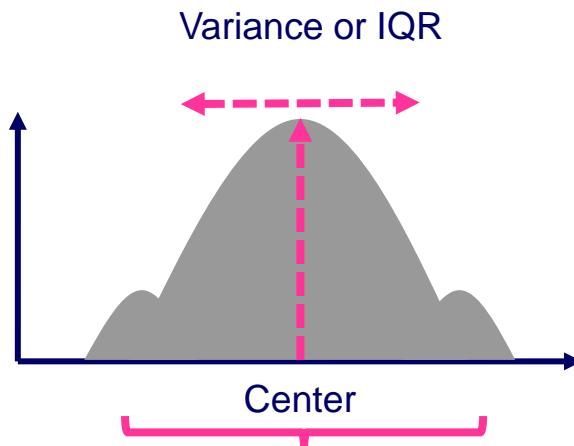
MEM accounts for variance and median of markers to identify enriched features on subsets of cells



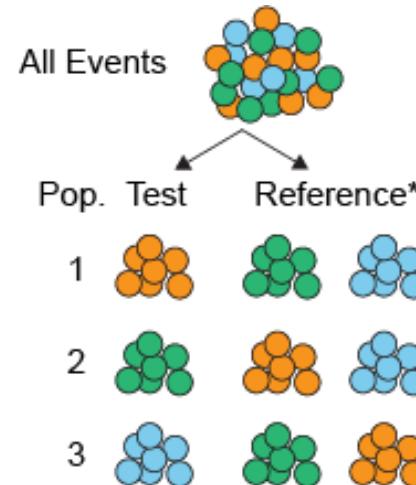
Shape (skewness, symmetry  
# peaks, outliers, etc.)



# MEM Quantifies Relative Enrichment by Combining Magnitude and Interquartile Range



Shape (skewness, symmetry  
# peaks, outliers, etc.)



MEM label

▲ HLADR<sup>+10</sup> CD20<sup>+9</sup> CD19<sup>+7</sup> IgM<sup>+5</sup> CD34<sup>+3</sup>  
CD45RA<sup>+3</sup> CXCR4<sup>+2</sup> CD47<sup>+2</sup> CD33<sup>+2</sup>  
▼ CD7<sup>-2</sup>

Linear transformation to -10 to +10

If  $MAG_{test} - MAG_{ref} < 0$ ,  $MEM = -MEM$

# 03\_blood\_cell\_PBMC\_example.html

## Run MEM

```
207 ````{r run_MEM}
208 # Time ~30 sec
209
210 # Run MEM on the FlowSOM clusters found from using UMAP axes
211 cluster = as.numeric(as.vector((FlowSOM.clusters)))
212 MEM.data = cbind(transformed.chosen.markers, cluster)
213
214 MEM.values = MEM(
215   MEM.data,           # input data (last column must contain cluster values)
216   transform = FALSE,  # data is already scaled in this case
217   cofactor = 1,
218   choose.markers = FALSE,
219   markers = "all",    # use all transformed, chosen markers from previous
220   selection
221   choose.ref = FALSE, # reference will be all other cells
222   zero.ref = FALSE,
223   rename.markers = FALSE,
224   new.marker.names =
225     "CD19,CD117,CD11b,CD4,CD8,CD20,CD34,CD61,CD123,CD45RA,CD45,CD10,CD33,CD68,CD69,CD15,CD16,CD44,CD38,CD25,CD3,IgM,HLA-DR,CD56", # rename channels
226   labels
227   file.is.clust = FALSE,
228   add.fileID = FALSE,
229   IQR.thresh = NULL
230 )
231
232 # build MEM heatmap and output enrichment scores
233 build.heatmaps(
234   MEM.values,          # input MEM values
235   cluster.MEM = "both", # dendrogram for columns and rows
236   display.thresh = 3,   # display threshold for MEM scores
237   newwindow.heatmaps = FALSE, # makes txt and PDF files for heatmap and MEM
238   output.files = TRUE,  # include labels in heatmap
239   scores
240   labels = TRUE,        # include labels in heatmap
241   only.MEMheatmap = FALSE
242 )````
```

This section performs MEM analysis on the identified FlowSOM clusters

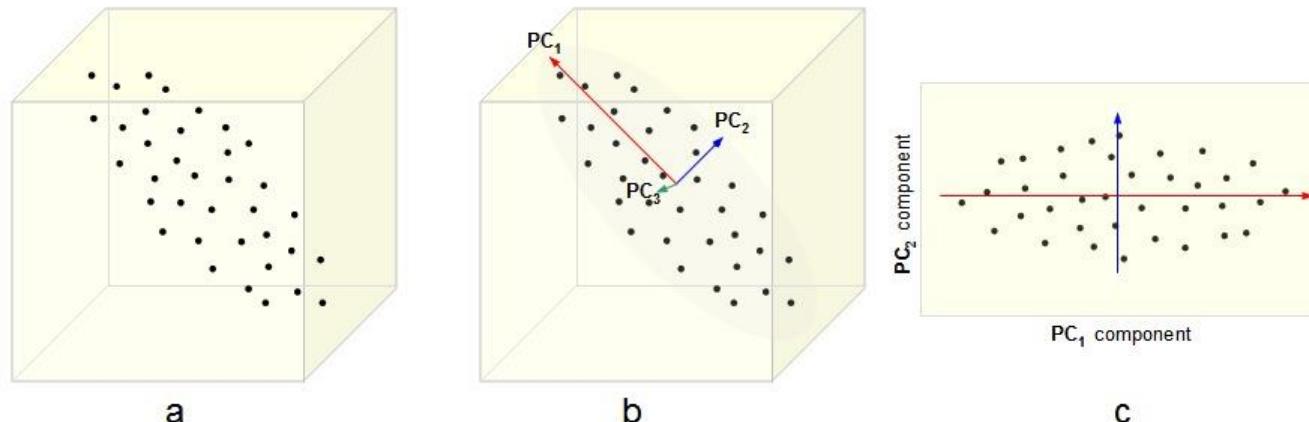
You can choose the markers for the MEM analysis as well as their names and the reference population

This section produces heatmaps and MEM (enrichment) scores

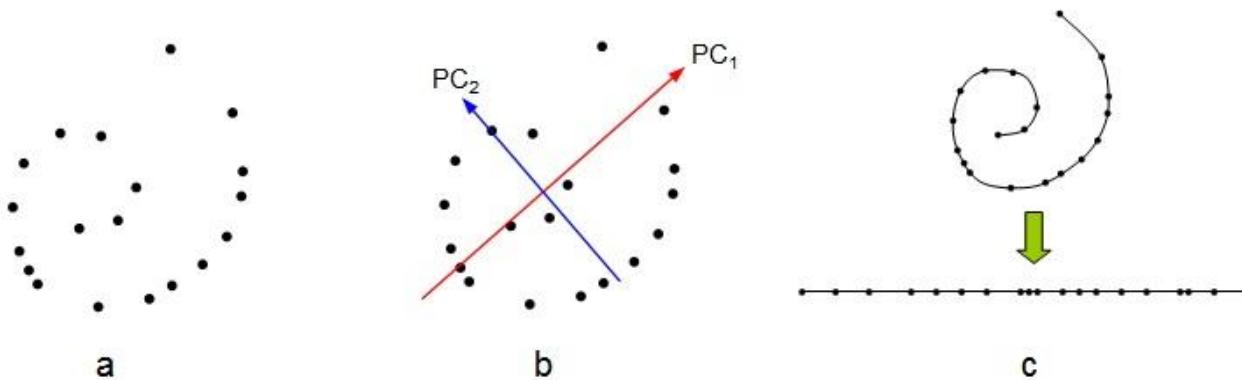
# Part 7: Many More Tools!

# PCA, Citrus, Trajectories

# PCA is a Linear Dimensionality Reduction Tool

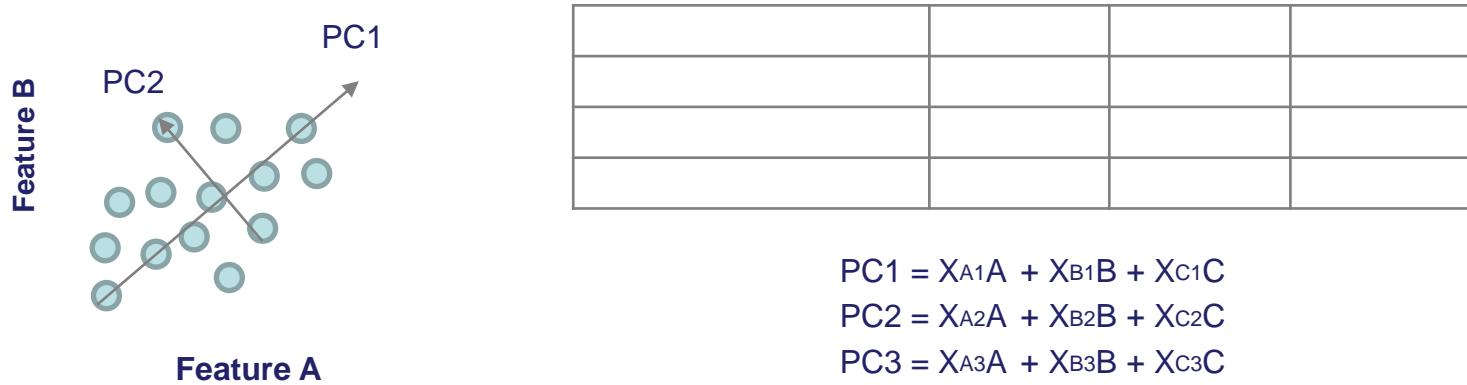


An illustration of PCA. **a)** A data set given as 3-dimensional points. **b)** The three orthogonal Principal Components (PCs) for the data, ordered by variance. **c)** The projection of the data set into the first two PCs, discarding the third one.



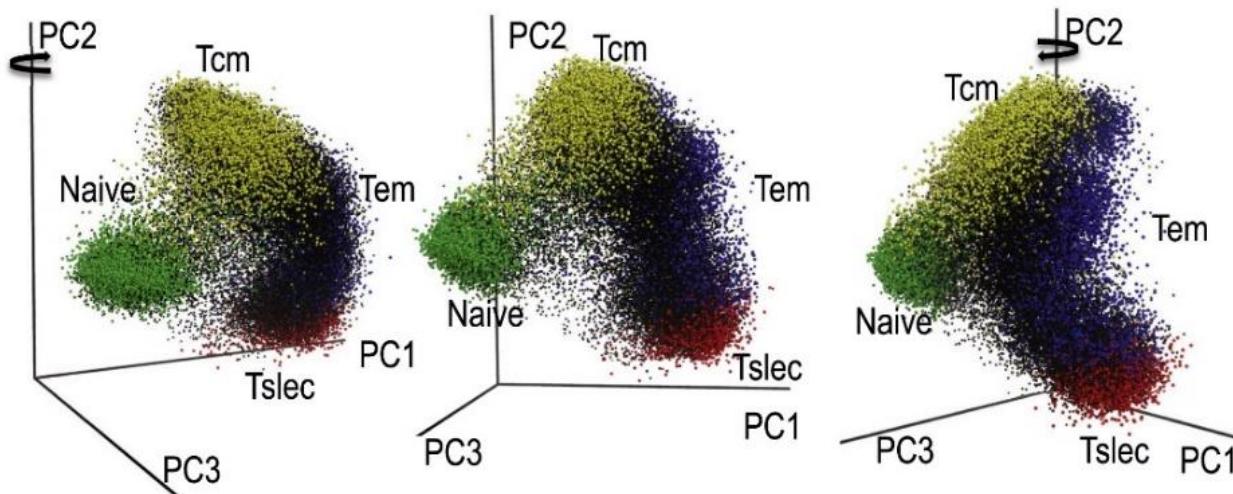
Effects of dimensionality reduction on an inherently non-linear data set. **a)** The original data given as a two-dimensional set. **b)** PCA identifies two PCs as contributing significantly to explain the data variance. **c)** However, the inherent topology (connectivity) of the data helps identify the set as being one-dimensional, but non-linear.

# Principal Component Analysis



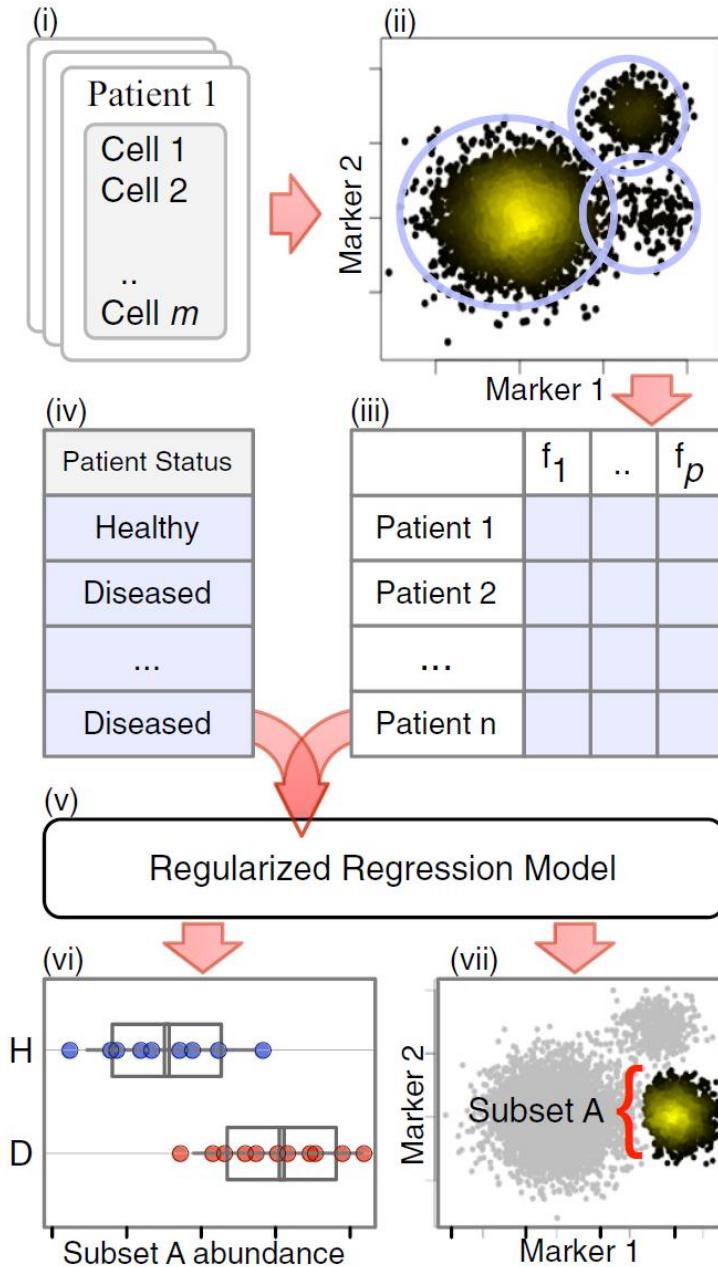
## PCA used to Reduce Dimensionality of CyTOF Data

A 3D-PCA view of CD8<sup>+</sup> T cell 25 parameter data



Newell et al 2012, *Immunity*

# Citrus: Supervised Population Finding

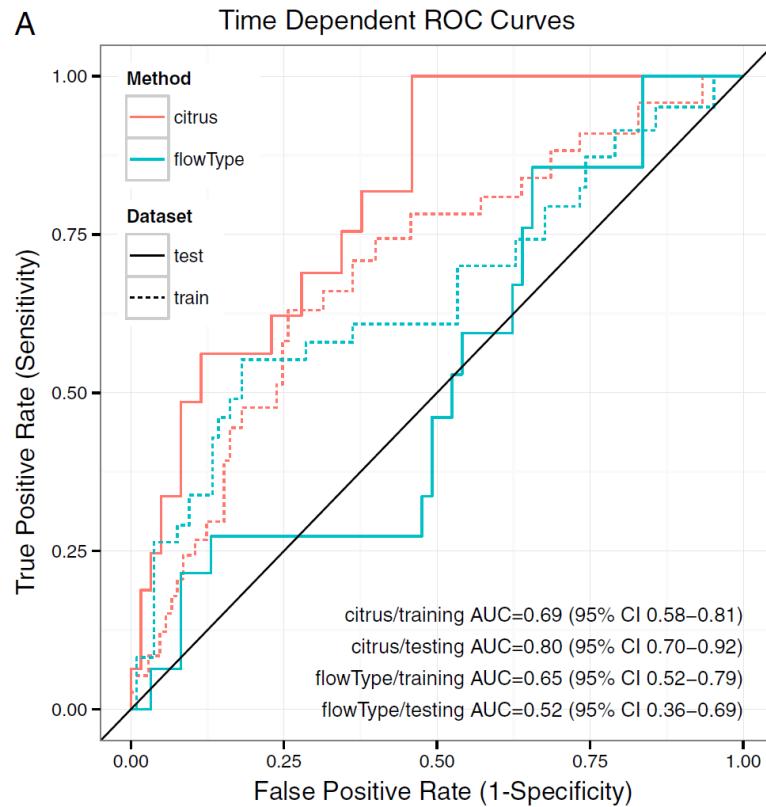


## Automated identification of stratifying signatures in cellular subpopulations

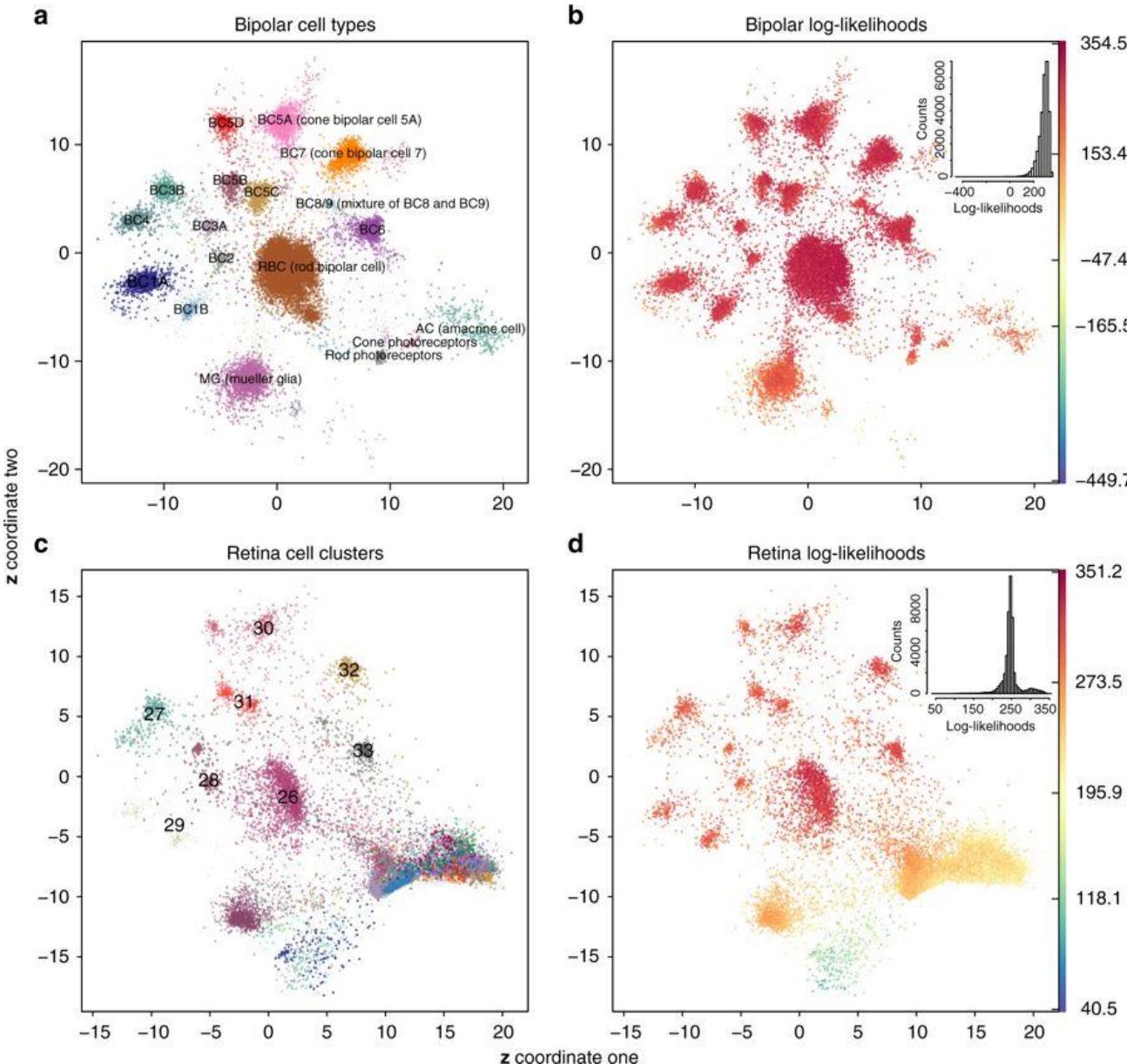
Robert V. Bruggner<sup>a,b</sup>, Bernd Bodenmiller<sup>c</sup>, David L. Dill<sup>d</sup>, Robert J. Tibshirani<sup>e,f,1</sup>, and Garry P. Nolan<sup>b,1</sup>

<sup>a</sup>Biomedical Informatics Training Program, Stanford University Medical School, Stanford, CA 94305; <sup>b</sup>Baxter Laboratory for Stem Cell Biology, Department of Microbiology and Immunology, and Departments of <sup>c</sup>Computer Science, <sup>d</sup>Health Research and Policy, and <sup>f</sup>Statistics, Stanford University, Stanford, CA 94305; and <sup>e</sup>Institute of Molecular Life Sciences, University of Zurich, CH-8057 Zurich, Switzerland

Contributed by Robert J. Tibshirani, May 14, 2014 (sent for review February 12, 2014)

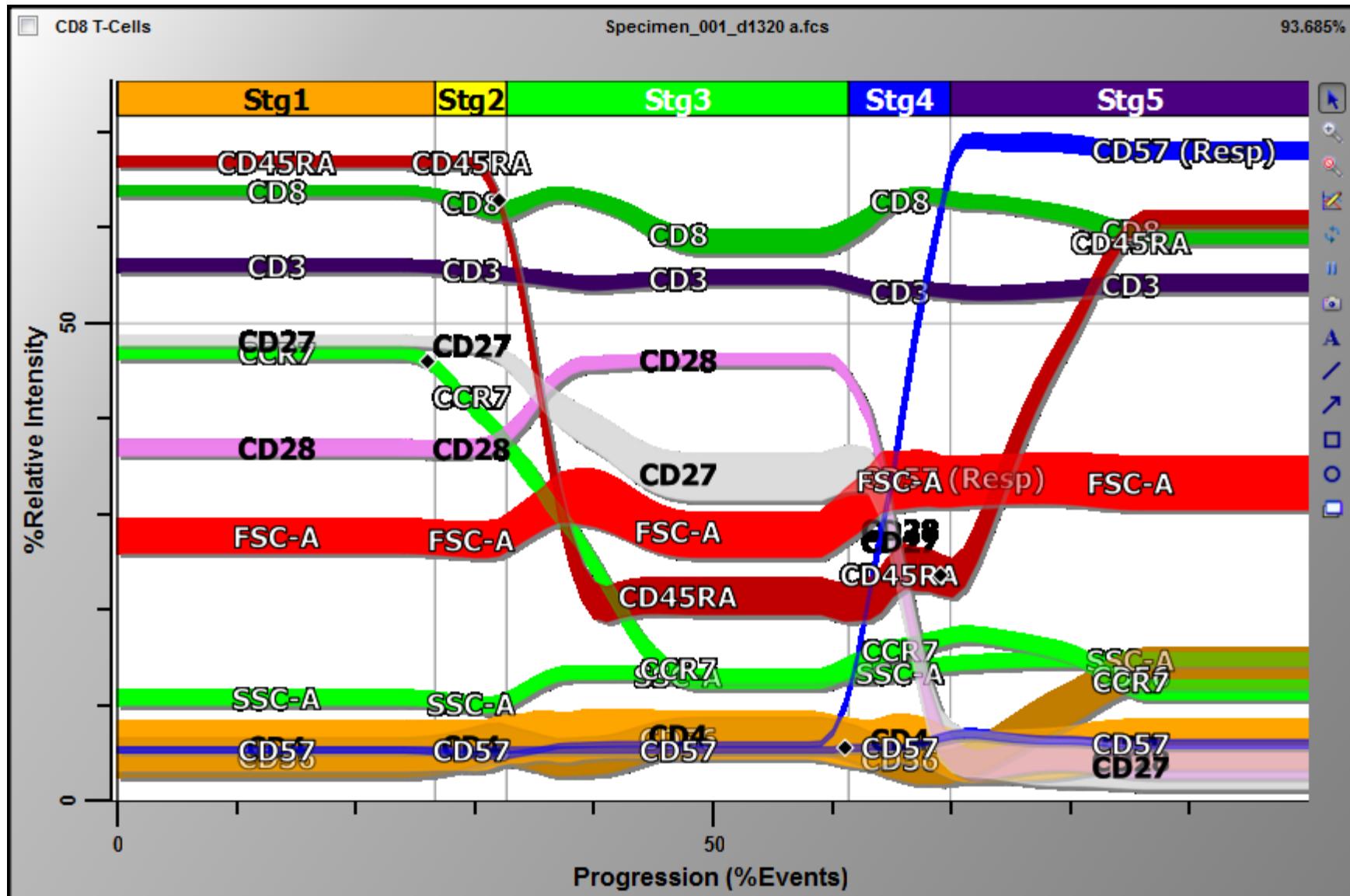


# In 2018, Ding et al. Used scvis & Probability to Characterize / Identify Cells (scRNA-seq)

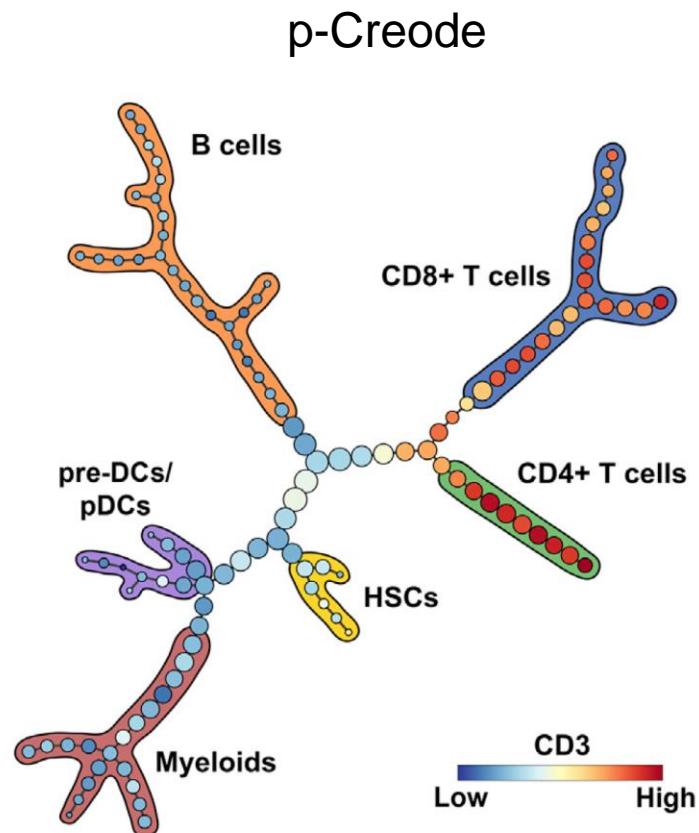
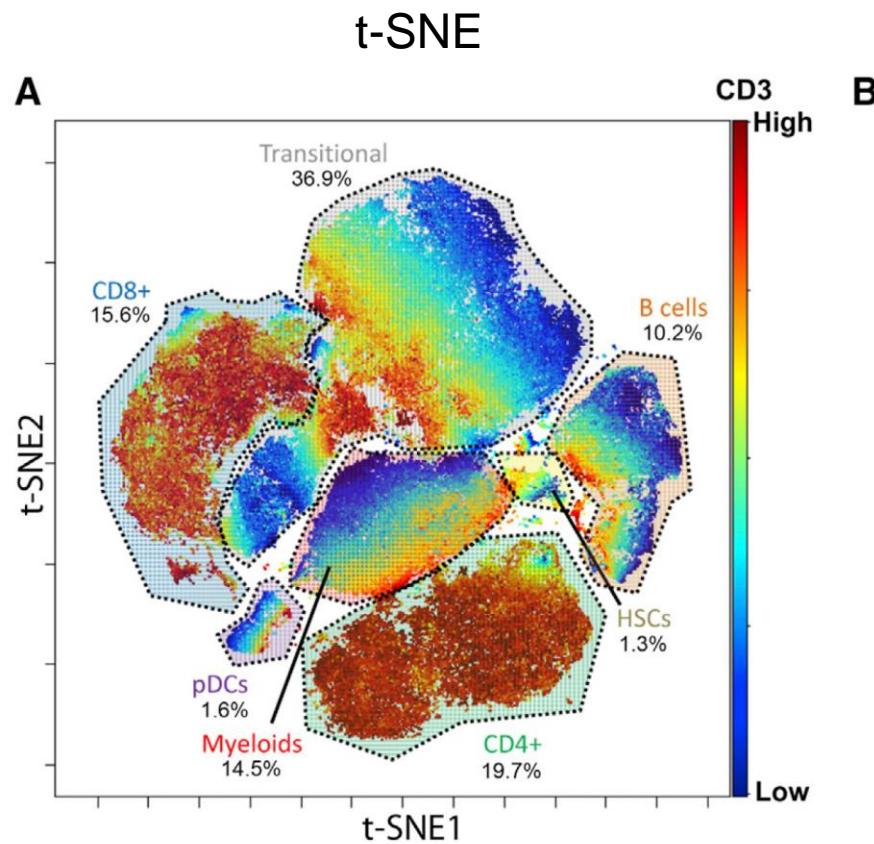


Learning a probabilistic mapping function from the bipolar data and applying the function to the independently generated mouse retina dataset. **a** scvis learned two-dimensional representations of the bipolar dataset, **b** coloring each point by the estimated log-likelihood, **c** the whole mouse retina dataset was directly projected to a two-dimensional space by the probabilistic mapping function learned from the bipolar data, and **d** coloring each point from the retina dataset by the estimated log-likelihood

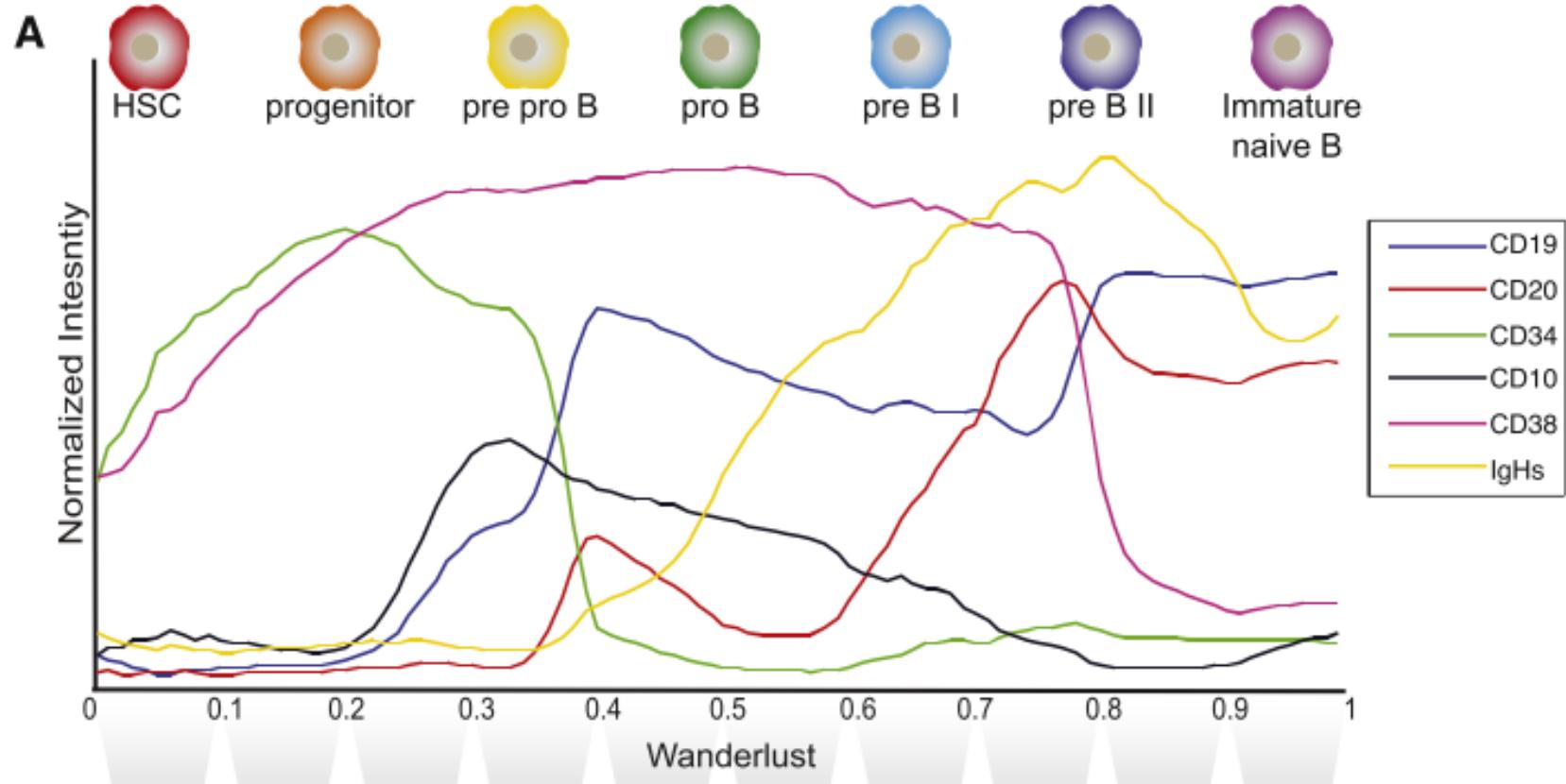
# Gemstone Uses Supervised Analysis to Identify Progressions



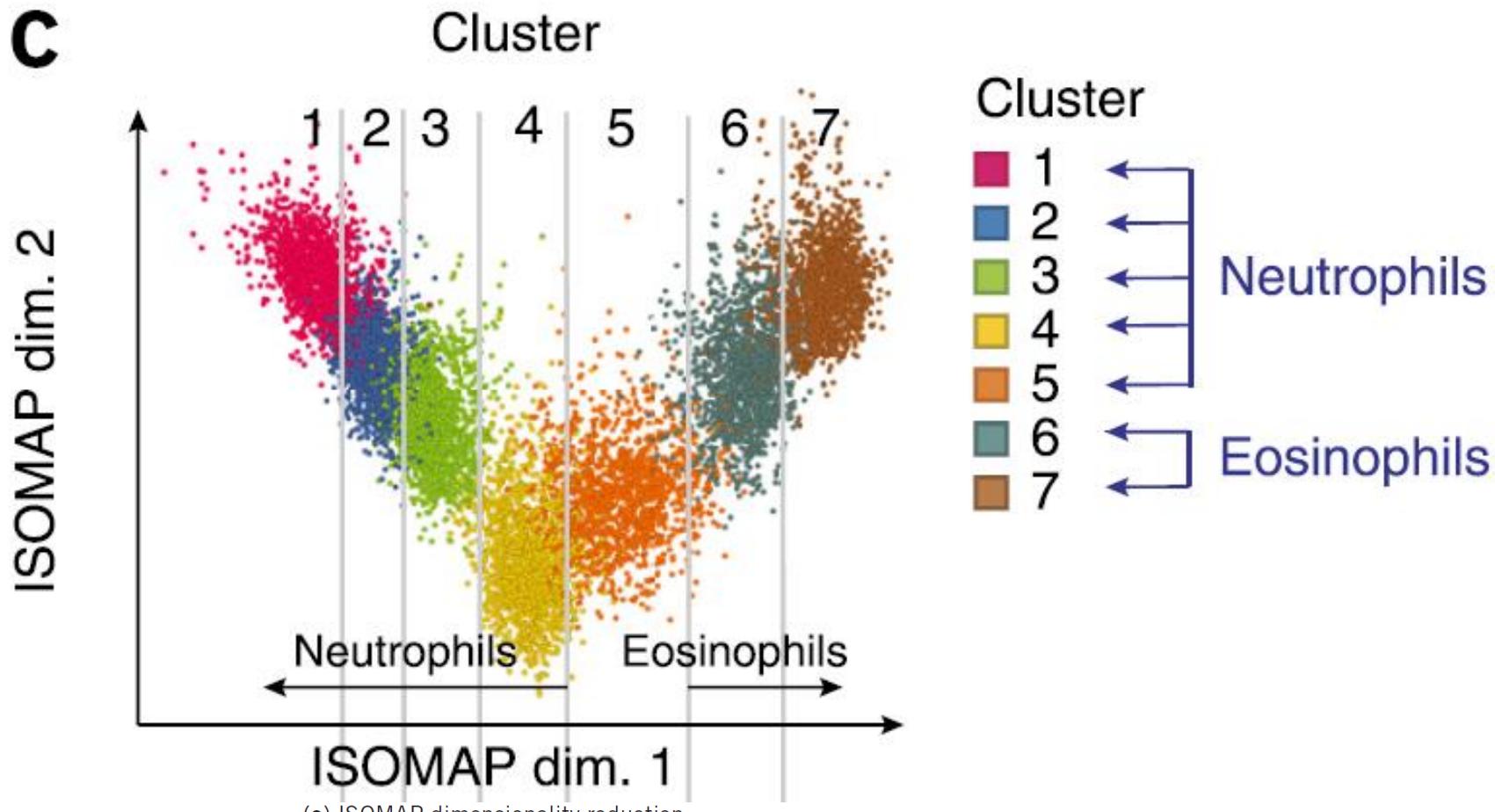
In 2018, Ken Lau's Group Developed p-Creode to Infer Continua in Single Cell Data (e.g., human bone marrow, CyTOF)



# Wanderlust Identifies Phenotypic Progression



# ISOMAP guided analysis

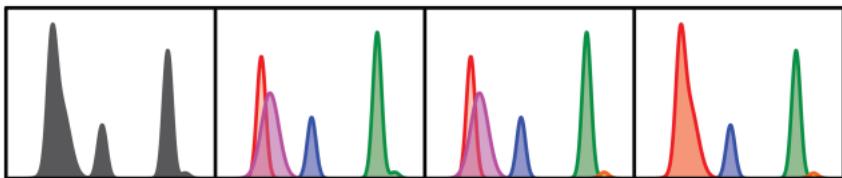


to compare overall phenotypic relatedness of populations of neutrophil-like and eosinophil-like cells<sup>31</sup>. Top, cells color-coded by DensVM cluster number are plotted by their scores for ISOMAP dimensions 1 and 2. Binned median expression of defining markers (middle) and the tissue composition (percentage of each cluster as a fraction of total granulocytes from each tissue, bottom) of cells along this phenotypic progression defined by ISOMAP dimension 1 and DensVM clusters 1–7 are plotted.

# Mixture Modeling

## SWIFT

A:



**Initial sub-populations:**

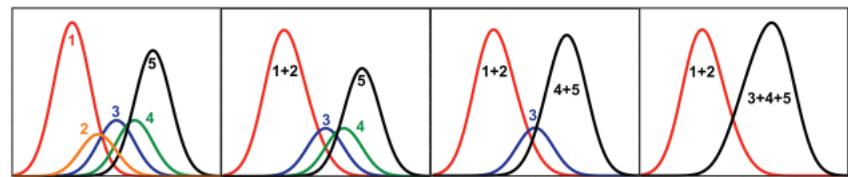
May be skewed;  
May overlap;  
May have a high  
dynamic range.

**1: EM fitting:** The EM algorithm fits data to a specified number of Gaussians, by weighted, iterative sampling. Large asymmetric peaks may be split, but rare peaks may not separate.

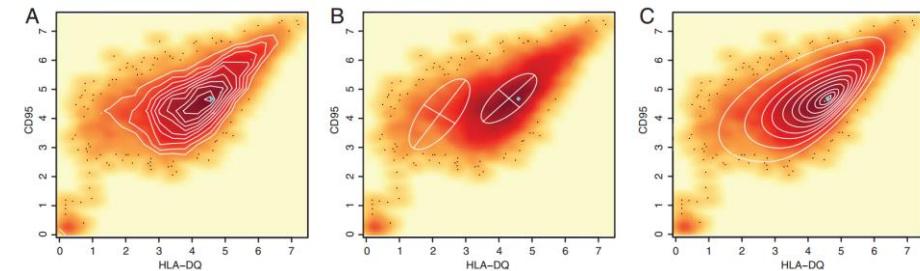
**2: Splitting:** Each cluster from Step 1 is tested by LDA for multiple modes in all combinations of dimensions. Clusters are split if necessary (using EM), until all are unimodal.

**3: Merging:** All cluster pairs are tested, and merged if the resulting cluster is unimodal in all dimensions. Agglomerative merging prevents over-merging due to 'bridging' Gaussians.

B:



## FLAME

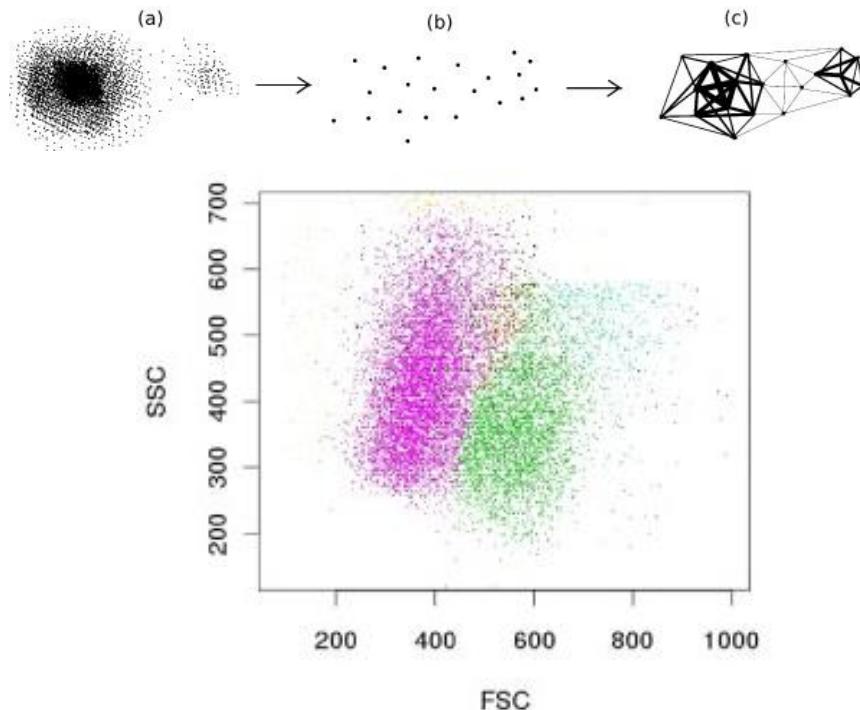


Pyne et al, 2009 *PNA*

Mosmann et al, 2014 *Cytometry A*

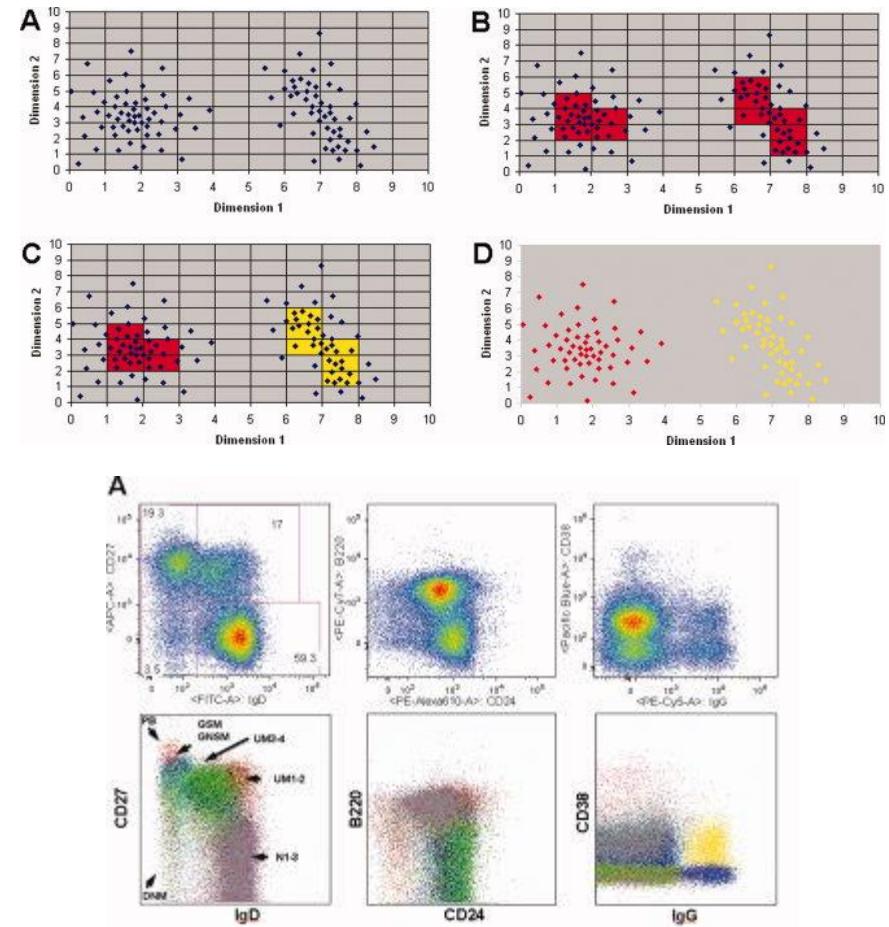
# Automated Clustering and Population Identification Methods Based on Density

## SamSpectral



Zare et al, 2010 *BMC Bioinformatics*

## FLOCK

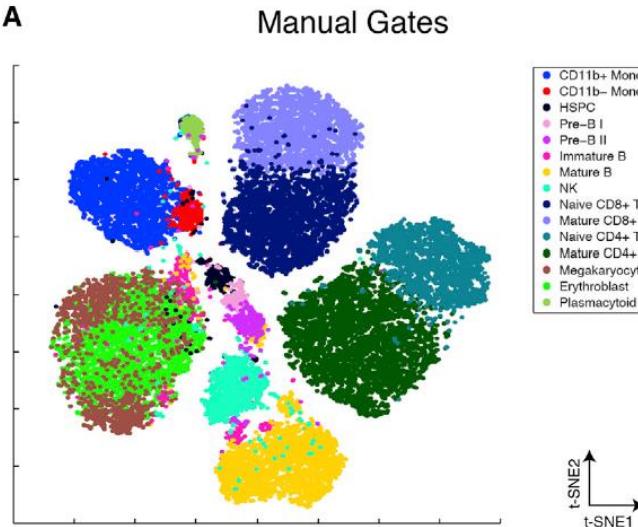


Qian et al, 2010 *Cytometry B Clin Cytom*

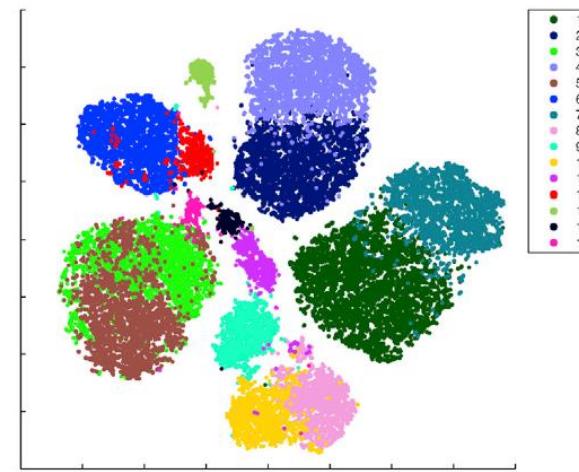
# Phenograph

# Phenograph Adds Fast Clustering & Meta-Analysis to viSNE

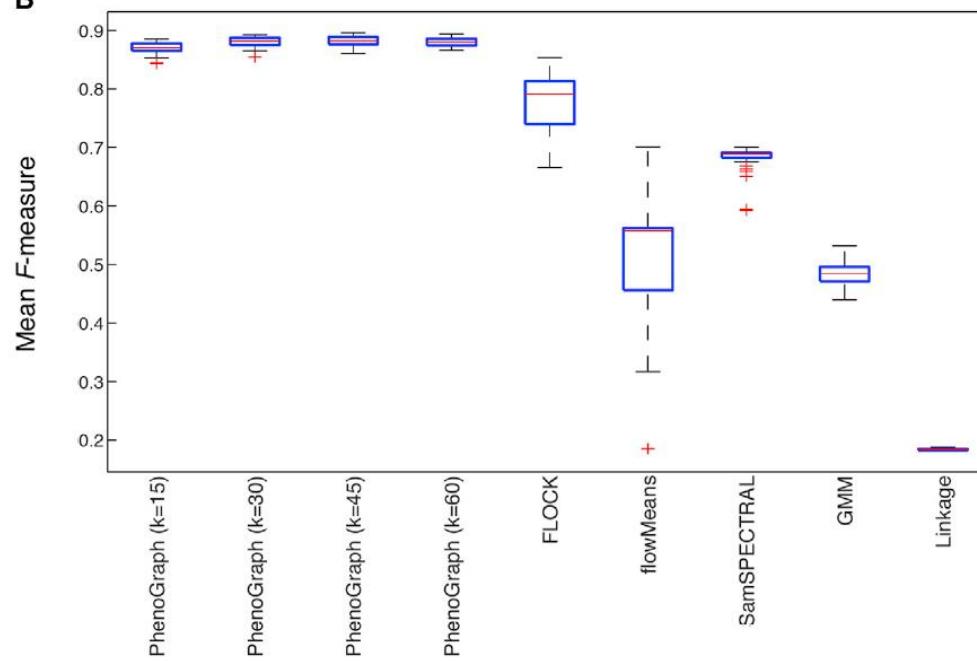
A



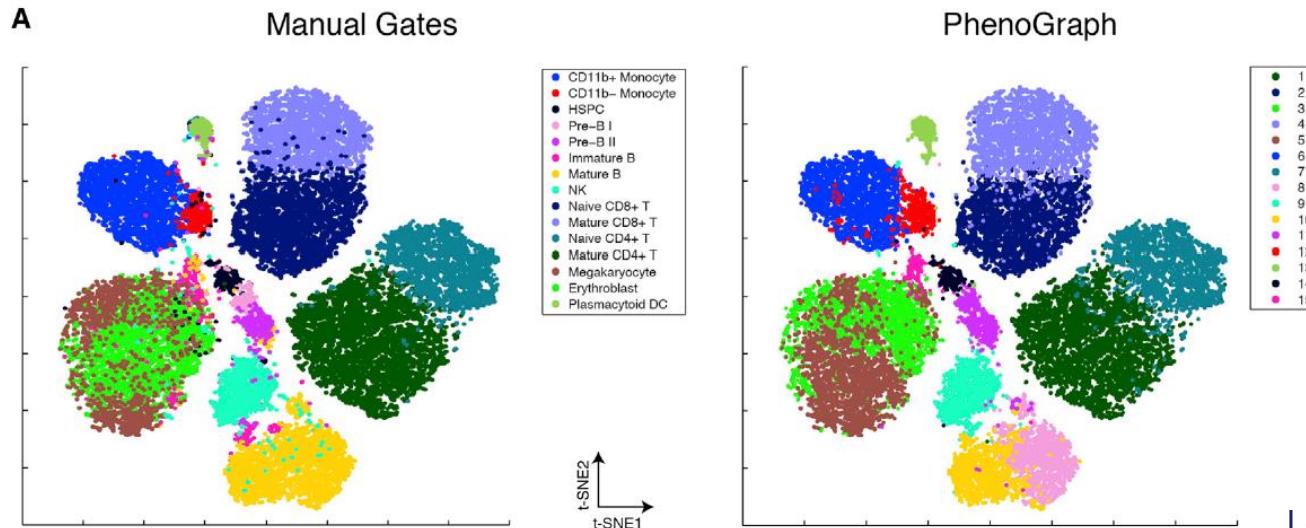
PhenoGraph



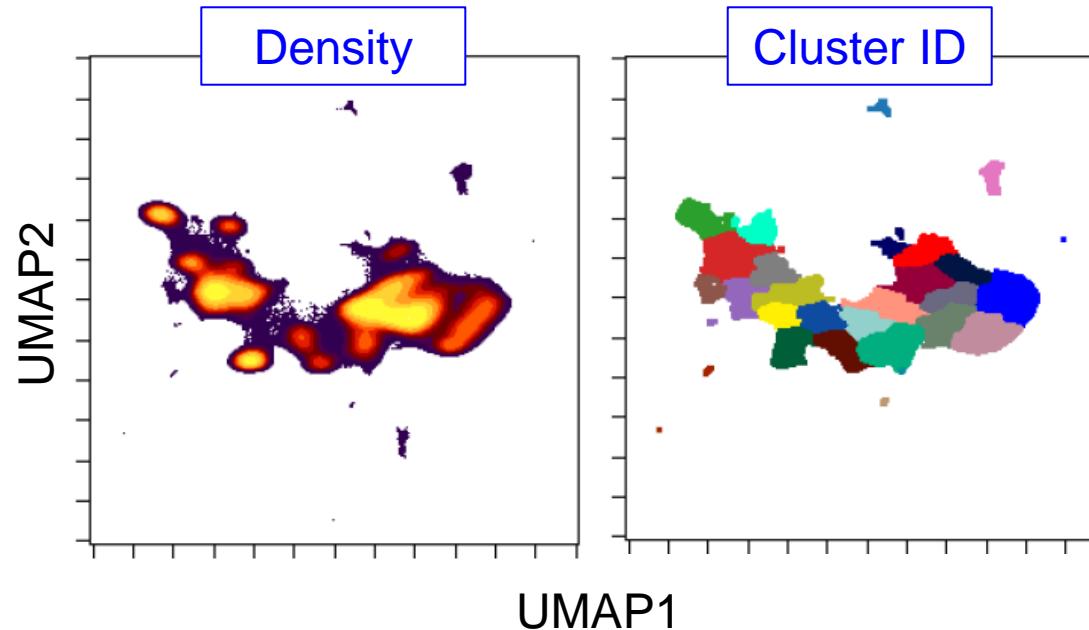
B



# Phenograph: Clustering 35 Features => t-SNE (Not the Reverse)



# Diggins: t-SNE or UMAP on Features => Clustering on 2 axes



UMAP

# Dimensionality reduction for visualizing single-cell data using UMAP

Etienne Becht<sup>1</sup>, Leland McInnes<sup>2</sup> , John Healy<sup>2</sup>, Charles-Antoine Dutertre<sup>1</sup>, Immanuel W H Kwok<sup>1</sup>, Lai Guan Ng<sup>1</sup>, Florent Ginhoux<sup>1</sup>  & Evan W Newell<sup>1,3</sup> 

Advances in single-cell technologies have enabled high-resolution dissection of tissue composition. Several tools for dimensionality reduction are available to analyze the large number of parameters generated in single-cell studies. Recently, a nonlinear dimensionality-reduction technique, uniform manifold approximation and projection (UMAP), was developed for the analysis of any type of high-dimensional data. Here we apply it to biological data, using three well-characterized mass cytometry and single-cell RNA sequencing datasets. Comparing the performance of UMAP with five other tools, we find that UMAP provides the fastest run times, highest reproducibility and the most meaningful organization of cell clusters. The work highlights the use of UMAP for improved visualization and interpretation of single-cell data.

# UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

Leland McInnes

Tutte Institute for Mathematics and Computing

[leland.mcinnes@gmail.com](mailto:leland.mcinnes@gmail.com)

John Healy

Tutte Institute for Mathematics and Computing

[jchealy@gmail.com](mailto:jchealy@gmail.com)

James Melville

[jlmelville@gmail.com](mailto:jlmelville@gmail.com)

December 7, 2018

<https://arxiv.org/abs/1802.03426>

## Abstract

UMAP (Uniform Manifold Approximation and Projection) is a novel manifold learning technique for dimension reduction. UMAP is constructed from a theoretical framework based in Riemannian geometry and algebraic topology. The result is a practical scalable algorithm that applies to real world data. The UMAP algorithm is competitive with t-SNE for visualization quality, and arguably preserves more of the global structure with superior run time performance. Furthermore, UMAP has no computational restrictions on embedding dimension, making it viable as a general purpose dimension reduction technique for machine learning.

Comments: Reference implementation available at [this http URL](http://this http URL)

Subjects: **Machine Learning (stat.ML)**; Computational Geometry (cs.CG); Machine Learning (cs.LG)

Cite as: [arXiv:1802.03426 \[stat.ML\]](https://arxiv.org/abs/1802.03426)

(or [arXiv:1802.03426v2 \[stat.ML\]](https://arxiv.org/abs/1802.03426v2) for this version)

## Submission history

From: Leland McInnes [[view email](#)]

[v1] Fri, 9 Feb 2018 19:39:33 UTC (958 KB)

[v2] Thu, 6 Dec 2018 18:54:07 UTC (7,966 KB)

# Flow Cytometry Data That Looks Like a Blob on a t-SNE Plot Appears to Have Structure on a UMAP Plot

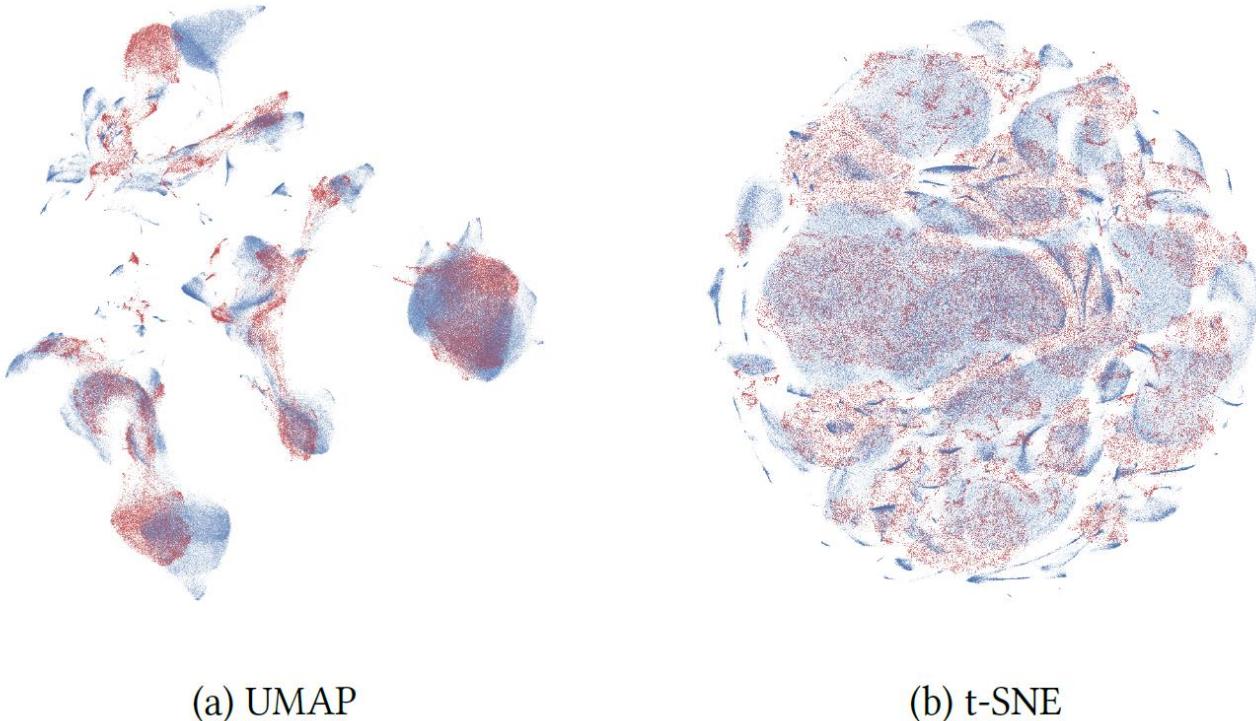
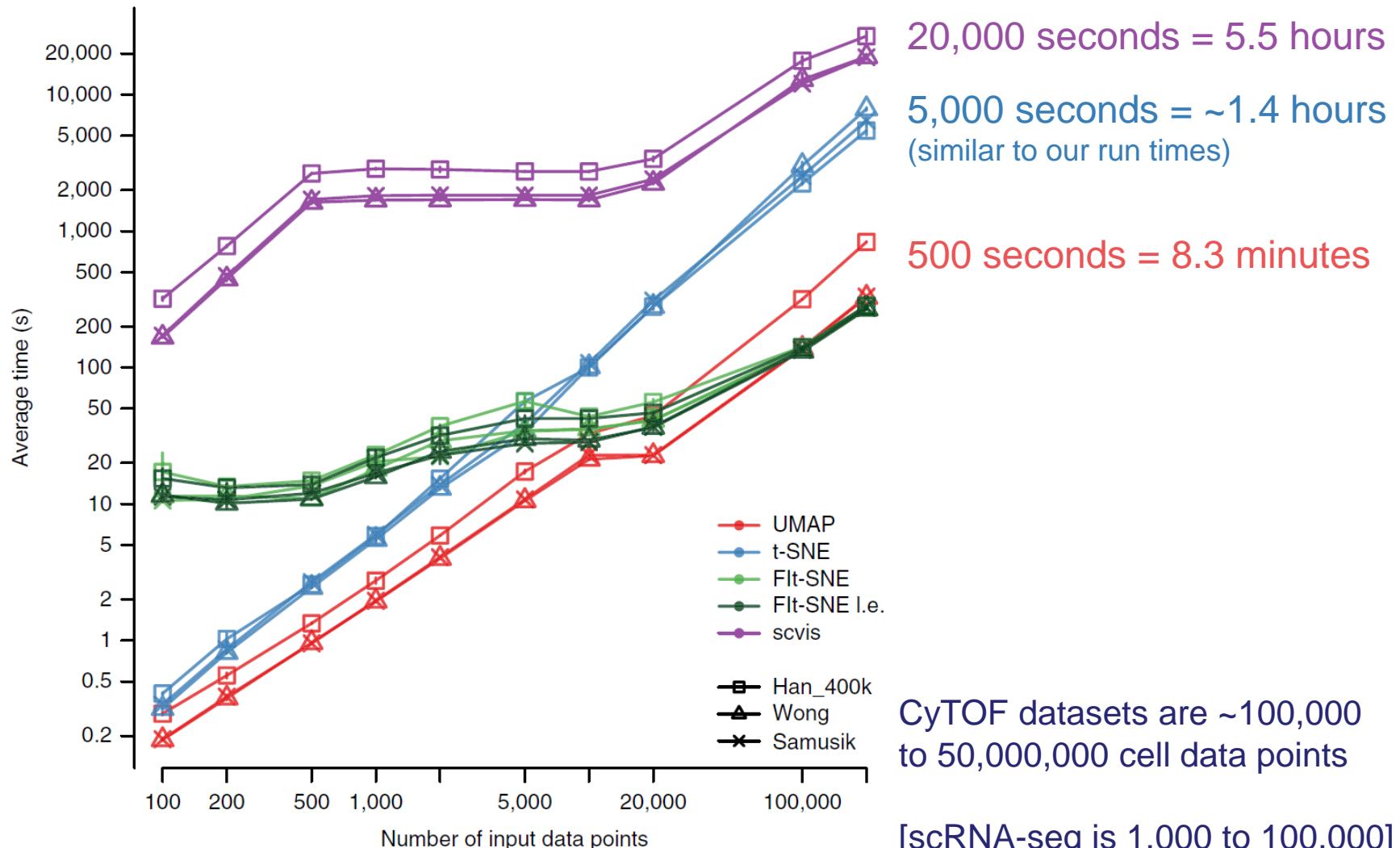


Figure 3: Procrustes based alignment of a 10% subsample (red) against the full dataset (blue) for the flow cytometry dataset for both UMAP and t-SNE.

In [Greek mythology](#), **Procrustes** ([Ancient Greek](#): Προκρούστης *Prokrōstēs*) or "the stretcher [who hammers out the metal]", also known as **Prokoptas** or **Damastes** (Δαμαστής, "subduer"), was a rogue smith and bandit from [Attica](#) who attacked people by stretching them or cutting off their legs, so as to force them to fit the size of an iron bed.

The word "Procrustean" is thus used to describe situations where different lengths or sizes or properties are fitted to an arbitrary standard.

# UMAP & Fit-SNE are Much Faster than Traditional t-SNE



**Figure 3** Run times of five dimensionality reduction methods for inputs of varying sizes. The average run time of three random subsamples is represented, with vertical bars representing s.d. after log-transforming the run times.

# Now, McInnes et al., UMAP Preserves Local and Global Structure

Datasets

COIL20      MNIST      Fashion MNIST      Word vectors

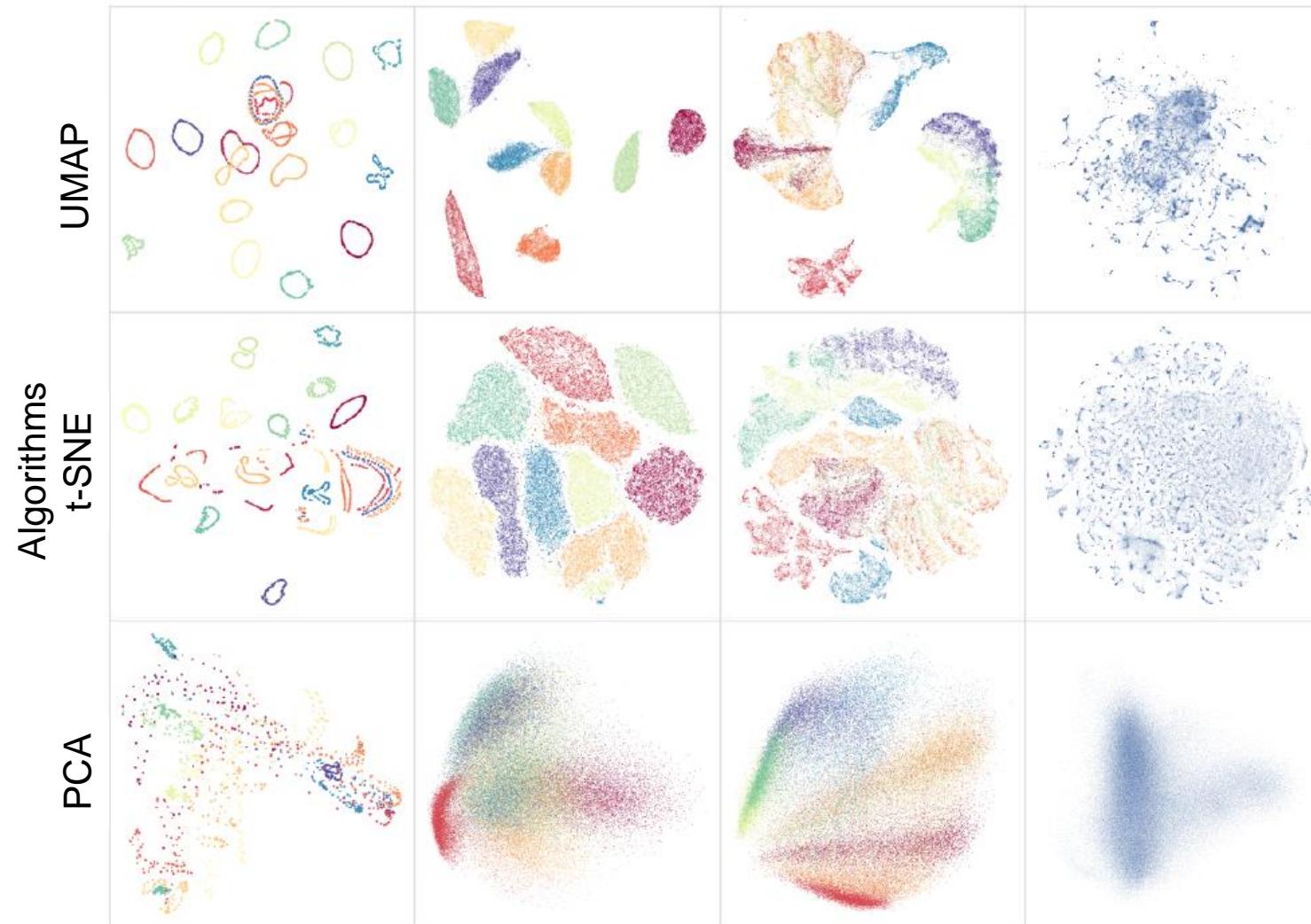


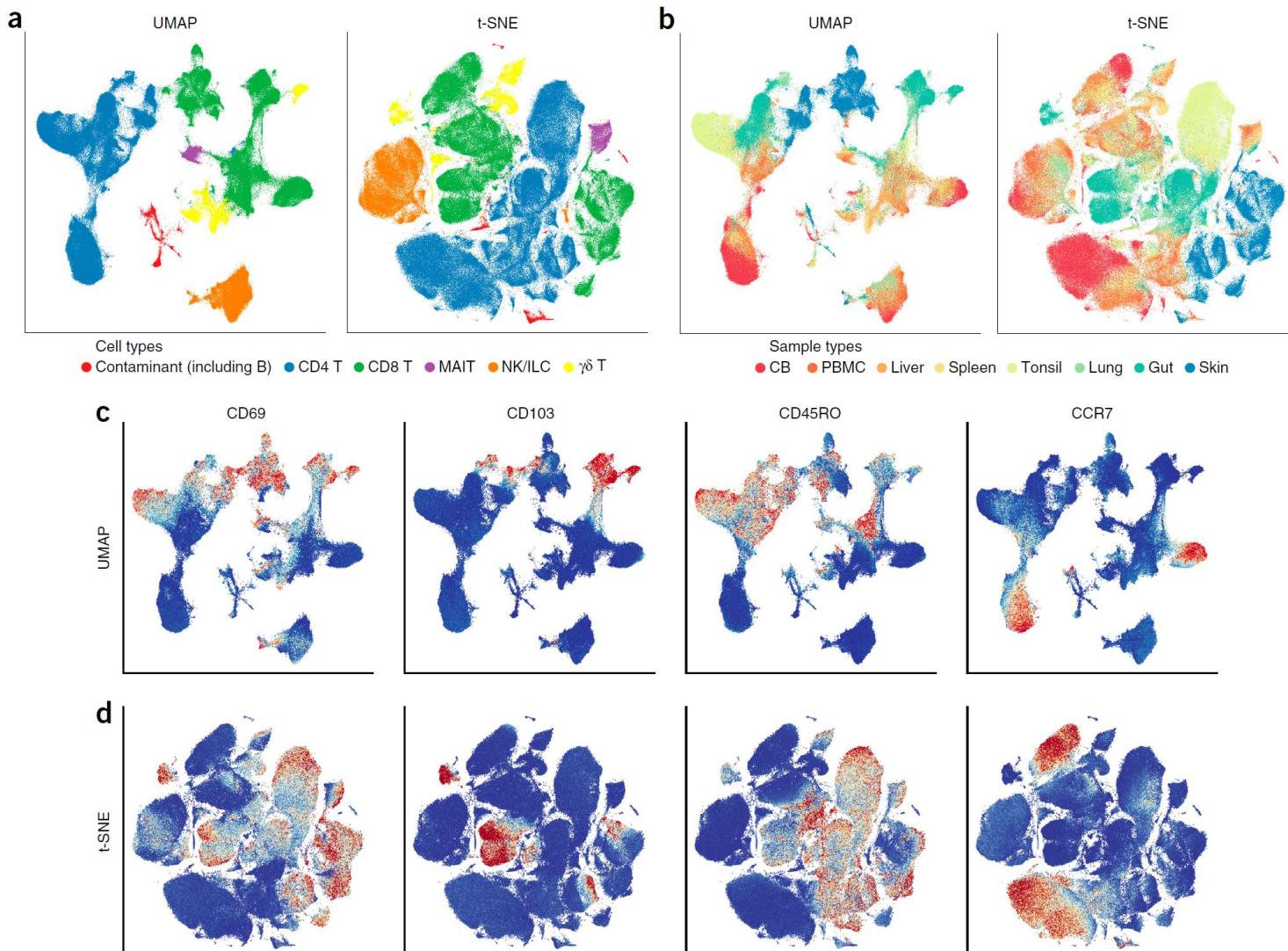
Figure 2: A comparison of several dimension reduction algorithms. We note that UMAP successfully reflects much of the large scale global structure that is well represented by Laplacian Eigenmaps and PCA (particularly for MNIST and Fashion-MNIST), while also preserving the local fine structure similar to t-SNE and LargeVis.

# Becht et al., UMAP Preserves Local and Global Structure (Analysis of Tissue T Cells; Color = Expert Knowledge / Source)

(a) UMAP better split CD8 T cells,  $\gamma\delta$  T cells, and contaminating cells

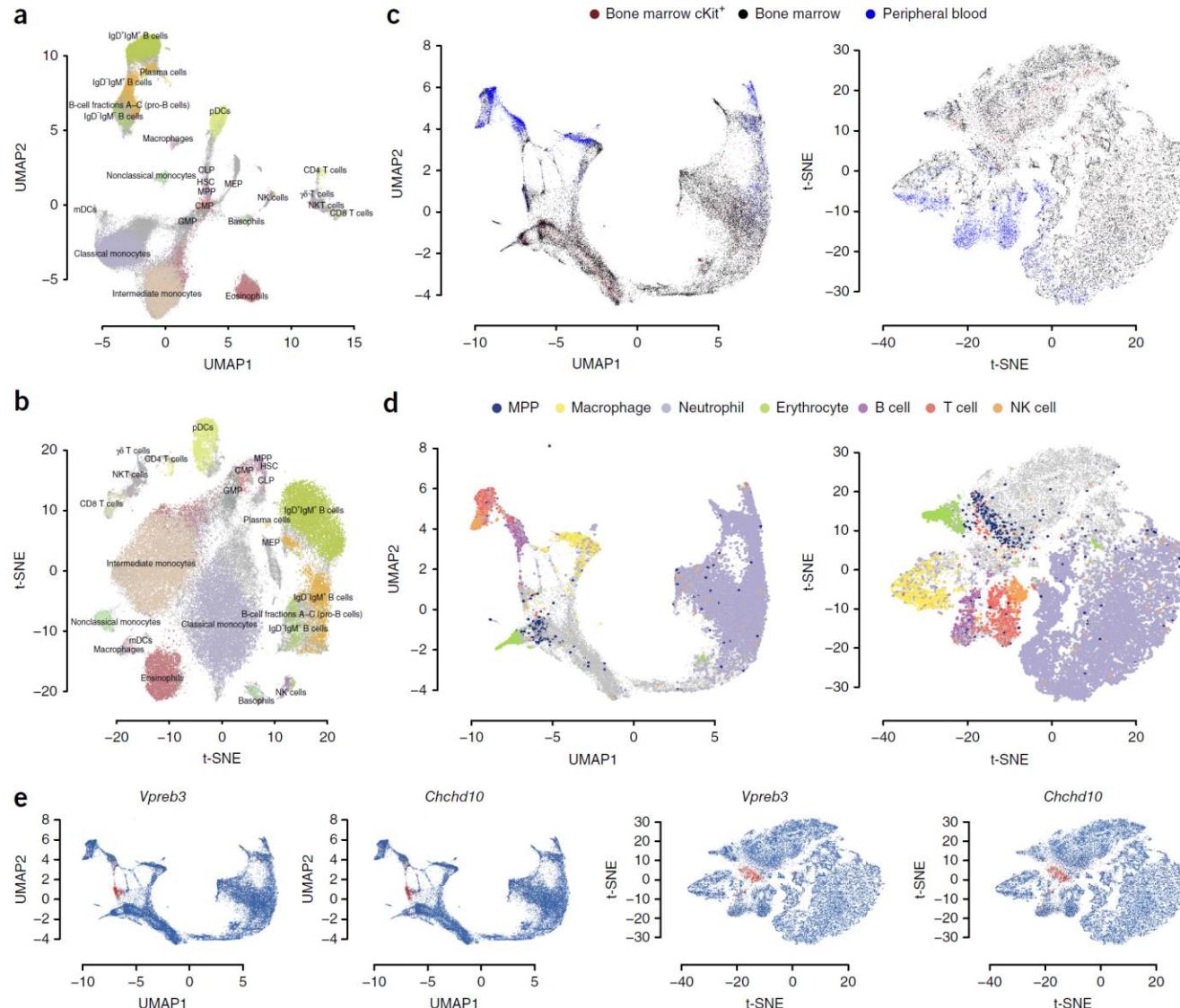
(b) color-coding the tissues of origin; t-SNE separated cell populations according to their tissue of origin more often than UMAP. UMAP instead ordered events according to their origin within each major cluster, roughly from 1) cord blood and peripheral blood mononuclear cells, to liver and spleen, and to tonsils on the one end to skin, gut and lung on the other.  
**Continua not apparent in t-SNE.**

Dataset covering 35 samples originating from 8 distinct human tissues enriched for T and natural killer (NK) cells, of more than >300,000 cell events with 39 protein targets (the Wong dataset).



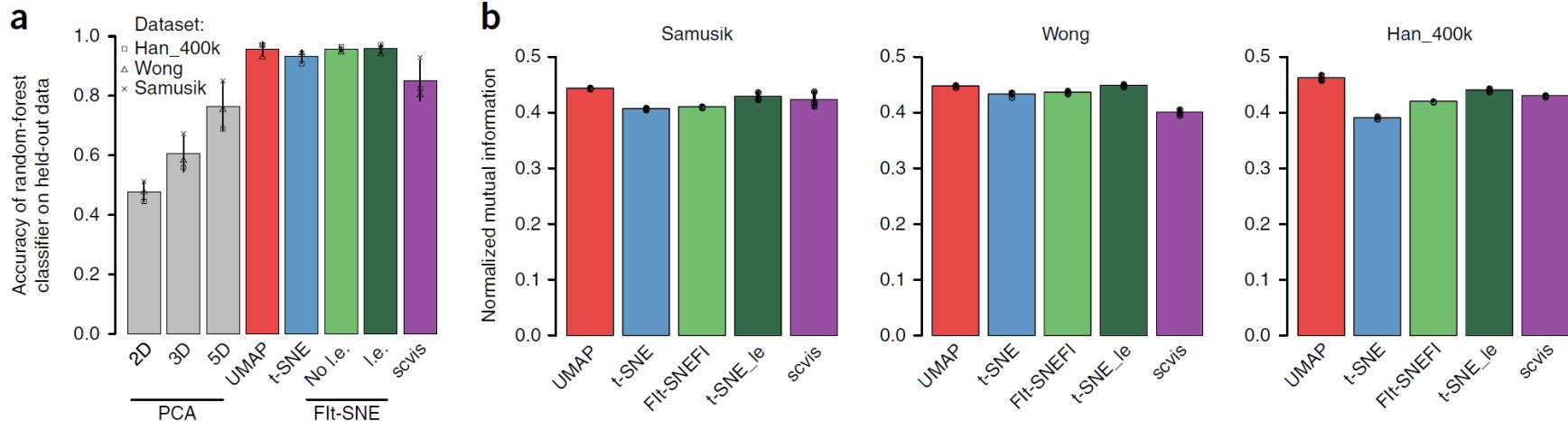
**Figure 1** UMAP embeds local and large-scale structure of the data. UMAP and t-SNE projections of the Wong *et al.* dataset colored according to (a) broad cell lineages, (b) tissue of origin, and for (c) UMAP and (d) t-SNE, the expression of CD69, CD103, CD45RO and CCR7. For c and d, blue denotes minimal expression, beige intermediate and red high. MAIT, mucosal-associated invariant T cell; ILC, innate lymphoid cell; CB, cord blood; PBMC, peripheral blood mononuclear cell.

# Becht et al., UMAP Captures Developmental Trajectories

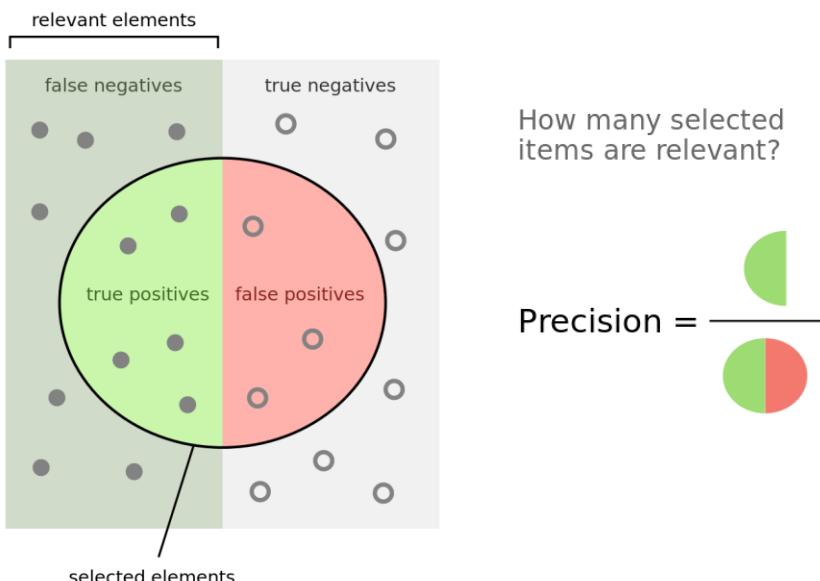


**Figure 2** UMAP embeddings of bone marrow and blood samples recapitulate hematopoiesis. **(a)** UMAP and **(b)** t-SNE projection of the Samusik\_01 dataset. Events are color-coded according to manual gates provided by the authors of the dataset. **(c,d)** UMAP and t-SNE projections of the Han dataset, color-coded by **(c)** tissue of origin or **(d)** cell populations. **(e)** Expression of the V-set pre-B cell surrogate light chain 3 (*Vpreb3*) and *Chchd10* genes on the UMAP and t-SNE projections of the Han dataset. Blue denotes minimal expression, beige intermediate and red high. pDC, plasmacytoid dendritic cell; mDC, myeloid dendritic cell; NKT, natural killer T.

# Speed is Nothing without Accuracy; UMAP is Also Accurate



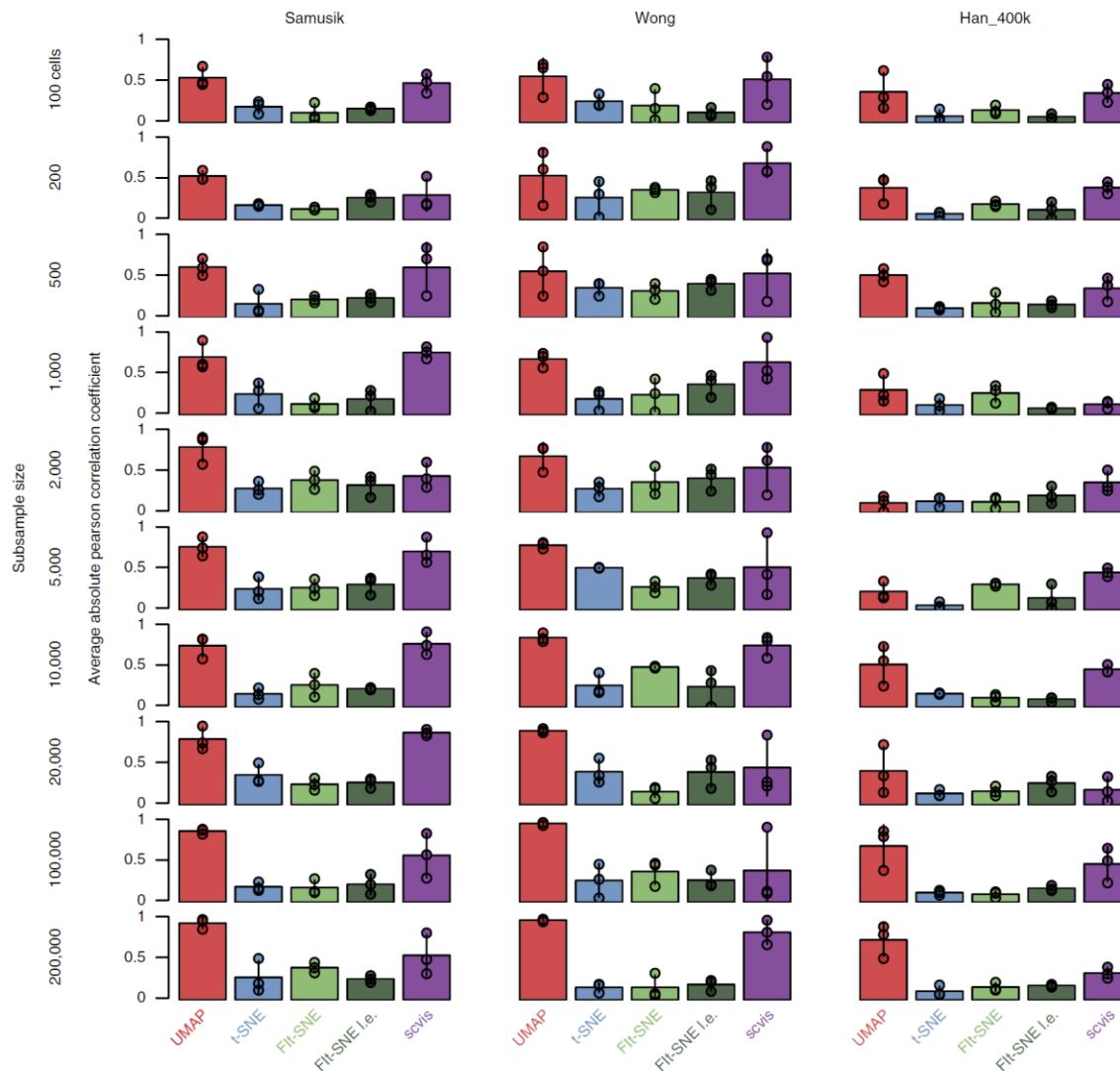
**Figure 4** Analysis of local data structure in embeddings produced by each algorithm. **(a)** Accurate classification rate on held-out data of random-forest classifiers predicting Phenograph cluster labels using embedded coordinates as input. The average across the three datasets is shown, with vertical bars representing s.d. **(b)** Average normalized mutual information of  $k$ -means clustering ( $k = 100$ ) performed on the embeddings of data subsamples and  $k$ -means clustering ( $k = 100$ ) performed on total datasets. The average across the three random subsamples of size 200,000 is shown, with vertical bars representing s.d.



$$F_1 \text{ score} = \text{accuracy}$$

$$2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$$

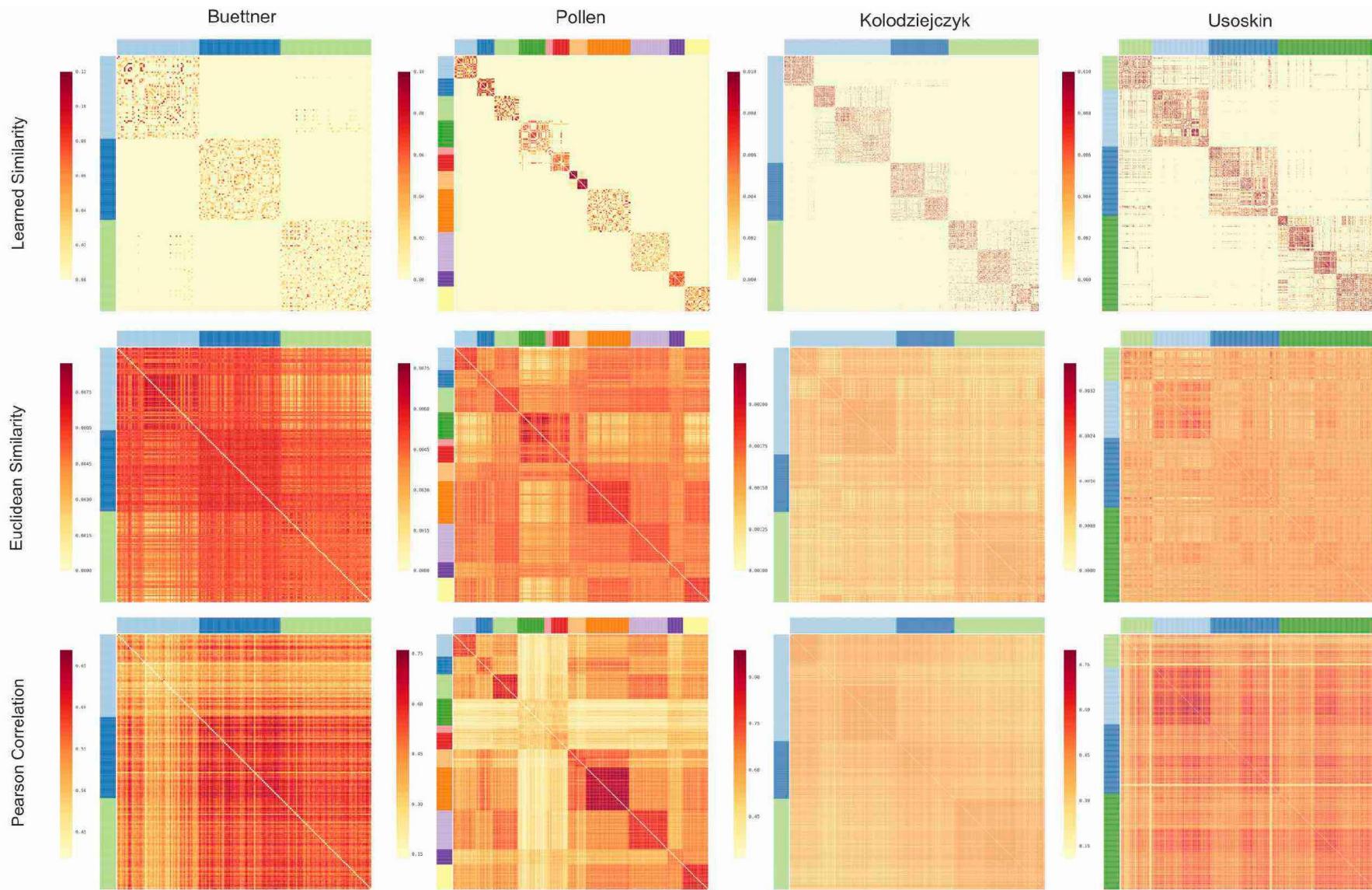
# UMAP Preserves “Large-Scale Structure” That t-SNE Ignores (Large = Position of Islands; Fine = Position of Cells in an Island)



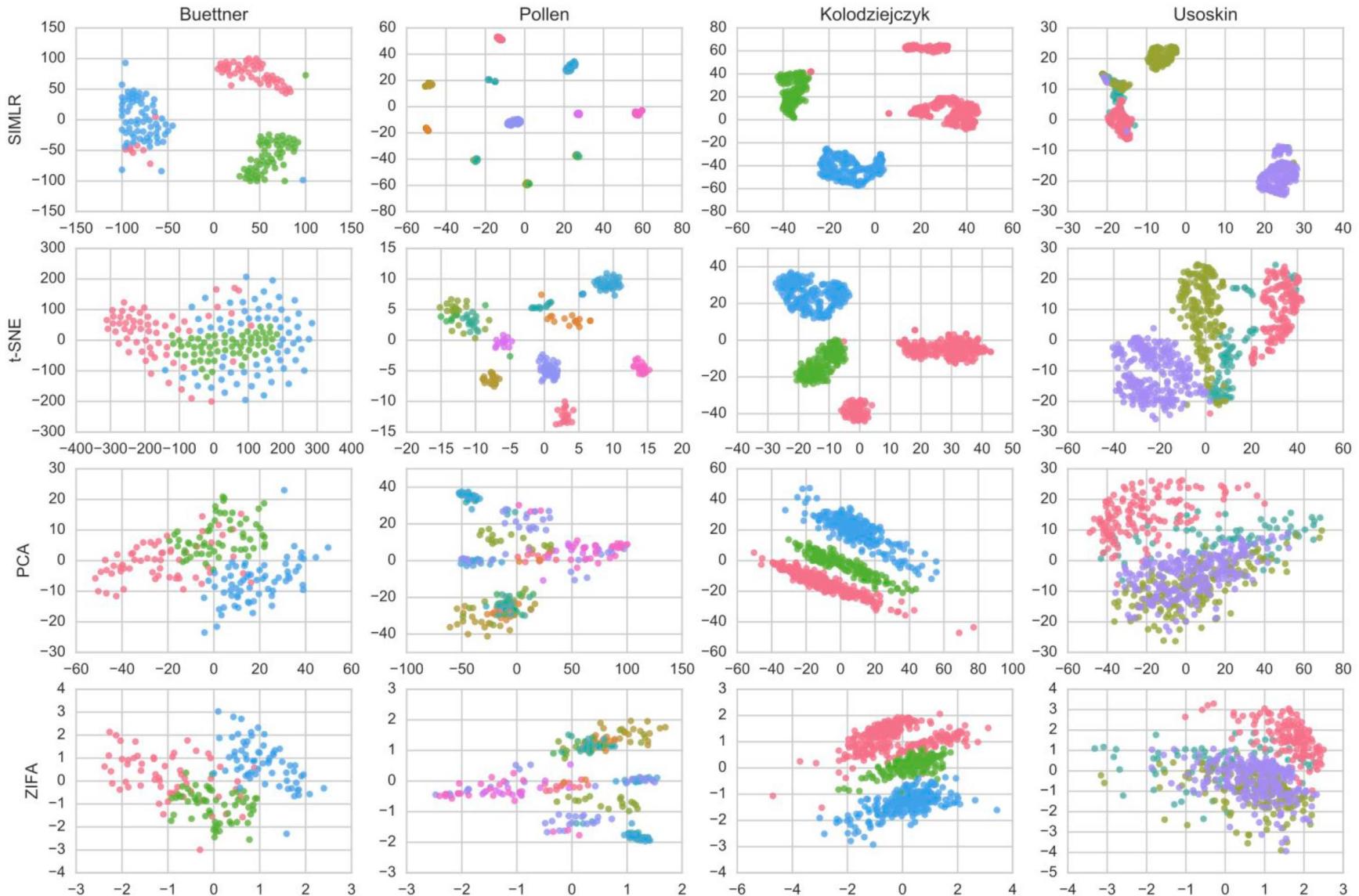
**Figure 6** Reproducibility of large-scale structures in embeddings. Bar plots represent the average unsigned Pearson correlation coefficient of the points' coordinates in the embedding of subsamples versus in the embedding of the full dataset, thus measuring the correlation of coordinates in subsamples versus in the embedding of the full dataset, up to symmetries along the graph axes. Bar heights represent the average across three replicates and vertical bars the corresponding s.d.

SIMLR

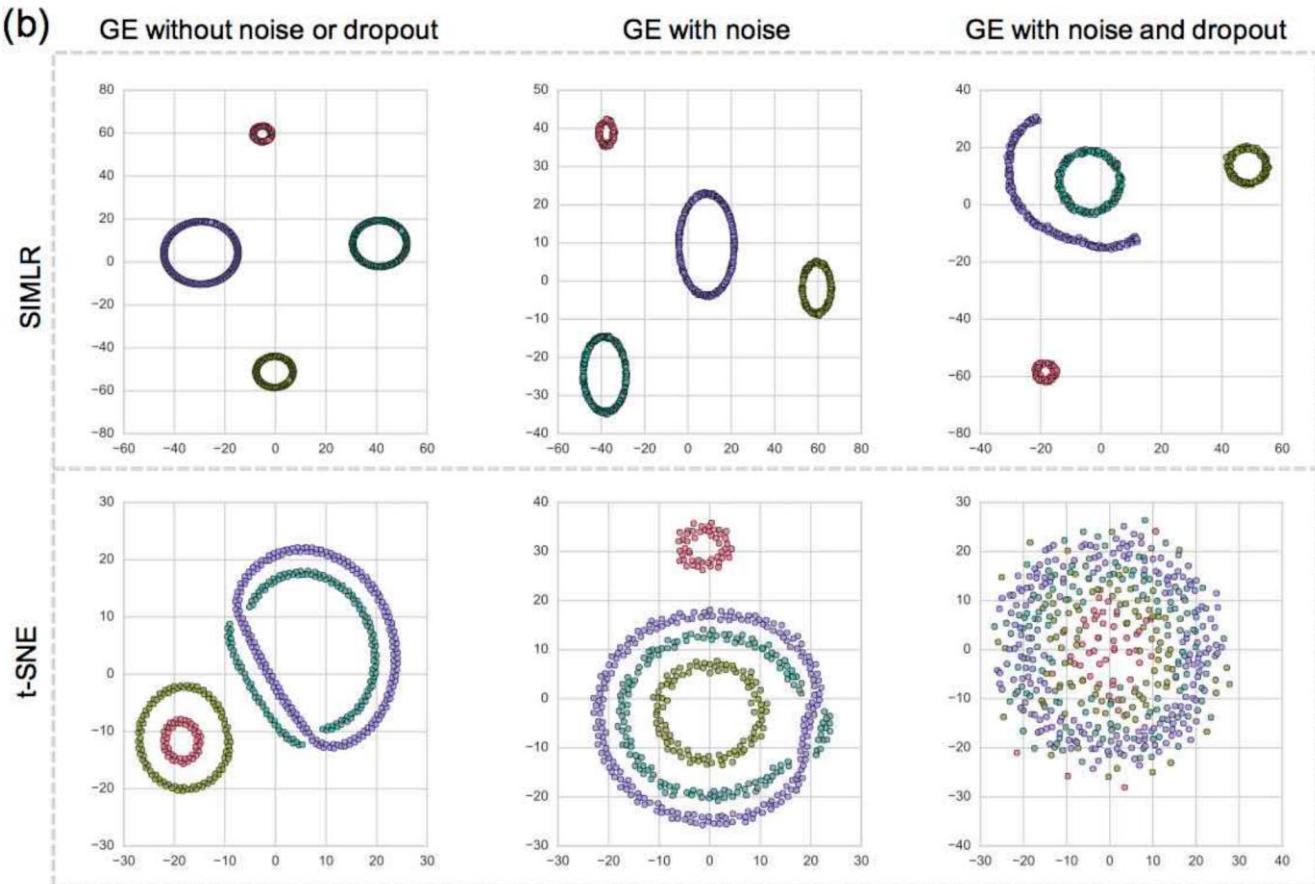
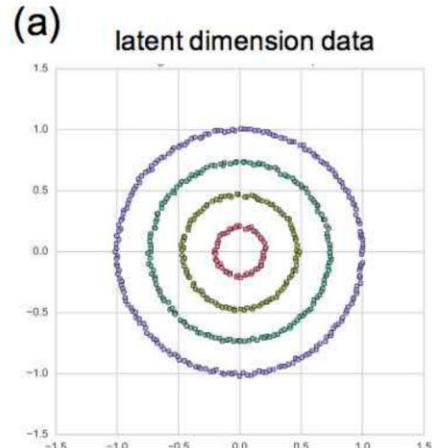
# Visualization and analysis of single-cell RNA-seq data by kernel-based similarity learning (SIMLR)



# SIMLR vs. t-SNE vs. PCA on Four scRNA-seq Datasets



# Analysis of “Toy” Gene Expression Data +/- Noise and Data Dropout



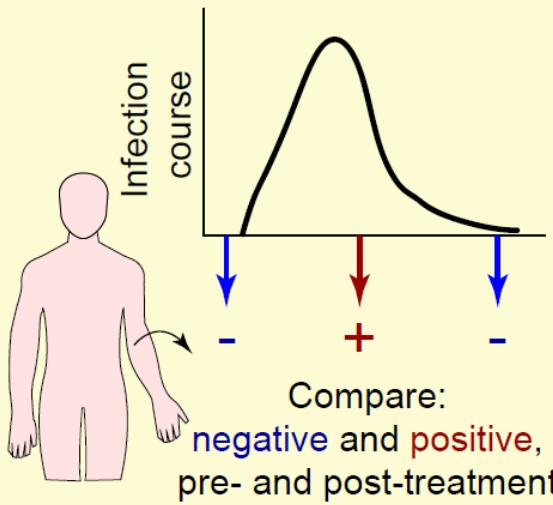
# Part 8: T-REX!



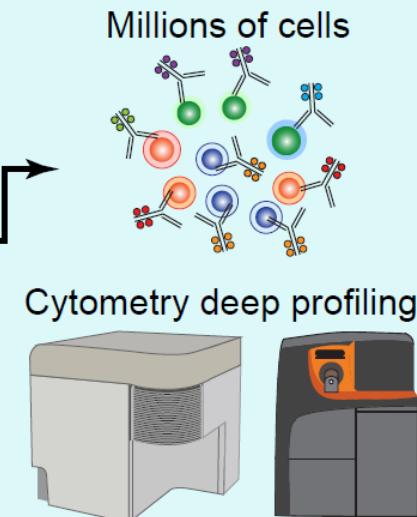
## Pinpointing Ultra-Rare Cells (e.g., virus-specific T cells)

# T-REX: Machine Learning Reveals Virus-Specific Cells in Longitudinal Immune Monitoring Cytometry

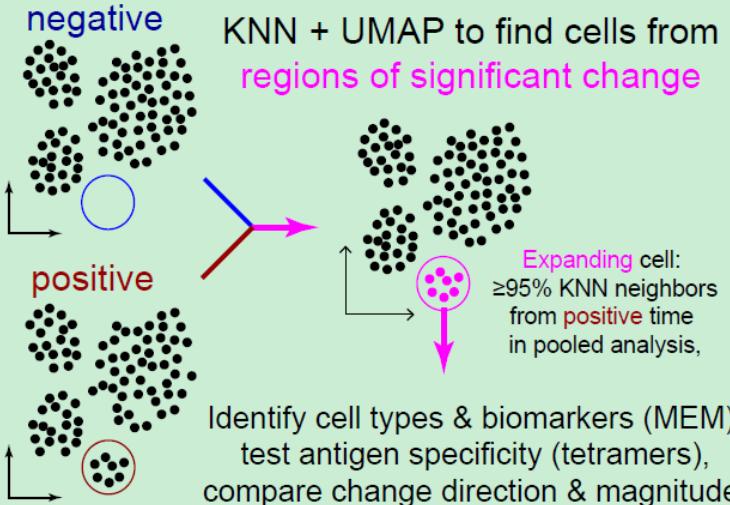
Sample individual over time



Track blood cell types



T-REX algorithm: reveal responding cells



New algorithm: T-REX (Tracking Responders EXpanding)

Code:

<https://github.com/cytolab/t-rex>

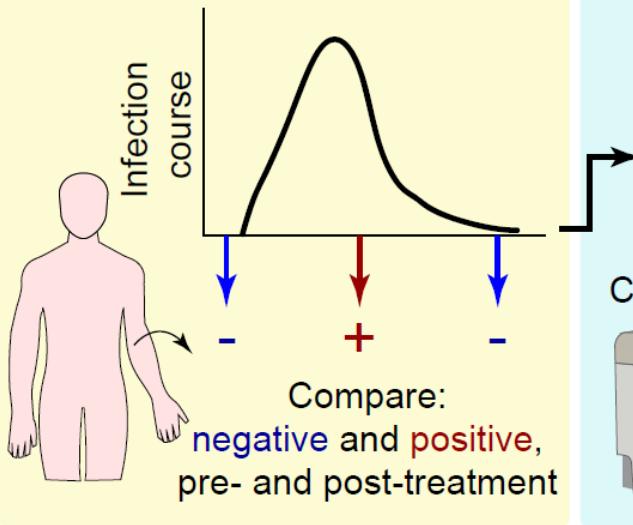
Manuscript:

<https://elifesciences.org/articles/64653>

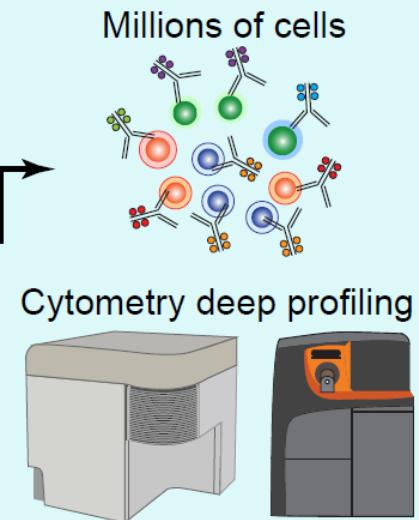


# Key Ideas & Findings in Today's Talk

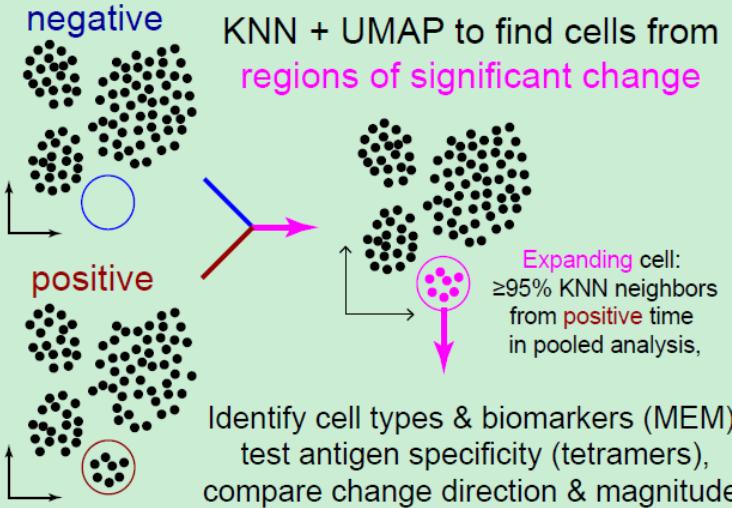
Sample individual over time



Track blood cell types



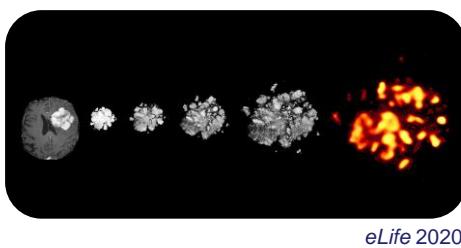
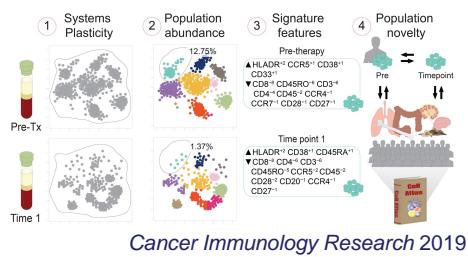
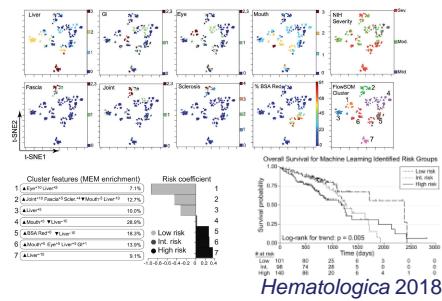
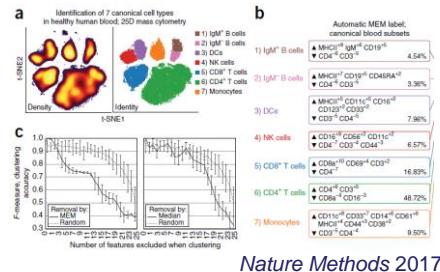
T-REX algorithm: reveal responding cells



- Idea 1: T-REX automatically reveals virus-specific T cells in rhinovirus & SARS-CoV-2 vaccine response (without the need for tetramers, sorting, or sequencing)
- Idea 2: Approach focuses on extreme change & can summarize disease, therapy, or perturbation response (direction & magnitude of change; rhinovirus, COVID-19, cancer therapy, compound screening)
- Finding: Mass cytometry + T-REX characterized SARS-CoV-2 vaccine-induced memory CD4 and CD8 T cells (phenotype: CD38++ ICOS++ CD45R0+ PD-1+ Ki-67+ CXCR5-)
- Finding: This phenotype of SARS-CoV-2 vaccine responding T cells closely matched rhinovirus-specific T cells



# T-REX Builds on Cytometry Data Science Tools



Diggins et al., PMC5330853

- MEM: machine labeling & identification of cell type clusters
- Enabled comparison across single cell platforms
- E.g., memory CD4+ T cells, 10-point enrichment scale:  
 $\text{ICOS}^{+8} \text{CD38}^{+8} \text{CD4}^{+7} \text{CD45R0}^{+6} \text{CD3}^{+5} \text{Ki-67}^{+4}$



Gandelman et al., PMC6312024

- Updated Diggins et al., 2015, precursor for RAPID
- Risk stratified cGVHD patients based on 8 organ scores
- Workflow: t-SNE => FlowSOM on embedding => MEM

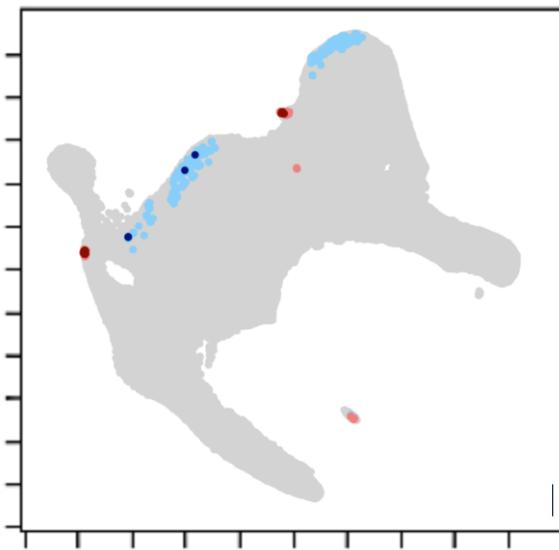
Greenplate et al., PMC6318034

- Set of tools for longitudinal single cell tumor immunology.
- Revealed abnormal immune cells in multiple tumor types.
- Includes datasets (AML, melanoma) used by T-REX
- See Greenplate et al.'s outstanding COVID-19 work (Science 2020)

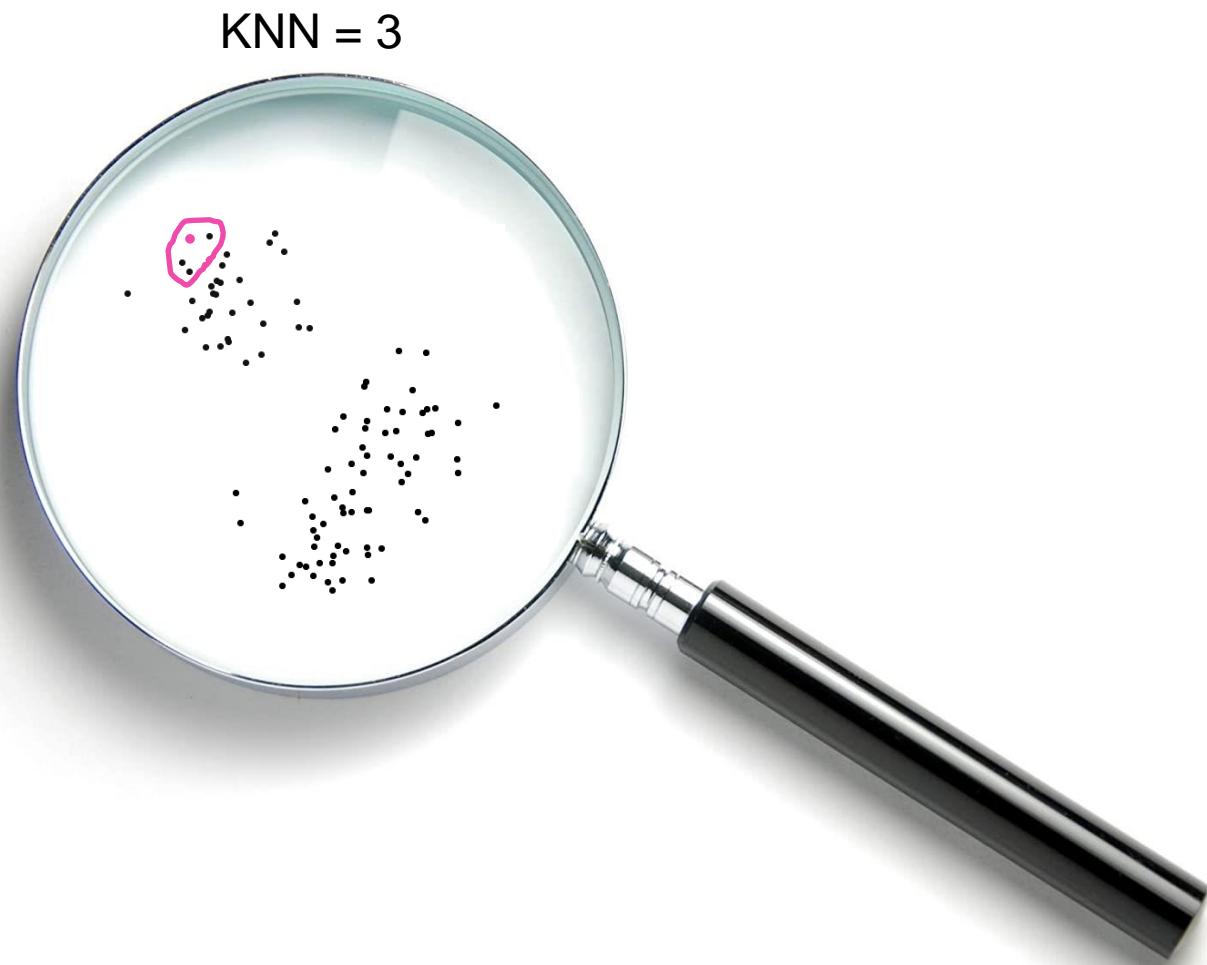
Leelatian and Sinnaeve et al., PMC7340505

- RAPID: maps probabilistic clinical outcomes on t-SNE & UMAP
- Reveals associations with extreme clinical outcomes
- Revealed JAK + AKT cooperation in glioblastoma cells

# T-REX Draws on Sconify's K-Nearest Neighbors (KNN) Approach To Make Each Cell A Phenotypic Neighborhood



T-REX will use  $k = 60^*$   
for cell neighborhoods

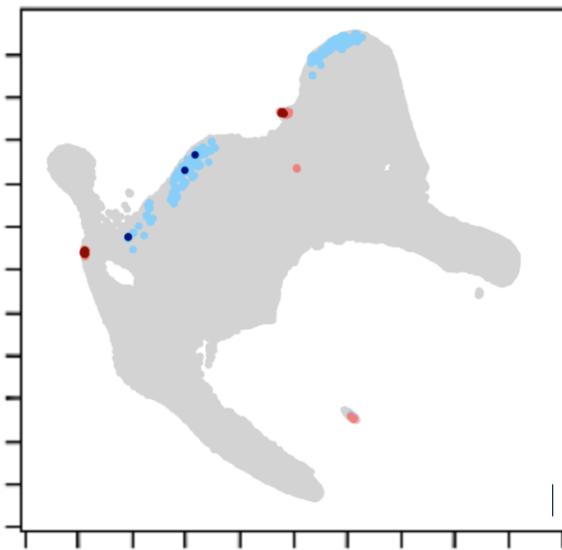


(\*shout out to ISAC's CYTO meeting,  
where questions & poster discussions  
in 2019 helped inspire T-REX)

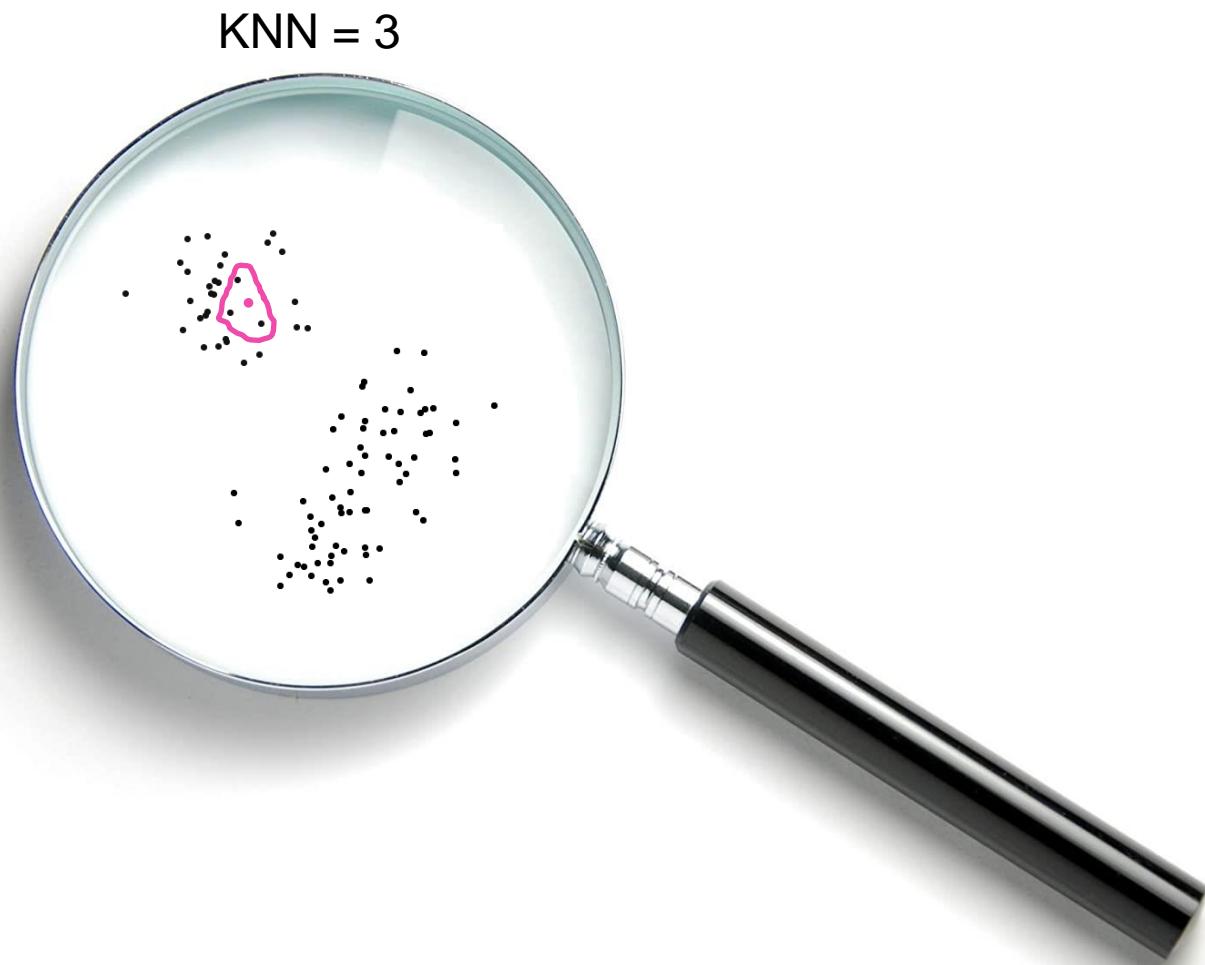


Sconify: Tyler Burns, Garry Nolan, and Nikolay Samusik, bioRxiv 2018; R/bioconductor 2020

# T-REX Draws on Sconify's K-Nearest Neighbors (KNN) Approach To Make Each Cell A Phenotypic Neighborhood



T-REX will use  $k = 60^*$   
for cell neighborhoods

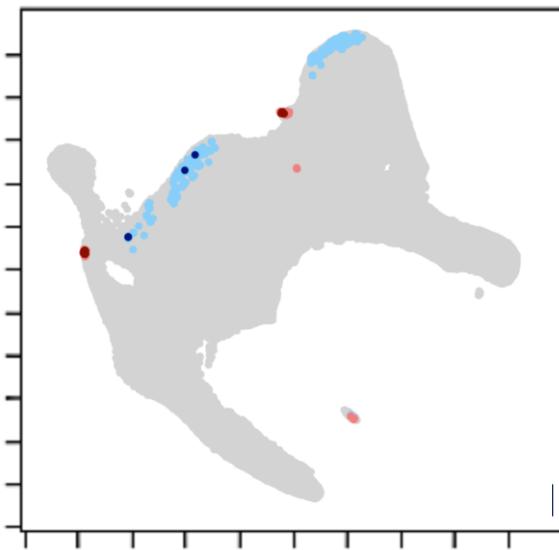


(\*shout out to ISAC's CYTO meeting,  
where questions & poster discussions  
in 2019 helped inspire T-REX)

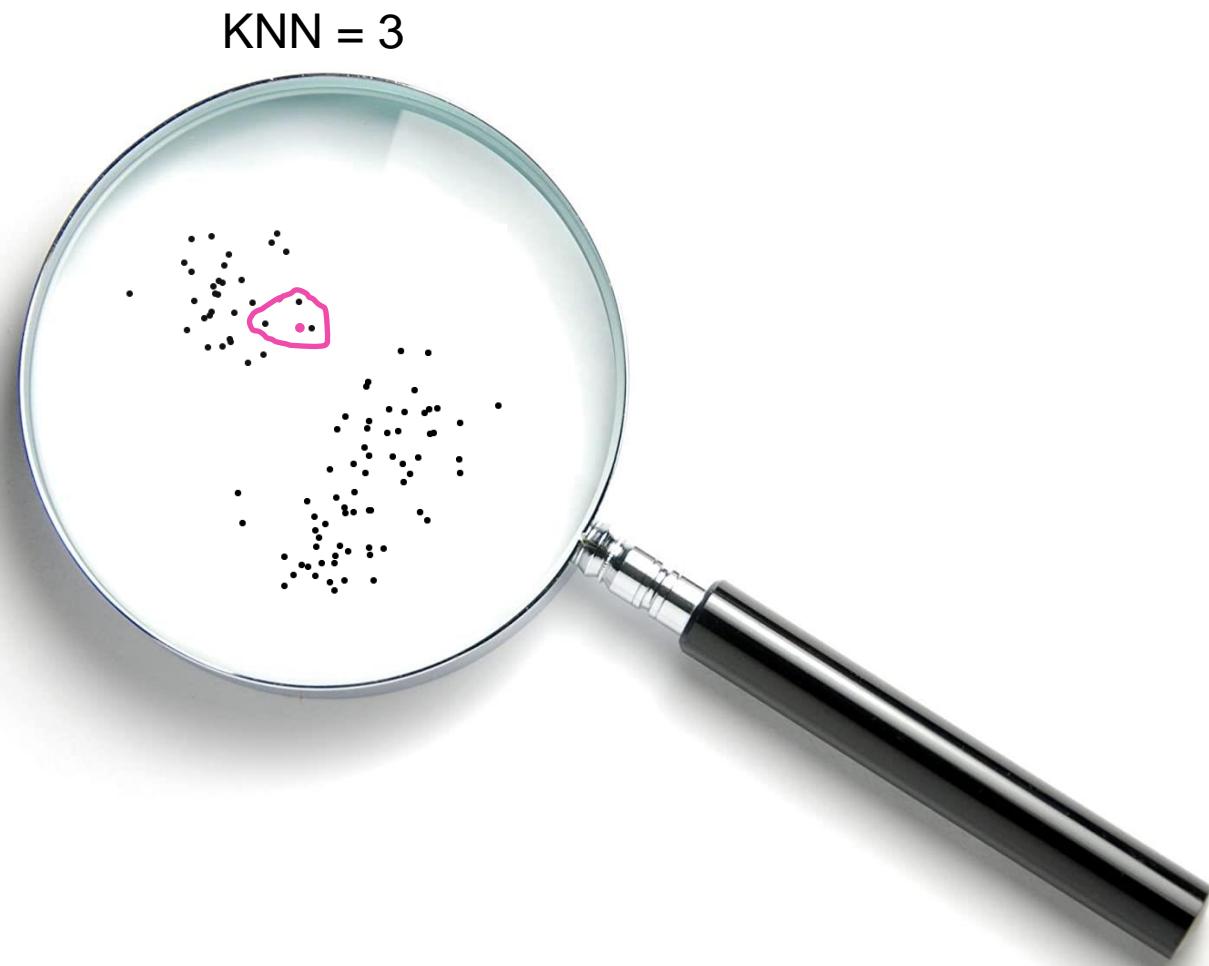


Sconify: Tyler Burns, Garry Nolan, and Nikolay Samusik, bioRxiv 2018; R/bioconductor 2020

# T-REX Draws on Sconify's K-Nearest Neighbors (KNN) Approach To Make Each Cell A Phenotypic Neighborhood



T-REX will use  $k = 60^*$   
for cell neighborhoods



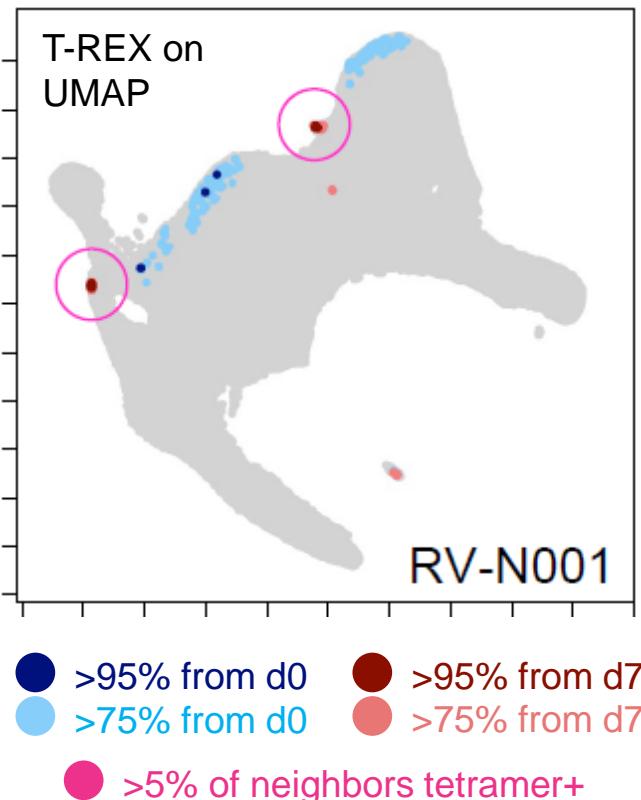
(\*shout out to ISAC's CYTO meeting,  
where questions & poster discussions  
in 2019 helped inspire T-REX)



Sconify: Tyler Burns, Garry Nolan, and Nikolay Samusik, bioRxiv 2018; R/bioconductor 2020

# T-REX: Tracking Responders Expanding, Every Cell Is Characterized in a Search for Hotspots of Change

Live CD4+ T cells



CD4 T cells, Day 0 vs. Day 7,  
individual infected with rhinovirus (RV-N001)  
no cell enrichment, Aurora data,  $\sim 3 \times 10^6$  cells

Color: cells in that phenotypic neighborhood  
are mostly from one sample

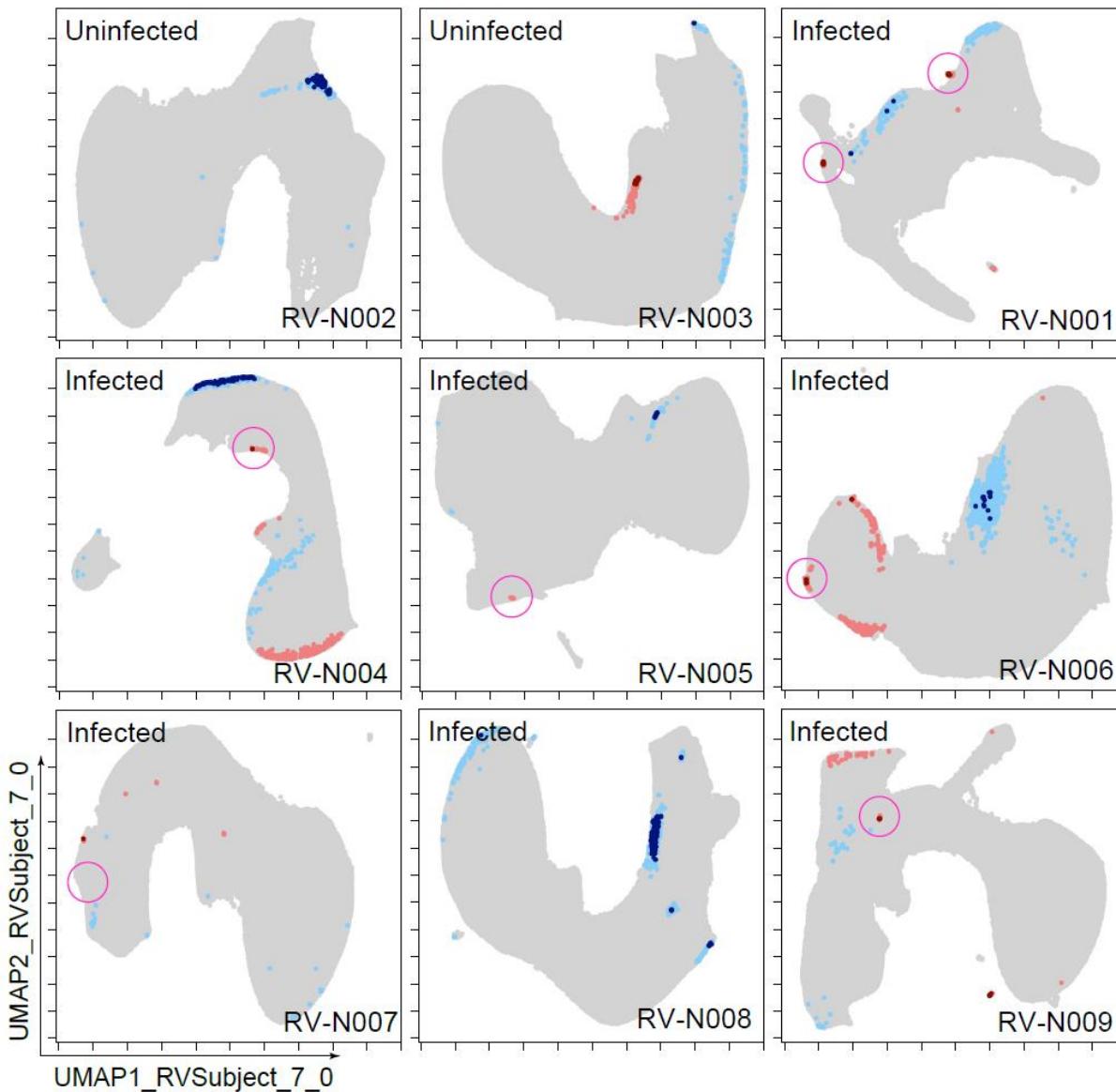
Dark red = cells mostly from day 7 (expanding)



MHCII tetramers marking rhinovirus specific CD4 T cells were not used to make the UMAP, instead used to show: Change hotspots were enriched for virus-specific T cells

# In Analysis of a Rhinovirus Challenge Cohort, T-REX Revealed Virus-Specific Cell Phenotypes

●  $\geq 95\%$  from day 0 ●  $> 85\%$  from day 0 ●  $> 85\%$  from day 7 ●  $\geq 95\%$  from day 7 ○ Tetramer+ hotspot



CD4 T cells, Day 0 vs. Day 7,  
individuals infected with rhinovirus  
no cell enrichment, Cytek Aurora data

In 5 of 7 infected individuals, expansion  
hotspots were enriched for virus-specific cells



The phenotype of rhinovirus-specific memory  
CD4+ T cells included:  
CCR5+ ICOS+ CD38+ PD-1+ CXCR5-

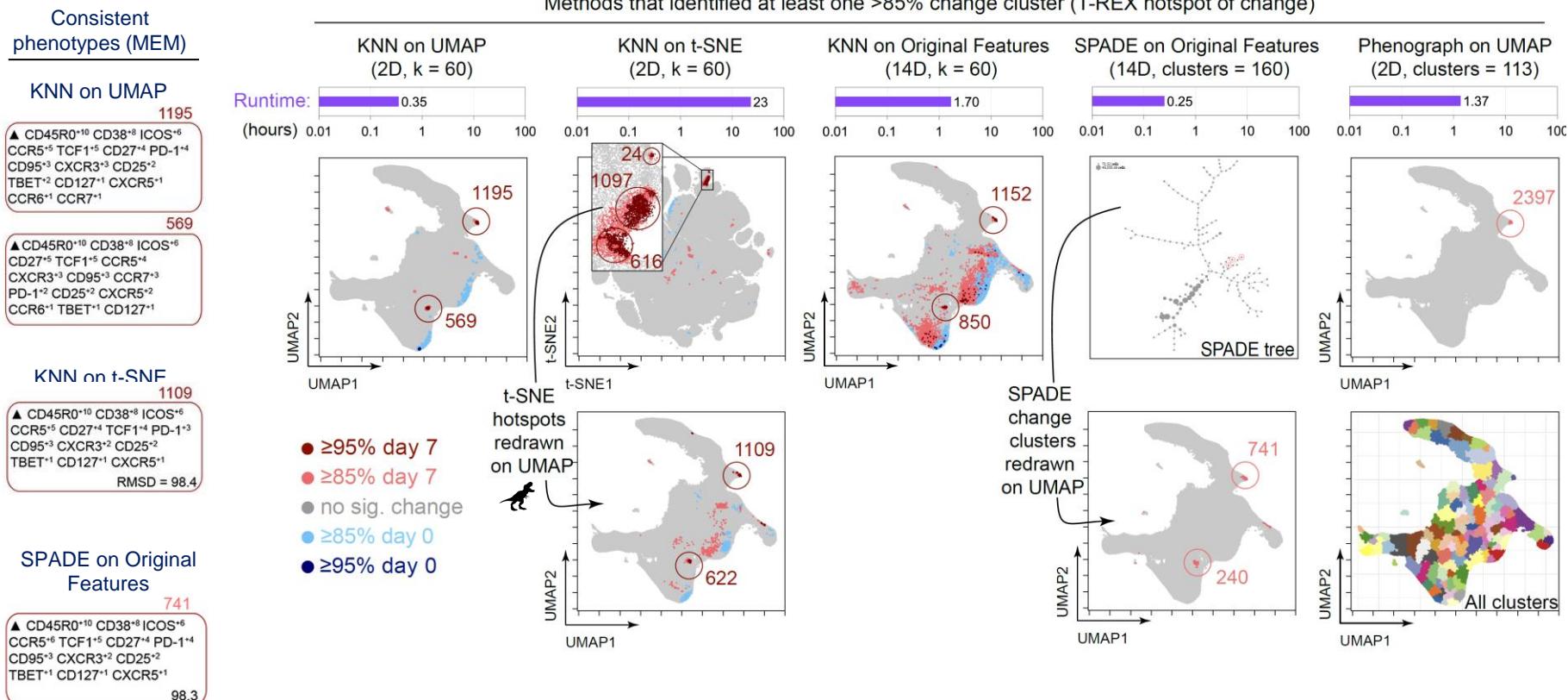
Gating on MEM enriched proteins =>  
identified tetramer+ cells  
(without gating on tetramers),  
suggesting: FACS sort based on MEM label

## T-REX revealed virus-specific T cells without tetramers



Would this approach work with other clustering algorithms?  
Is it 'OK' to do KNN on UMAP axes as parameters?  
(Perhaps: all embeddings are wrong, but some are useful...)

# T-REX Worked with Other Algorithms to Identify Comparable Cells, But KNN on UMAP or t-SNE Outperformed KNN on Original Features



---

## T-REX revealed virus-specific T cells without tetramers

Also found to work for:

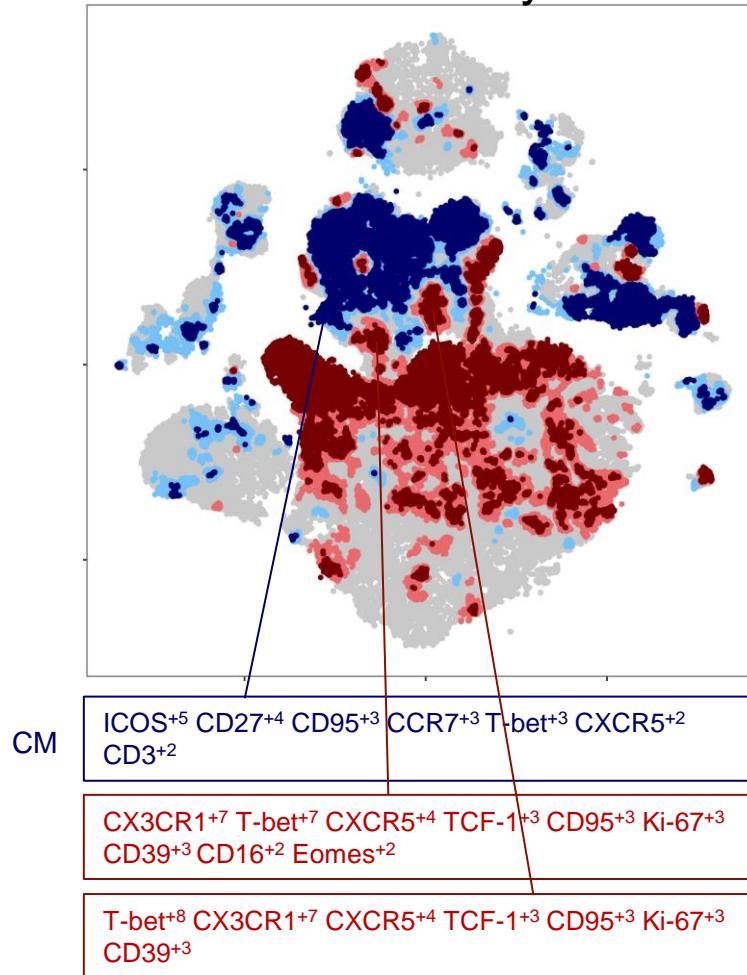
- a range of k-values ( $k = 60$  was optimal)
- post-infection as the comparison point to day 7
- data from a range of cytometers, studies, and labs
- COVID-19, melanoma immunotherapy response, AML

(see the manuscript for this & more!)

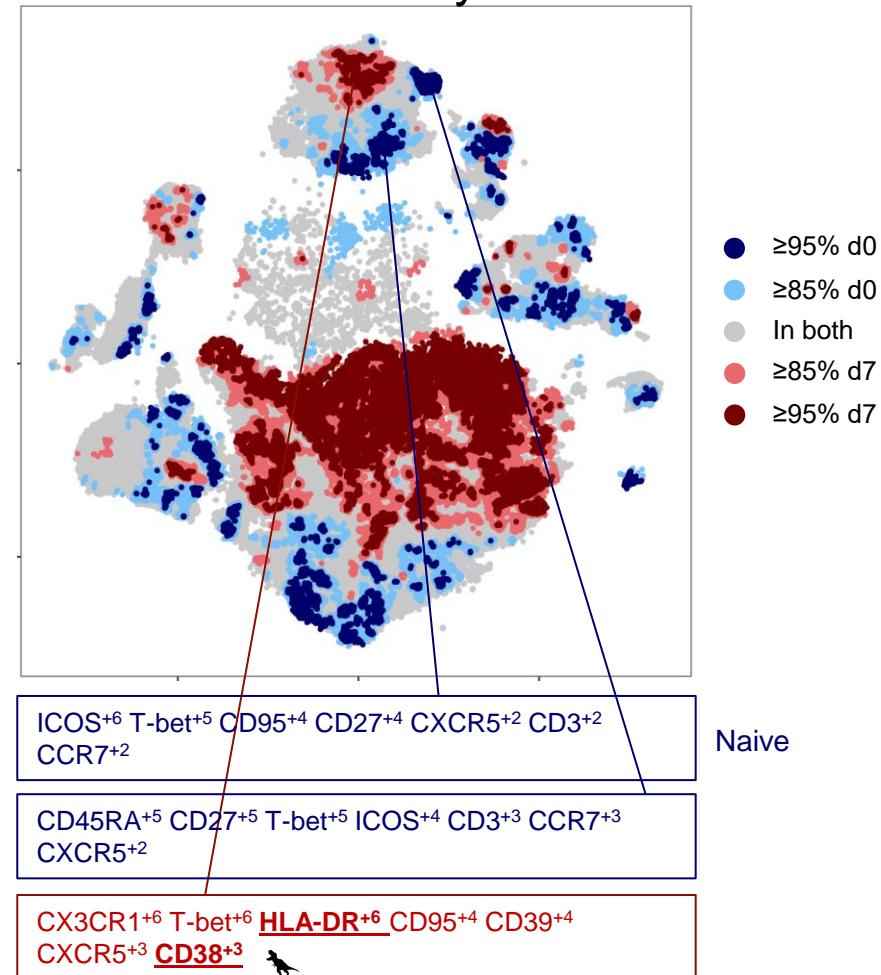


# Massive Immune Change, Common Shifts in Expanding Cell Subsets Observed Between Day 0 and Day 7 in COVID-19

COV-994535 Day 0 vs. 7

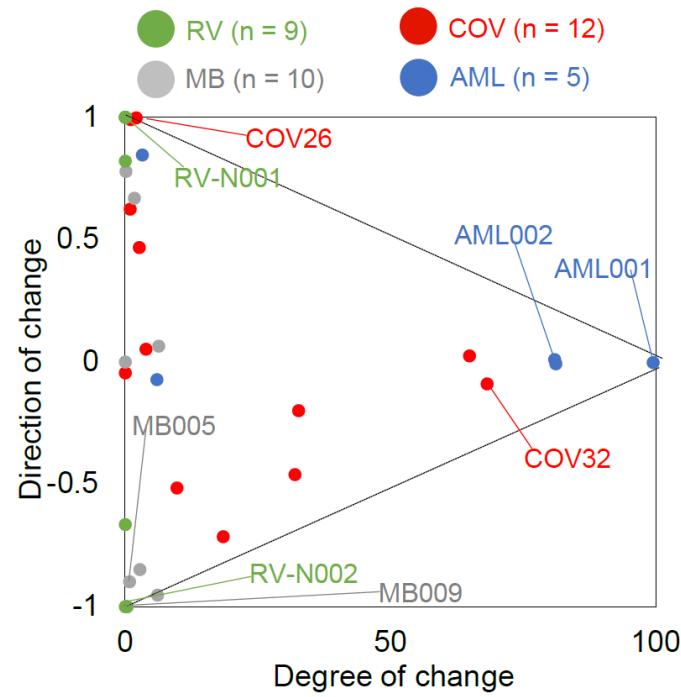


COV-994536 Day 0 vs. 7

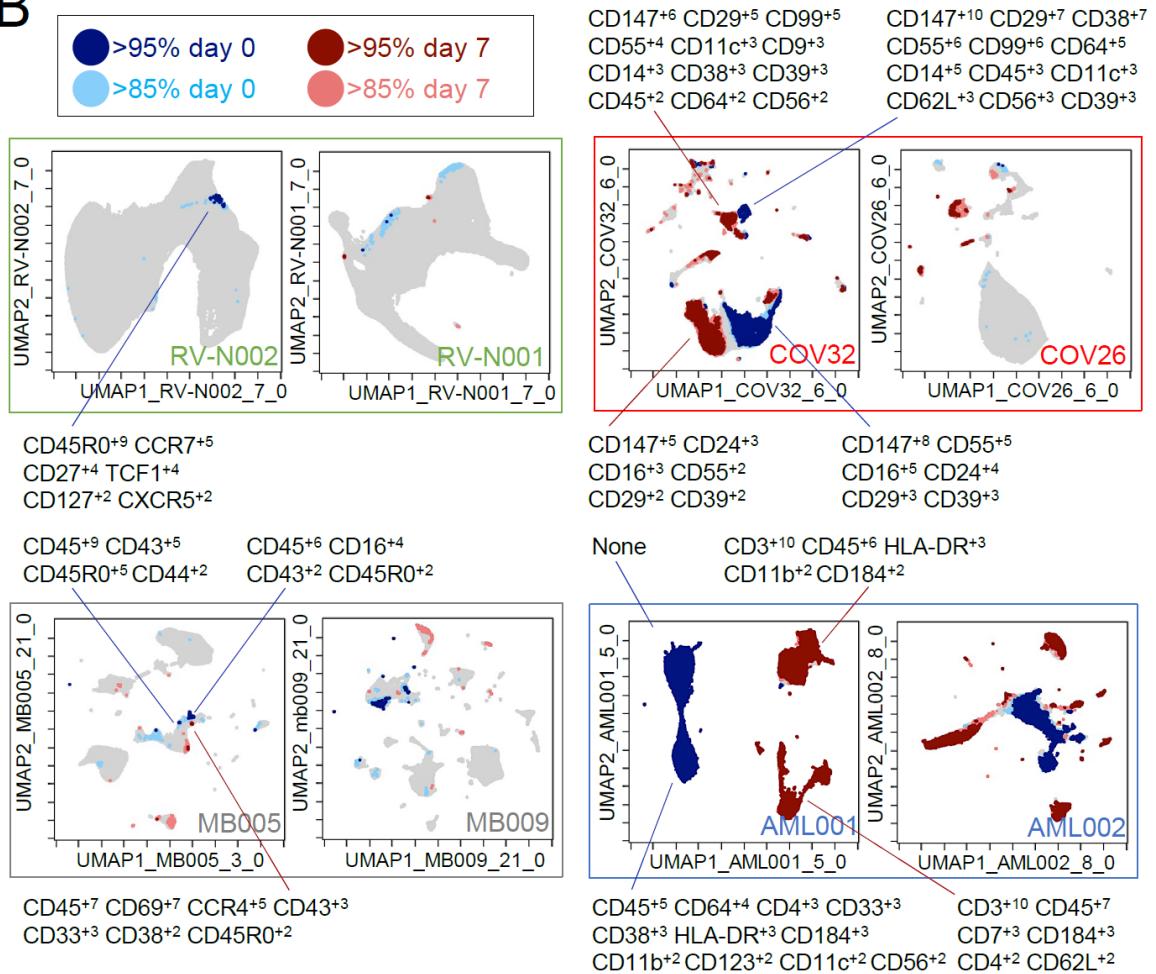


# Half of COVID-19 Patients Displayed Immune Changes Comparable to AML Patients with a Complete Response to Chemotherapy

A



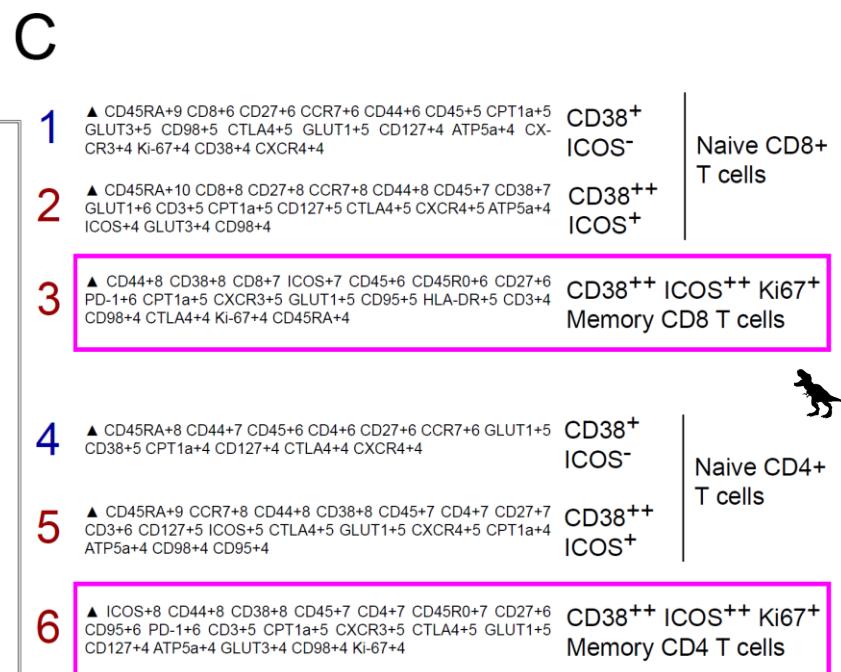
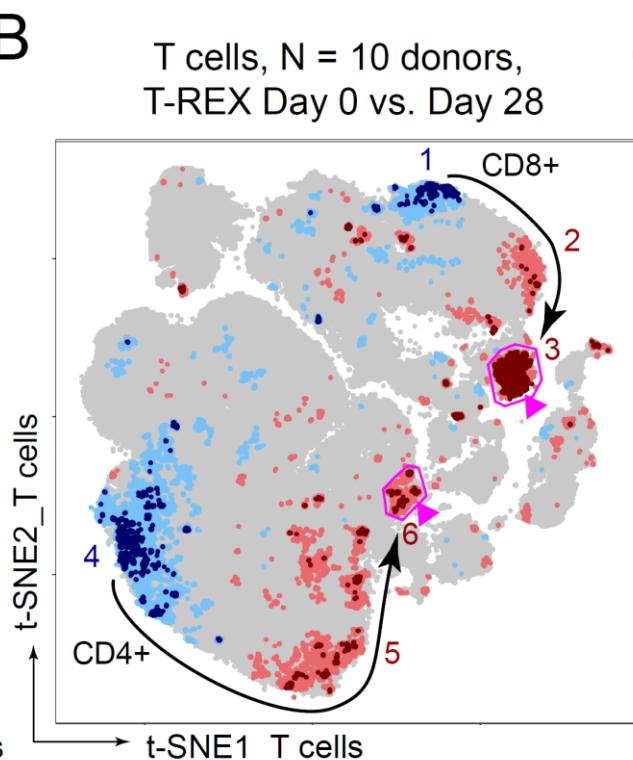
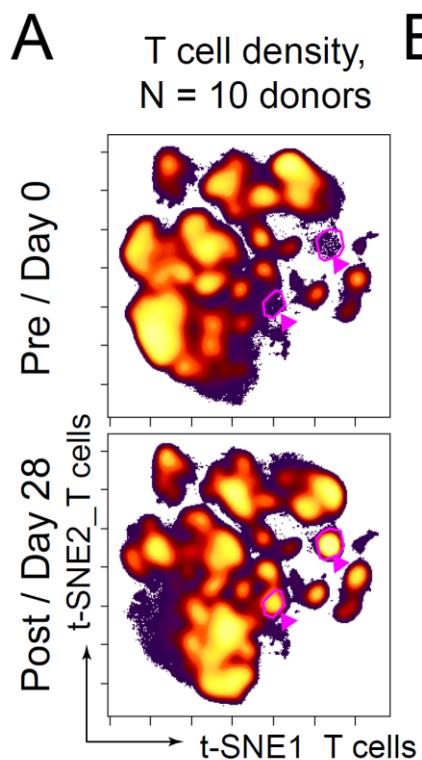
B



T-REX revealed virus-specific T cells without tetramers  
& characterized massive immune changes in COVID-19

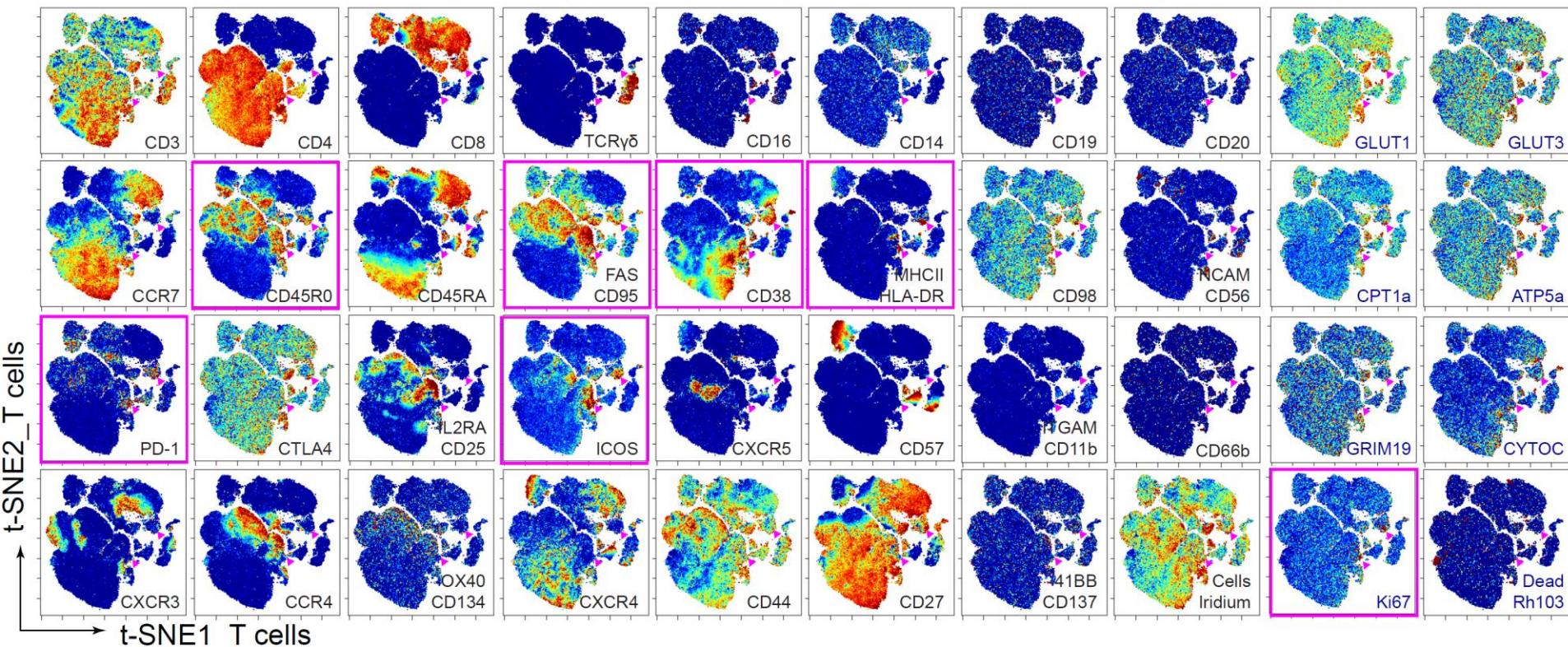
Would it also work to characterize SARS-CoV-2 vaccine response? 

# T-REX Reveals Memory CD4 & CD8 T Cell Phenotypes Expanding following BNT162b2 SARS-CoV-2 RNA Vaccine

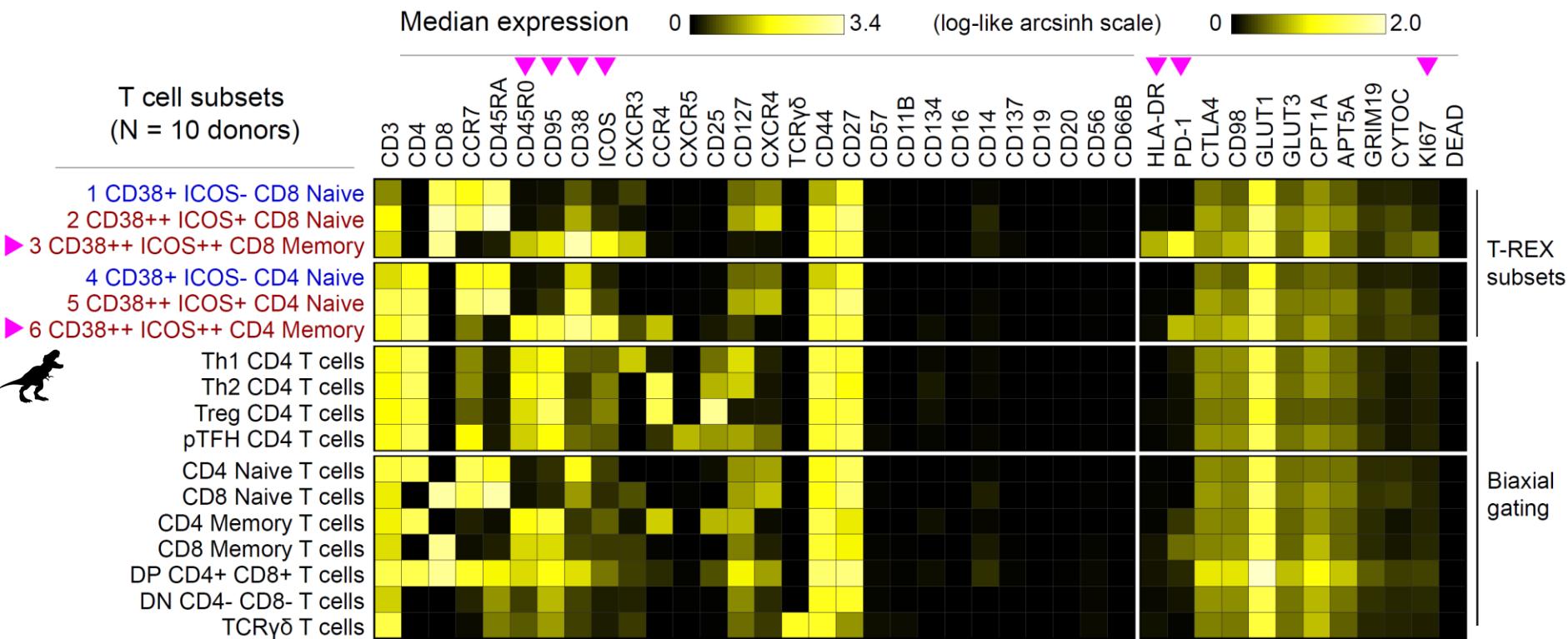


# Mass Cytometry Phenotyping of ICOS+ CD38+ PD-1+ Ki-67+ CXCR5- Memory CD4 & CD8 T Cells following SARS-CoV-2 Vaccination

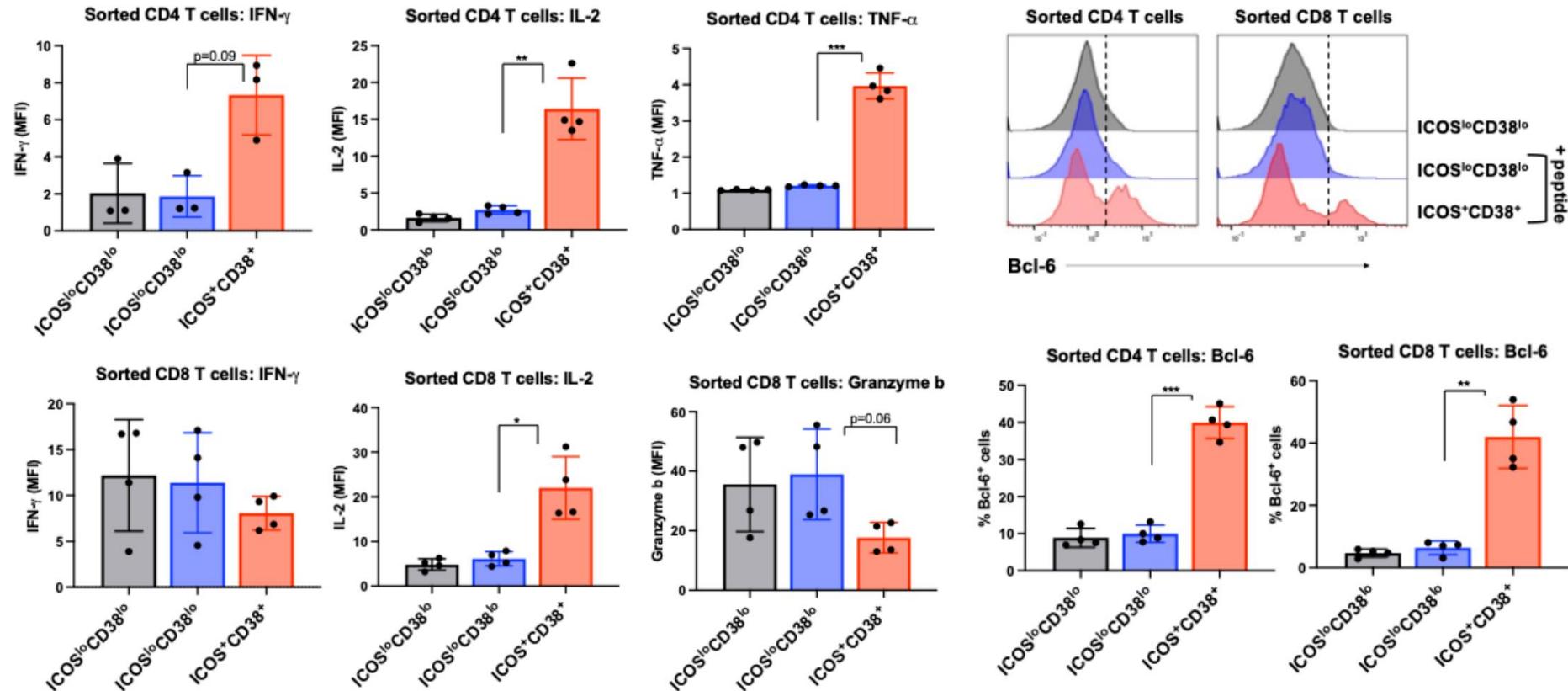
T cell mass cytometry panel on merged post-vaccine data (Day 28, N = 10)



# Mass Cytometry Phenotyping of ICOS+ CD38+ PD-1+ CXCR5- Memory CD4 & CD8 T Cells following SARS-CoV-2 Vaccination



# ICOS++ CD38++ T Cells from Day 28 Post-Vaccine Responded to SARS-CoV-2 Spike Peptide (Autologous PBMCs)



T<sub>FH</sub>/T<sub>FC</sub>? Only half of these cells were BCL-6+, and the cells from T-REX were CXCR5-



T-REX revealed virus-specific T cells without tetramers, characterized massive immune changes in COVID-19, & identified a SARS-CoV-2 reactive non-canonical memory T cell that expands by day 28 following RNA vaccination

Check out the pre-print for more, including plasmablasts, B cell LIBRA-seq, and a breakthrough case who did NOT generate the ICOS+ CD38+ T cells.



# T-REX & COVID-19 Acknowledgements

## Irish Lab at Vanderbilt University + Cancer & Immunology Core



Stephanie Medina  
PhD Student



Amanda Kouaho  
VU Undergraduate



Sierra Barone Lima  
Data Science  
Program Coordinator



Todd Bartkowiak  
PhD Postdoc  
K00 Fellow



Caroline Roe  
CIC/MCCE  
Senior Research  
Specialist



Madeline Hayes  
Lab Development  
Program  
Coordinator

## University of Virginia Collaborators (Rhinovirus T-REX)



Alberta Paul



Lyndsey Muehling



Judith Woodfolk



Kevin Kramer



Erin Wilfong



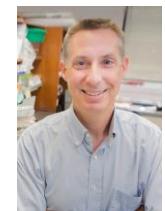
Kelsey Voss



Rachel Bonami



Ivelin Georgiev



Jeff Rathmell

Thank you for the invitation!



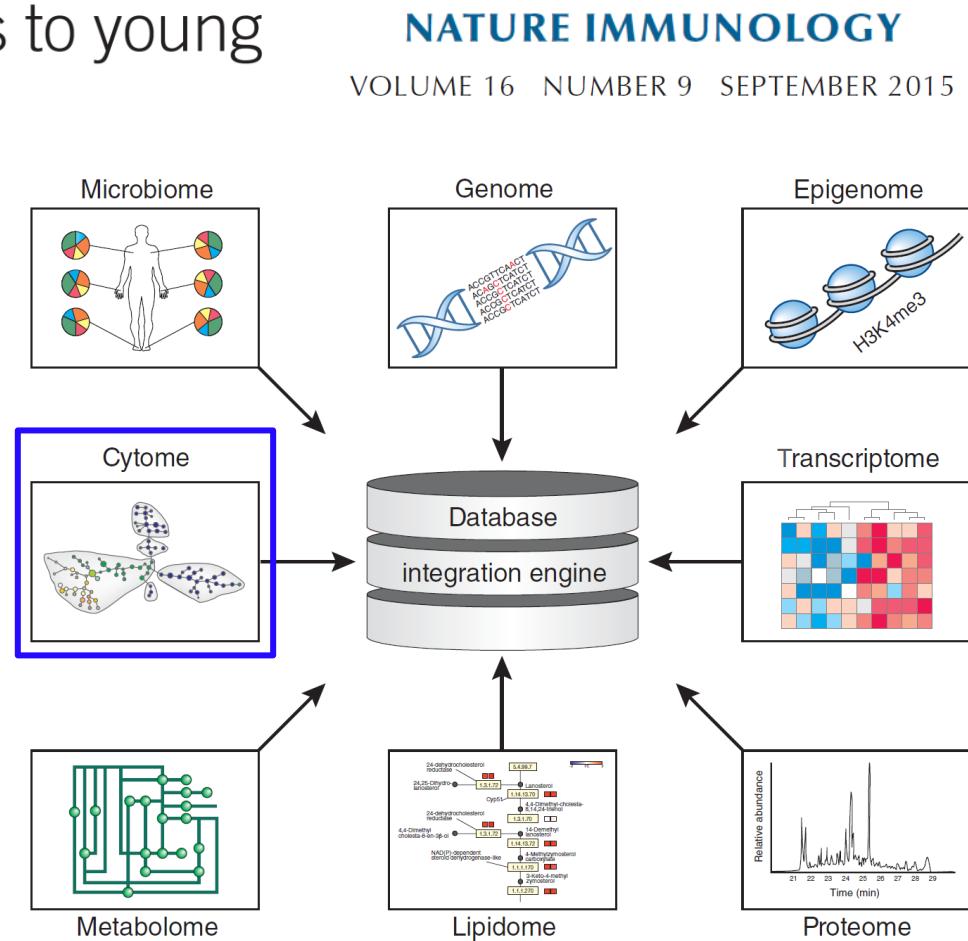
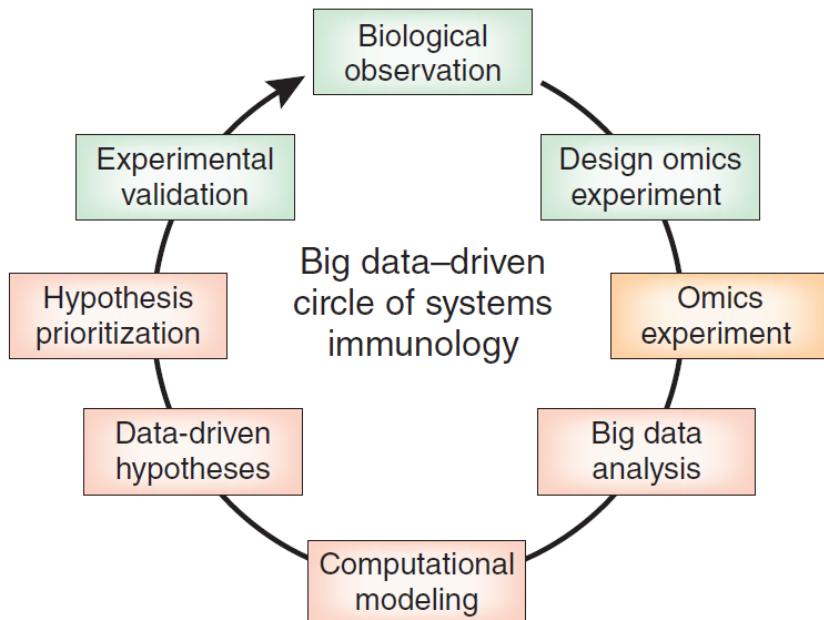
# Part 2: Quantifying Cell Biology & Cytometry Tools

The data type for today: **cytometry**  
(quantitative single cell measurements)

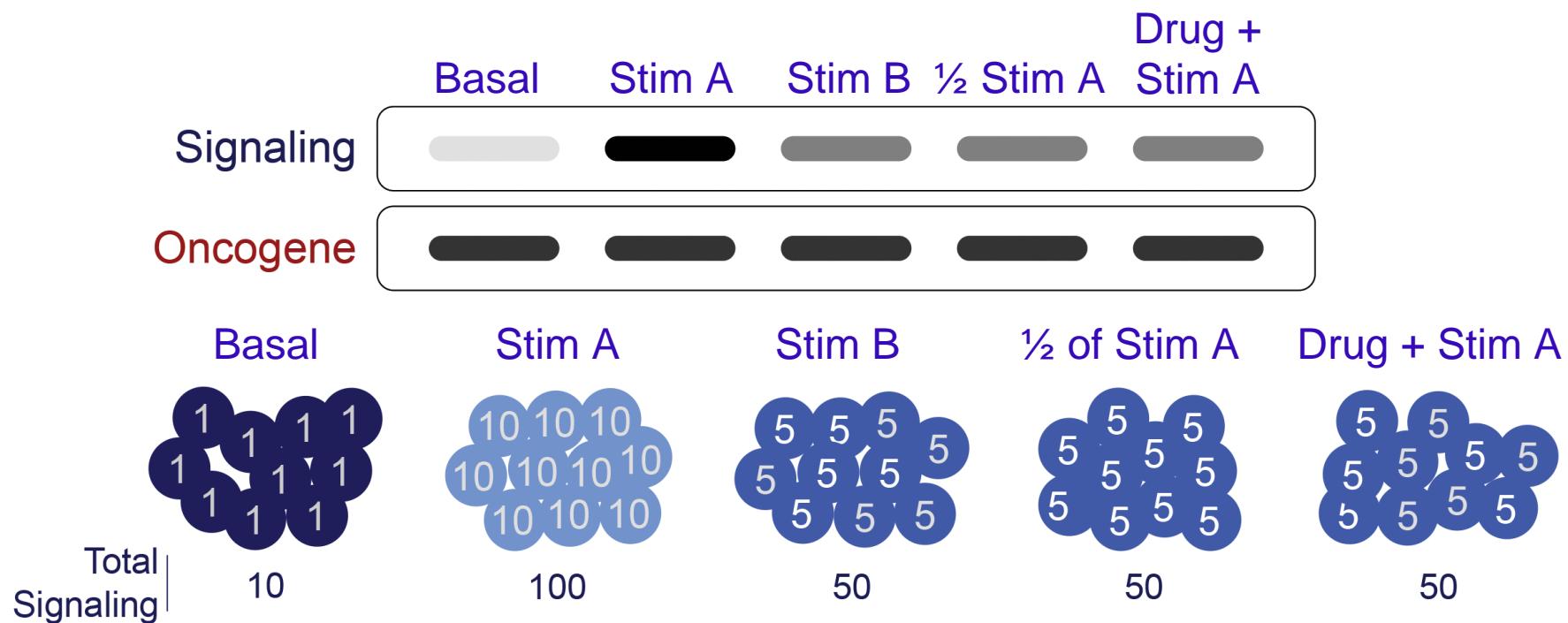
# Cytomics: The ‘Omics of Cells & Cell Identity

Teaching ‘big data’ analysis to young immunologists

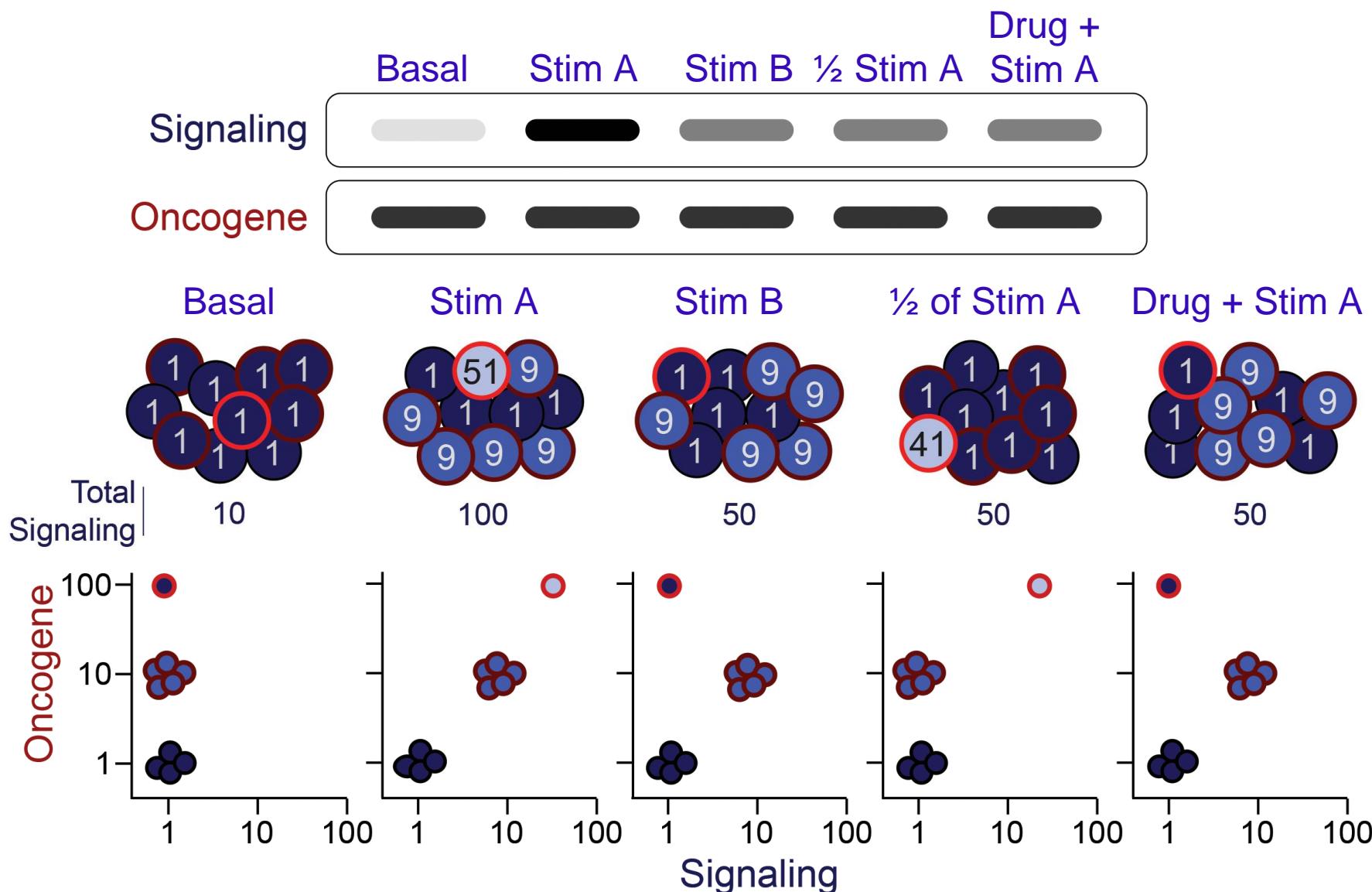
Joachim L Schultze



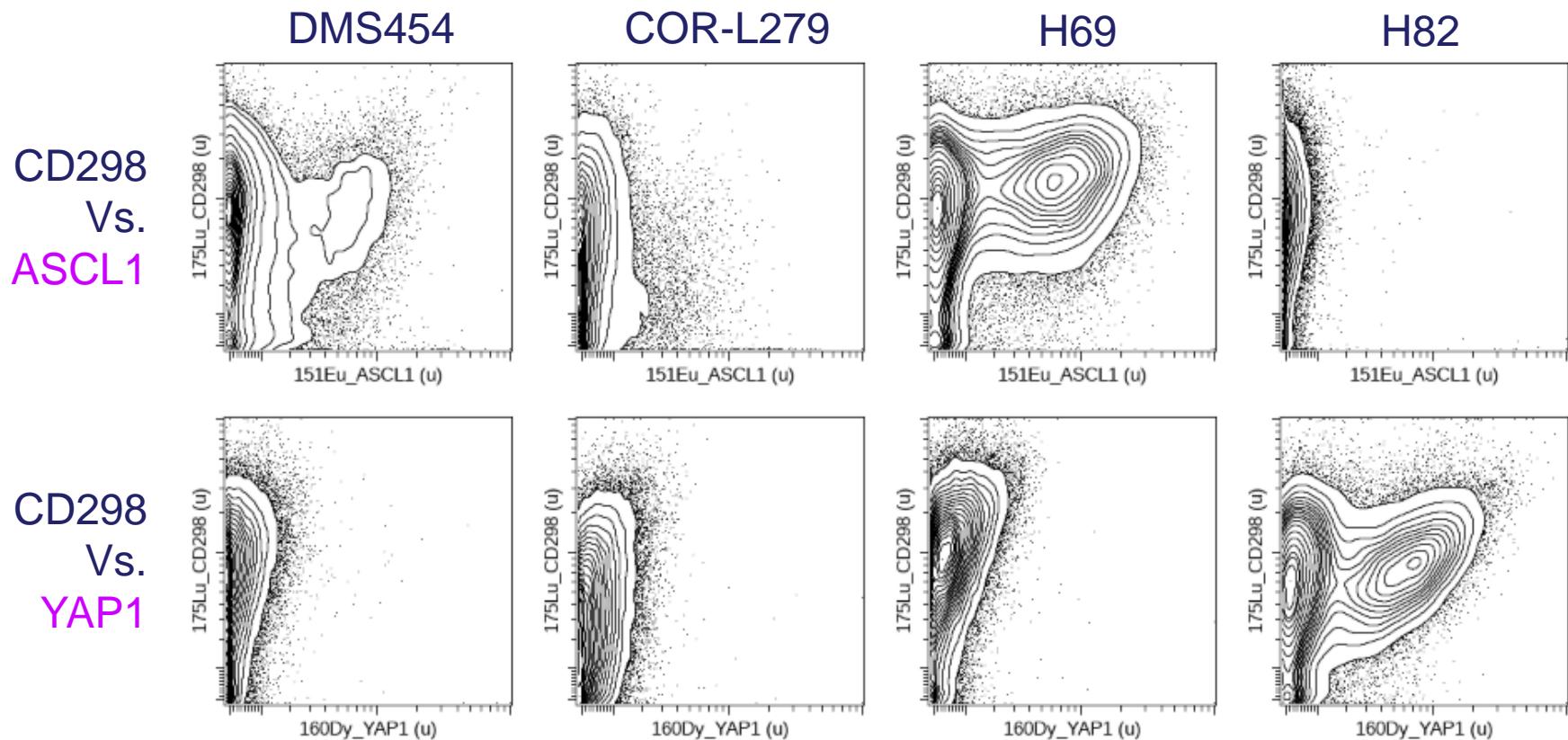
# Single Cell Biology: Which Cell, How Much?



# Single Cell Biology: Which Cell, How Much?



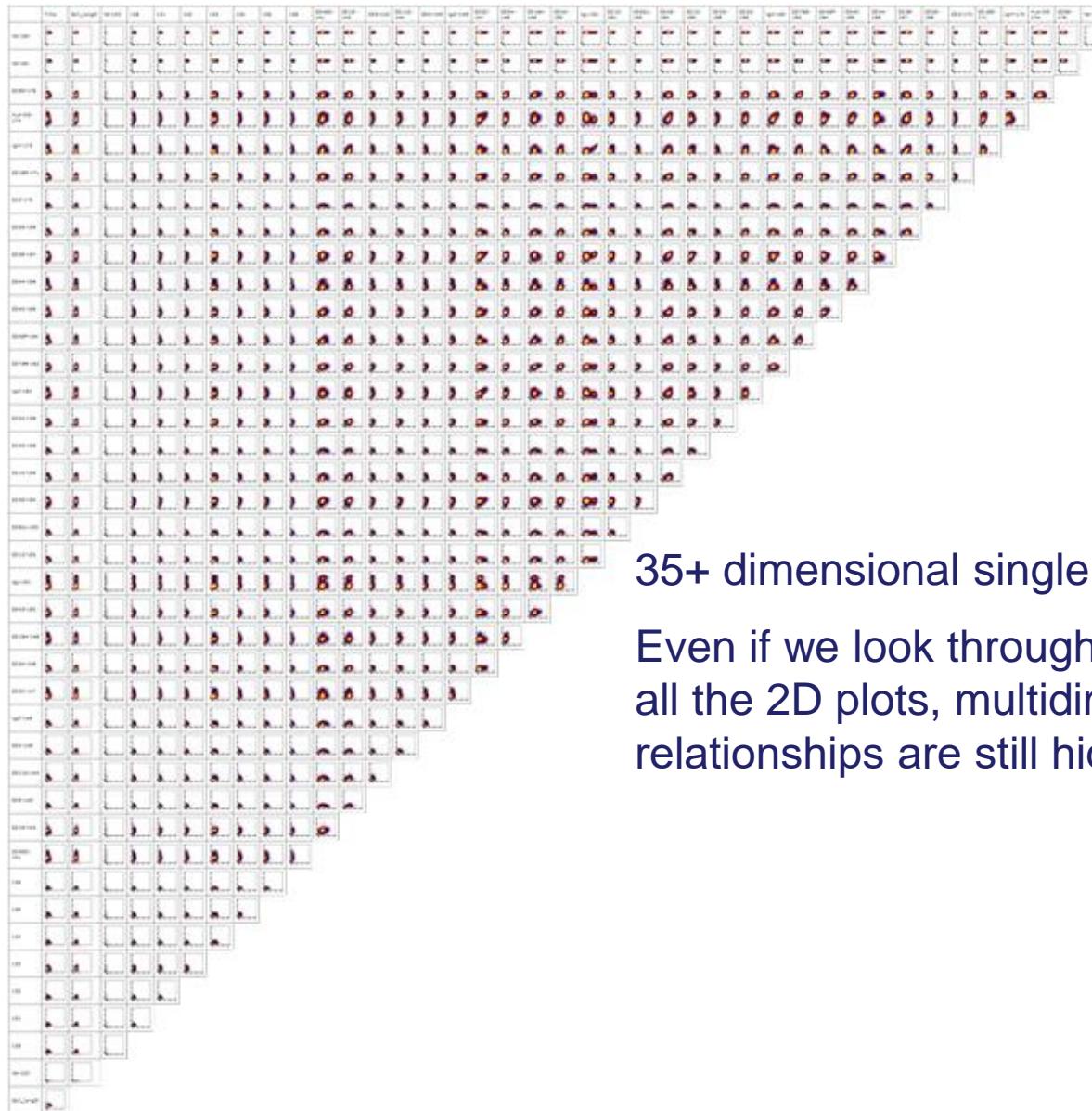
# Cytometry Analysis Reveals Stable Subpopulations with Key Subtype-Defining Features Even in ‘Clonal’ Cell Lines (!)



CD298 (y-axis, Marker for All cells) vs. **ASCL1** & **YAP1** (x-axis, SCLC Subtypes),  
Staining on CyTOF (mass cytometry)

CD298-ATP1B3 Is a “Universal” Human Cell Marker Developed for Tracking PDXs (Lawson et al., *Nature* 2015 Werb) & Cellular Barcoding (Hartmann et al., *Scientific Reports* 2018)

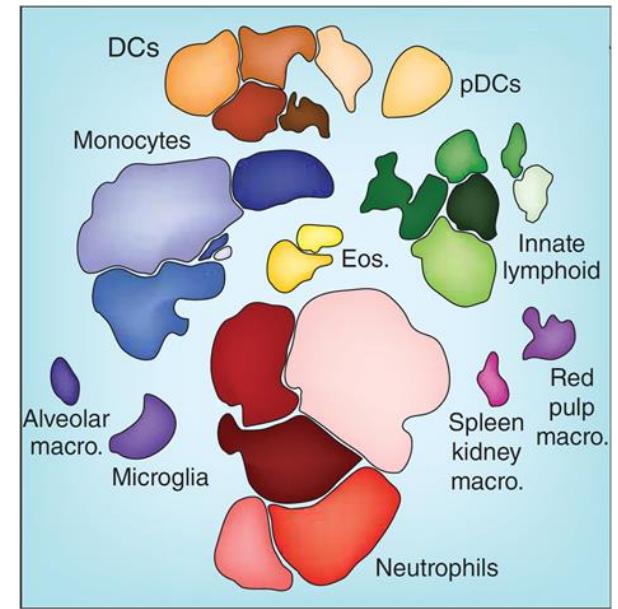
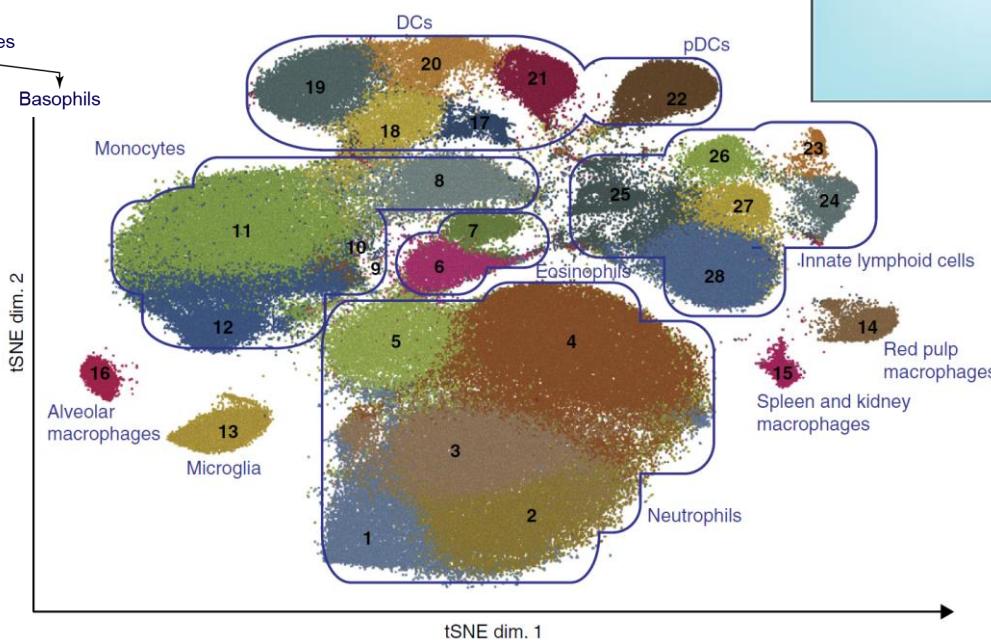
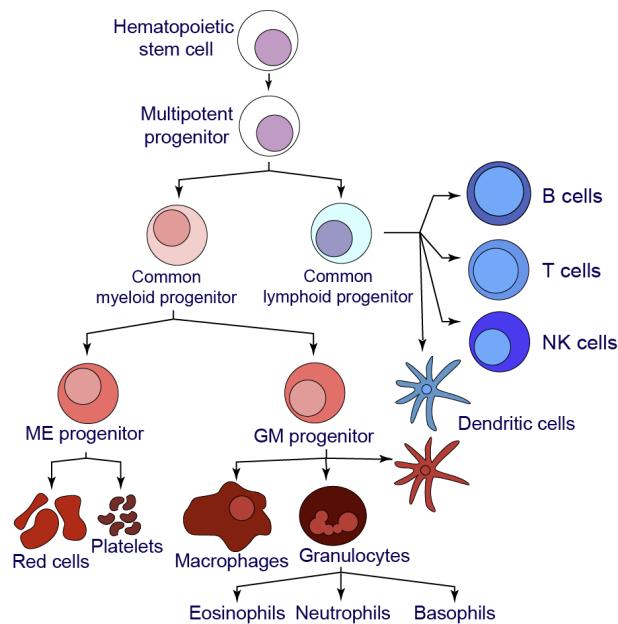
We Now Make Billions of Multi-D Single Cell Measurements  
=> Need for Machine Learning Tools & Human Readable Views



35+ dimensional single cell data:

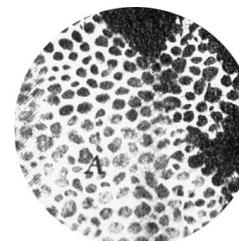
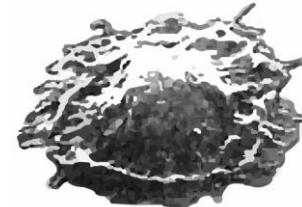
Even if we look through  
all the 2D plots, multidimensional  
relationships are still hidden...

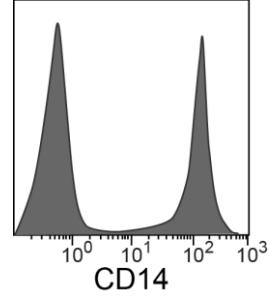
# The Big Idea: Automatically Identify All Cell Types in Primary Tissues, Create Reference Models to Study Impact of Disease, Genetic Changes, etc.



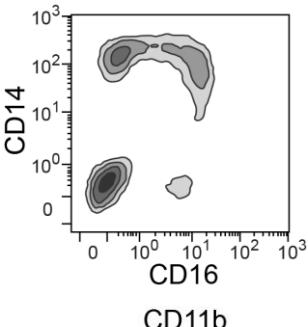
# New Technology Reveals & Characterizes New Cells

Date	Approach	Dimensions (D) Per Cell & Speed	
1665*	Light microscopy	Low	Low
1908**	Light microscopy	Low	Low
1946	Scanning EM	Low	Low
1989	Flow cytometry identification	Low	1K cells/s
2001	Flow cytometry subsetting	4D	2 – 50K cell/s
2011	Mass cytometry + SPADE	32D	500 cell/s
2014	Mass cytometry + t-SNE / viSNE	38D	500 cell/s
(now)	Flow or Imaging MC + UMAP, FlowSOM, MEM	38D	500 cell/s



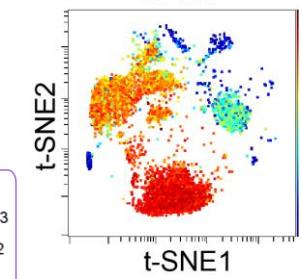
CD14



CD14

CD16

CD11b



t-SNE2

t-SNE1

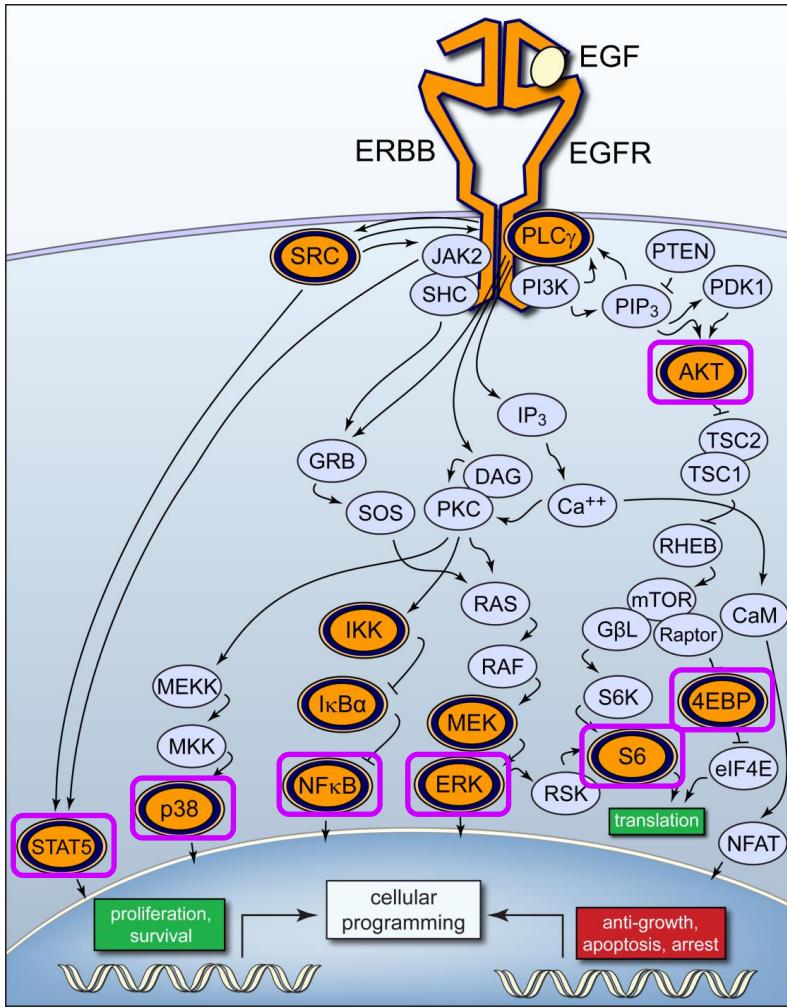
▲ CD206<sup>+3</sup> CD33<sup>+2</sup> CD32<sup>+2</sup>  
 ▼ CD163<sup>-4</sup> CD86<sup>-4</sup> HLA-DR<sup>-3</sup>  
 MerTK<sup>-2</sup> CD14<sup>-2</sup> S100A9<sup>-2</sup>  
 8) MDSC\_b (40%)

\* Robert Hooke describes 'cells' in *Micrographia: or Some Physiological Descriptions of Miniature Bodies Made by Magnifying Glasses*

\*\* Élie Metchnikoff characterizes mononuclear phagocytes: Lectures on the Comparative Pathology of Inflammation, Pasteur Institute in 1891, Nobel Prize in 1908 w/ Ehrlich.

# An Updated Mass Cytometry Panel for Glioblastoma Cell Identity, Stemness, Function, & Signaling

## Phospho-protein effectors of key mutant proteins

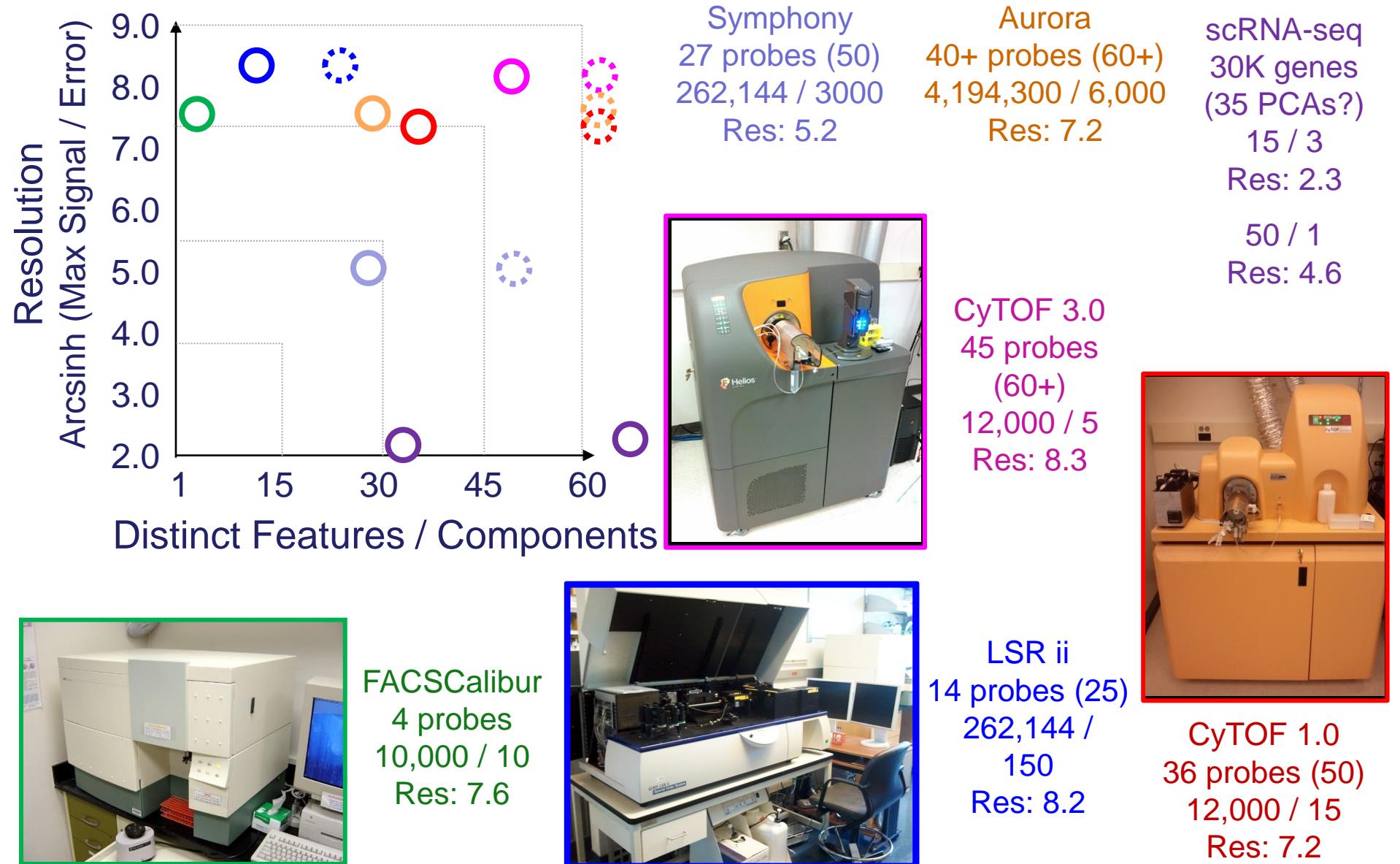


Phospho-flow on lymphoma tumors: Irish et al., PNAS 2010

Target	Mass	Clone	Signaling & proteins		Stain		
			Panel	t-SNE	Live	Sap	MeOH
Rhodium	103	-	●		✓		
Cyclin B1	139	GNS-1	●				✓
TUJ1	141	TUJ1	●	■			✓
cCasp3	142	5A1E	●				✓
CD117	143	104D2	●		✓		
S100B	144	19/S100B	●	■			✓
CD31	145	WM59	●	■*	✓		
$\gamma$ H2AX	147	JBW301	●				✓
CD34	148	581	●	■	✓		
p-4E-BP1 (T37/T46)	149	236B4	●				✓
p-STAT5 (Y694)	150	47	●	■			✓
BMX	151	40/BMX	●				✓
p-AKT (S473)	152	D9E	●	■			✓
p-STAT1 (Y701)	153	58D6	●	■			✓
CD45	154	HI30	●	■*	✓		
NCAM/CD56	155	HCD56	●	■	✓		
p-p38 (T180/Y182)	156	D3F9	●				✓
p-STAT3 (Y705)	158	4/P-STAT3	●	■			✓
ITG $\alpha$ 6/CD49F	159	GoH3	●	■	✓		
CD133	160	AC133	●	■	✓		
PDGFR $\alpha$	161	16A1	●	■	✓		
SOX2	163	O30-678	●	■			✓
SSEA-1/CD15	164	W6D3	●	■	✓		
EGFR	165	AY13	●	■	✓		
p-NF $\kappa$ B p65 (S529)	166	K10-895.12.50	●	■			✓
L1CAM	167	5G3	●	■	✓		
Nestin	168	10C2	●	■			✓
CD44	169	BJ18	●	■	✓		
GFAP	170	1B4	●	■			✓
p-ERK1/2 (T202/Y204)	171	D13.14.4E	●				✓
p-S6 (S235/S236)	172	N7-548	●	■			✓
SOX10	173	A-2	●	■			✓
HLA-DR	174	L243	●	■	✓		
p-HH3	175	HTA28	●				✓
Histone H3	176	D1H2	●				✓

GBM panel: Leelatian & Sinnaeve et al., bioRxiv 2019  
Tumor methods: Leelatian & Doxie et al., Cytometry B 2017

# Cytometry Probes & Instruments Keep Improving (ca. 2019)



# HD Cytometry Balances Signal, Throughput, & Cost

Mass cytometry: Imaging or flow cytometry method of multiplexed single cell analysis. Standard mass cytometry panels detect **37+ features** per cell using pre-validated antibodies. The dynamic range is **>10,000 intensity units** per feature and a small flow-based mass cytometry dataset might include **1.2 million cells** from 12 samples collected at a rate of **500 cells/second** (~40 min instrument time) for a total cost of ~\$4,500 (\$0.004 per cell), including personnel time/effort.

Diggins et al., *Bench to Bedside to Bytes, in review*

Mass cytometry vs. other single cell technologies: Mistry et al., *FEBS J* 2018

Some of the literature from “big seq” is hilariously inaccurate when it comes to flow cytometry...

From: “Single cell RNA sequencing to explore immune cell heterogeneity”, *Nat Rev Immunology* 2017

	FACS	CyTOF	qPCR	Plate-based protocols (STRT-seq, SMART-seq, SMART-seq2)	Fluidigm C1	Pooled approaches (CEL-seq, MARS-seq, SCRB-seq, CEL-seq2)	Massively parallel approaches (Drop-seq, InDrop)
Cell capture method	Laser	Mass cytometry	Micropipettes	FACS	Microfluidics	FACS	Microdroplets
Number of cells per experiment	Millions	Millions	300–1,000	50–500	48–96	500–2,000	5,000–10,000
Cost	\$0.05 per cell	\$35 per cell	\$1 per cell	\$3–6 per well	\$35 per cell	\$3–6 per well	\$0.05 per cell
Sensitivity	Up to 17	Up to 40	10–30 genes	7,000–10,000 genes	6,000–9,000 genes for cell 000–5,000 per cell for cells	7,000–10,000 genes per cell for cell lines; 2,000–6,000 genes per cell for primary cells	5,000 genes per cell for cell lines; 1,000–3,000 genes per cell for primary cells

Wait, so this review says each experiment costs... ?!

FACS: \$50,000 (\$0.05 x 1,000,000 cells)

CyTOF: \$35,000,000 (\$35.00 x 1,000,000 cells)

and sequencing; STRT-seq, single-cell tagged reverse transcription sequencing.

(mass cytometry); FACS, fluorescence-activated cell sorting; qPCR, quantitative PCR; SCRB-seq, single-cell RNA barcoding

# HD Cytometry Dissects Cellular Mechanisms of Cancer Immune Response

Cell

Article

Spitzer et al.,  
*Cell* 2017

## Systemic Immunity Is Required for Effective Cancer Immunotherapy

Uses mass cytometry to characterize essential role of peripheral blood CD4<sup>+</sup> T cells in immunotherapy response

ARTICLE

doi:10.1038/nature22079

Huang et al.,  
*Nature* 2017

## T-cell invigoration to tumour burden ratio associated with anti-PD-1 response

Uses mass cytometry to reveal peripheral blood CD8 T cells associated with anti-PD-1 immunotherapy responses

Cell

Article

Wei et al.,  
*Cell* 2017

## Distinct Cellular Mechanisms Underlie Anti-CTLA-4 and Anti-PD-1 Checkpoint Blockade

Uses mass cytometry to characterize similar & distinct tumor-infiltrating immune cell subsets (mostly T cells) following immunotherapies

# Referenced Comparisons of Single Cell Techniques

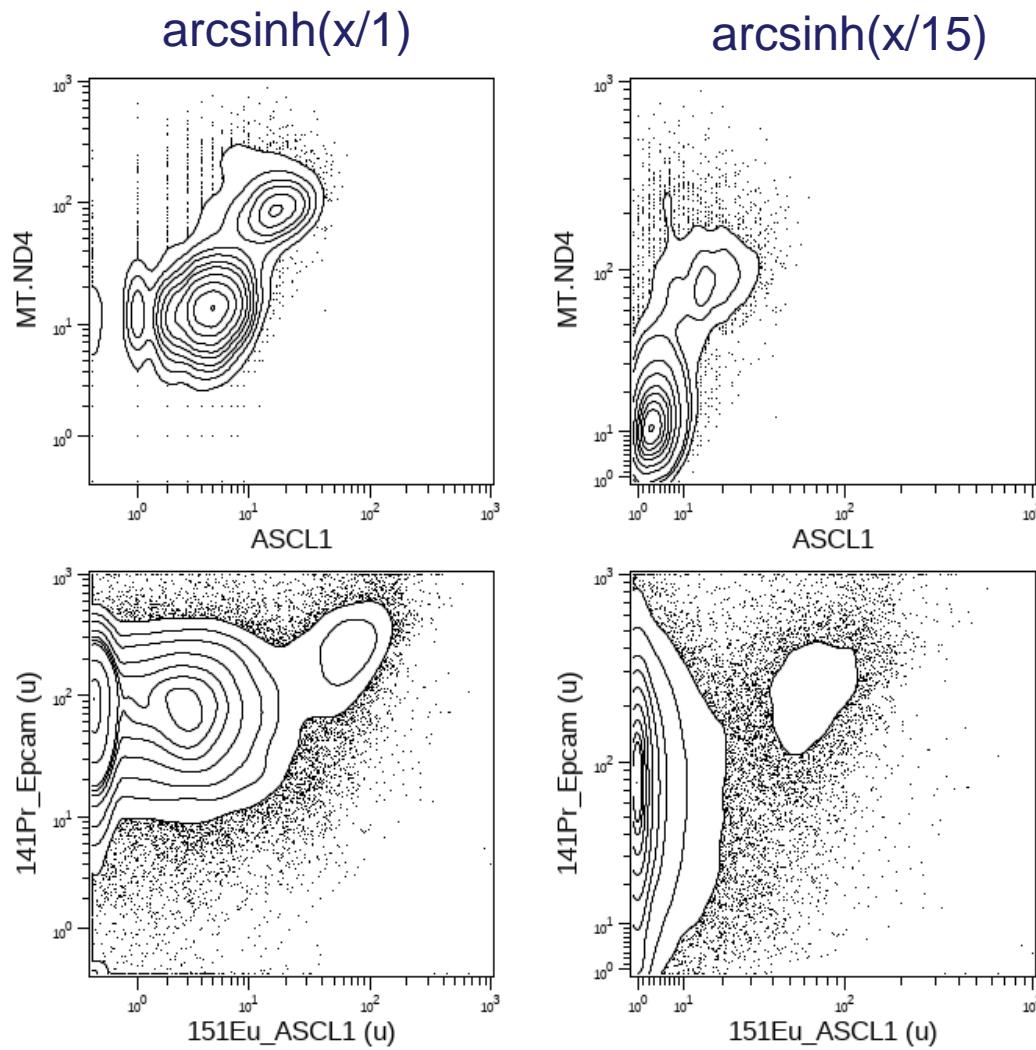
	Mass cytometry	Fluorescence cytometry	scRNA-seq
Detection method	Time-of-flight of ionized heavy elemental isotopes	Emitted light and light scatter	Nucleic acid sequence
Example probes	Metal-conjugated antibodies, metal-conjugated small molecules	Fluorochrome-conjugated antibodies, fluorescent molecules	Poly-A targeted oligonucleotides, antibody-conjugated oligonucleotides [144]
Example cellular features and targets measured	Proteins, phospho-proteins, chromatin modifications [124], RNA transcripts [131–133], platinum drug uptake [128,129], cell cycle status and DNA synthesis [130], cell size [125], apoptosis [48], and viability [27]	Proteins, phospho-proteins, chromatin modifications [145], RNA transcripts [146], fluorescent drug uptake [147], metabolism and redox state [148,149], cell cycle status and DNA synthesis [150], cell size and granularity, apoptosis [151] and viability [152]	Poly-adenylated RNA transcripts, CITE-seq probes [144]
Minimum cells per sample needed at start of protocol	50 000 cells	50 000 cells	200 cells [153] <sup>a</sup>
Cell capture rate	30–60% [31,32]	> 95%	5–65% [154]
Target capture	> 95%	> 95%	10–40% [154–156]
Example of analyzed cell events per study	2 000 000 cells per study	20 000 000 cells per study	5000 cells per study [157]
Resolution, arcsinh (max/error) <sup>b</sup>	~ 6–9 e.g., arcsinh (12 000/5)	~ 5–9 e.g., arcsinh (262 144/150), arcsinh (4 194 300/6000) <sup>c</sup>	~ 2–3 e.g., arcsinh (10/1)
Cell throughput	Up to $10^4$ cells·s <sup>-1</sup> [33]	Up to $10^5$ cells·s <sup>-1</sup> [31]	Up to $10^4$ cells·s <sup>-1</sup> [155]
Features/available channels	50/200 channels	30/64 channels <sup>c</sup>	~ 30 000/N/A
Degree of crosstalk between parameters	3% (range: 0–8.6%) [55]	19.5% (4.9–51.1%) [56]	N/A
Accessibility	Major research institutions	> 20 features: major research institutions; ~ 4 color common	Major research institutions
Total cost	\$0.004–\$0.01/cell [31]	< \$0.001/cell	\$0.05–\$3.00/cell [158]

# Example Challenge, Scaling to Compare scRNA-seq vs. CyTOF (Here with one example cell line, DMS454, and ASCL1 on the x-axis)

DMS454  
SCLC cell line

scRNA-seq

CyTOF



Qualitative: Is there a subset of cells that is ASCL1<sup>+</sup>?

Quantitative: How much ASCL1 is in each cell? How much does ASCL1 change after Tx?

# Spectral Flow Cytometry Can Separate ‘Overlapping’ Probes

40 Colors

35 Colors

28 Colors

24 Colors

Small Particles

AF Extraction

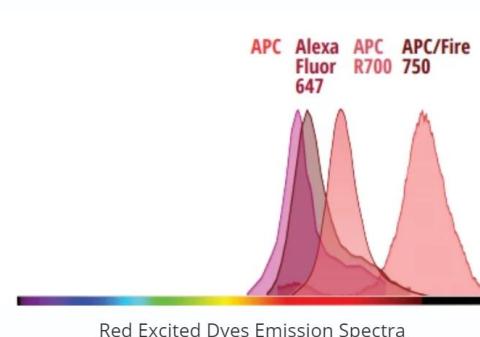
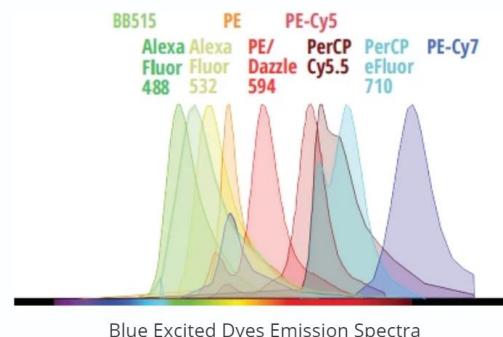
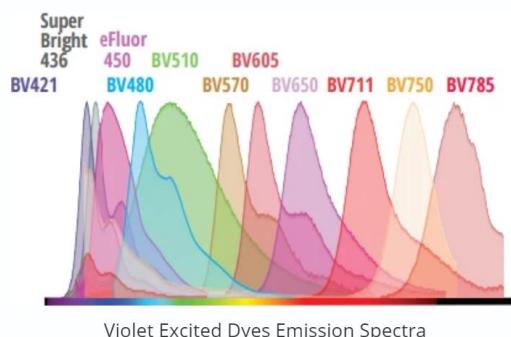
Overlapping Dyes

## More Choice, Greater Flexibility, Easier Setup

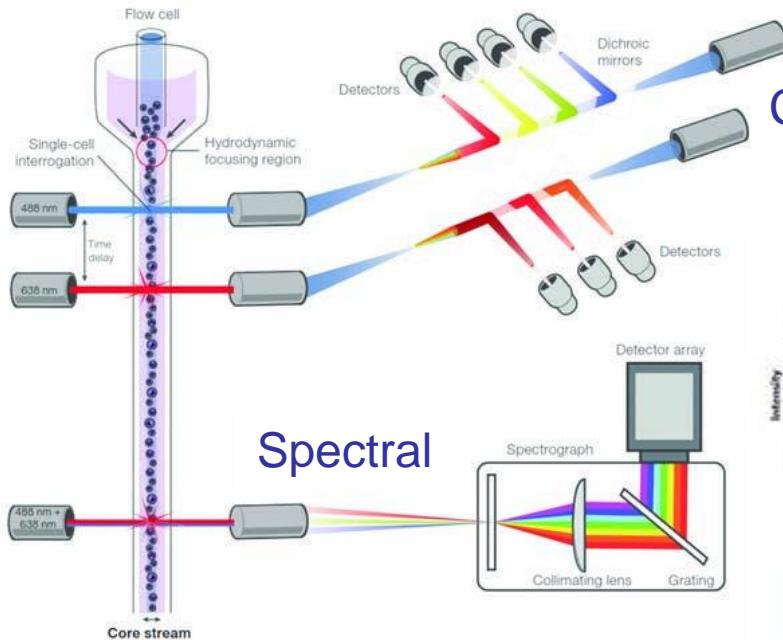
The optical design combined with the unmixing capability in SpectroFlo® software allows greater fluorochrome choice, panel flexibility, and easy setup without having to change filters. The three laser configuration provides outstanding multi-parametric data for a wide array of applications. Markers and fluorochromes in a 24-color panel designed for identification of circulating cell subsets in human peripheral blood are summarized in the table below:

SPECIFICITY	FLUOROCHROME	SPECIFICITY	FLUOROCHROME	SPECIFICITY	FLUOROCHROME
CCR7	Brilliant Violet 421™	CD11c	BD Horizon™ BB515	CD27	APC
CD19	Super Bright 436	CD45RA	Alexa Fluor® 488	CD123	Alexa Fluor® 647
CD16	eFluor® 450	CD3	Alexa Fluor® 532	CD127	BD Horizon™ APC R700
TCR γ/δ	BD Horizon™ BV480	CD25	PE	HLA DR	APC/Fire™ 750
CD14	Brilliant Violet 510™	IgD	PE/Dazzle™ 594		
CD8	Brilliant Violet 570™	CD95	PE-Cy™5		
CD1c	Brilliant Violet 605™	CD11b	PerCP-Cy™5.5		
PD-1	Brilliant Violet 650™	CD38	PerCP-eFluor® 710		
CD56	Brilliant Violet 711™	CD57	PE-Cy™7		
CD4	Brilliant Violet 750™				
CD28	Brilliant Violet 785™				

## The 24-Color Panel Includes Many Highly Overlapping Dyes:



# By Detecting “Across the Spectrum” (vs. “Channels” of Light) Spectral Flow Distinguishes Probes with Overlapping Peaks

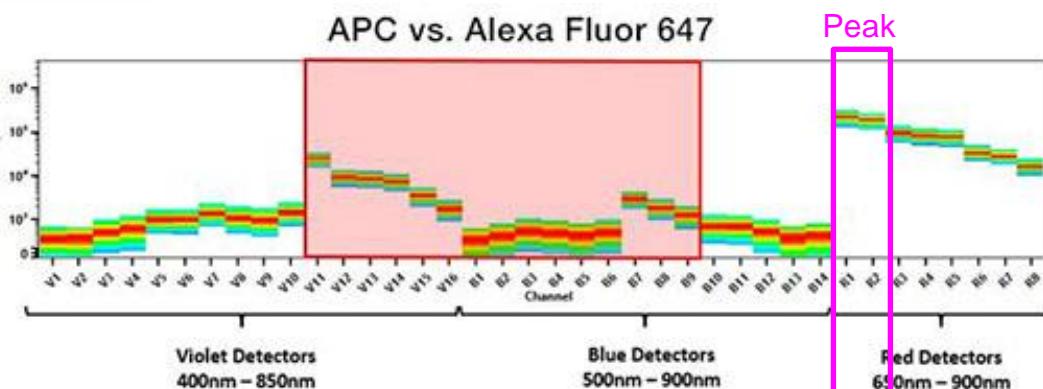


Themo Fisher, adapted from John Nolan et al. 2013

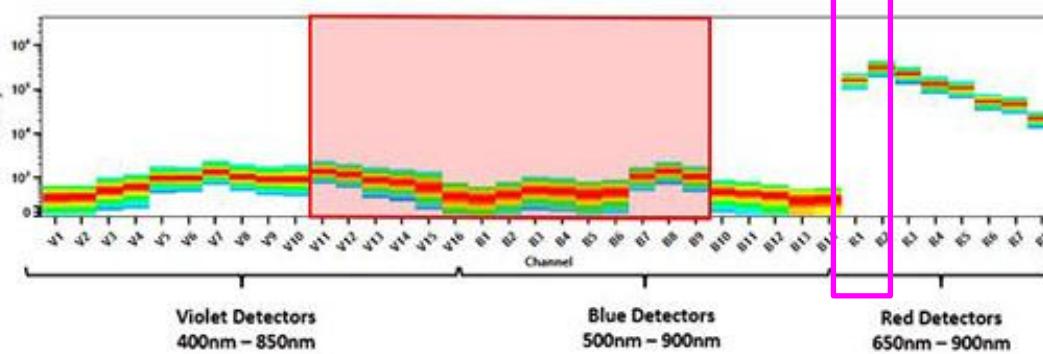
Conventional

Probes with similar “peaks” are resolved

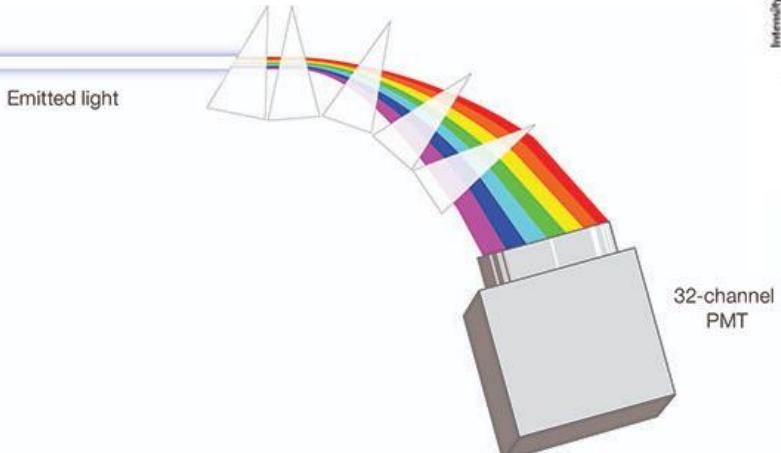
APC vs. Alexa Fluor 647



Peak

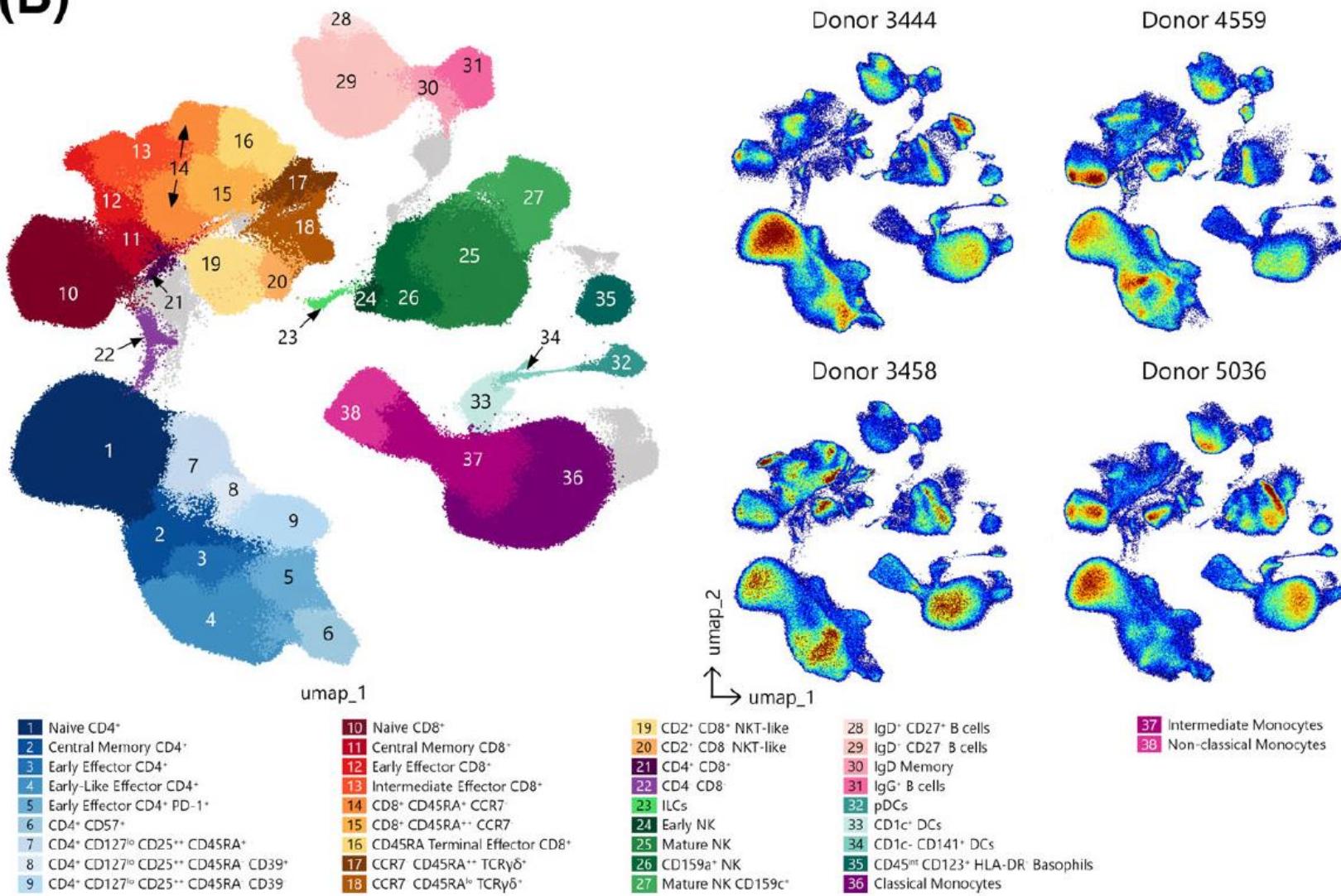


Conventional detector measures only peaks,  
uses dyes with distinct peaks

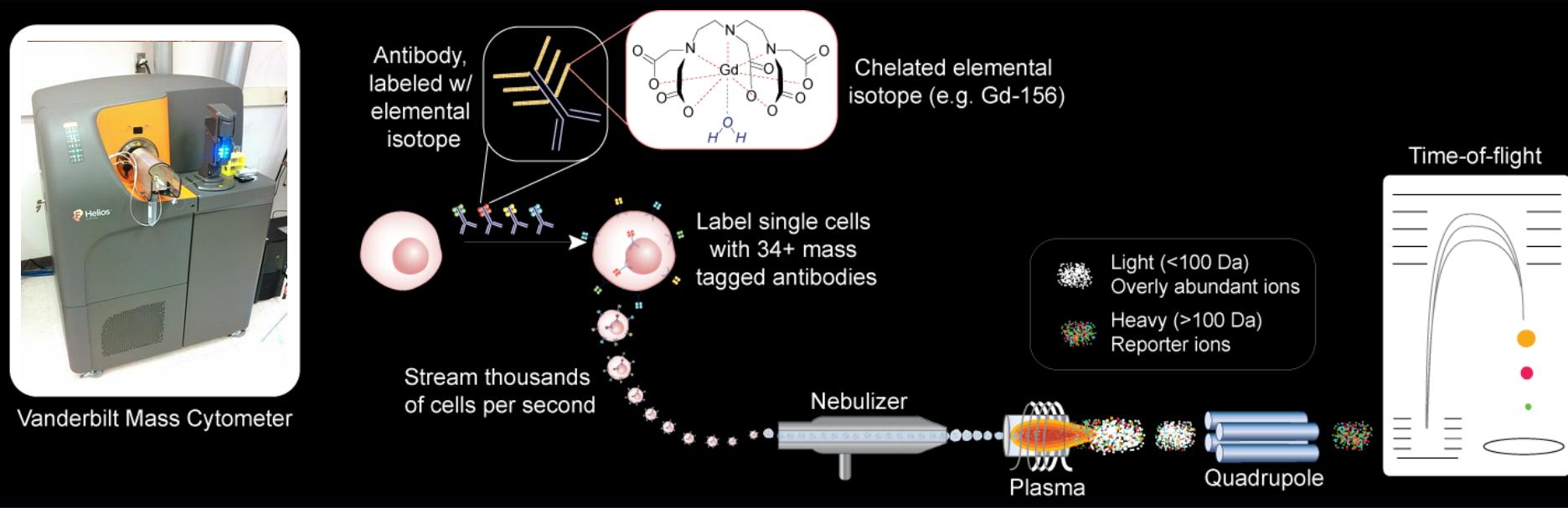


# Spectral Flow Cytometry's State of the Art in 2020: Excellent Resolution of 40 Antibodies (Rivals CyTOF)

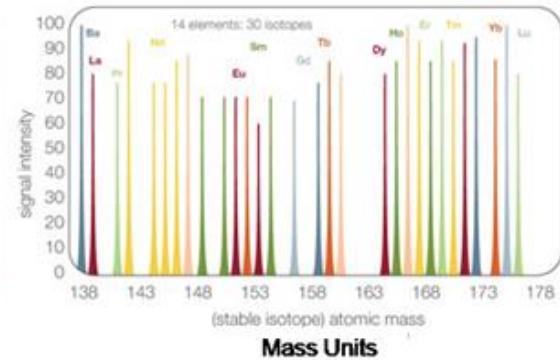
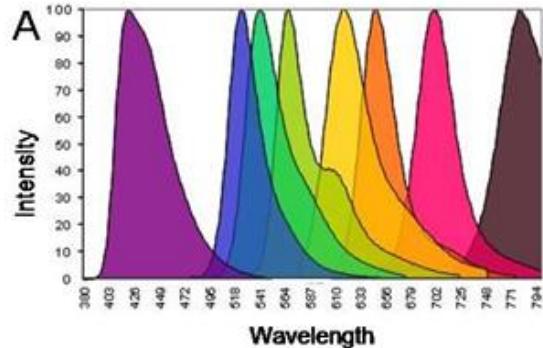
(B)



# Mass Cytometry Uses Isotopically Pure Metals As Probes



Mass cytometry: standard panels detect 35+ features per cell using pre-validated antibodies. The dynamic range is >10,000 intensity units per feature. A small dataset might include 1.2 million cells from 12 samples collected at a rate of 500 cells/second (~40 min instrument time) for a total cost of ~\$4,500 (\$0.004 per cell), including reagents & personnel.

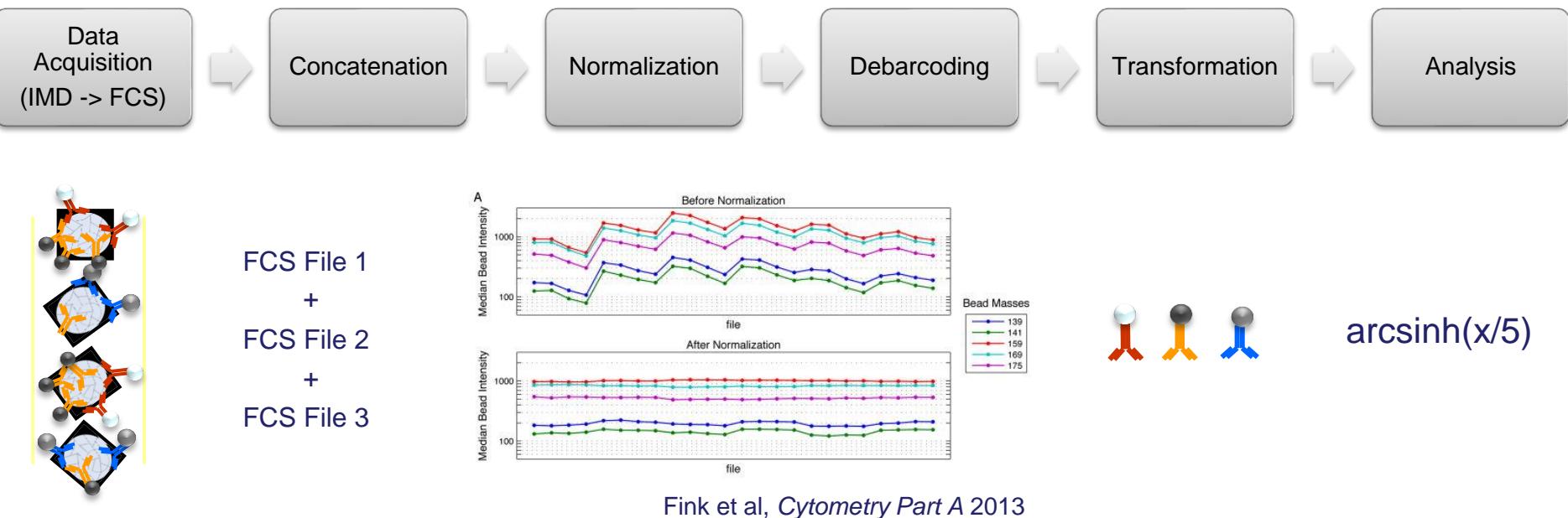


Experiment Protocols: Leelatian et al., *Methods in Molecular Biology* 2015  
Leelatian et al., *Current Protocols in Molecular Biology* 2017

Analysis Protocols: Diggins et al., *Methods* 2015  
Diggins et al., *Current Protocols in Cytometry* 2018

Reviewed in: Spitzer & Nolan, *Cell* 2016  
Adapted from Bendall et al., *Science* 2011

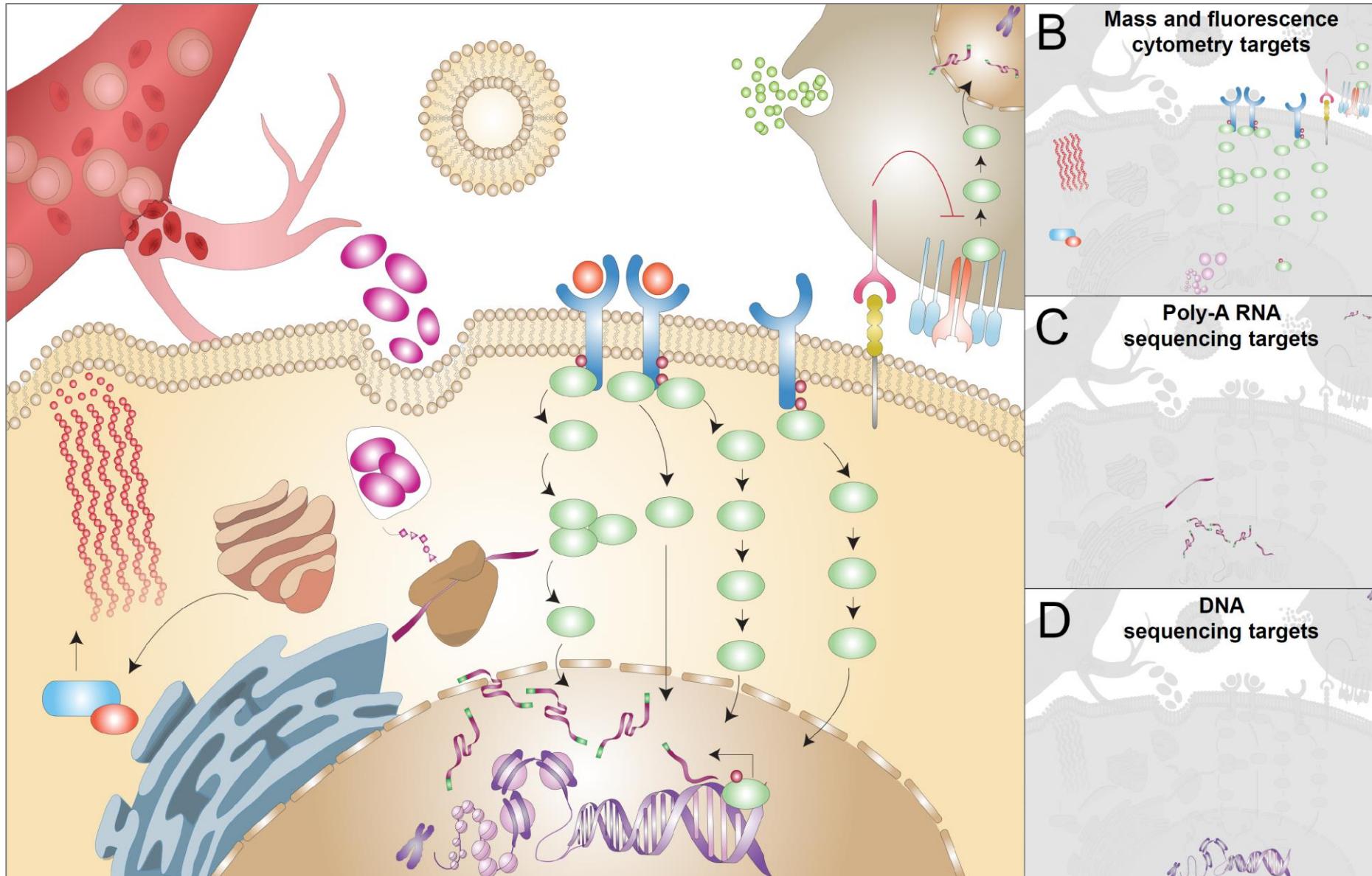
# Modern Cytometry Includes Internal Batch Controls



## Resources:

- Concatenation: downloadable tool from Cytobank (<http://support.cytobank.org/help/kb/cytobank-utilities/concatenating-fcs-files>)
- Normalization: Cytometry Part A [Volume 83A, Issue 5, pages 483-494, 19 MAR 2013 DOI: 10.1002/cyto.a.22271](http://onlinelibrary.wiley.com/doi/10.1002/cyto.a.22271/full#fig6) <http://onlinelibrary.wiley.com/doi/10.1002/cyto.a.22271/full#fig6>
- Barcoding: Bodenmiller et al, *Nature Biotechnology* 2012 (<http://www.nature.com/nbt/journal/v30/n9/full/nbt.2317.html>)

# Challenge: Get a Full Picture of Cell Function (Multiple Data Types)



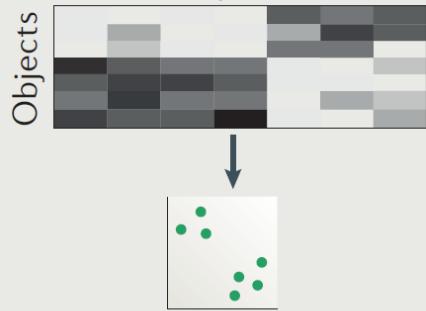
Beyond the Message: Advantages of Snapshot Proteomics with Single-Cell Mass Cytometry in Solid Tumors  
Mistry et al., *FEBS Journal* 2018

# Computational Flow Cytometry: There Is Help for Data Analysis

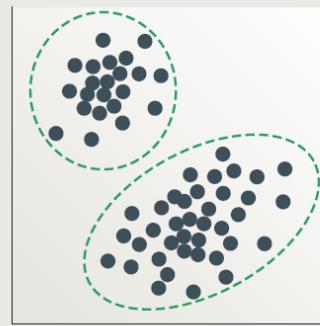
## a Unsupervised machine learning: learning structures

Dimensionality reduction

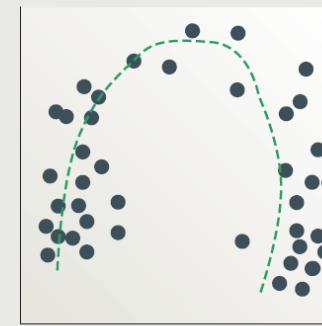
Properties



Clustering

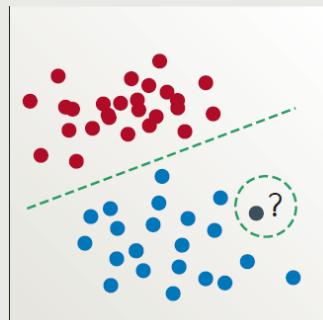


Seriation

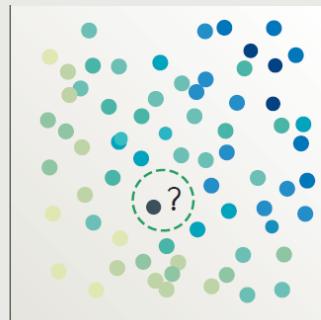


## b Supervised machine learning: learning from examples

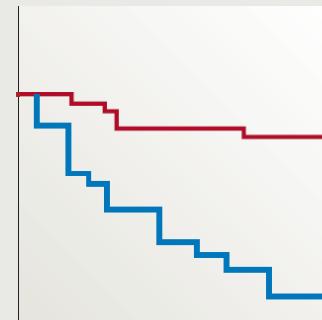
Classification



Regression



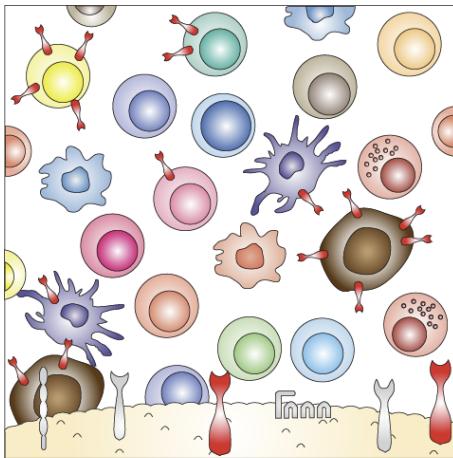
Survival analysis



# Part 9: Samples Over Time Reveal Immune System Dynamics

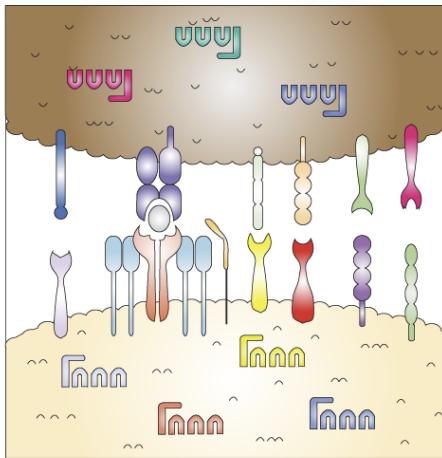
# Applications of HD Cytometry in Cancer & Cell Biology

Microenvironment



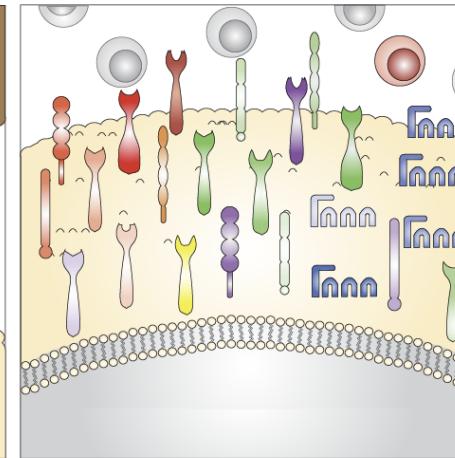
Track key biomarkers  
on all cell types;  
find cytokine producers

Cell:cell interactions



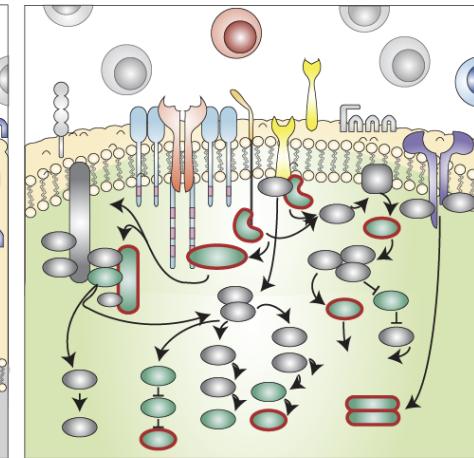
Monitor immune checkpoint  
proteins on T cells, APCs,  
& cancer cells

Immunophenotype



Measure differentiation;  
deep phenotype  
using fewer cells

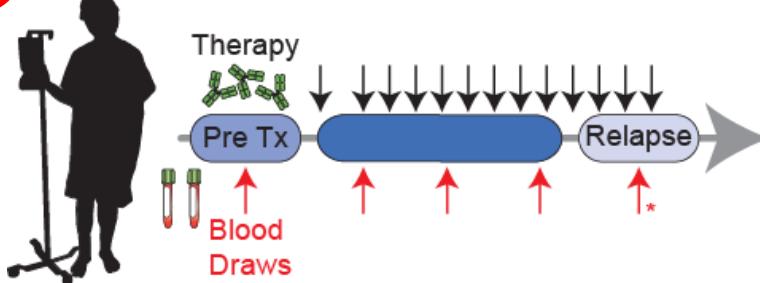
Signaling & function



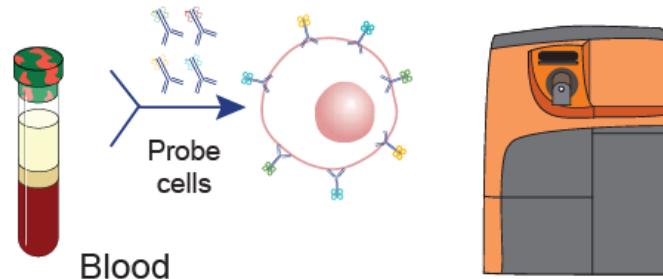
Dissect signaling changes;  
characterize mechanisms  
of treatment response

# Human Immune Monitoring: T cell Focus

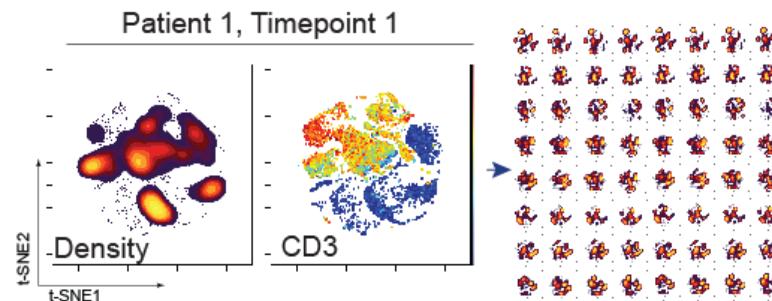
1 Live cells from patients



2 High dimensional, single cell measurements



3 Analyze systems immune changes



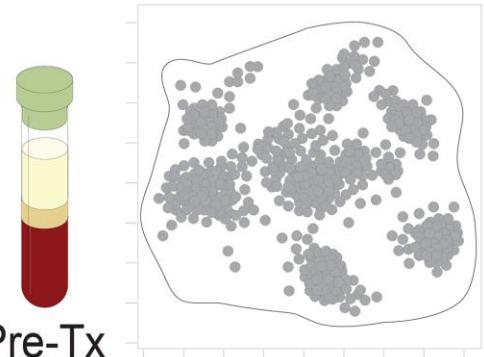
Tag Isotope #	Target
141 Pr	ICOS
142 Nd	CD19
<b>143 Nd</b>	<b>TIM3</b>
144 Nd	CCR5
145 Nd	CD4
146 Nd	CD64
147 Sm	CD20
<b>148 Nd</b>	<b>CD38</b>
149 Sm	CCR4
150 Ns	CD43
151 Eu	CD14
152 Sm	TCRgd
153 Eu	CD45RA
154 Sm	CD45
156 Gd	CXCR3
158 Gd	CD33
159 Tb	CCR7
160 Gd	CD28
<b>161 Dy</b>	<b>CD32</b>
162 Dy	CD69
<b>163 Dy</b>	<b>HLA-DR</b>
164 Dy	CD45RO
165 Ho	CD16
166 Er	CD44
167 Er	CD27
168 Er	CD8
169 Tm	CD25
170 Er	CD3
<b>171 Yb</b>	<b>CXCR5</b>
172 Yb	CD57
<b>174 Yb</b>	<b>PD-1</b>
175 Lu	PD-L1
176 Yb	CD56
191/193 Ir	Nucleic Acid
195 Pt	Viability

Custom conjugates  
Commercially available

# Clinical Trial Monitoring: What Do We Need to Know? Automate Four Key Readouts vs. Clinical Outcomes

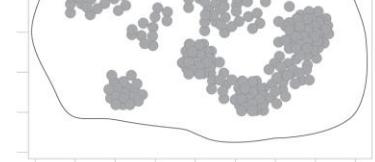
## Features of Dynamic Populations

### 1 Systems Plasticity



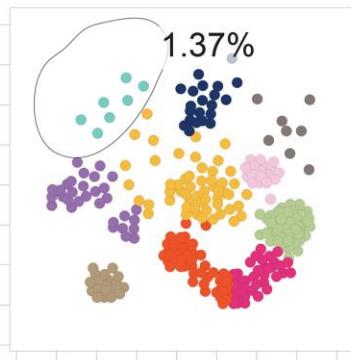
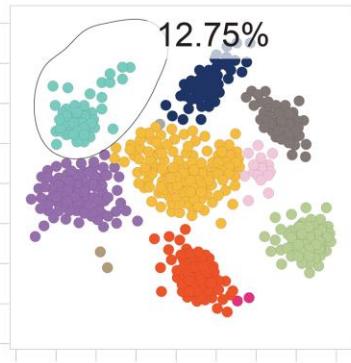
Pre-Tx

Time 1



Earth Mover's Distance  
on t-SNE or UMAP

### 2 Population abundance



Traditional gating  
or cluster frequency

### 3 Signature features

#### Pre-therapy

- ▲ HLADR<sup>+2</sup> CCR5<sup>+1</sup> CD38<sup>+1</sup>  
CD33<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD45RO<sup>-6</sup> CD3<sup>-6</sup>  
CD4<sup>-4</sup> CD45<sup>-2</sup> CCR4<sup>-1</sup>  
CCR7<sup>-1</sup> CD28<sup>-1</sup> CD27<sup>-1</sup>

#### Time point 1

- ▲ HLADR<sup>+2</sup> CD38<sup>+1</sup> CD45RA<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD4<sup>-6</sup> CD3<sup>-6</sup>  
CD45RO<sup>-5</sup> CCR5<sup>-2</sup> CD45<sup>-2</sup>  
CD28<sup>-2</sup> CD20<sup>-1</sup> CCR4<sup>-1</sup>  
CD27<sup>-1</sup>

Marker Enrichment  
Modeling (MEM)

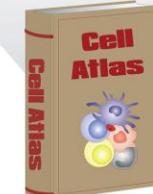
### 4 Population novelty



Pre



Timepoint

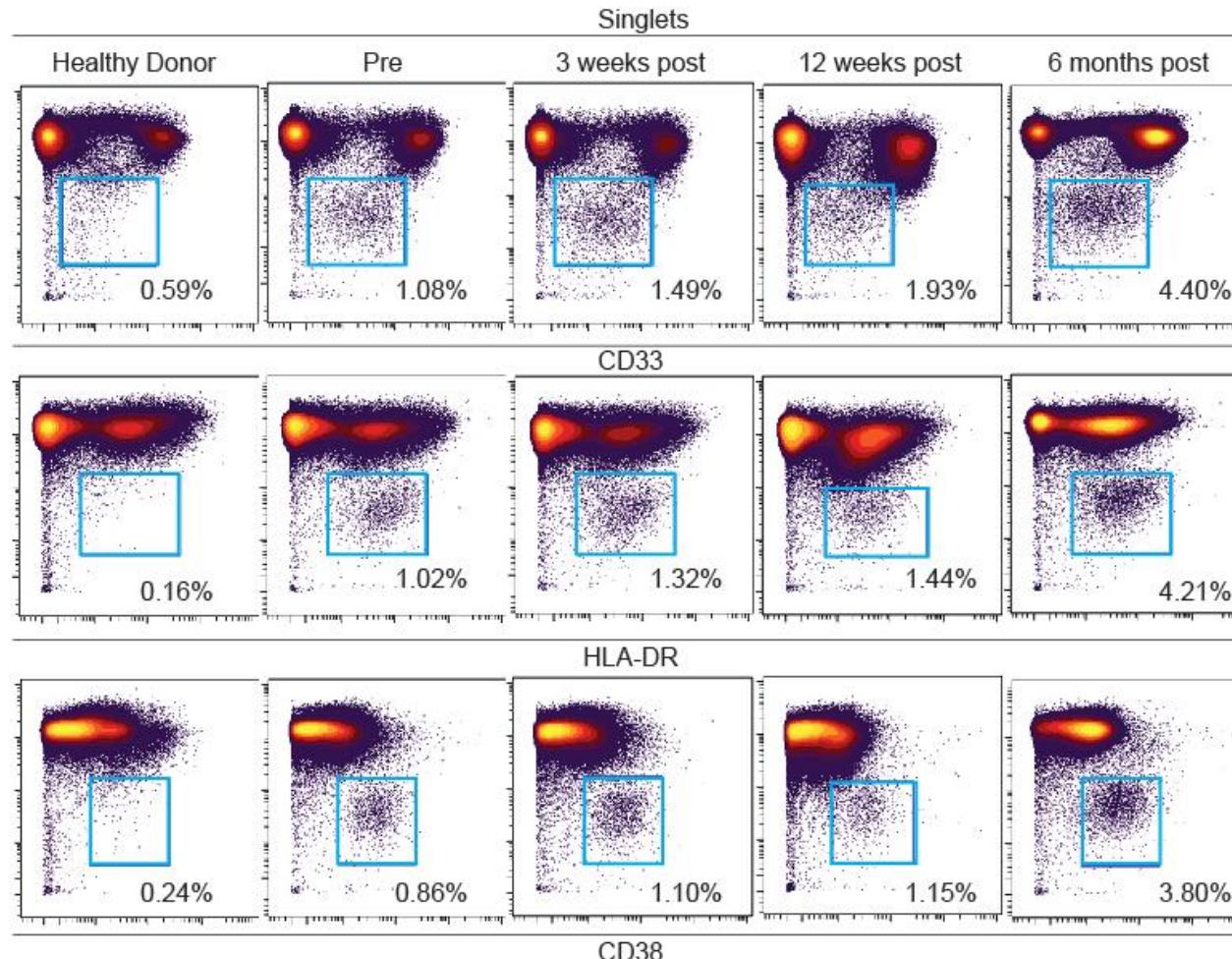


ΔMEM vs. Timepoint  
or Cell Atlas

How we quantified

# A Case Study: Systems Immune Monitoring with Mass Cytometry Reveals A Clinically Significant Rare Cell Subset

## MDS in Melanoma Patient Revealed During $\alpha$ -PD-1 Therapy



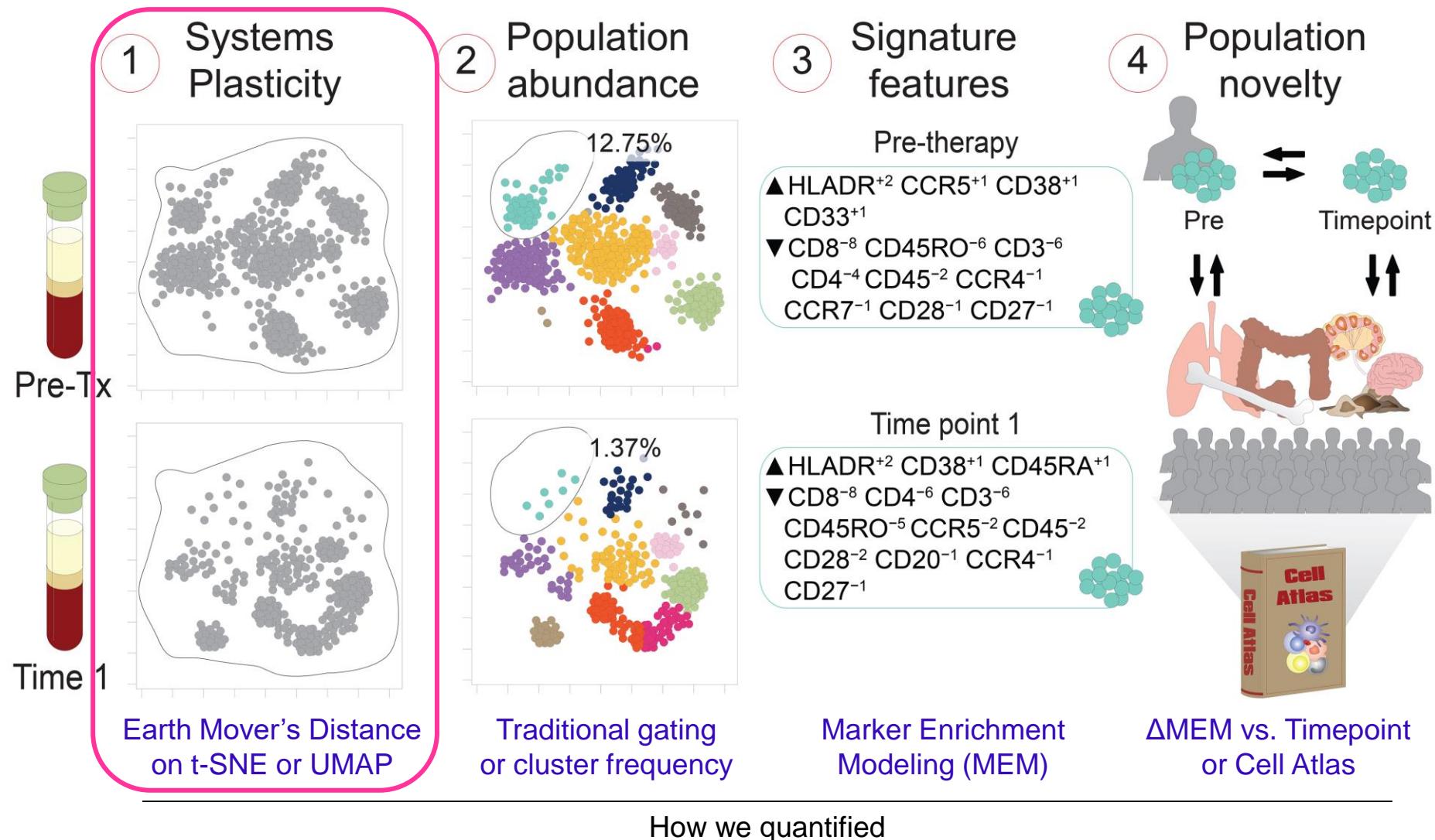
Healthy donor looks similar to melanoma in 2D views

At Pre-Tx, MDS blasts were not detected by standard CBC

High dimensional panel allowed review of PD-1 on MDS blasts w/ existing data

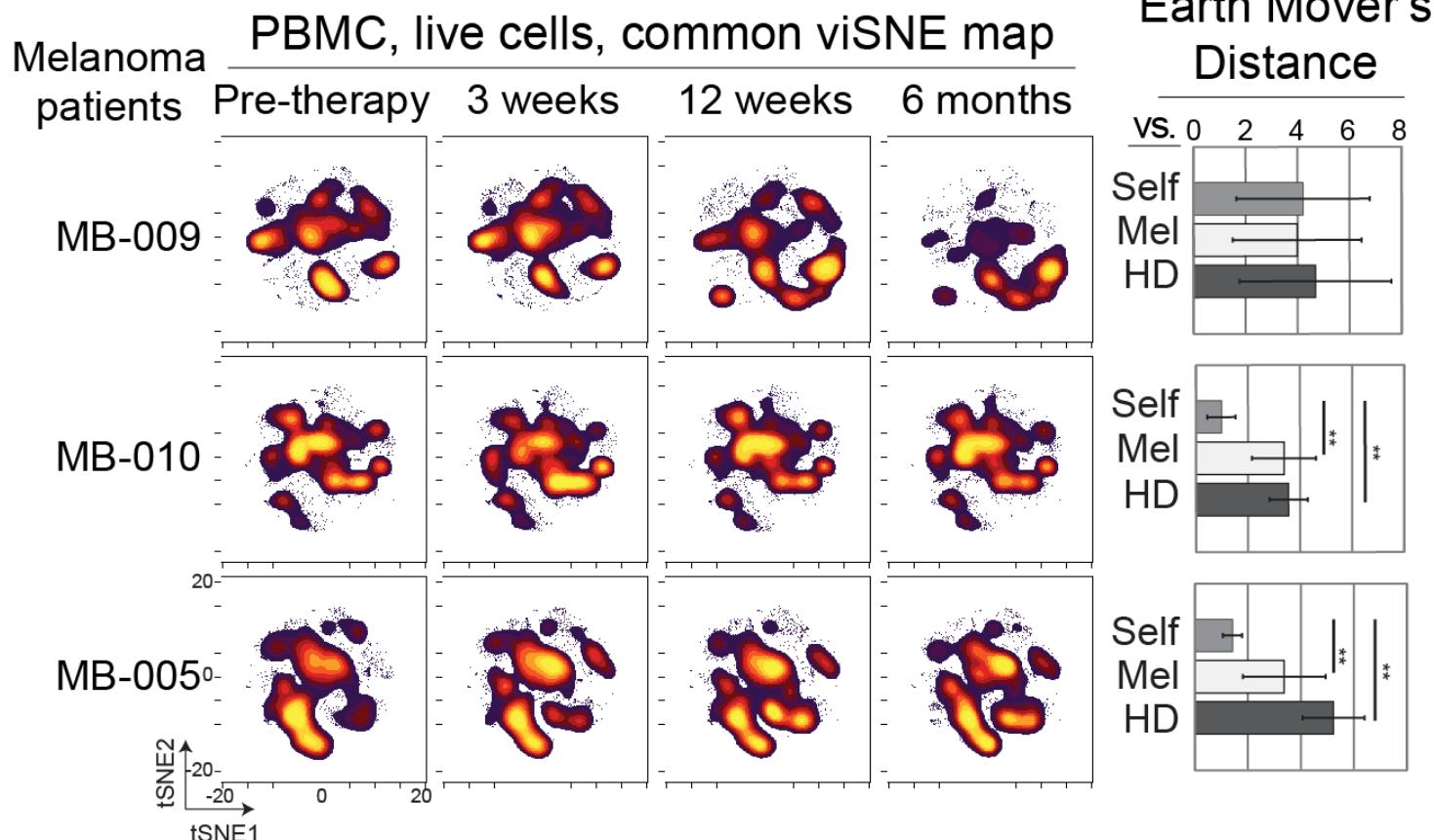
# Clinical Trial Monitoring: What Do We Need to Know? Automate Four Key Readouts vs. Clinical Outcomes

## Features of Dynamic Populations



# Plasticity / Stability: Earth Mover's Distance Quantifies Change Over Time Within a t-SNE Analysis

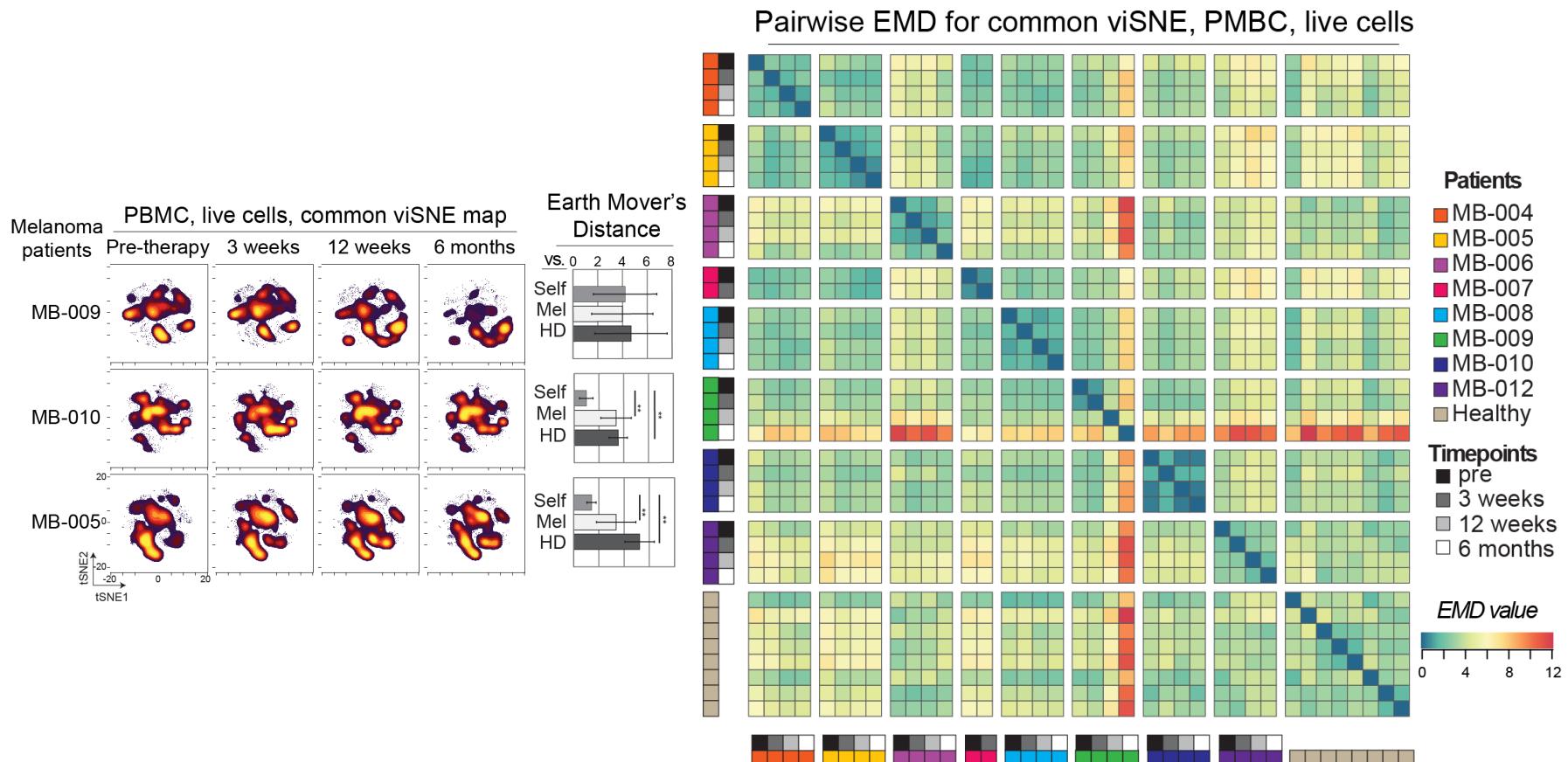
Melanoma Patients Treated with  $\alpha$ -PD-1 Therapy, Monitored by Mass Cytometry



Systems immune monitoring reveals an unexpected pattern in MB-009  
Individuals can be their own significantly stable baseline

# Plasticity / Stability: Earth Mover's Distance Quantifies Change Over Time Within a t-SNE Analysis

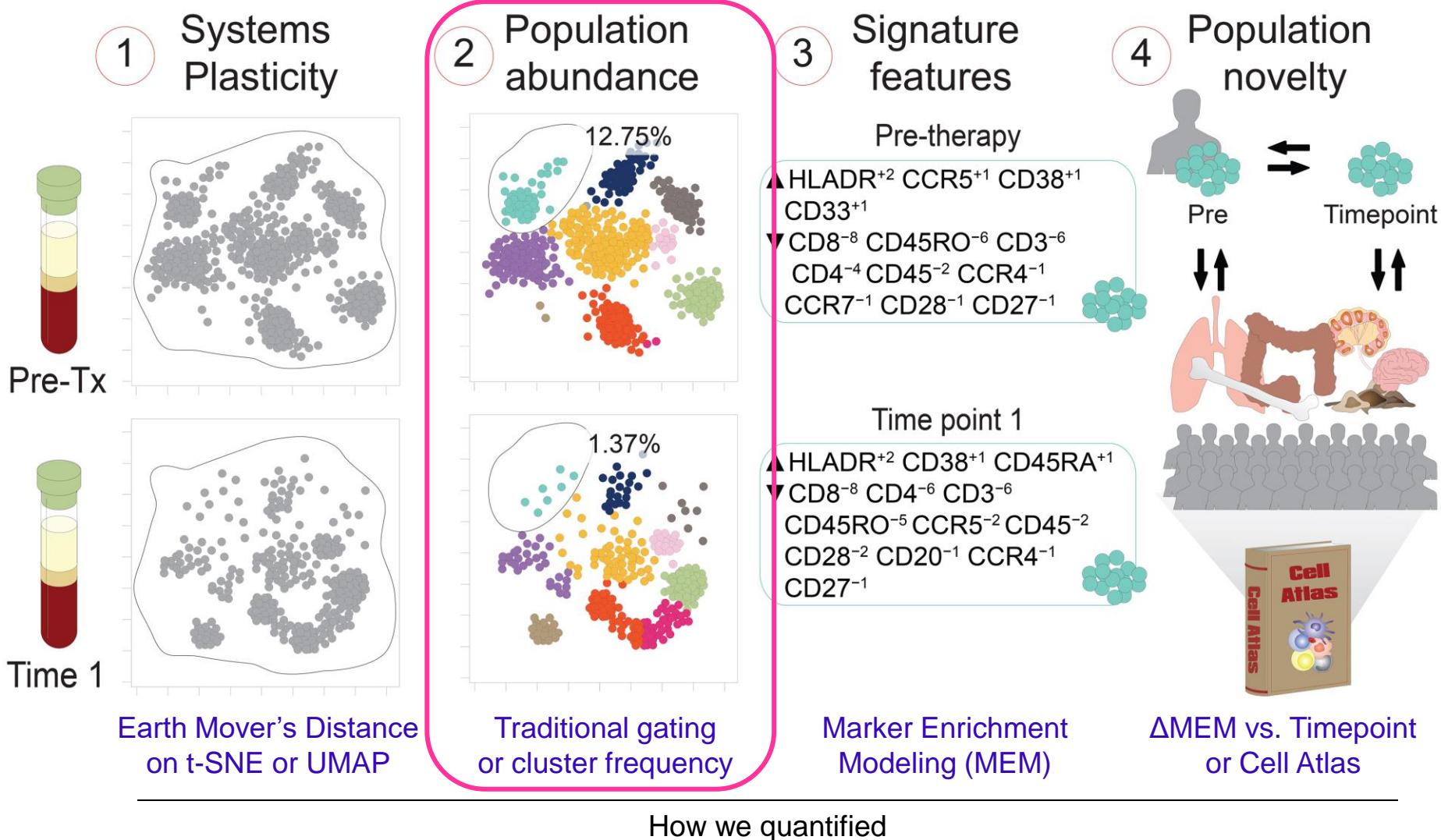
Melanoma Patients Treated with  $\alpha$ -PD-1 Therapy, Monitored by Mass Cytometry



Systems immune monitoring reveals an unexpected pattern in MB-009

# Clinical Trial Monitoring: What Do We Need to Know? Automate Four Key Readouts vs. Clinical Outcomes

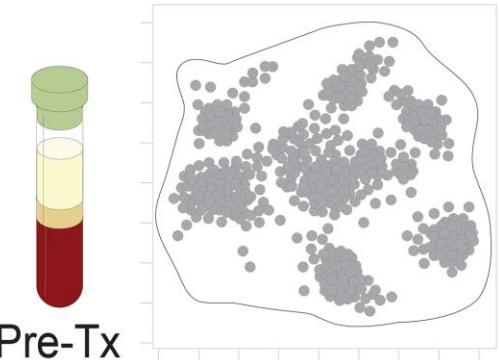
## Features of Dynamic Populations



# Clinical Trial Monitoring: What Do We Need to Know? Automate Four Key Readouts vs. Clinical Outcomes

## Features of Dynamic Populations

### 1 Systems Plasticity

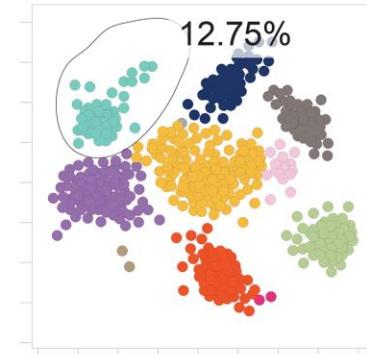


Pre-Tx

Time 1

Earth Mover's Distance  
on t-SNE or UMAP

### 2 Population abundance



Traditional gating  
or cluster frequency

### 3 Signature features

#### Pre-therapy

- ▲ HLA<sup>DR</sup><sup>+2</sup> CCR5<sup>+1</sup> CD38<sup>+1</sup>  
CD33<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD45RO<sup>-6</sup> CD3<sup>-6</sup>  
CD4<sup>-4</sup> CD45<sup>-2</sup> CCR4<sup>-1</sup>  
CCR7<sup>-1</sup> CD28<sup>-1</sup> CD27<sup>-1</sup>

#### Time point 1

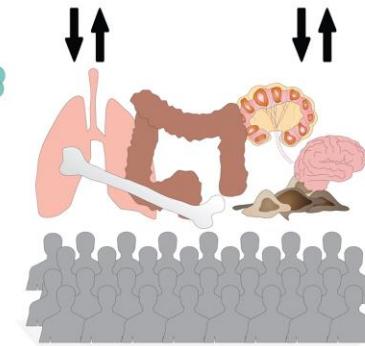
- ▲ HLA<sup>DR</sup><sup>+2</sup> CD38<sup>+1</sup> CD45RA<sup>+1</sup>
- ▼ CD8<sup>-8</sup> CD4<sup>-6</sup> CD3<sup>-6</sup>  
CD45RO<sup>-5</sup> CCR5<sup>-2</sup> CD45<sup>-2</sup>  
CD28<sup>-2</sup> CD20<sup>-1</sup> CCR4<sup>-1</sup>  
CD27<sup>-1</sup>

Marker Enrichment  
Modeling (MEM)

### 4 Population novelty

Pre

Timepoint

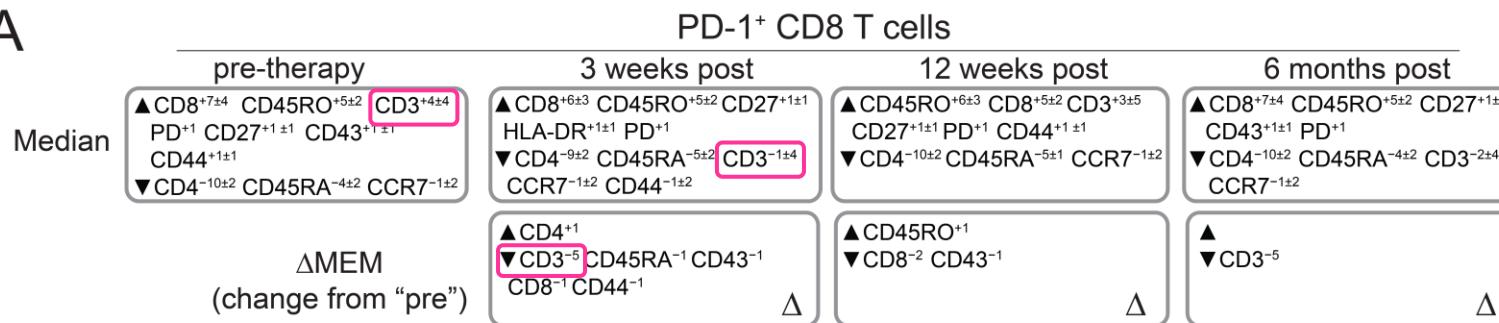


ΔMEM vs. Timepoint  
or Cell Atlas

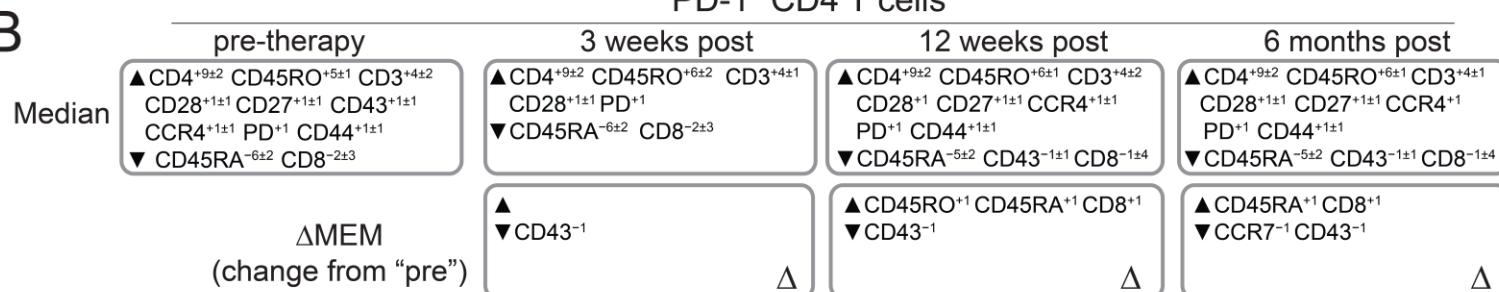
How we quantified

# $\Delta$ MEM Reveals CD8 $^{+}$ Specific Decrease in Per-Cell CD3 in Melanoma Patient PBMC at 3 Weeks after $\alpha$ -PD-1 Therapy

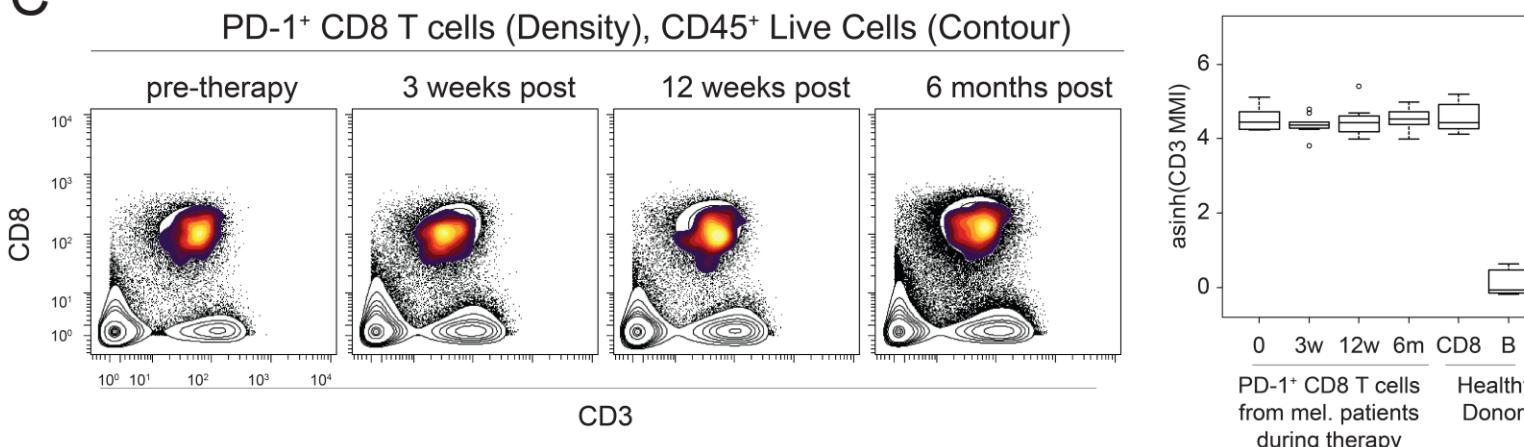
A



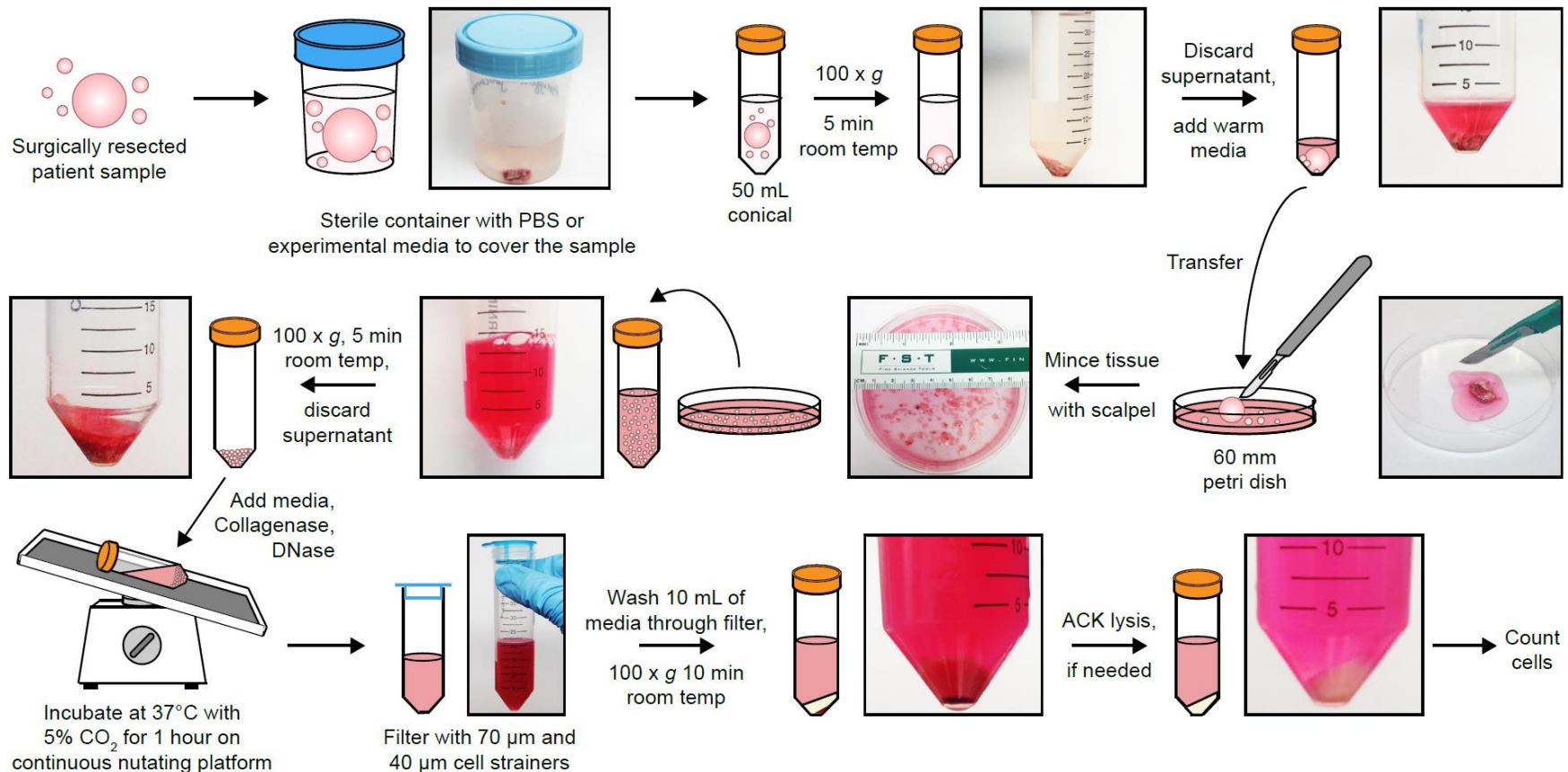
B



C

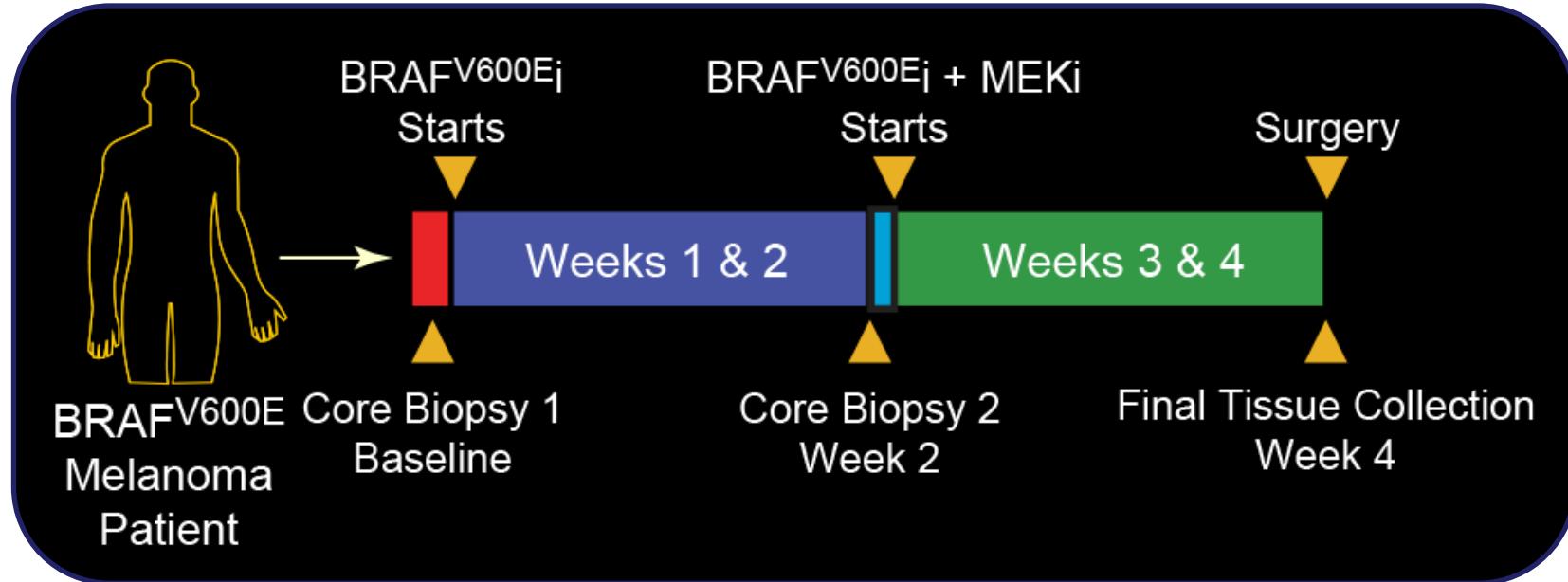


# Tissue Biopsies Go from the Operating Room to the Lab and Viable Single Cells Isolated for Functional Studies



Protocol: Leelatian et al., *Current Protocols in Molecular Biology* 2017  
Original research: Based on Leelatian and Doxie et al., *Cytometry B* 2017  
CIC services: <https://my.vanderbilt.edu/cancerimmunology/>

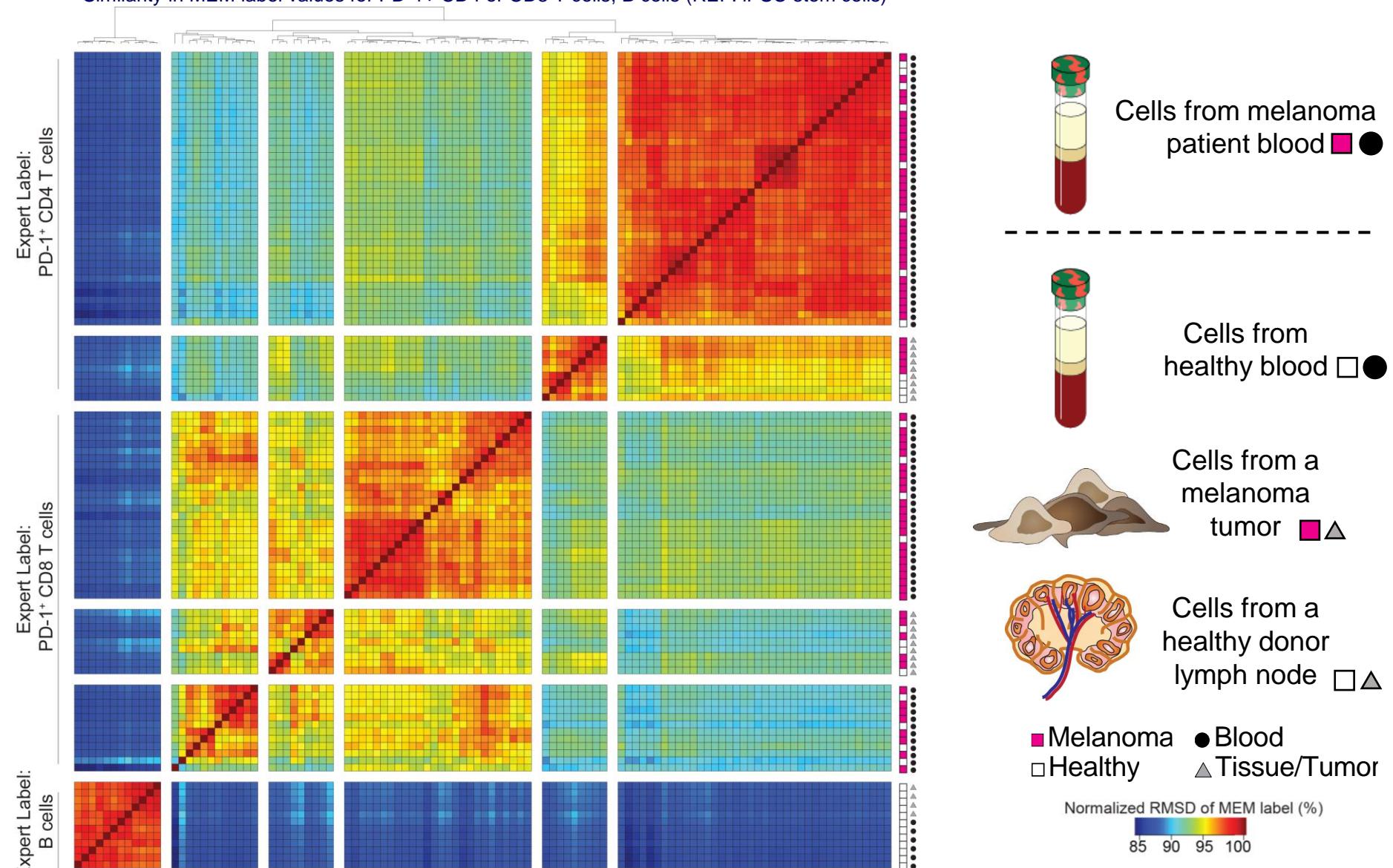
# Identify and Characterize Tumor & Immune Cells That Regress or Persist During Clinical Trial Therapies



Melanoma patient tissues are brought directly from surgery to the lab for mass cytometry + machine learning single cell analysis

# Distinct Phenotypes of PD-1<sup>+</sup> CD8<sup>+</sup> T cells in Melanoma Tumors Revealed by Quantitatively Comparing MEM Text Labels

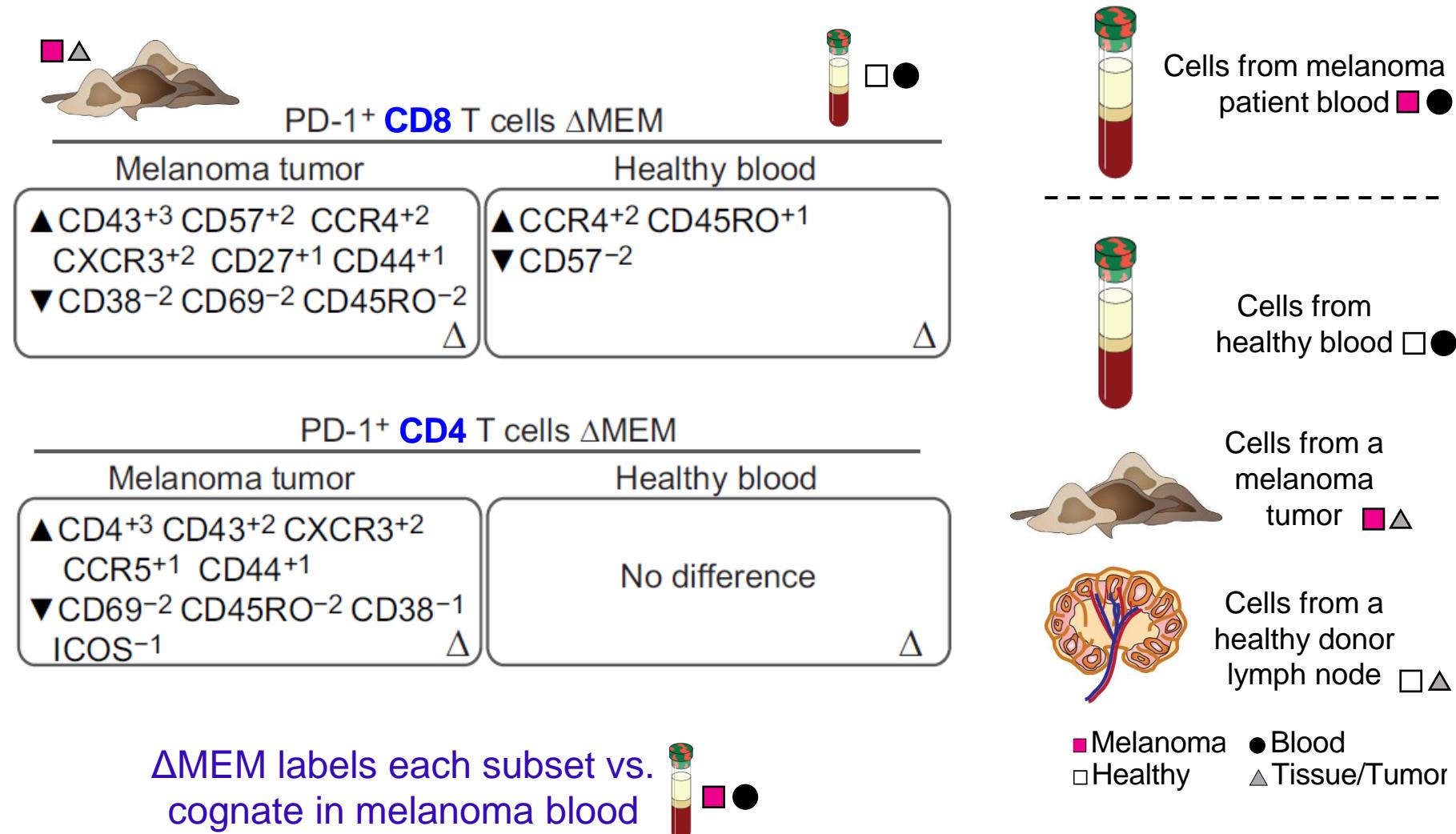
Similarity in MEM label values for PD-1<sup>+</sup> CD4 or CD8 T cells, B cells (REF: iPSC stem cells)



Greenplate et al., *Cancer Immunology Research* 2019

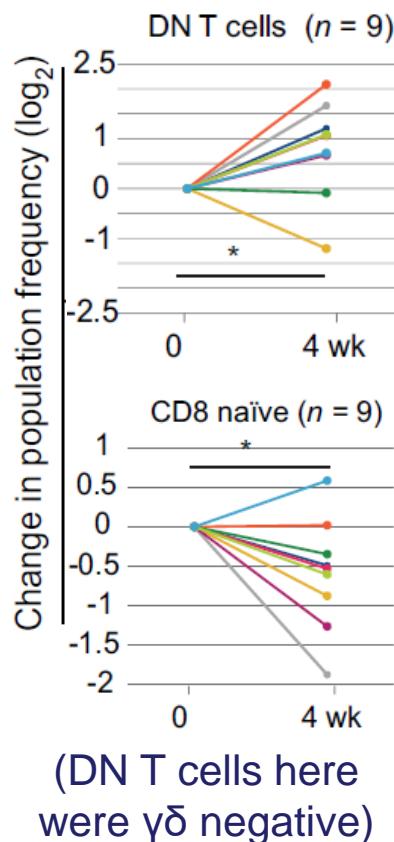
Methods: Diggins et al., *Nature Methods* 2017; *Curr Prot Cyt* 2018

# Blood vs. Tumor: CXCR3 Is Enriched on PD-1<sup>+</sup> CD4 and CD8 Melanoma TIL; CD57 Is Gained on Tumor Infiltrating PD-1<sup>+</sup> CD8

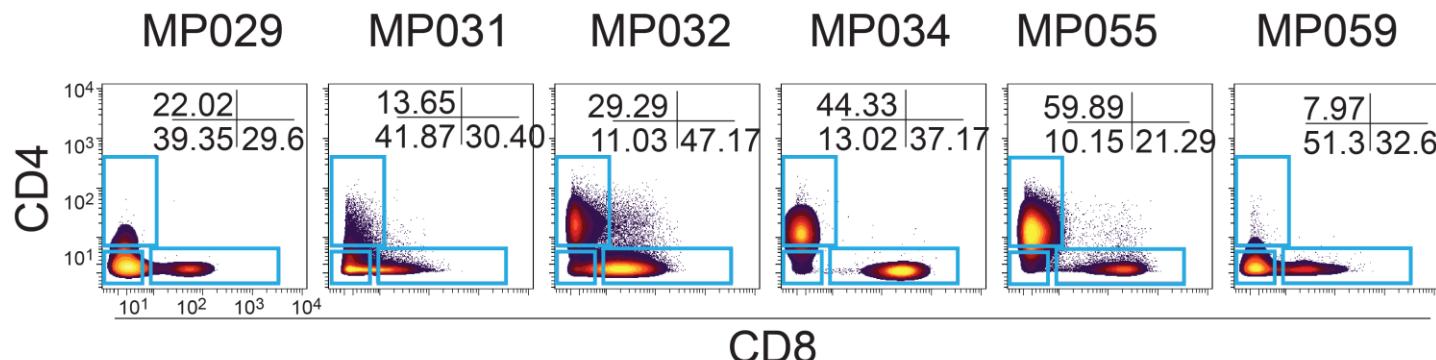


Greenplate et al., *Cancer Immunology Research* 2019  
Methods: Diggins et al., *Nature Methods* 2017; *Curr Prot Cyt* 2018

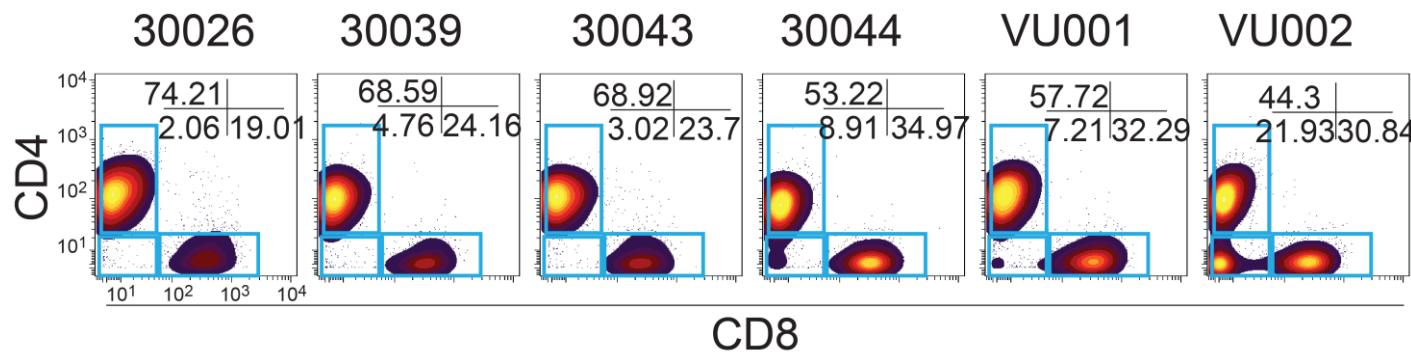
# Abnormal CD4<sup>-</sup> CD8<sup>-</sup> Double Negative (DN) T cells Are Enriched in Melanoma Tumors Following BRAFi + MEKi



Melanoma patient biopsy, all CD3+ tumor-infiltrating T cells  
(4 weeks after BRAF + MEK inhibitor therapy start)

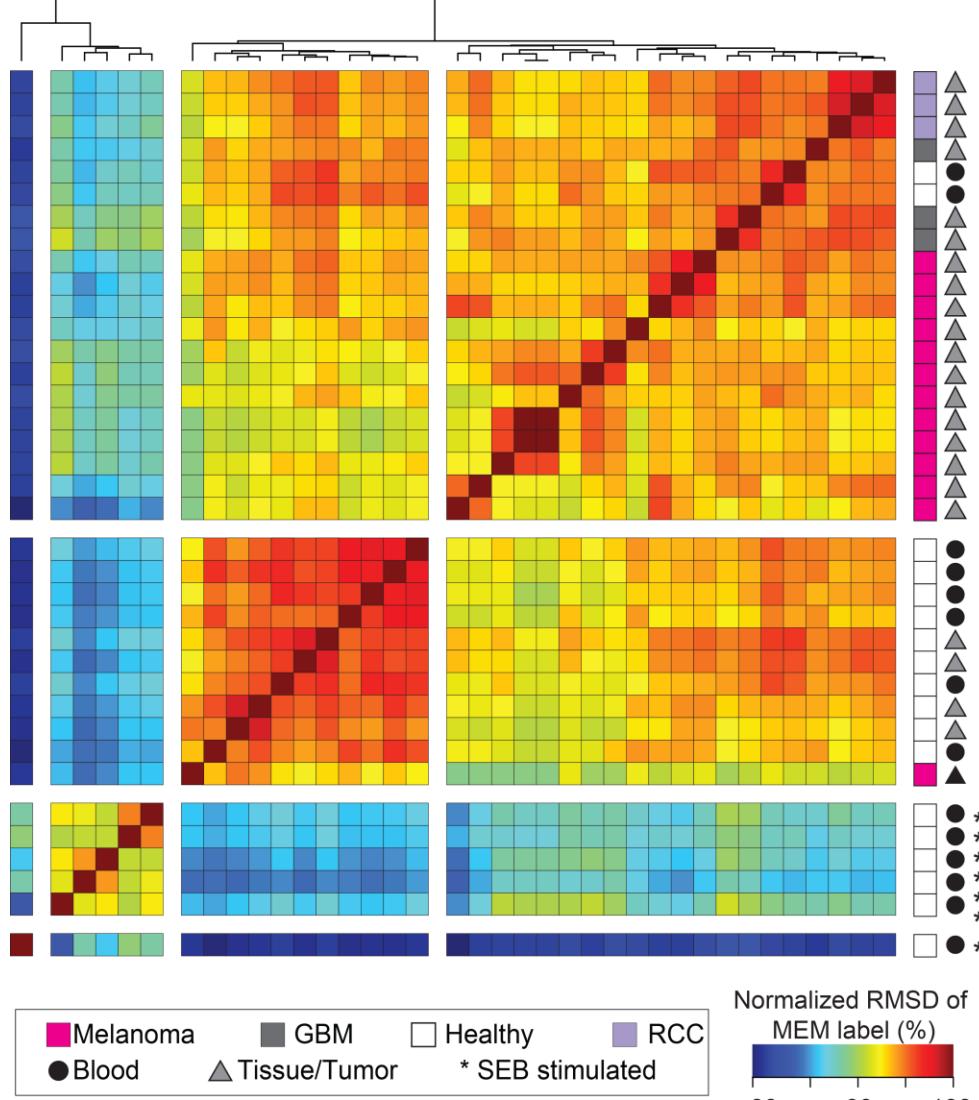


Comparison point: activated T cells from healthy blood  
(Staphylococcal enterotoxin B stimulated)



# Similar Abnormal Tumor-infiltrating CD4<sup>-</sup> CD8<sup>-</sup> DN T cells Are Observed Across Diverse Tumor Types

Comparing MEM labels for DN T cells across human tissues



Median MEM Labels iPSCs as reference

- RCC Tumor**
  - ▲ CD45<sup>+7±1</sup> CD45RO<sup>+2±0</sup> CD3<sup>+2±0</sup> CD44<sup>+2±1</sup>
  - ▼ CD57<sup>-5±0</sup> CD56<sup>-3±0</sup> CCR7<sup>-1±0</sup> PD1<sup>-1±0</sup>
  - CD45RA<sup>-1±1</sup>
- GBM Tumor**
  - ▲ CD45<sup>+8±1</sup> CD45RO<sup>+3±1</sup> CD3<sup>+2±0</sup> CD69<sup>+2±1</sup>
  - CD44<sup>+2±1</sup>
  - ▼ CD57<sup>-5±2</sup> CD56<sup>-3±1</sup> PDL1<sup>-2±0</sup> CCR5<sup>-1±1</sup>
- Melanoma Tumor**
  - ▲ CD45<sup>+7±3</sup> CD3<sup>+3±1</sup> CD45RO<sup>+2±3</sup> CD44<sup>+2±3</sup>
  - ▼ CD57<sup>-6±2</sup> CCR4<sup>-3±1</sup> CD56<sup>-2±1</sup> PDL1<sup>-2±1</sup>
- Non-malignant LN**
  - ▲ CD45<sup>+8±1</sup> CD3<sup>+4±0</sup> CD44<sup>+2±1</sup> CD45RA<sup>+1±2</sup>
  - ▼ CD57<sup>-5±1</sup> CCR4<sup>-3±0</sup> CD56<sup>-3±0</sup> PDL1<sup>-2±0</sup>
  - CCR7<sup>-1±0</sup> CXCR3<sup>-2±0</sup>
- Healthy Donor PBMC**
  - ▲ CD45<sup>+10±1</sup> CD3<sup>+4±1</sup> CD44<sup>+2±1</sup> CD45RA<sup>+1±2</sup>
  - ▼ CD57<sup>-5±2</sup> CCR4<sup>-3±0</sup> CD56<sup>-3±0</sup> PDL1<sup>-2±0</sup>
  - CCR7<sup>-1±0</sup> CD28<sup>-1±1</sup> CXCR3<sup>-1±1</sup> CCR5<sup>1±1</sup>
- Activated T cells PBMC**
  - ▲ CD45<sup>+10±1</sup> CD45RO<sup>+6±2</sup> CD69<sup>+5±1</sup> CD4<sup>+4±2</sup>
  - CD3<sup>+3±1</sup> CCR7<sup>+3±2</sup> HLADR<sup>+2±1</sup> CD8<sup>+2±1</sup>
  - CD27<sup>+2±2</sup> CD25<sup>+2±3</sup>
  - ▼ CD57<sup>-5±1</sup> CD16<sup>-1±0</sup>

Greenplate et al., *Cancer Immunology Research* 2019

Methods: Diggins et al., *Nature Methods* 2017; *Curr Prot Cyt* 2018