



Status: Lecture 11

25.11.2019	Helena	hands-on session #2: cytometry	cytof null comparison	ICA-Based Clustering of Genes from Microarray Expression Data (LK, MP, RZ)	X
02.12.2019	Mark	single-cell 2: cell type definition, differential state	scRNA exercise 2	Capturing Heterogeneity in Gene Expression Studies by Surrogate Variable Analysis (C.B, T.F)	Molecular Cross-Validation for Single-Cell RNA-seq (AY, SM, GH)
09.12.2019	Pierre-Luc	hands-on session #3: single-cell RNA-seq	full scRNA-seq pipeline	X	X
16.12.2019	Mark	loose ends: HMM, EM, robustness	segmentation, peak finding	Shrinkage estimation of dispersion in Negative Binomial models for RNA-seq experiments with small sample size (AS, CP, IP)	Empirical Bayes Analysis of a Microarray Experiment (JS, CW, DS)



Projects

Reminders:

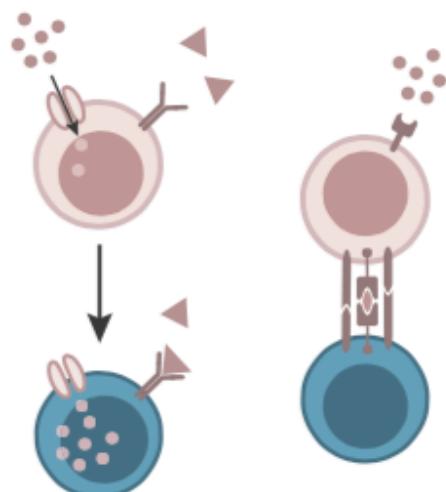
- project plan by 02.12. (a few bullet points)
- due 10.01.2020

When you have plan ready, let me know topic, group members and I will make a private github repo for it and invite you (via github classroom) to it.

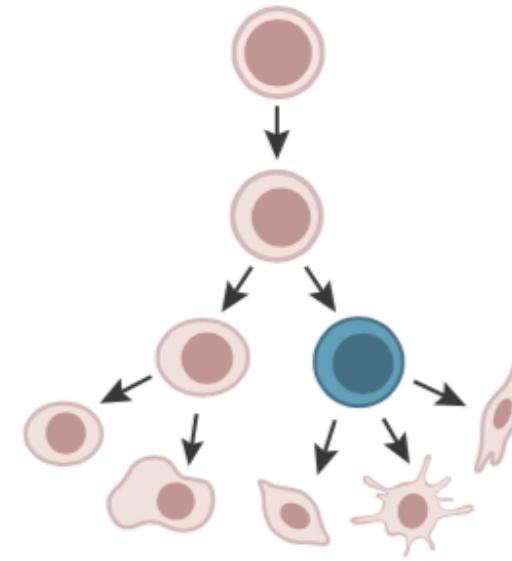
One topic idea: *compare fitting subpopulation-level models (muscat) versus fitting models for all subpopulations simultaneously and exploring interaction terms*

a

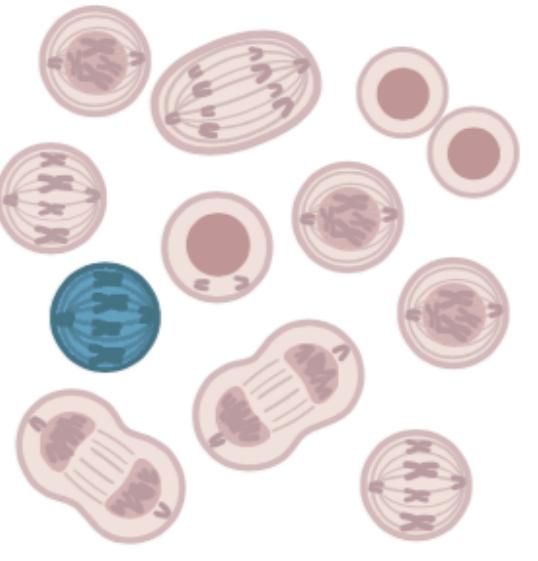
Environmental stimuli



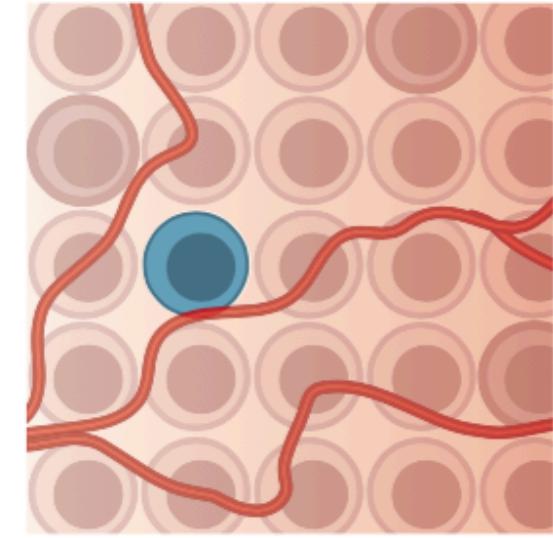
Cell development



Cell cycle



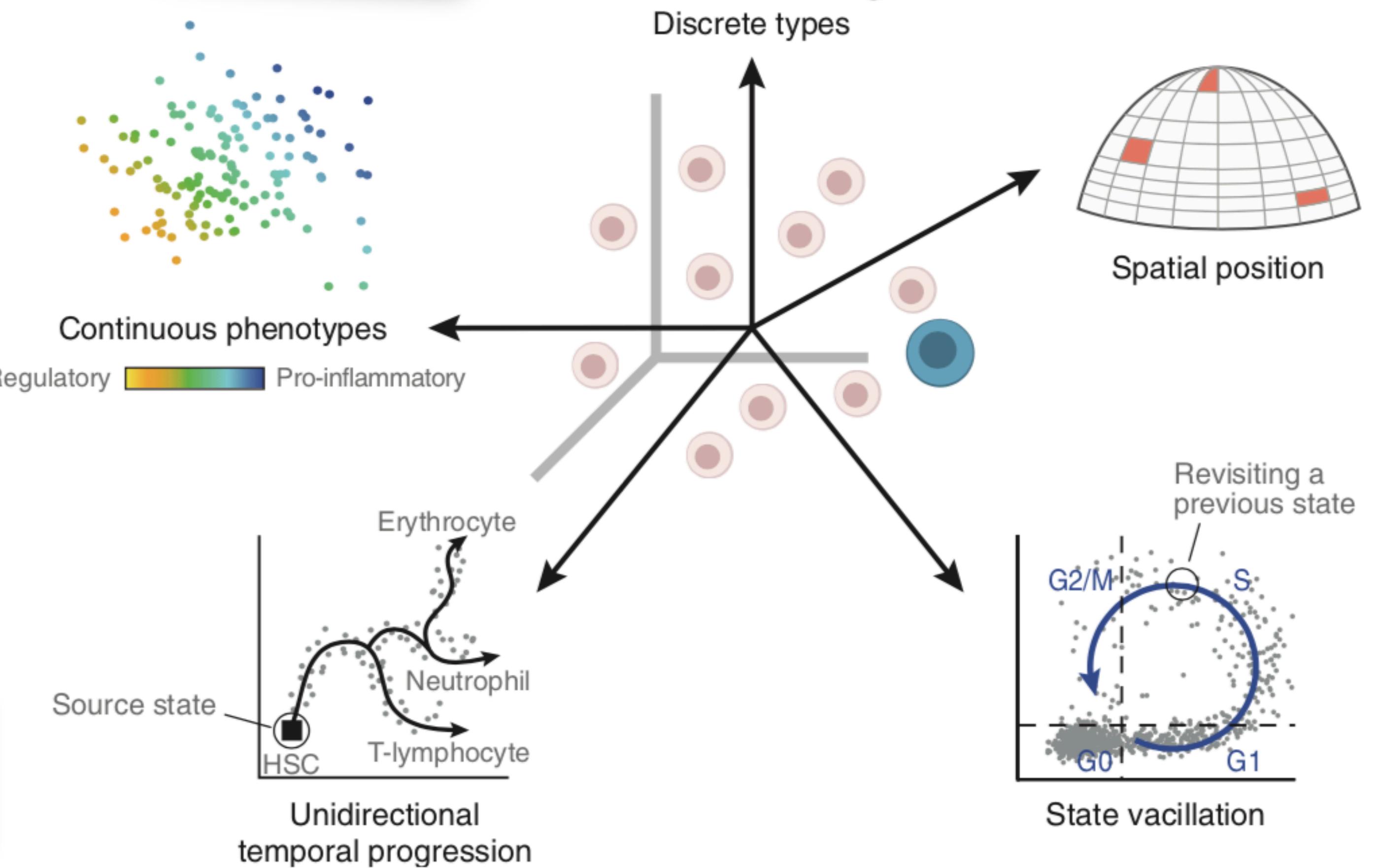
Spatial context



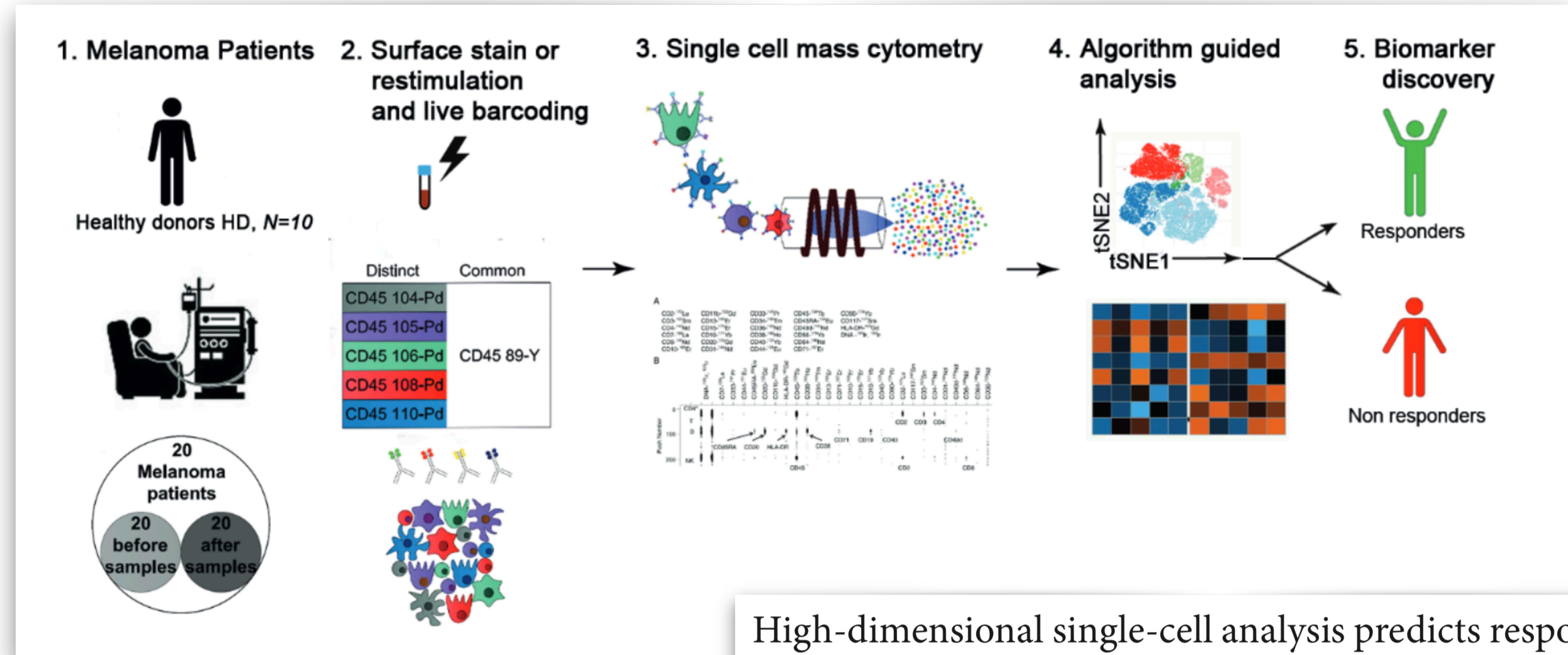
Applications

Revealing the vectors of cellular identity with single-cell genomics

Allon Wagner¹, Aviv Regev^{2,3,5} & Nir Yosef^{1,4,5}



Motivation for differential analysis of single cell (cytometry) data: *finding cancer biomarkers*

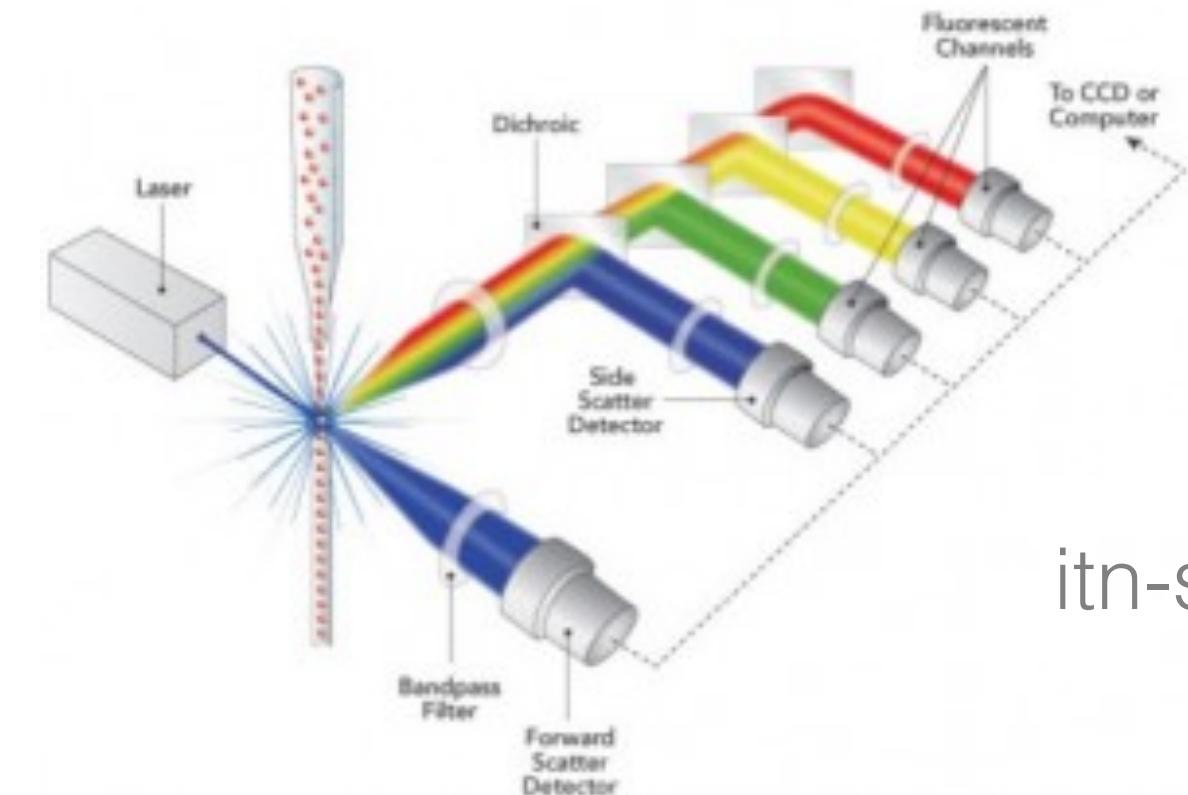


High-dimensional single-cell analysis predicts response to anti-PD-1 immunotherapy

Carsten Krieg^{1,6} , Malgorzata Nowicka^{2,3}, Silvia Guglietta⁴, Sabrina Schindler⁵, Felix J Hartmann¹ , Lukas M Weber^{2,3} , Reinhard Dummer⁵, Mark D Robinson^{2,3} , Mitchell P Levesque^{5,7} & Burkhard Becher^{1,7}

High-dimensional cytometry

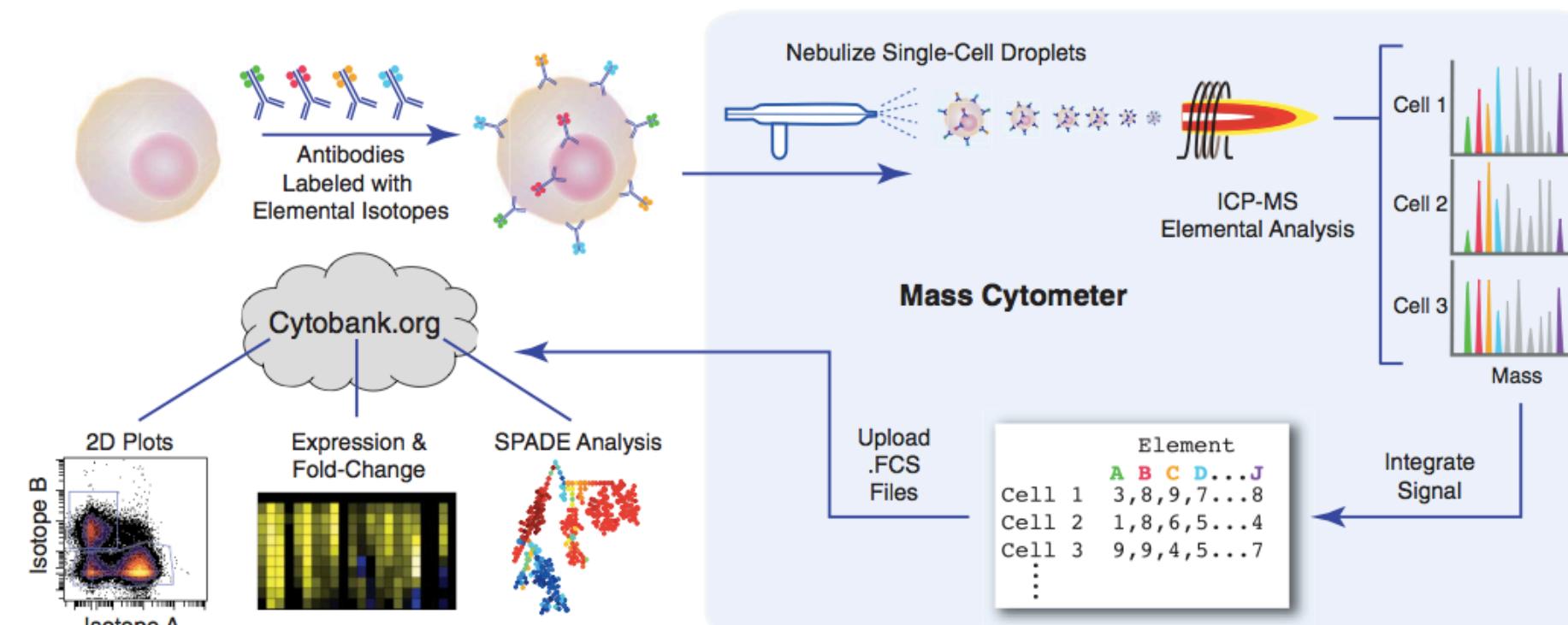
Measure **targeted protein expression levels**
in single cells using **antibodies**



itn-snal.net

(A) Fluorescent flow cytometry / FACS

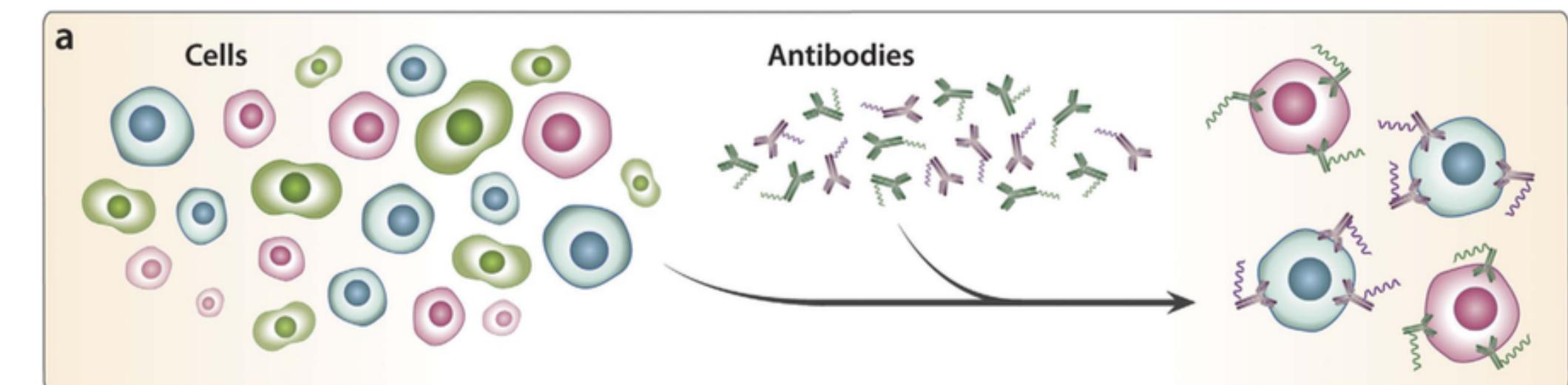
>20 proteins/cell; 1000s cells/sec,
non-destructive



Bendall et al.
(2011), Fig. 1A

(B) Mass cytometry / CyTOF

>40 proteins/cell; 100s cells/sec,
destructive



(C) sequence-based cytometry

>100 proteins/cell; destructive

Shahi et al. (2017), Fig. 1A

Flow cytometry

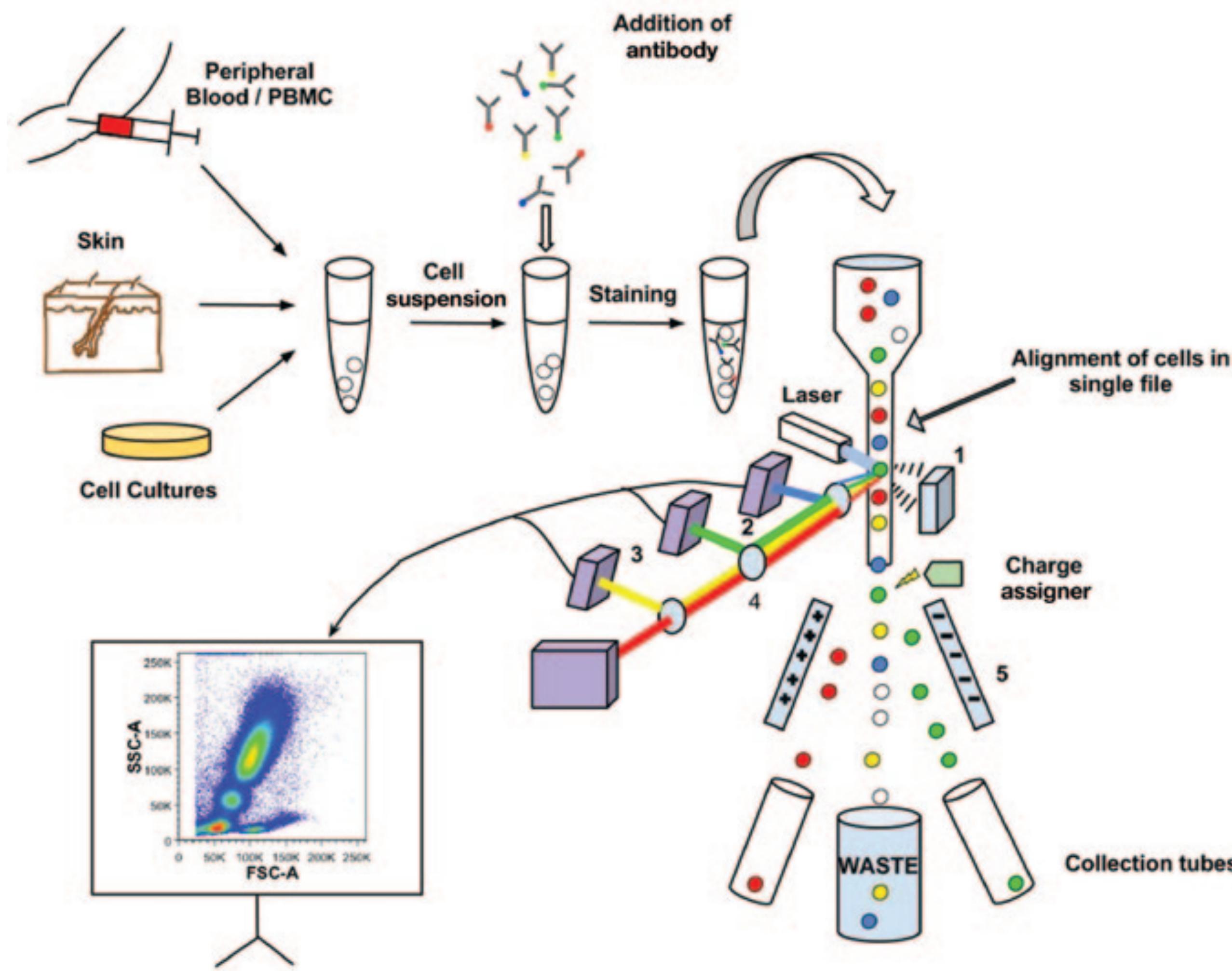
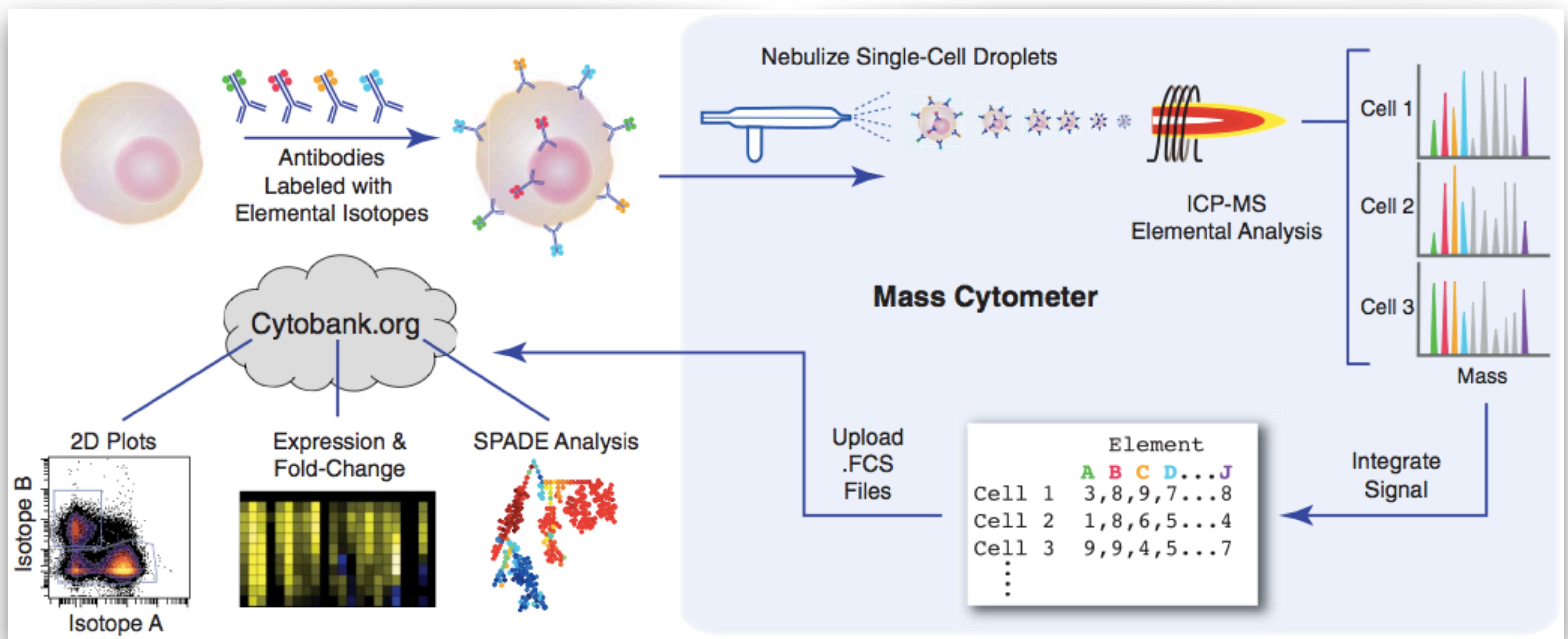
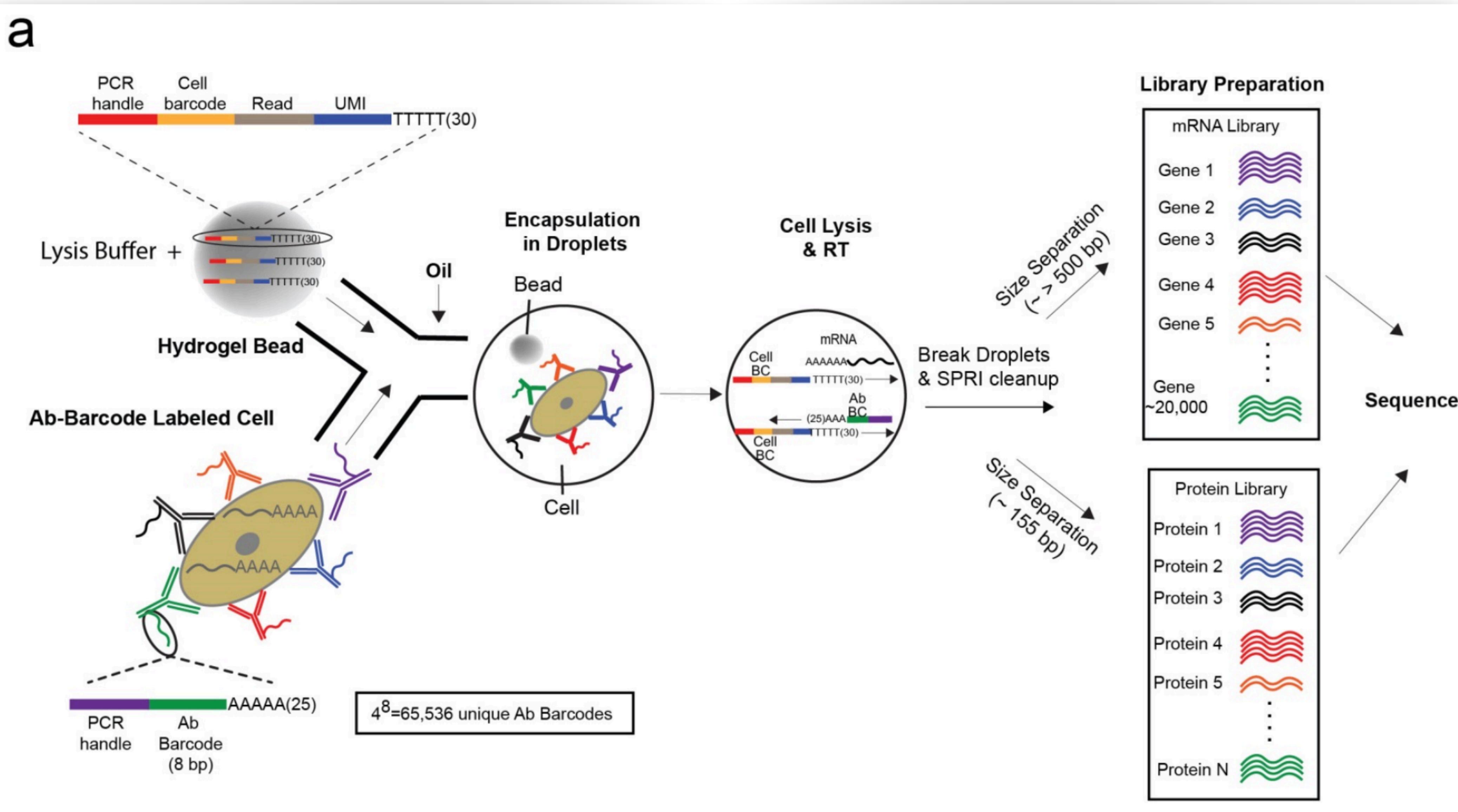


Figure 1. Schematic representation of a flow cytometer. For details please see text. (1) Forward-scatter detector, (2) side-scatter detector, (3) fluorescence detector, (4) filters and mirrors, and (5) charged deflection plates.

Mass cytometry



REAP-seq / CITE-seq



Spectral overlap vs. spillover

- CyTOF = increase in the number of parameters + massive decrease in spectral overlap

- but, still three sources of signal overlap:

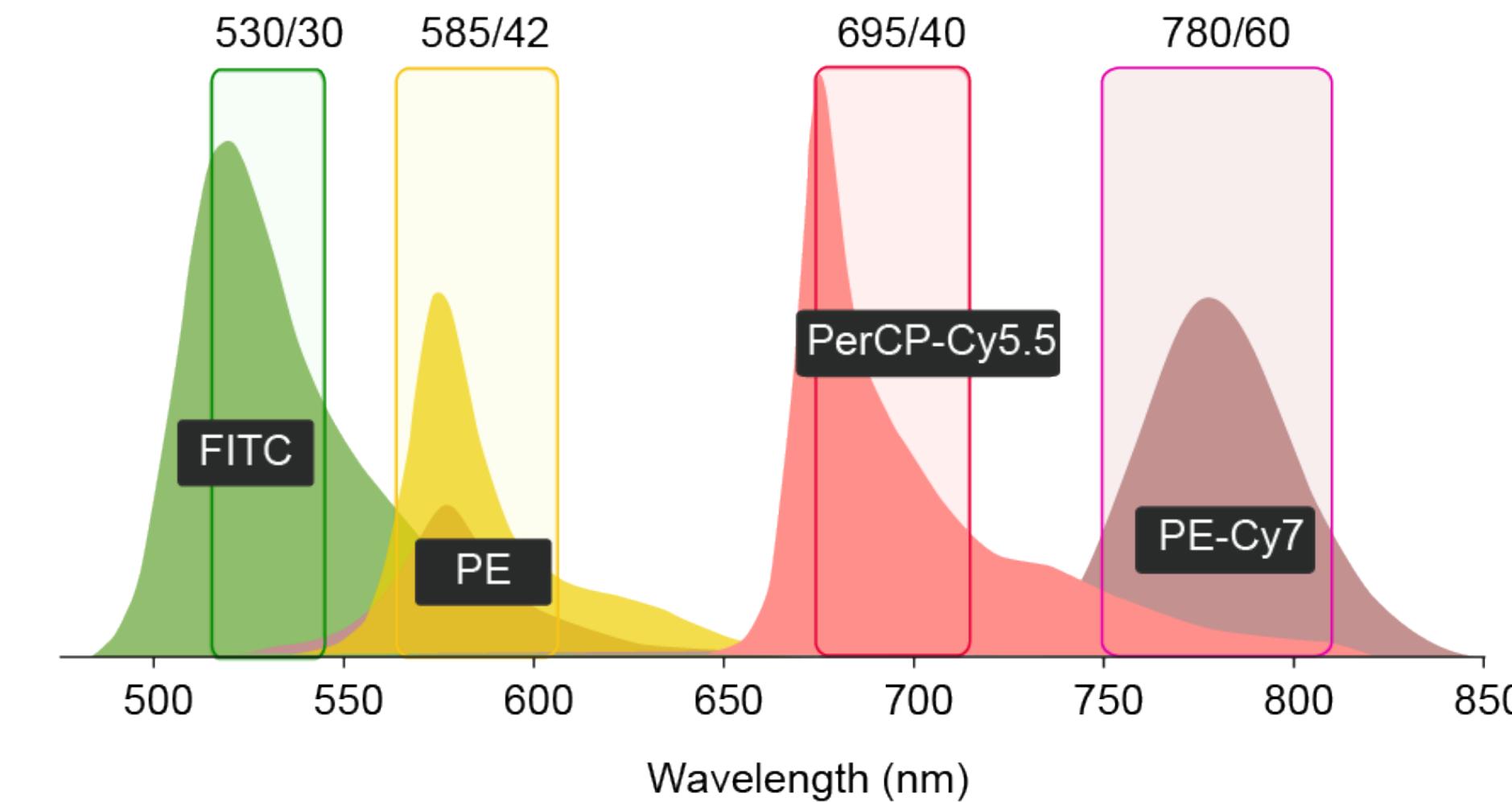
1. abundance

sensitivity := $(M \pm 1) / M$

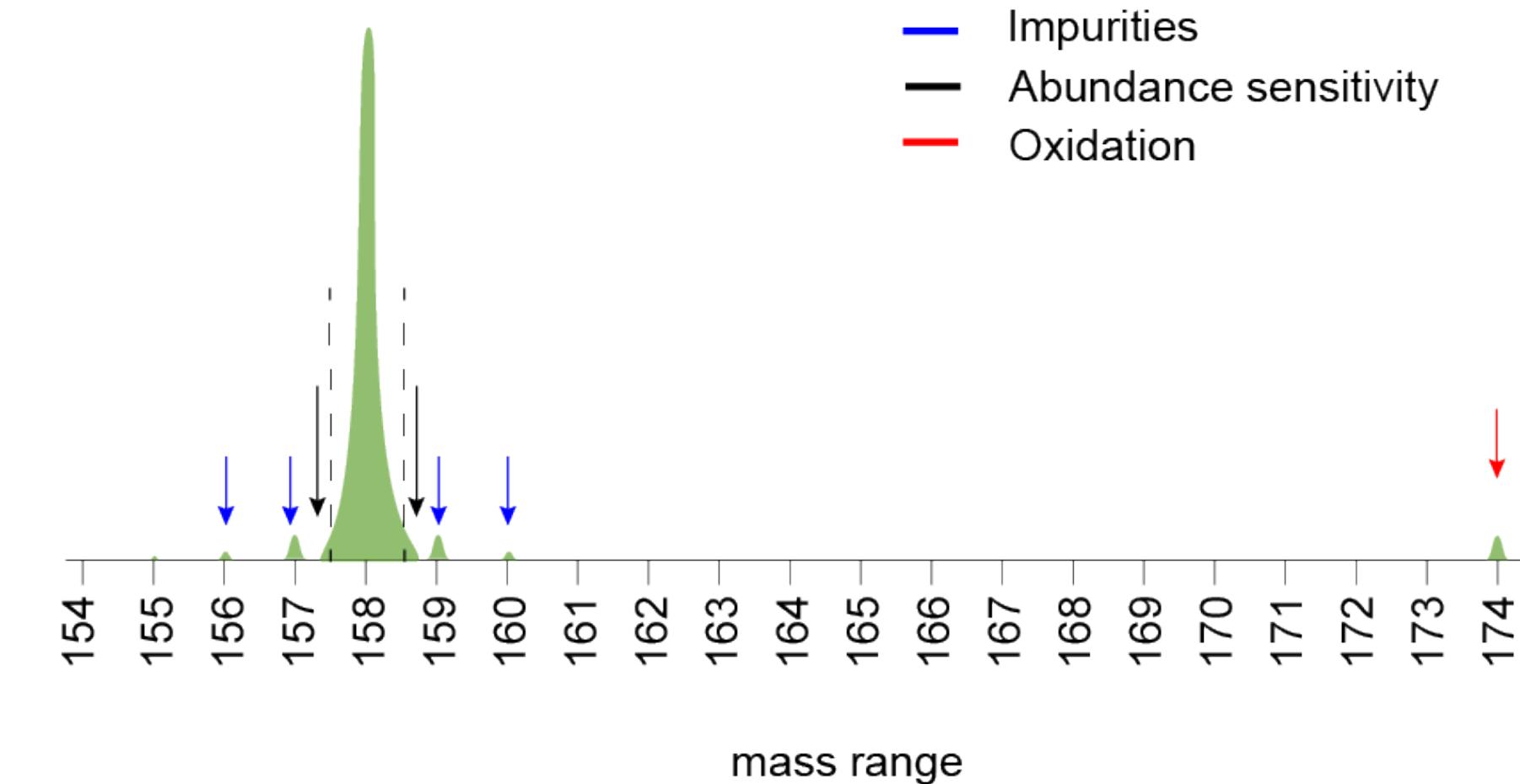
2. oxide formation: $+16M$

3. isotopic impurities: up to $\pm 6M$

FACS



Mass cytometry



Do we need to compensate CyTOF data?

The ability to multiplex up to 40 cellular subset markers in mass cytometry, without a requirement for compensation for overlap in fluorescence signals as needed in conventional flow cytometry, makes mass cytometry an ideal technology to deeply phenotype cells in complex cell populations. This feature was elegantly demonstrated by

Atkuri et al. 2015 Drug Metabolism and Disposition

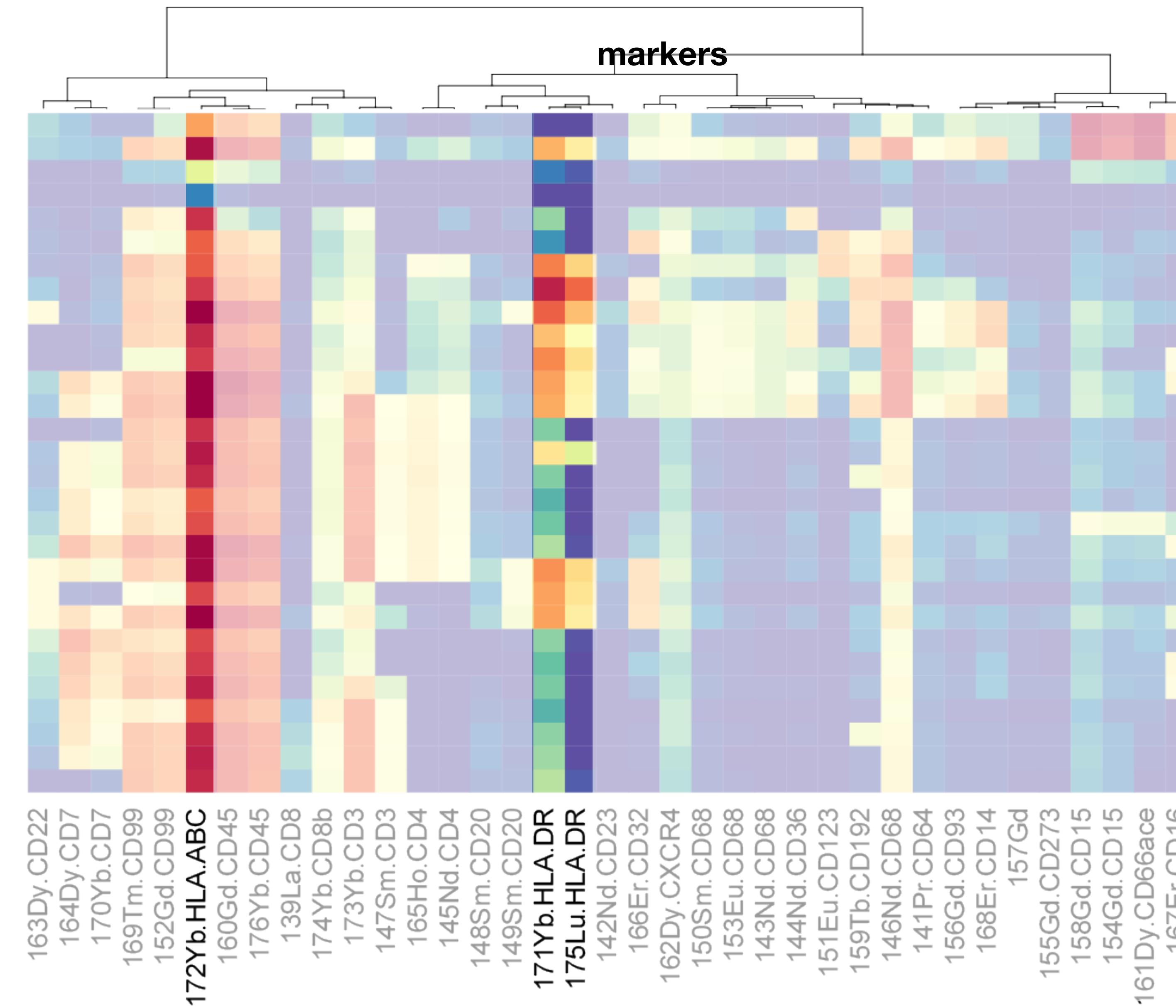
The metals that are sold as part of antibody labeling kits are of very high purity (98% and higher in most cases). As a practical matter, this means that “compensation” analogous to fluorescent antibodies is not needed, as most of the signal will be of the specified mass, with little to no signal at “M+1” or another contaminating mass. However, metal salts from other commercial sources may be of lesser purity. For example, the

Leipold al. 2015 Immunosenescence: Methods and Protocols, Methods in Molecular Biology

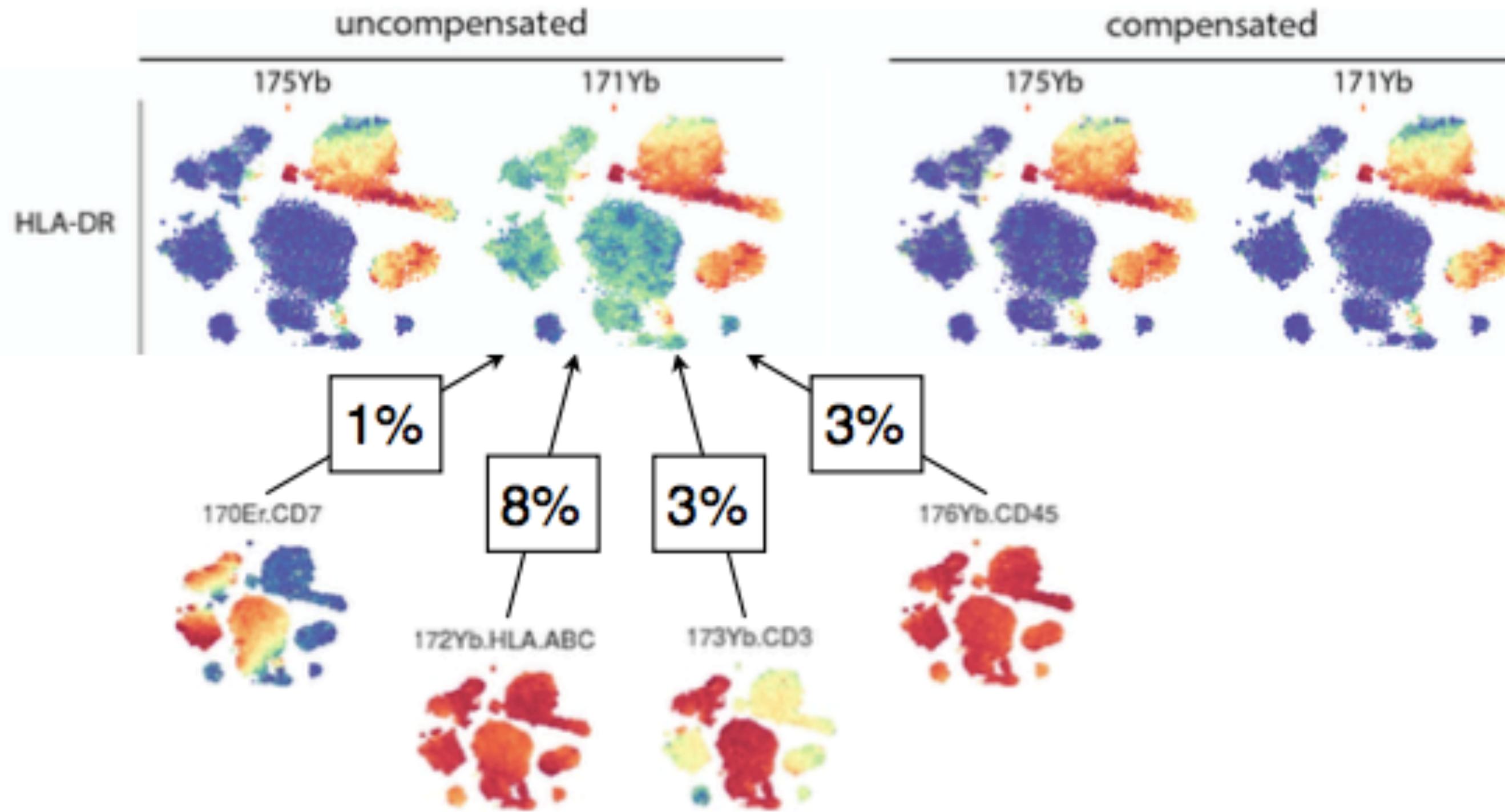
It should be made clear, though, that “minimal spillover” does not equal “no spillover.”

Short answer: Yes, even ~2-5% of a high signal can be significant

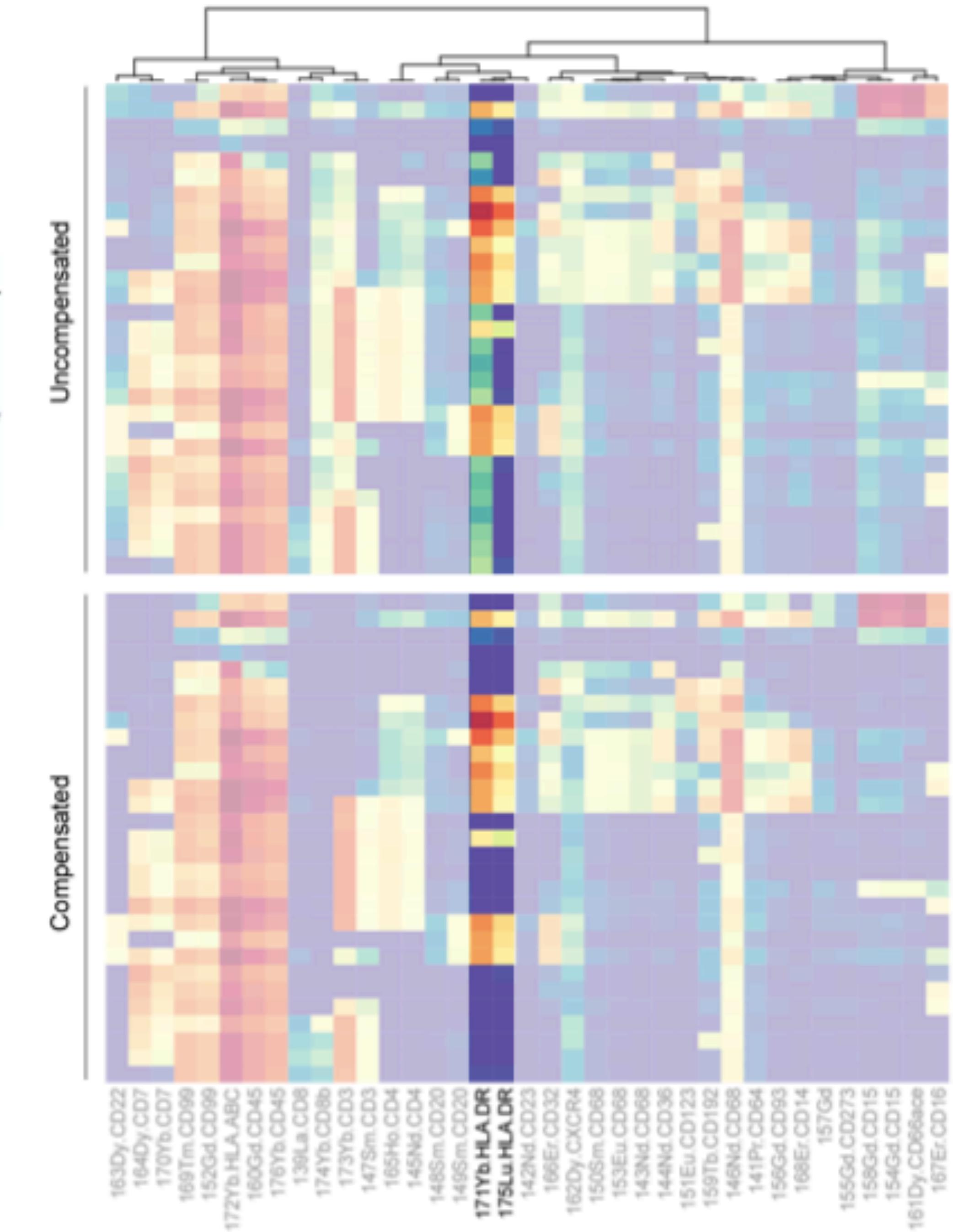
PBMCs measured,
clustered with
Phenograph; several
antibodies used twice
with different metals



Correction of spillover artefacts on a 36ab panel

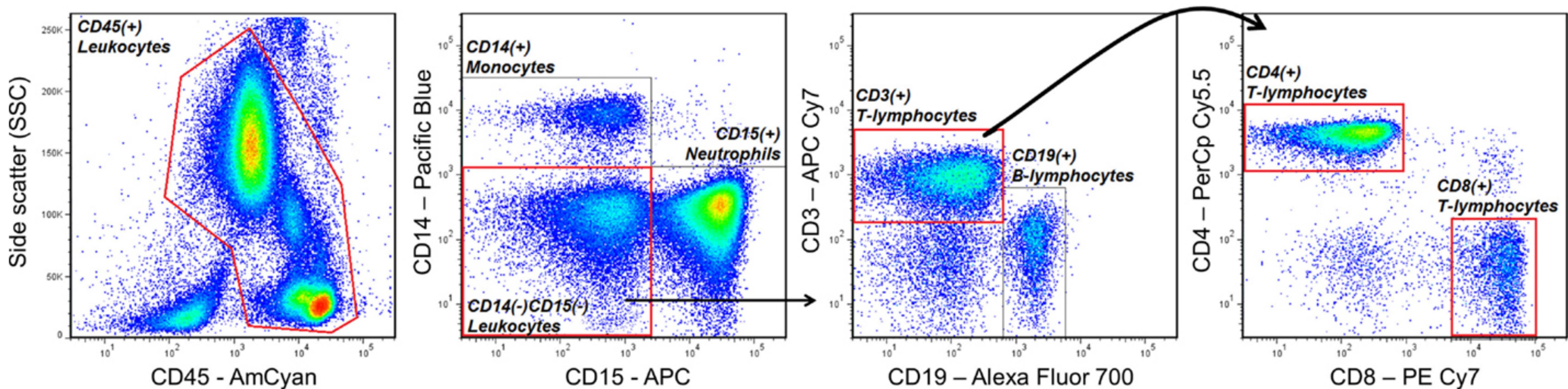


Spillover matrix estimated via single-stain beads:
non-negative least squares



Manual gating

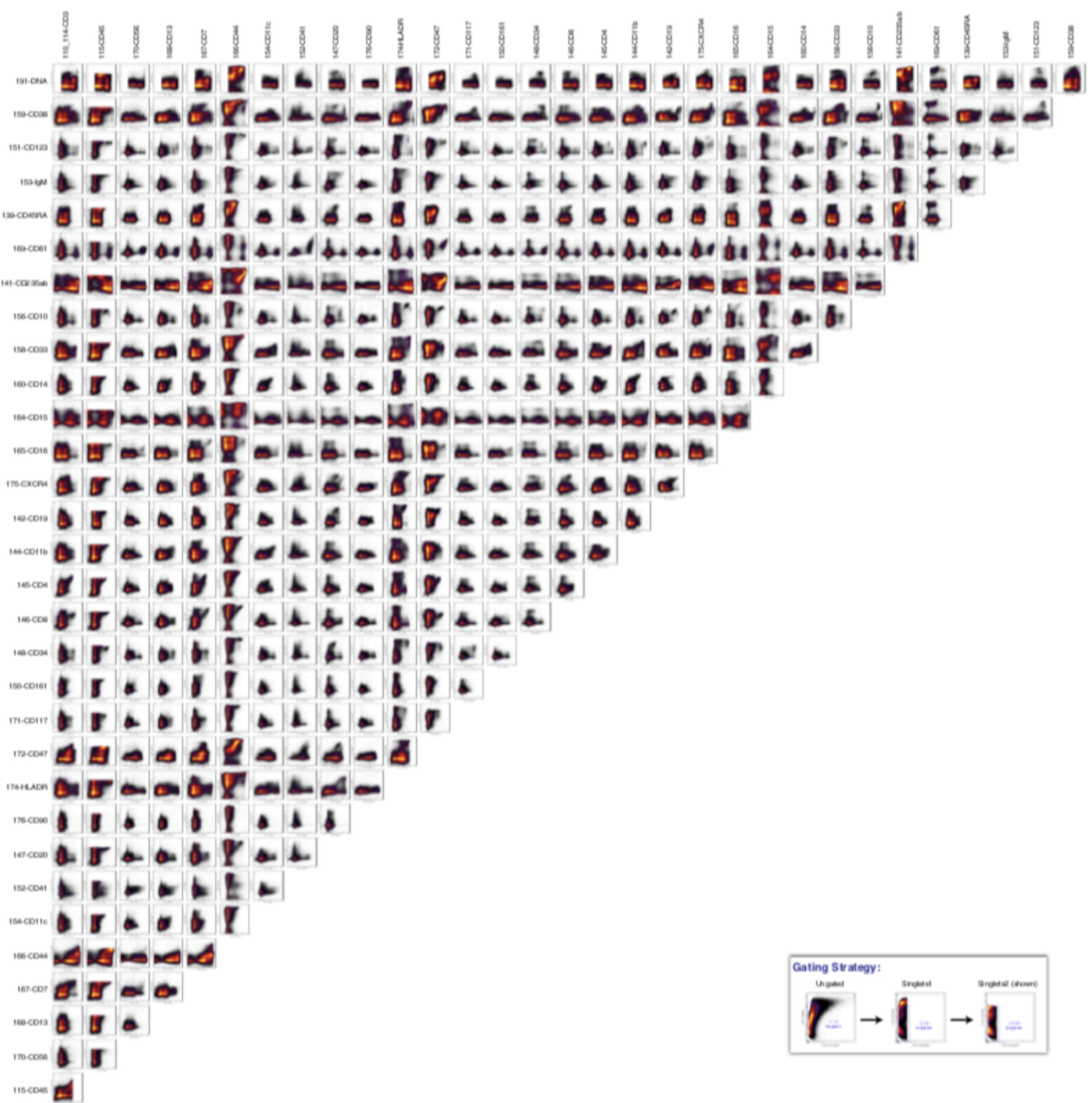
Identification of cell populations



Verschoor et al. (2015)

Manual gating

Not feasible in high-dimensional data



Bendall et al. (2011), Supp.

Clustering high-dimensional flow and mass cytometry

Motivation: Many new computational methods, explosion in the number of dimensions (both FACS and CyTOF) — what works “best”?



Lukas

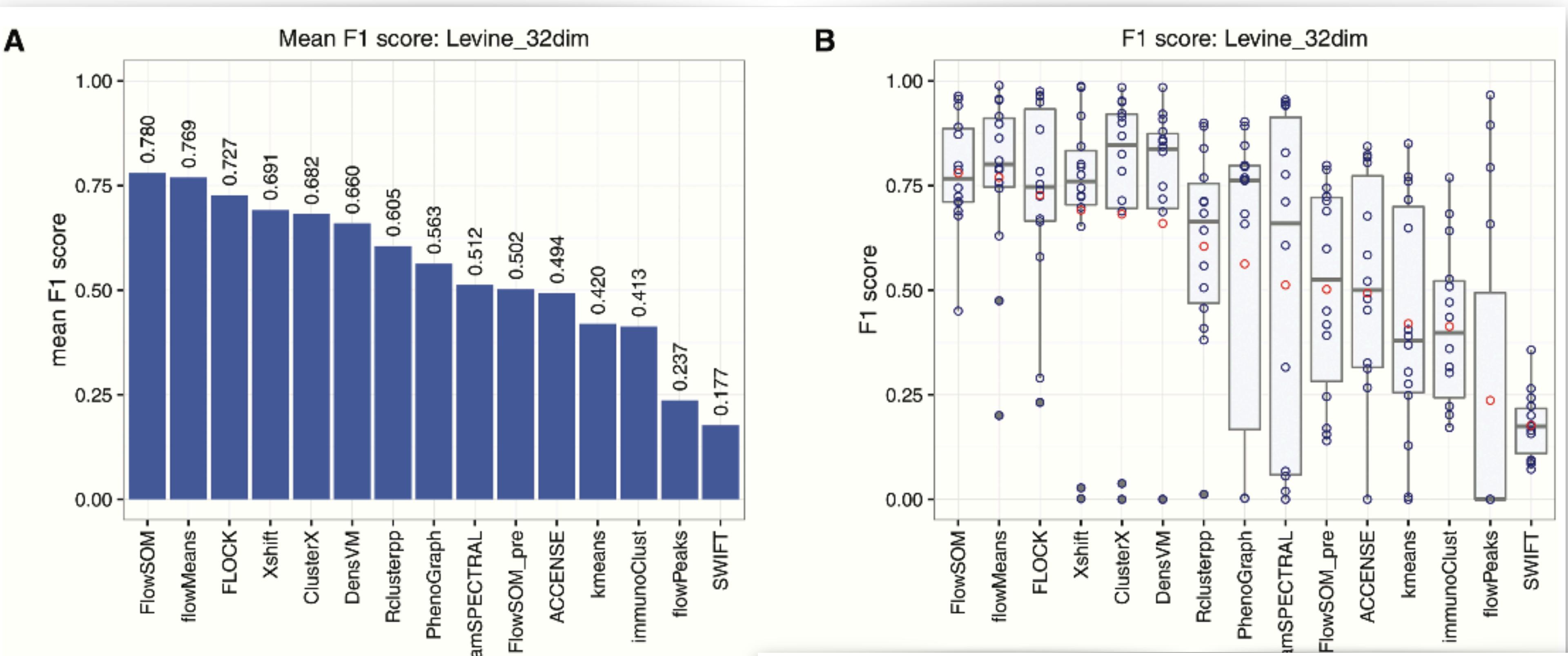
EDITOR'S CHOICE

Cytometry
PART A
Journal of the International Society for Advancement of Cytometry

Comparison of Clustering Methods for High-Dimensional Single-Cell Flow and Mass Cytometry Data

Lukas M. Weber,^{1,2} Mark D. Robinson^{1,2*}

Comparison of clustering methods



F1 score

From Wikipedia, the free encyclopedia

"F score" redirects here. For the significance test, see [F-test](#).

In statistical analysis of [binary classification](#), the **F₁ score** (also **F-score** or **F-measure**) is a measure of a test's accuracy. It considers both the [precision](#) p and the [recall](#) r of the test to compute the score: p is the number of correct positive results divided by the number of all positive results, and r is the number of correct positive results divided by the number of positive results that should have been returned. The F₁ score can be interpreted as a weighted average of the [precision](#) and [recall](#), where an F₁ score reaches its best value at 1 and worst at 0.

The traditional F-measure or balanced F-score (**F₁ score**) is the [harmonic mean](#) of precision and recall — multiplying the constant of 2 scales the score to 1 when both recall and precision are 1:

$$F_1 = 2 \cdot \frac{1}{\frac{1}{\text{recall}} + \frac{1}{\text{precision}}} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}.$$

EDITOR'S CHOICE



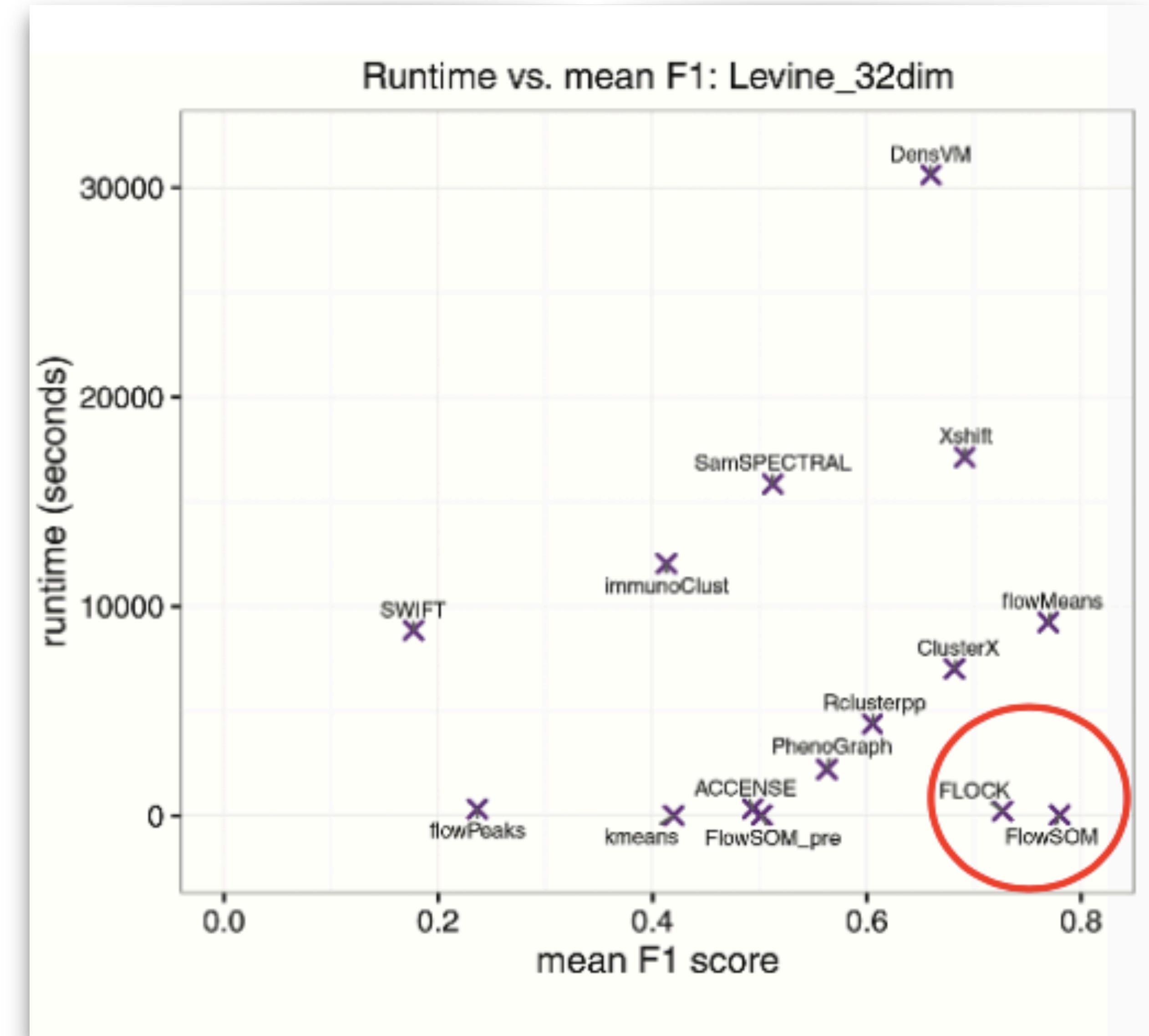
Comparison of Clustering Methods for High-Dimensional Single-Cell Flow and Mass Cytometry Data

Lukas M. Weber^{1,2} Mark D. Robinson^{1,2*}

Hungarian algorithm to match clusters to populations

Comparison of clustering methods

- several methods performed well:
FlowSOM, X-shift, PhenoGraph,
Rclusterpp, flowMeans
- **FlowSOM** gave best performance (for several data sets) and was fast
- **X-shift** gave best performance for rare cell populations
- several methods sensitive to random starts (rare populations)
- code, data freely available



Revealing the vectors of cellular identity with single-cell genomics

Allon Wagner¹, Aviv Regev^{2,3,5} & Nir Yosef^{1,4,5}

Box 1 The many facets of a cell's identity

We define a cell's identity as the outcome of the instantaneous intersection of all factors that affect it. We refer to the more permanent aspects in a cell's identity as its type (e.g., a hepatocyte typically cannot turn into a neuron) and to the more transient elements as its state. Cell types are often organized in a hierarchical

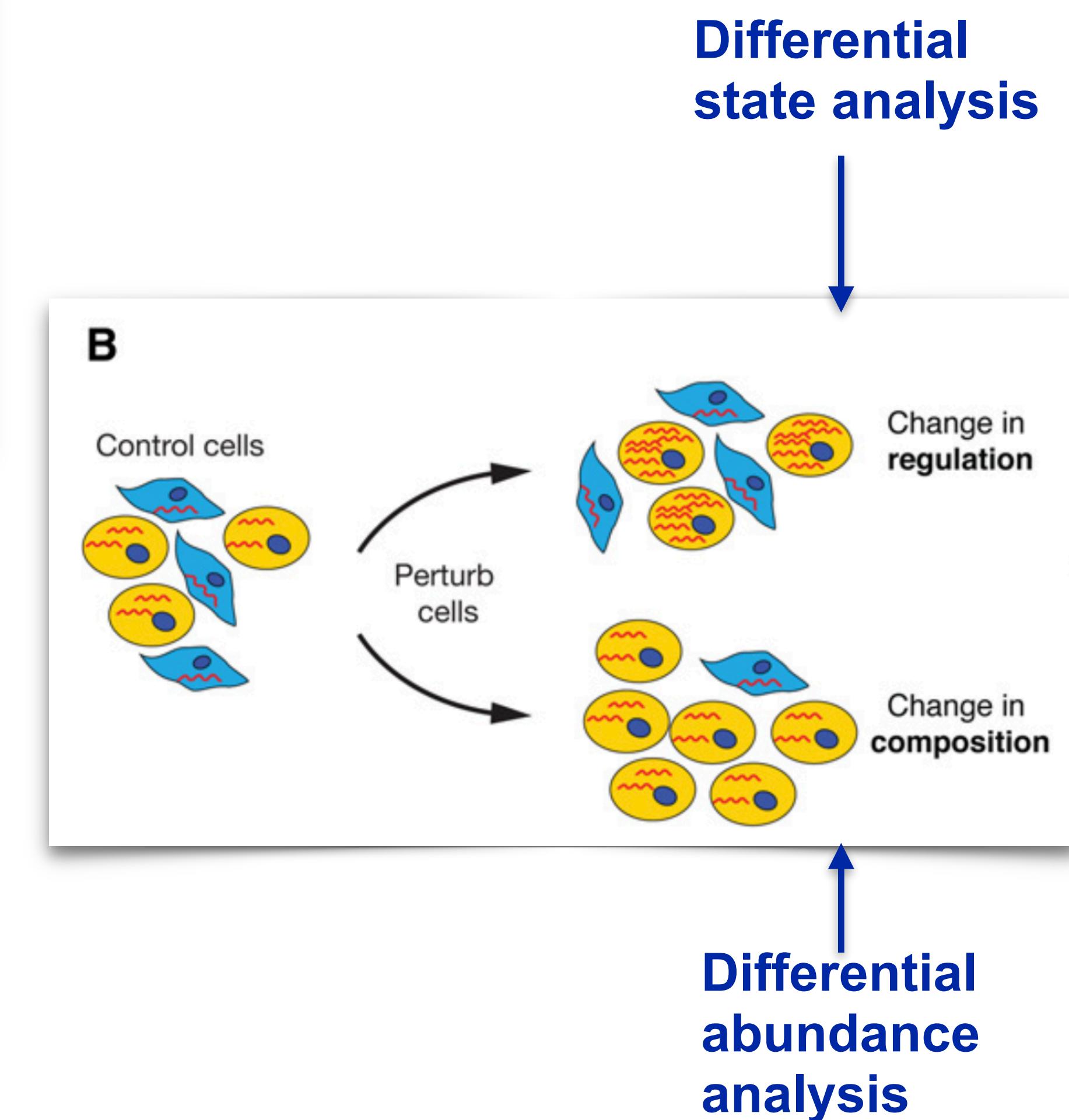
Type: more permanent
State: more transient

Perspective

Defining cell types and states with single-cell genomics

Cole Trapnell

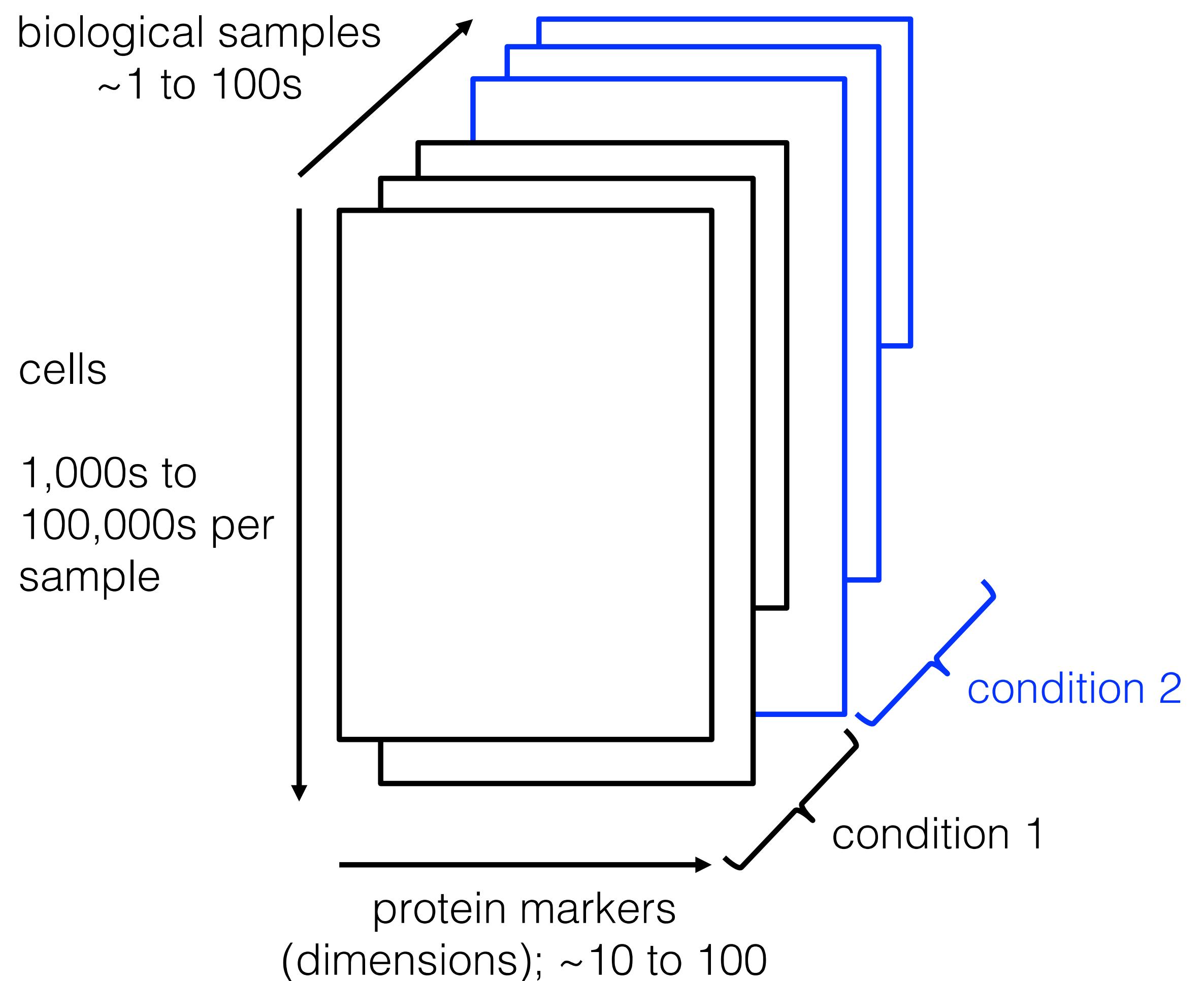
Department of Genome Sciences, University of Washington, Seattle, Washington 98105, USA



Data structure and differential analysis

Two types of differential analysis

- **differential abundance** (DA) of cell populations
- **differential states**
 - e.g., differential expression of functional proteins (e.g., signaling) within cell populations



Key elements of CyTOF workflow

- Exploration of various data aspects at each step
- Separation of **type** and **state** markers
- Put all samples together and cluster (FlowSOM or other)
- Optional: manually merge clusters (via visualizations: heatmaps, low dimensional projections)
- Differential abundance analysis (count-based model, somewhat similar to RNA-seq)
- For **state** markers, differential state analysis (aggregate and use linear model)

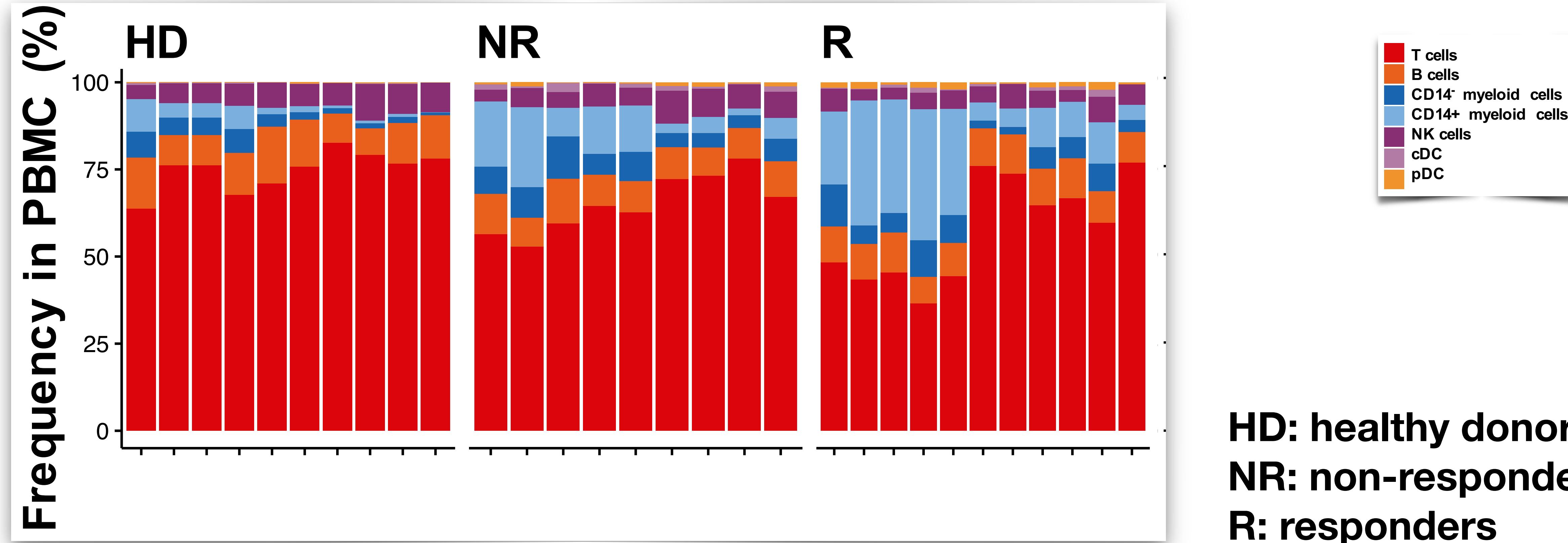
Good / Bad news: large batch effect, but nice experimental design (all conditions in every batch) so can be separated in statistical models.

High-dimensional single-cell analysis predicts response to anti-PD-1 immunotherapy

Carsten Krieg^{1,6} , Małgorzata Nowicka^{2,3}, Silvia Guglietta⁴, Sabrina Schindler⁵, Felix J Hartmann¹ , Lukas M Weber^{2,3} , Reinhard Dummer⁵, Mark D Robinson^{2,3} , Mitchell P Levesque^{5,7}  & Burkhard Becher^{1,7} 

Part 1:

Differential abundance of cell populations



After clustering (and manual merging), *generalized linear mixed model* is applied to cell count table to find differential abundance (n.b.: similar to RNA-seq differential expression).

Models for differential abundance similar to those for RNA-seq, but lower dimension



Manual merging of cell populations based on phenotypes

Generalized linear mixed models (differential abundance)

$$E(Y_{ij} | \beta_0, \beta_1, \gamma_i, \xi_{ij}) = \text{logit}^{-1}(\beta_0 + \beta_1 x_{ij} + \gamma_i + \xi_{ij}),$$

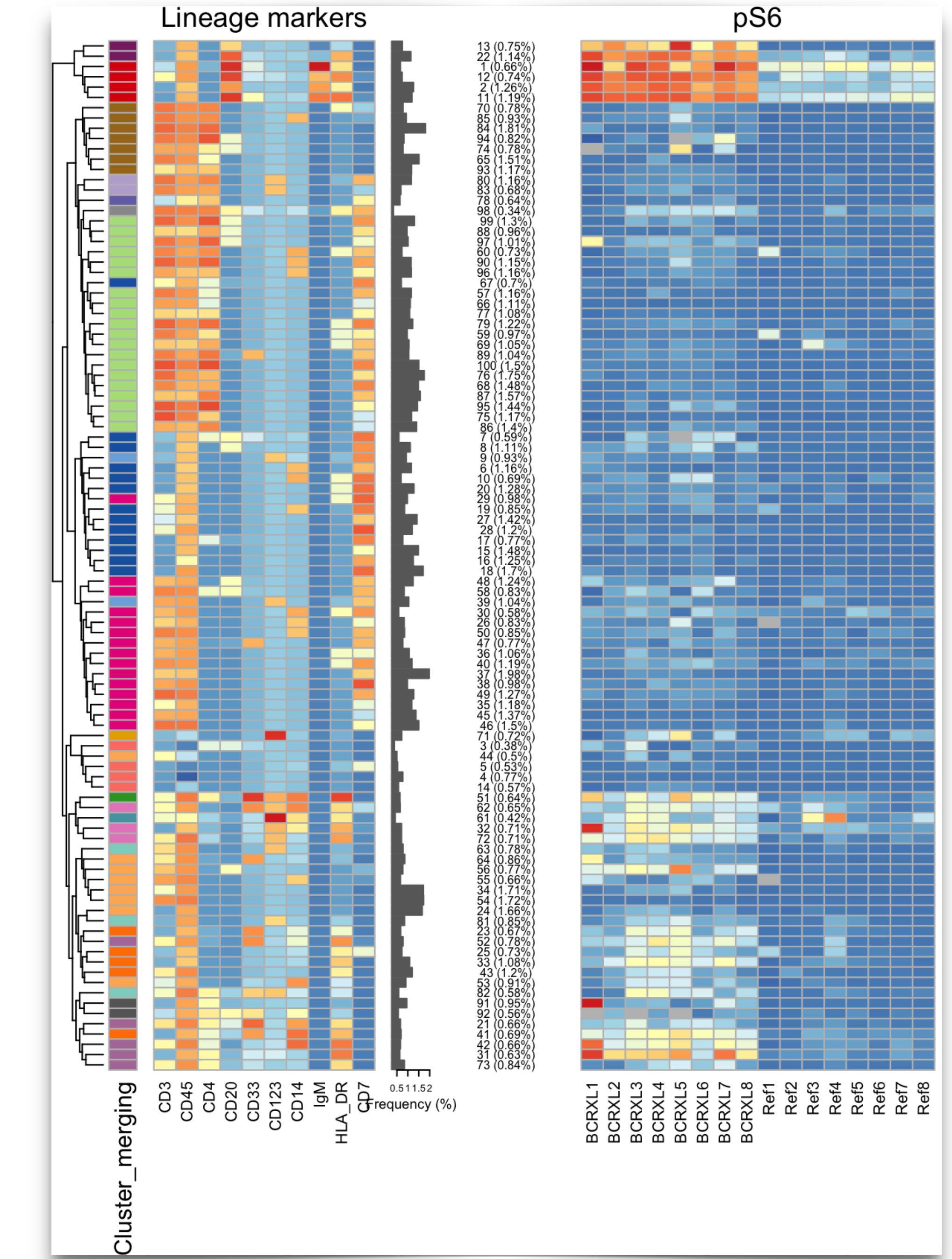
Linear mixed models (differential expression within populations)

$$Y_{ij} = \beta_0 + \beta_1 x_{ij} + \gamma_i + \epsilon_{ij},$$

Part 2: subpopulation-specific differential analyses

Cluster to some number of groups based on lineage/type markers; look across samples in functional marker

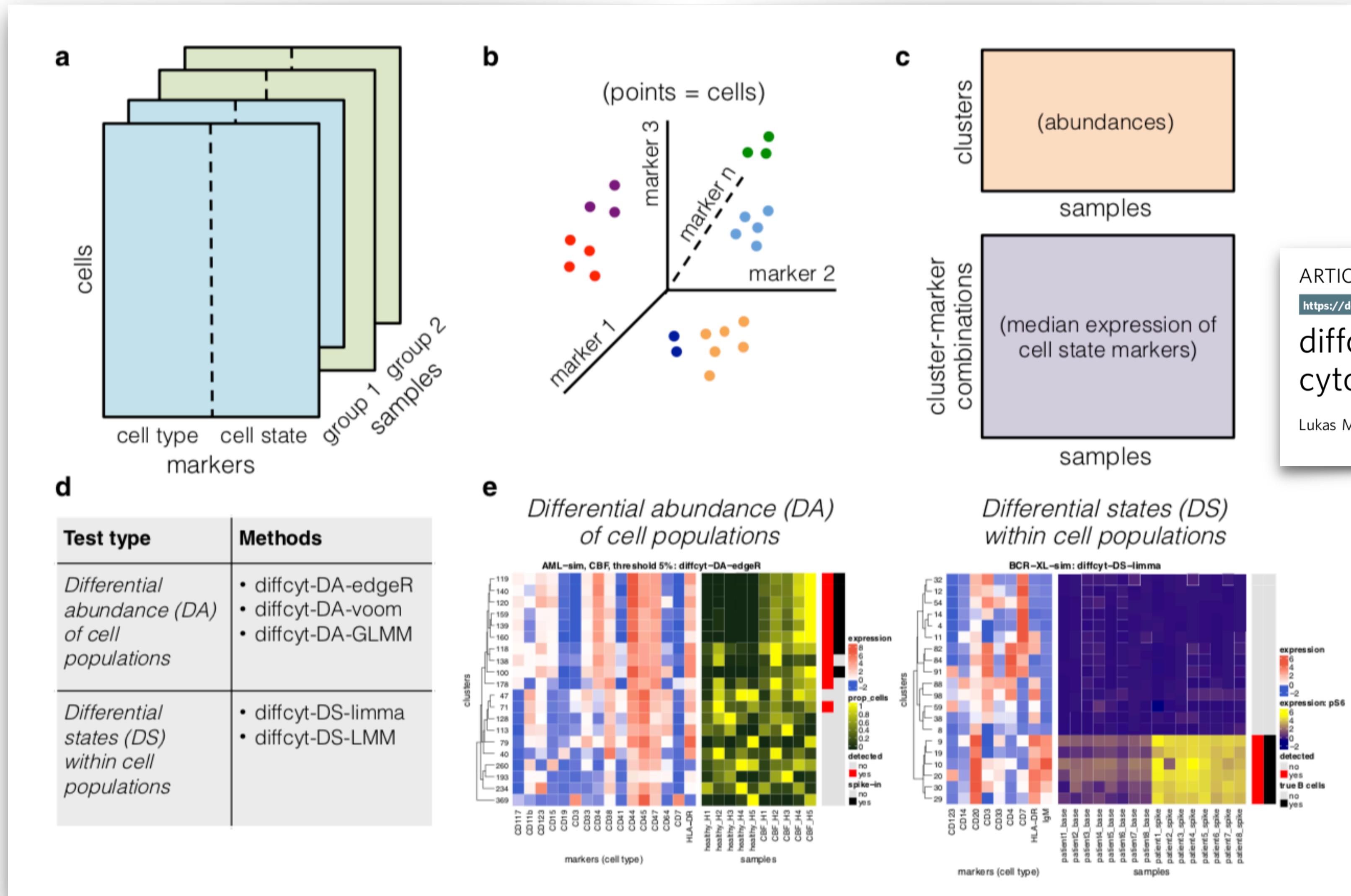
→ median lineage marker signal by cluster



diffcyt: differential tests more formalised



Lukas



ARTICLE

<https://doi.org/10.1038/s42003-019-0415-5>

OPEN

diffcyt: Differential discovery in high-dimensional cytometry via high-resolution clustering

Lukas M. Weber^{1,2}, Małgorzata Nowicka^{1,2,3}, Charlotte Soneson^{1,2,4} & Mark D. Robinson^{1,2}

Note: for differential state analysis, aggregates are always taken. We are testing this now with scRNA-seq data