



University of  
Zurich<sup>UZH</sup>



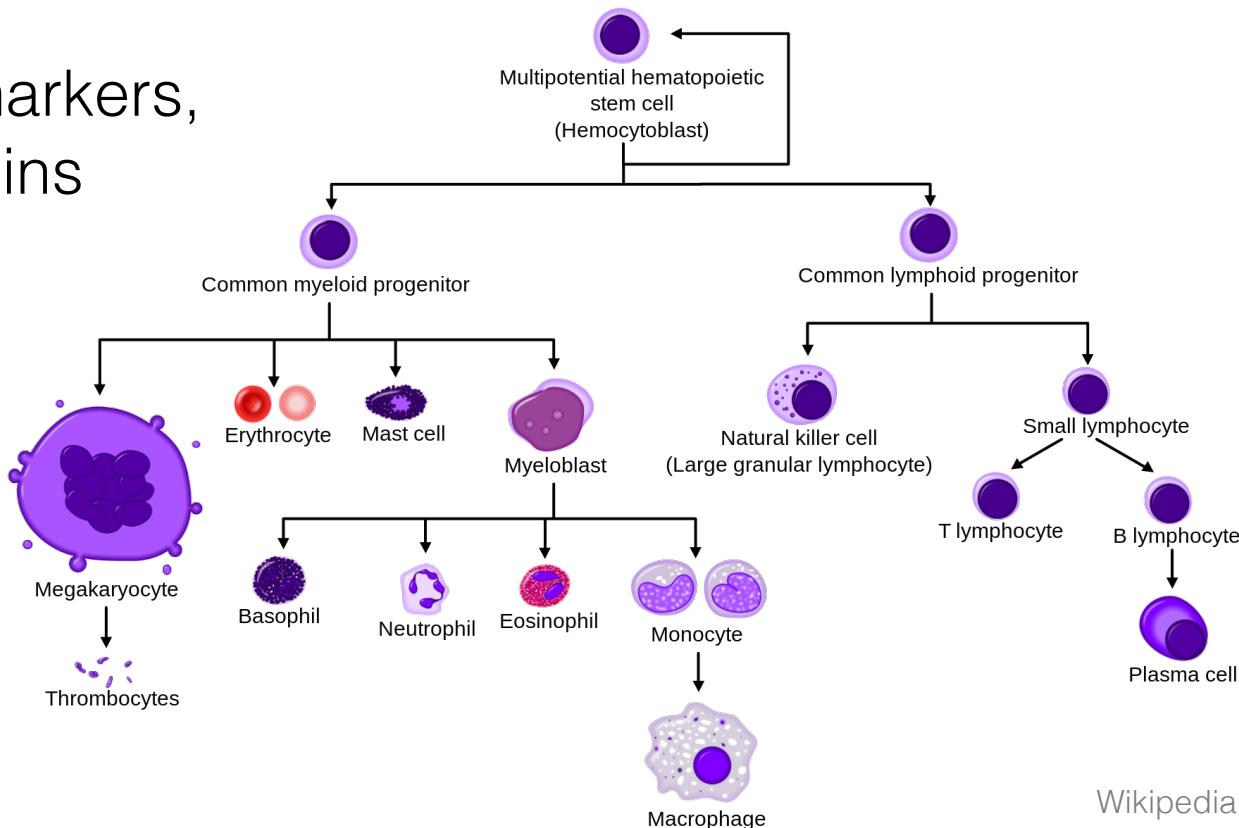
# Comparison of clustering methods for high-dimensional single-cell flow and mass cytometry data

C1omics 2015 conference  
Single-cell omics methods and applications  
Manchester, UK

Lukas Weber, 25 November 2015

# Background

- single cell analysis: cytometry
- protein expression levels to characterize cell populations
- surface protein markers, intracellular proteins
- flow cytometry, mass cytometry



# Applications

- Analysis of immune cell populations in biological samples
  - disease diagnostics, e.g. leukemia
  - vaccine research
- Differential analysis: comparisons between groups of samples
- Focus may be on major populations or rare cell types
  - e.g. analysis of proportions of major populations in a sample
  - or detection of a single rare cell type

# Technologies

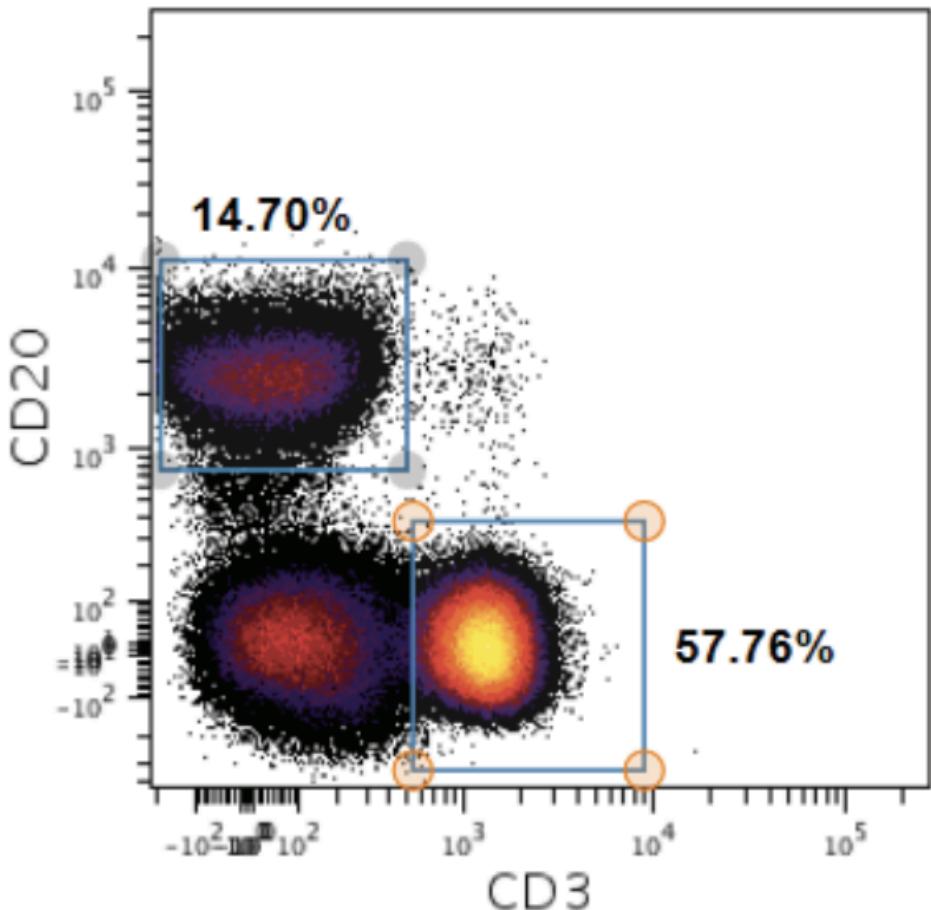
- **Flow cytometry:** 1000s cells/sec, 5-15 proteins/cell
- **Mass cytometry:** 100s cells/sec, 40 proteins/cell
- FCS file format: n cells, p protein expression levels

> head(data)

| n cells | p proteins (dimensions) |                |               |               |                |
|---------|-------------------------|----------------|---------------|---------------|----------------|
|         | CD45RA(La139)Di         | CD133(Pr141)Di | CD19(Nd142)Di | CD22(Nd143)Di | CD11b(Nd144)Di |
| [1,]    | 0.81704921              | -0.1479468     | -0.033481941  | 0.3321835     | -0.045922440   |
| [2,]    | 3.80138493              | -0.1914464     | -0.083273850  | 0.3723878     | 4.494378567    |
| [3,]    | 3.20443869              | -0.1611056     | 0.369612783   | -0.2149521    | -0.009404267   |
| [4,]    | 2.23738217              | -0.1380714     | -0.088311136  | -0.2204303    | 4.006597996    |
| [5,]    | -0.04404699             | -0.1515095     | 0.402548134   | 2.5817690     | 6.742060184    |
| [6,]    | 1.15033615              | -0.1475202     | -0.001792617  | -0.1497730    | 1.529571056    |

# Data analysis

- “gating”
  - visual inspection of series of 2D scatter plots
- problems
  - too many dimensions
  - multidimensional structure
  - unknown populations
  - subjective
- automated methods



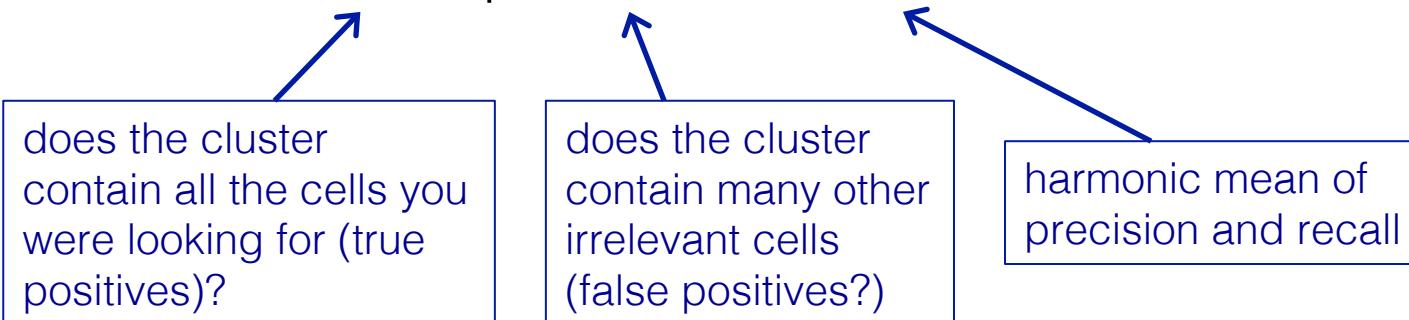
gating example:  
<http://www.cytobank.org>

# Our project

- Compare performance of clustering methods for automated detection of cell populations
  - specialized clustering methods for high-dimensional flow/mass cytometry data
  - active area of research, many new methods
  - ACCENSE, DensVM, FLOCK, flowMeans, FlowSOM, immunoClust, PhenoGraph, Rclusterpp, SamSPECTRAL, SWIFT, ...

# Our project

- 11 (13) methods
- 2 publicly available data sets, with manually gated cell populations as “truth”
  - Data set 1: 14 major populations
  - Data set 2: a single rare population of interest
- Metrics: recall, precision, F1 score



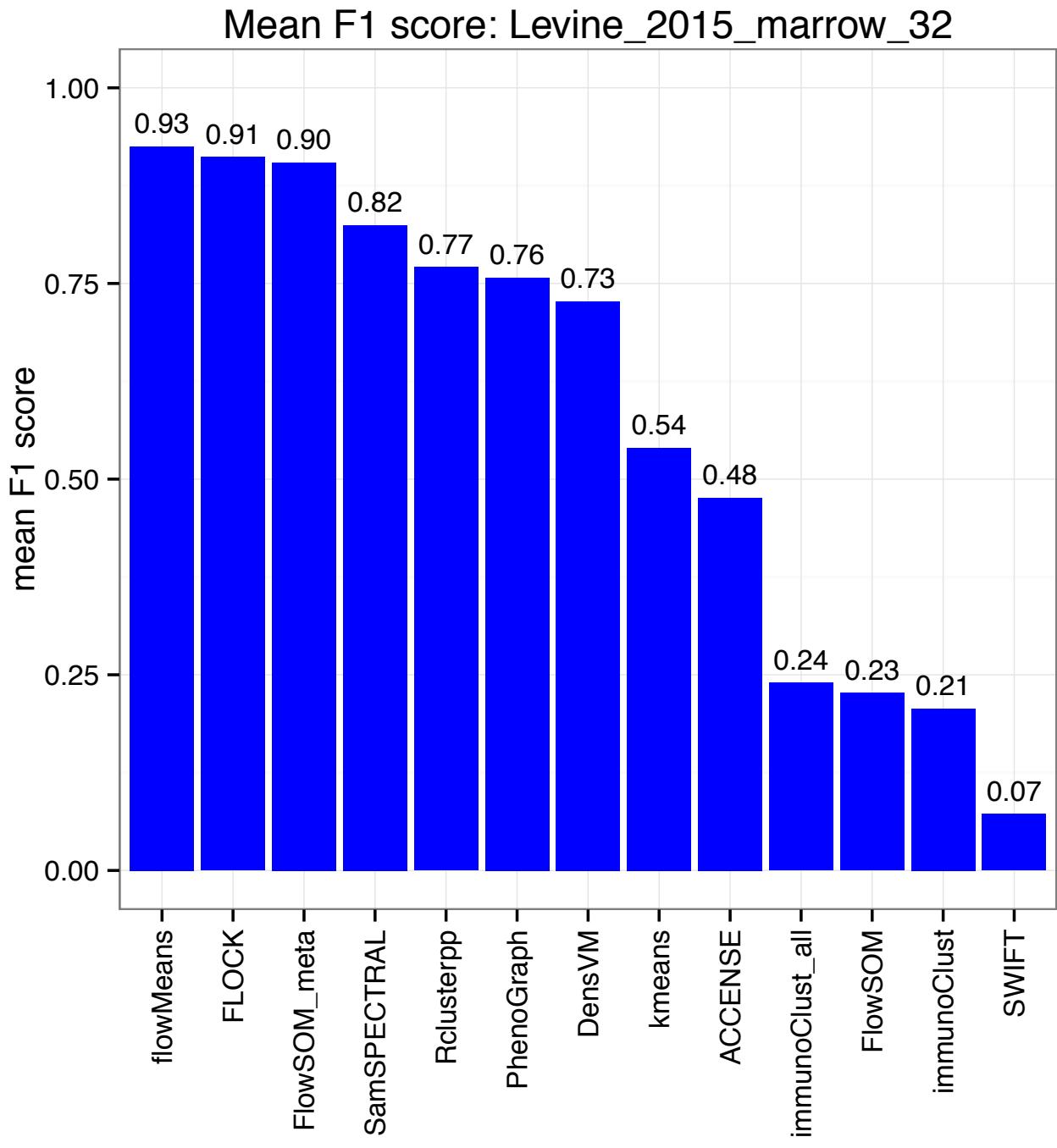
# Results: Data set 1

- Task: detection of 14 major immune cell populations
- Data set source:
  - Levine et al. (2015), *Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis*, Cell.
- 32-dimensional mass cytometry data set
  - healthy human bone marrow mononuclear cells (BMMCs) from 2 individuals
  - 104,184 cells (size  $n$ )
  - expression levels of 32 surface marker proteins (dimensionality  $p$ )
  - 14 manually gated major cell populations (truth)

# Results

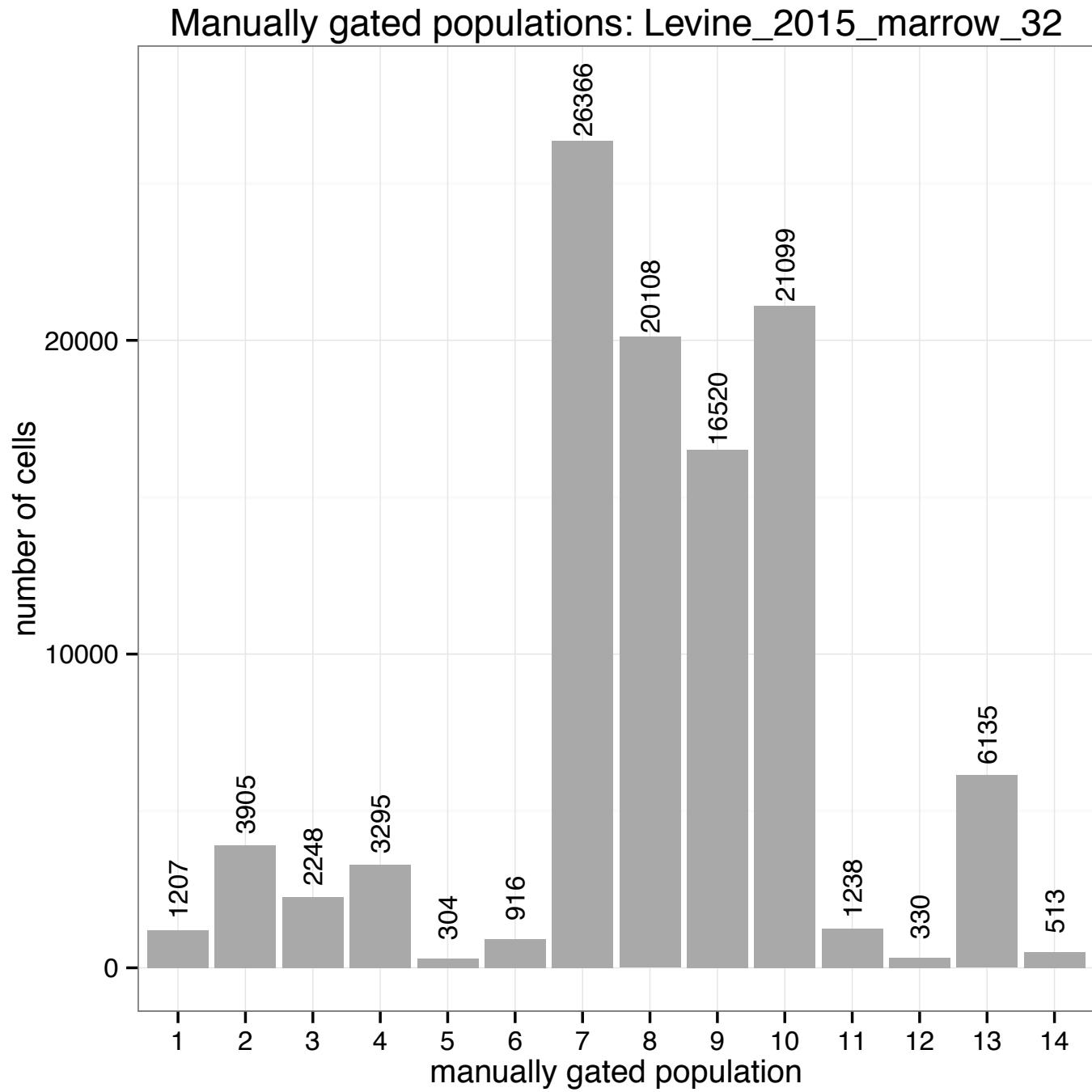
Mean F1 score

averaged over  
14 true clusters,  
weighted by  
number of cells



# Results

number of  
cells per true  
(manually  
gated) cluster



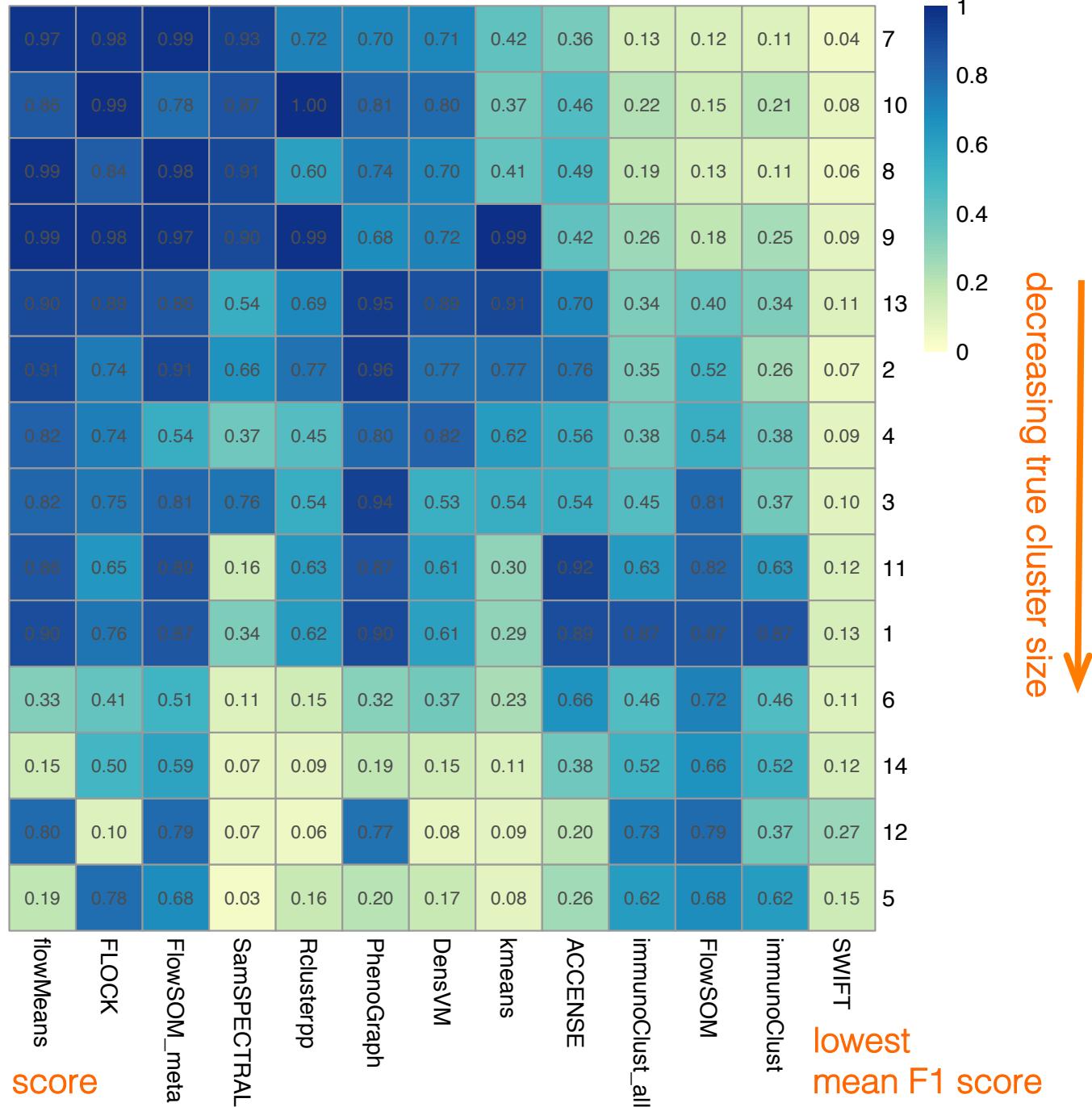
F1 score: Levine\_2015\_marrow\_32

# Results

F1 score by  
individual  
true cluster

clusters  
ranked by  
size

highest  
mean F1 score



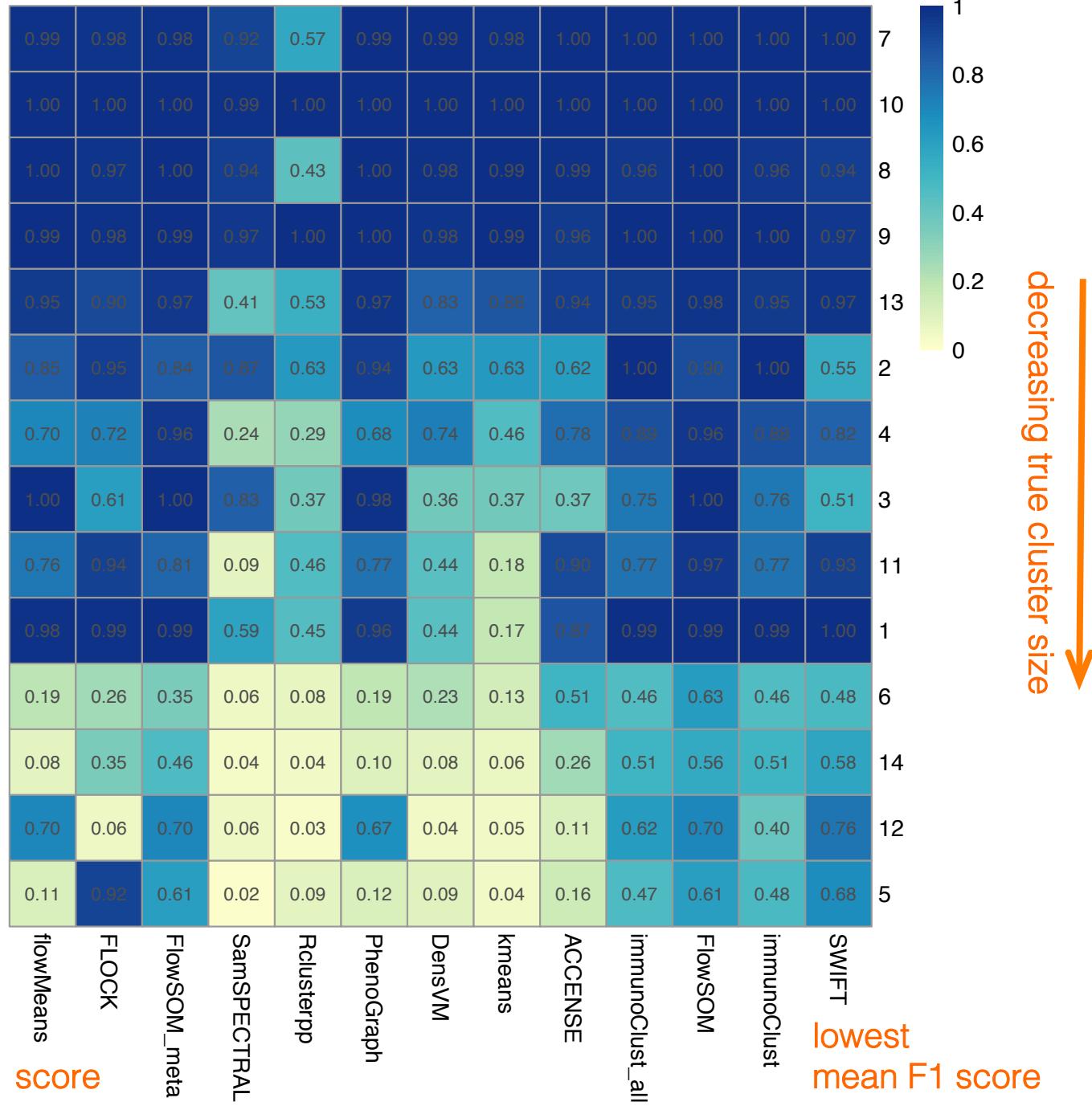
lowest  
mean F1 score

Precision: Levine\_2015\_marrow\_32

# Results

precision by  
individual  
true cluster

clusters  
ranked by  
size

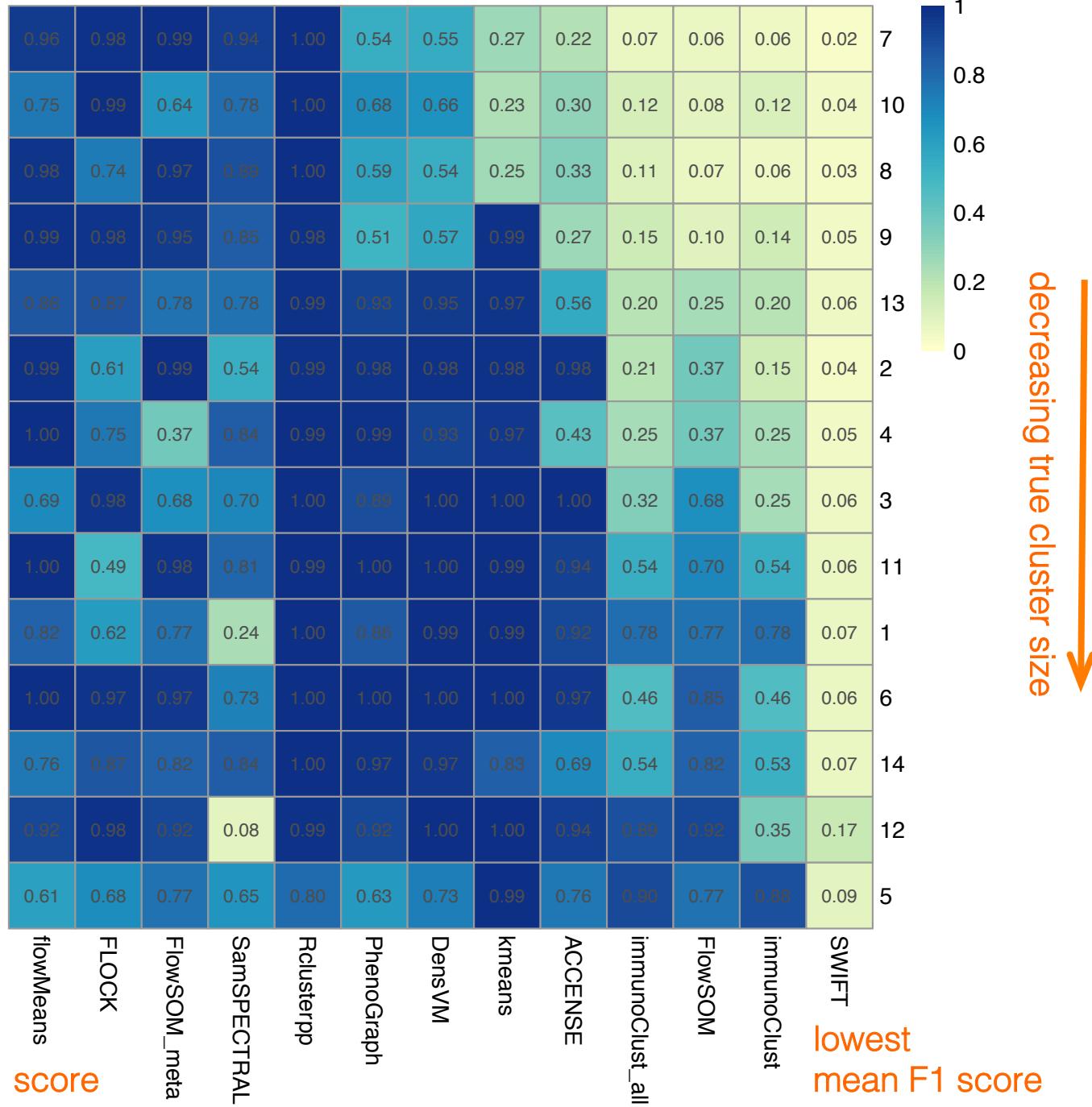


Recall: Levine\_2015\_marrow\_32

# Results

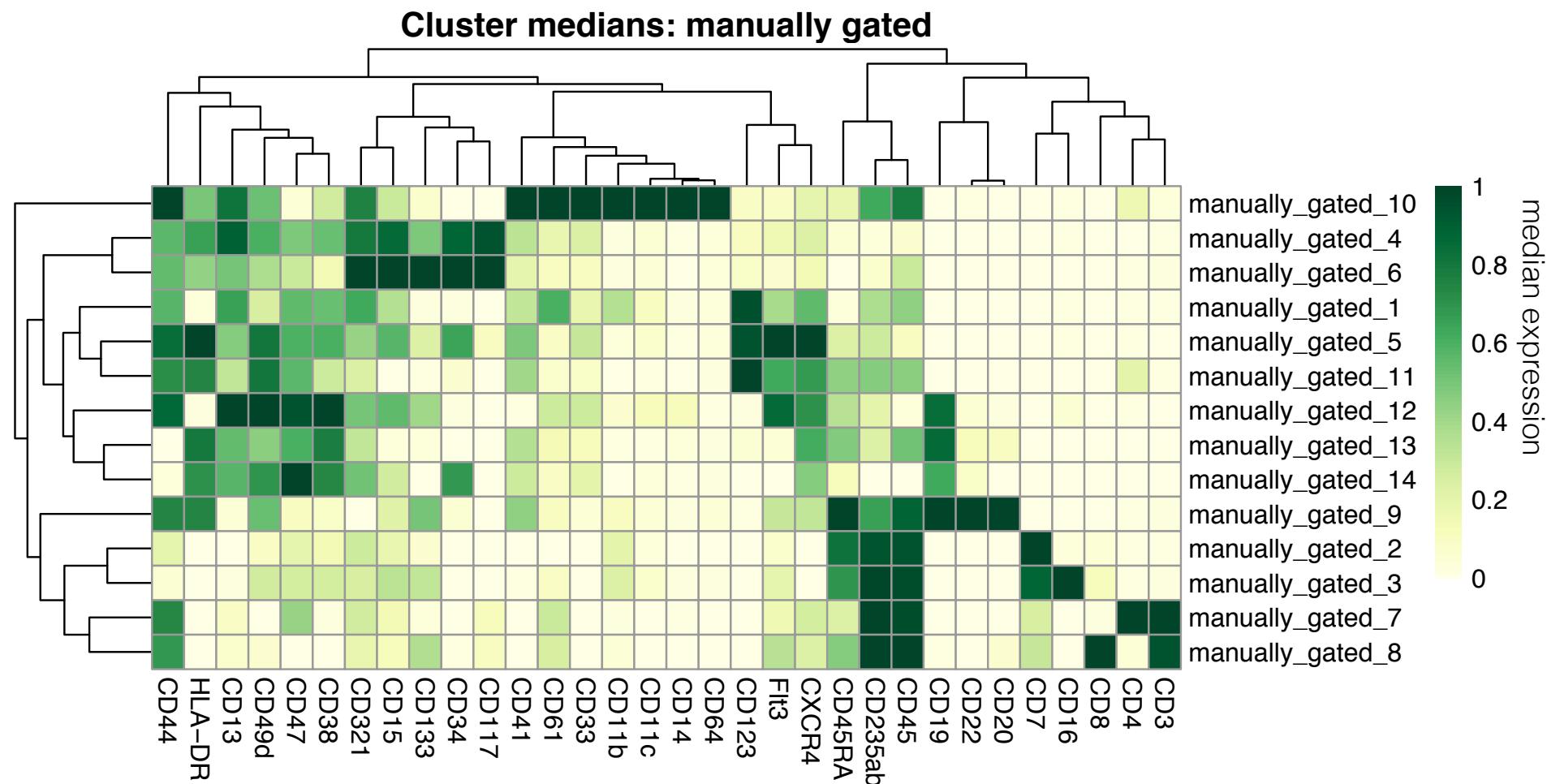
recall by individual true cluster

clusters ranked by size



# Results: cluster medians

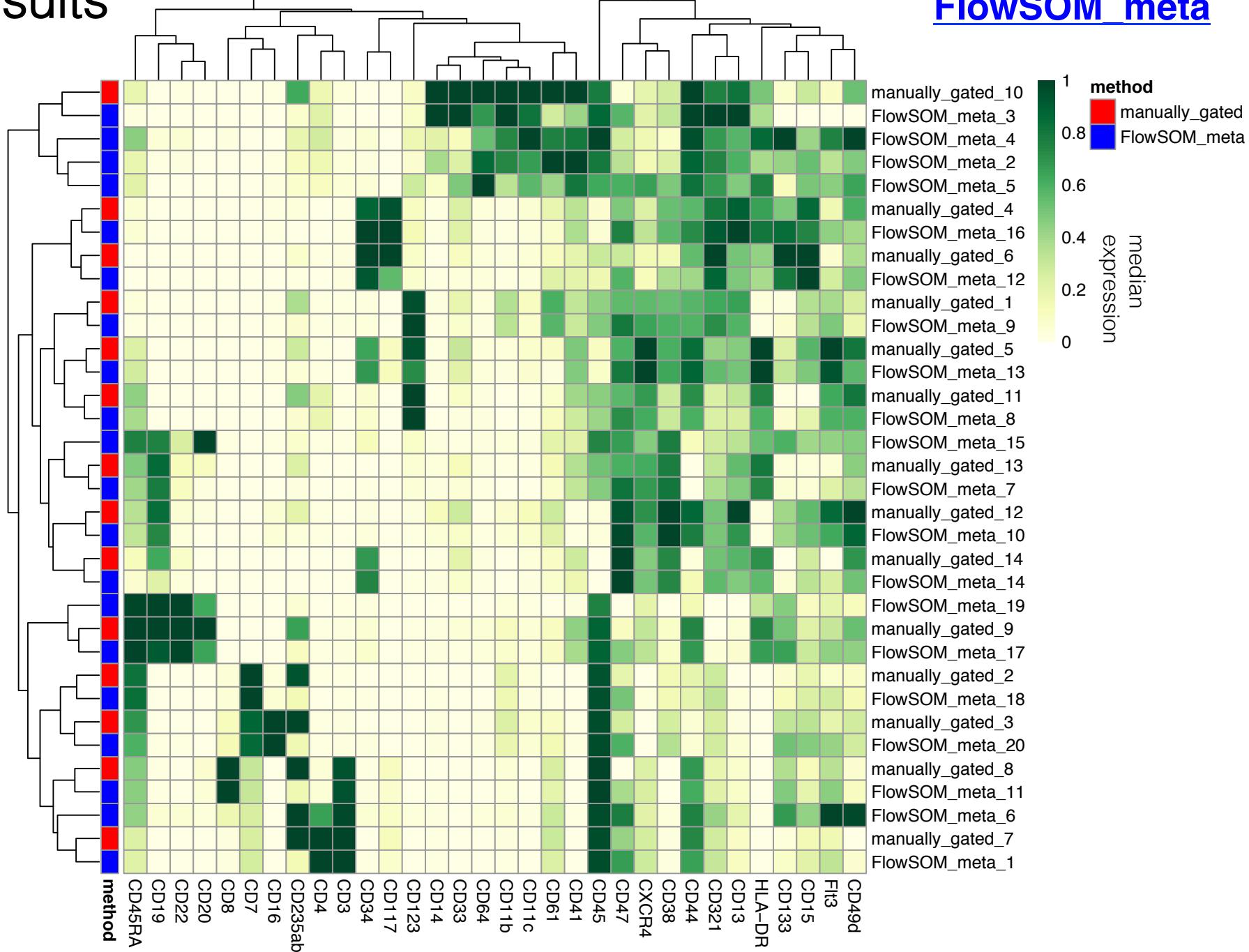
### manually gated (truth)



# Results

Cluster medians: manually gated and FlowSOM\_meta

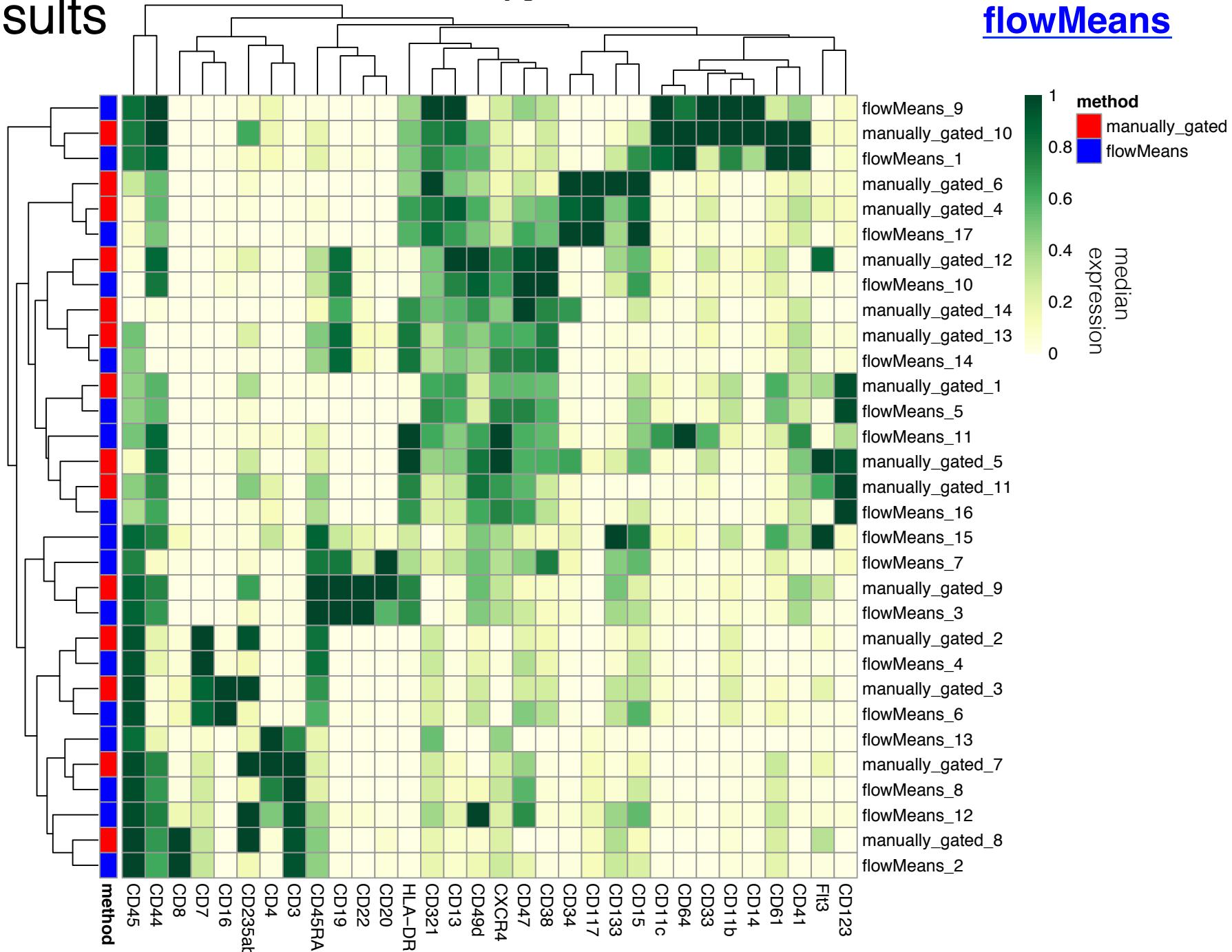
**FlowSOM meta**



# Results

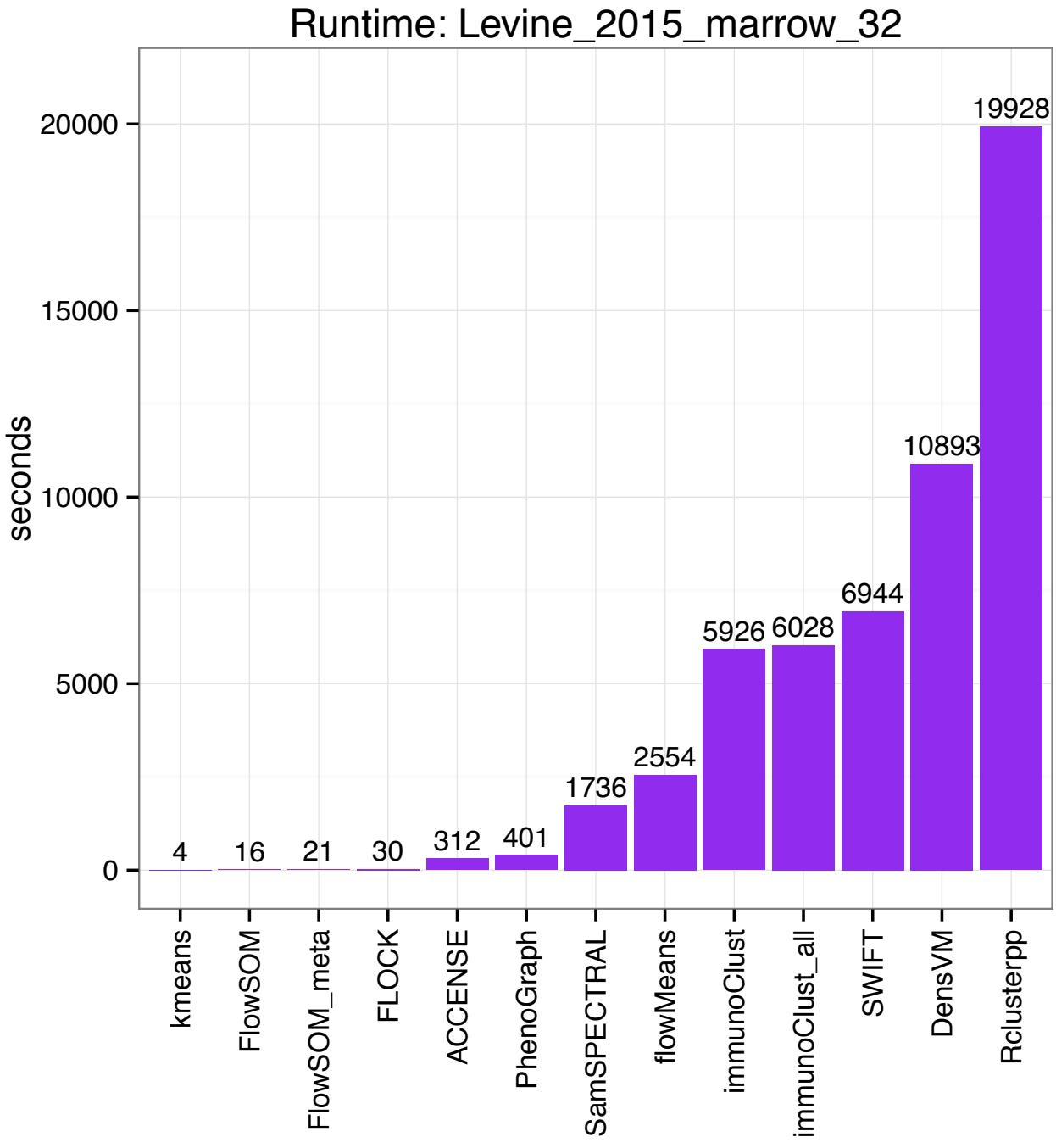
Cluster medians: manually gated and flowMeans

**flowMeans**



# Results

runtime in  
seconds



# Results: Data set 2

- Task: detection of a single rare cell population of interest
- Data set source:
  - Mosmann et al. (2014), *SWIFT – Scalable Clustering for Automated Identification of Rare Cell Populations in Large, High-Dimensional Flow Cytometry Datasets, Part 2: Biological Evaluation*, Cytometry Part A.
- 15-dimensional flow cytometry data set
  - healthy human peripheral blood mononuclear cells (PBMCs) exposed to influenza antigens
  - rare population of live, activated (cytokine-producing) memory CD4 T cells
  - 396,460 cells (size  $n$ ); 109 cells from rare population
  - expression levels of 15 proteins (dimensionality  $p$ )

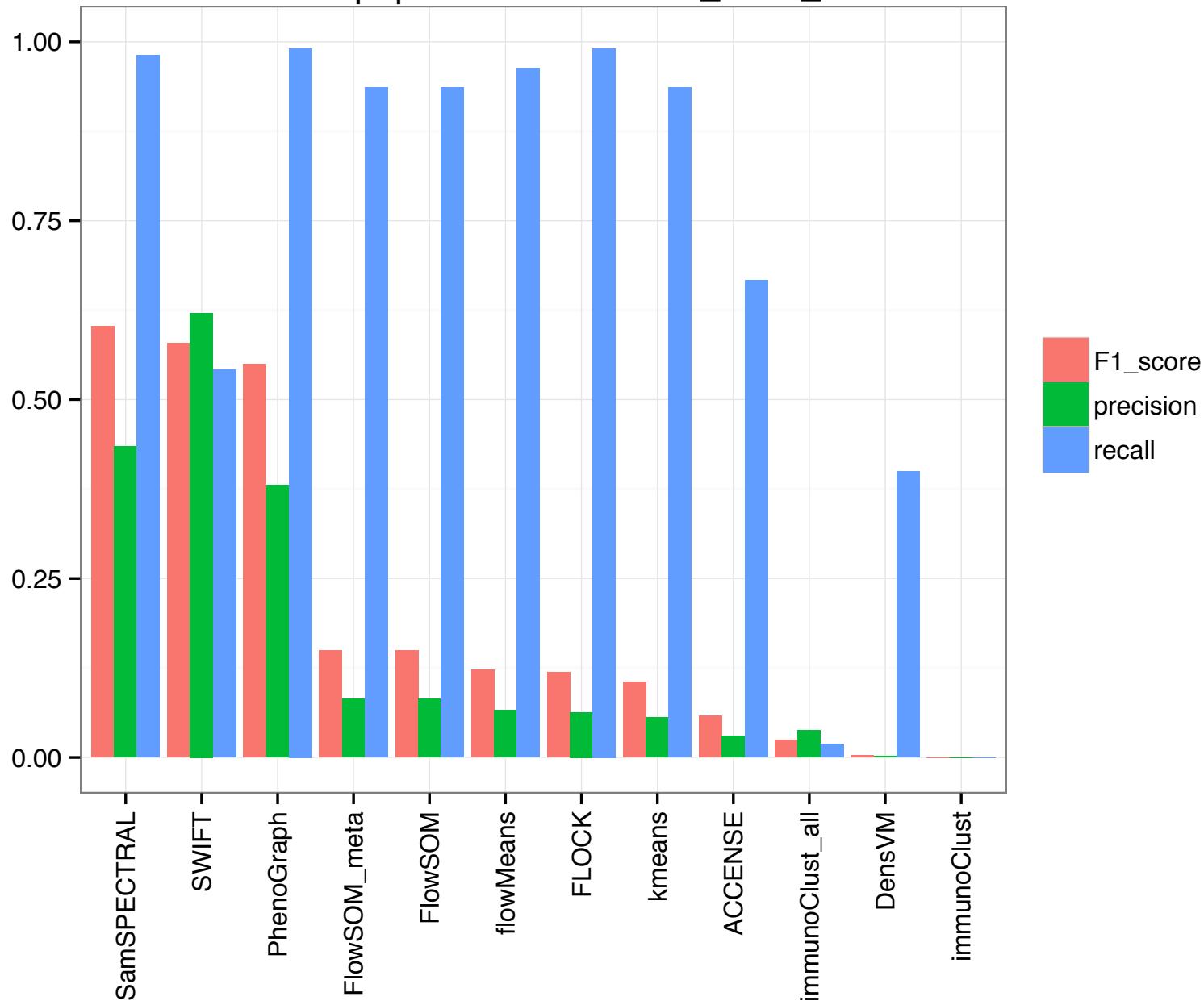
# Results

detection  
of rare cell  
population

F1 score,  
precision,  
recall

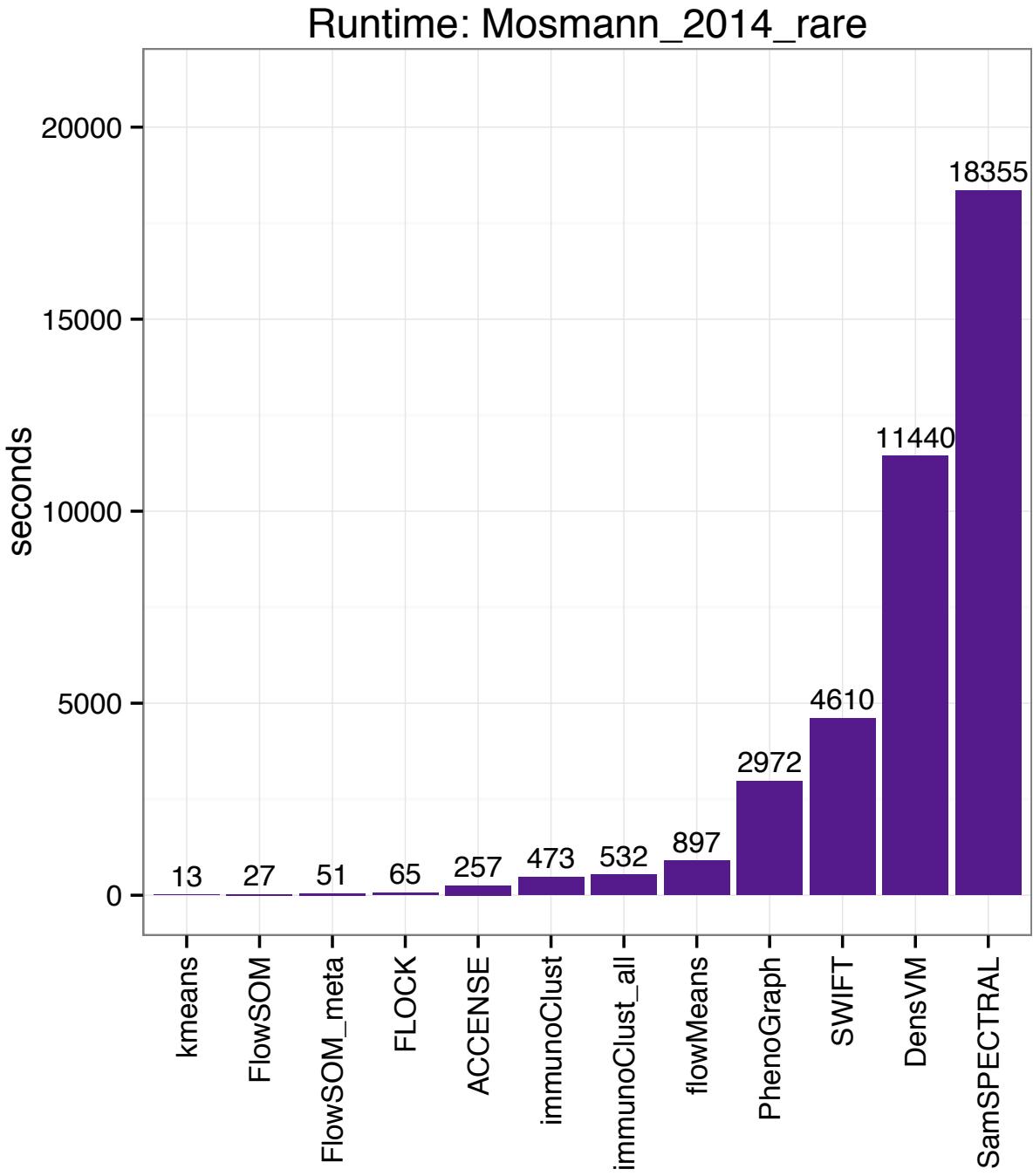
109 cells  
(out of  
396,460)

Rare cell population: Mosmann\_2014\_rare



# Results

runtime in  
seconds



# Conclusions

- Data set 1: 14 major cell populations
  - FlowSOM\_meta performs very well; fast runtime
- Data set 2: single rare population
  - PhenoGraph, SamSPECTRAL, SWIFT
- Runtime
- Next steps: more data sets; robustness

# Acknowledgments

Mark Robinson

Robinson lab (UZH)

Alexander Roth (UZH)

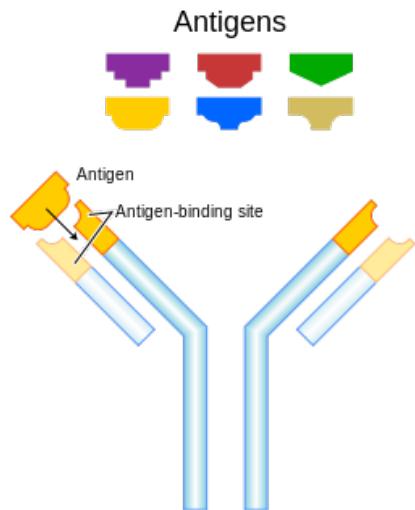
Michael Baudis (UZH)

RADIANT

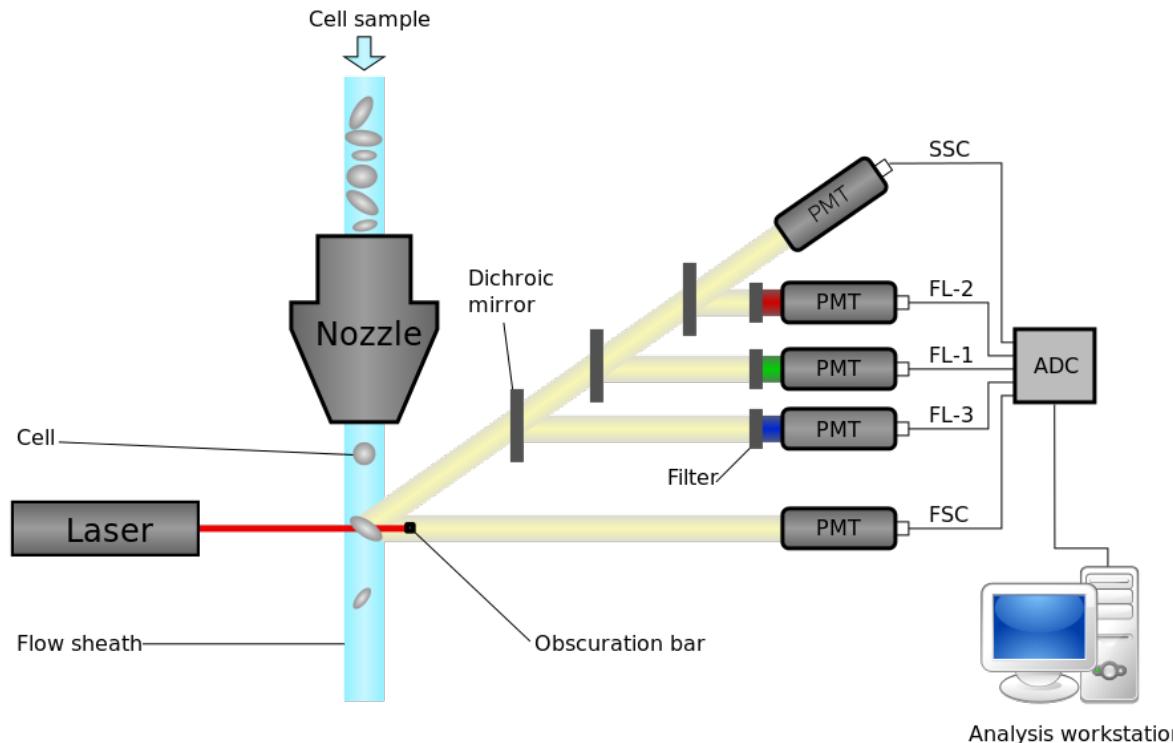
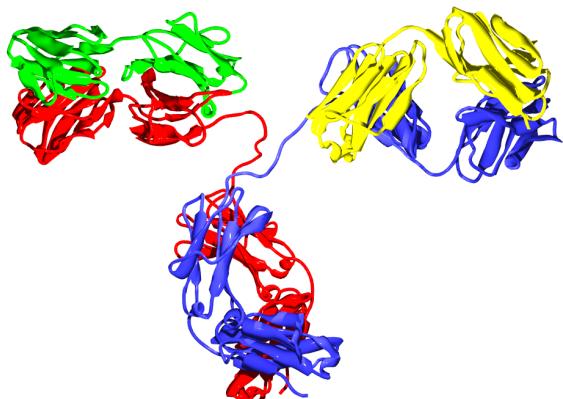


# Additional slides

# Flow cytometry (FACS)

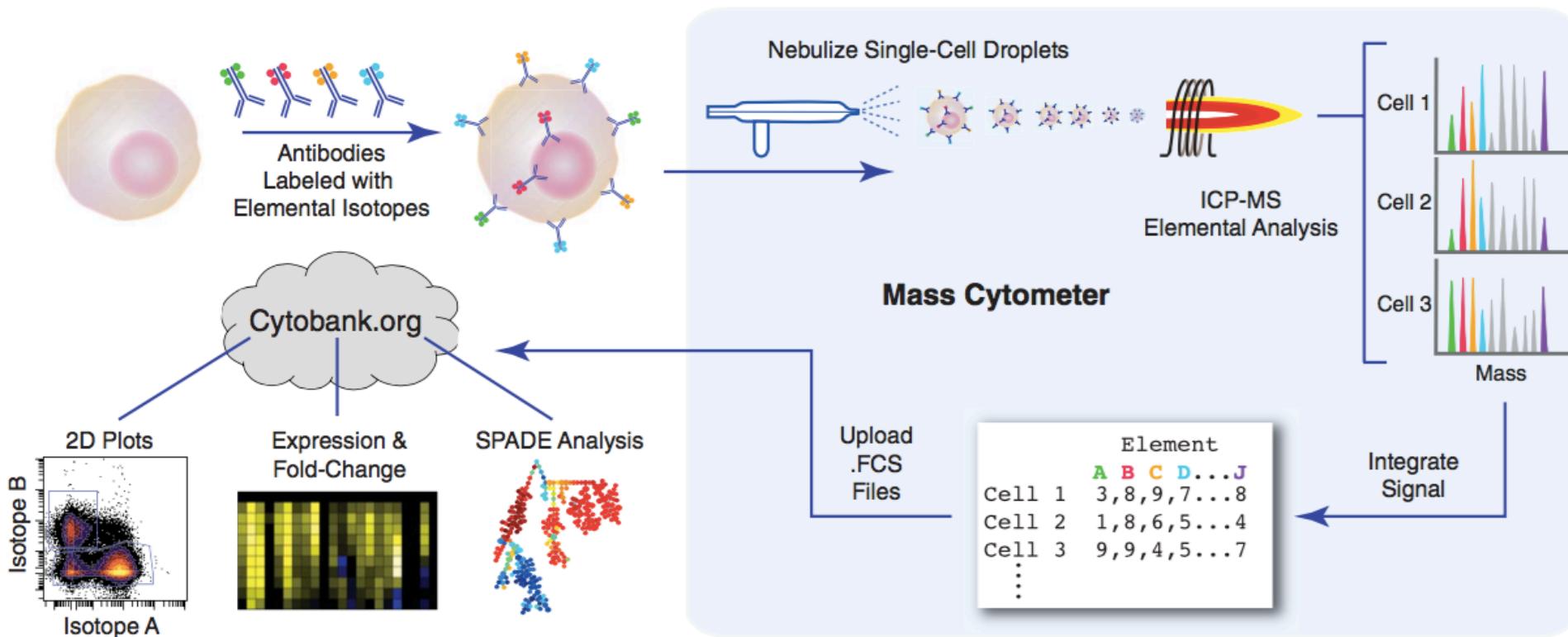


Antibodies



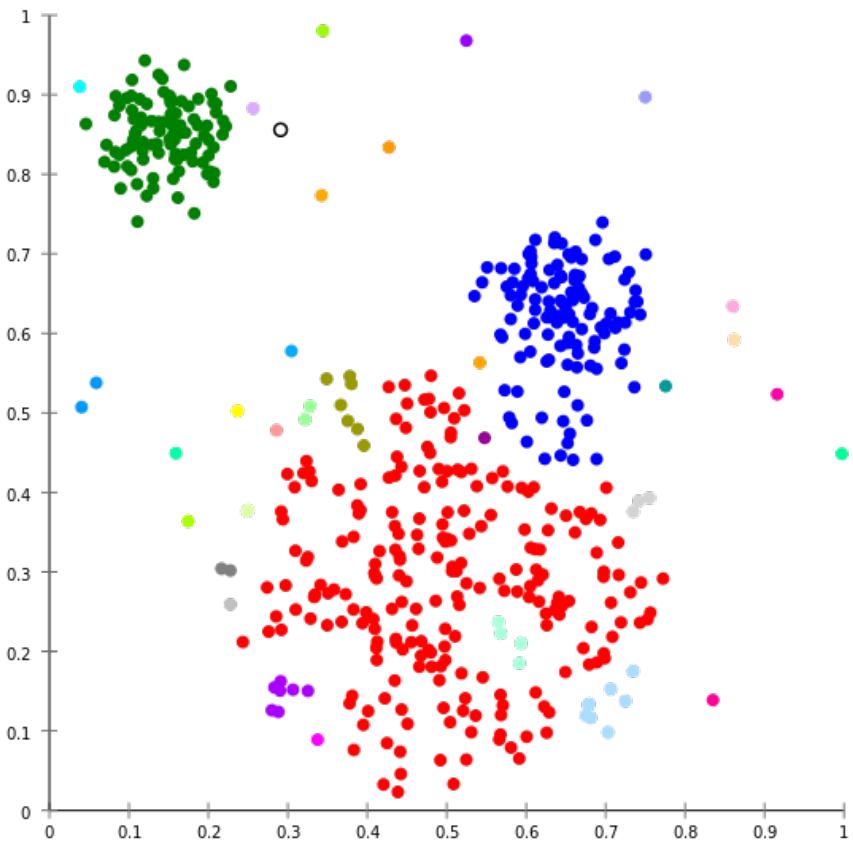
Flow cytometry schematic

# Mass cytometry

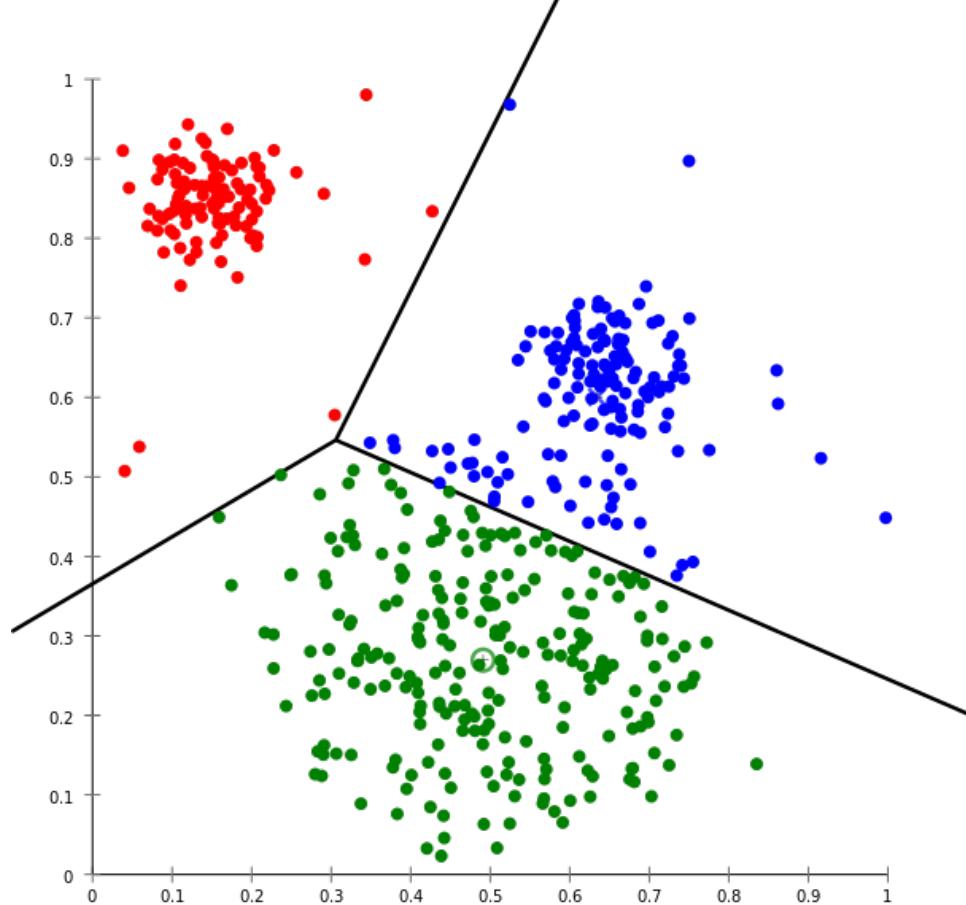


Bendall et al. (2011), Fig. 1A

# Cluster analysis



hierarchical clustering



k-means

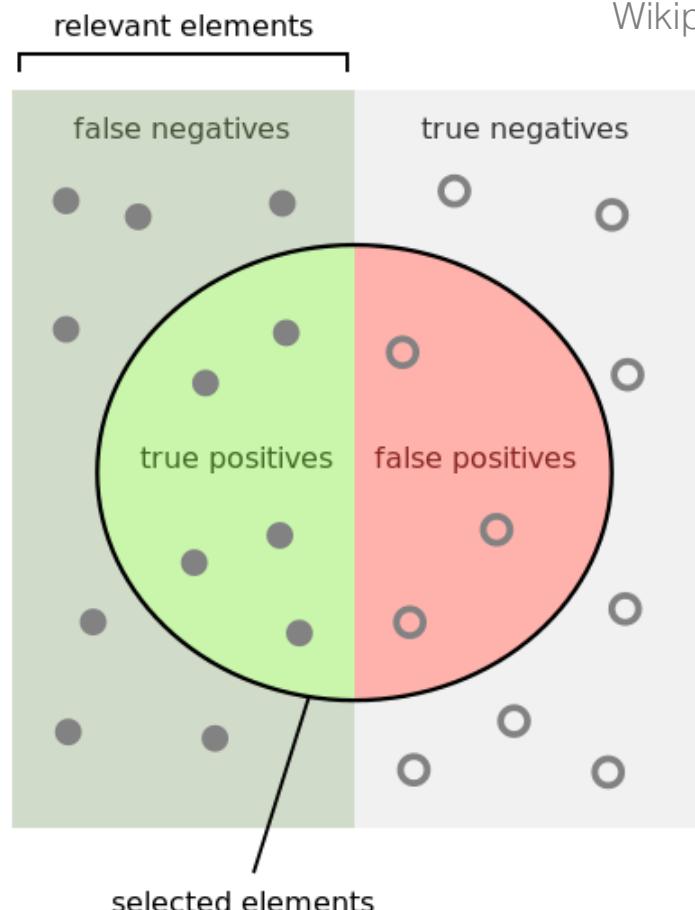
# F1 score

harmonic mean of precision and recall

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

precision = positive predictive value

recall = sensitivity = true positive rate (TPR)



How many selected items are relevant?

$$\text{Precision} = \frac{\text{green}}{\text{red} + \text{green}}$$

How many relevant items are selected?

$$\text{Recall} = \frac{\text{green}}{\text{green} + \text{grey}}$$

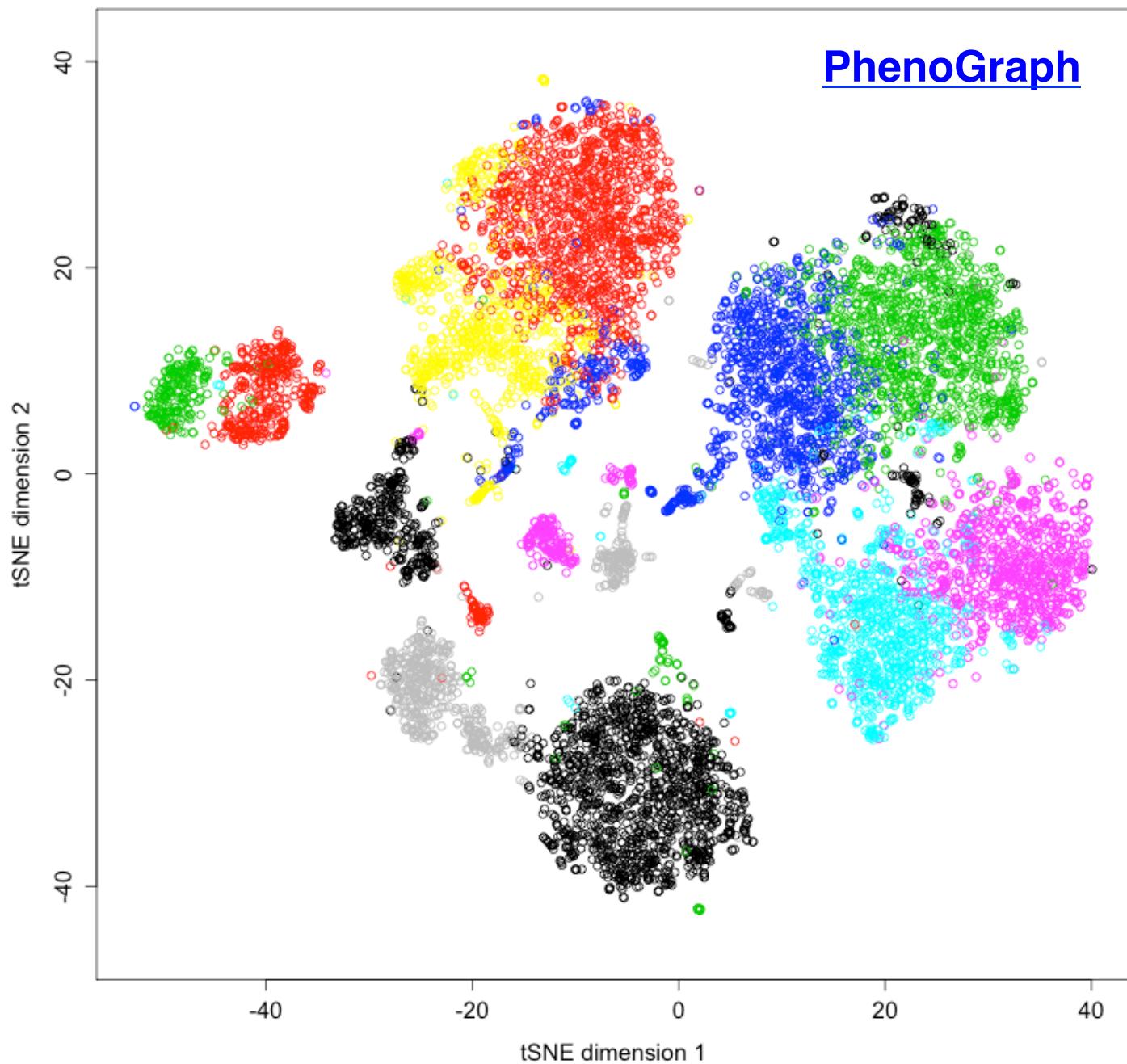
2D tSNE plot: PhenoGraph cluster labels

# Results

visualization  
with tSNE /  
viSNE

nonlinear  
2D  
projection

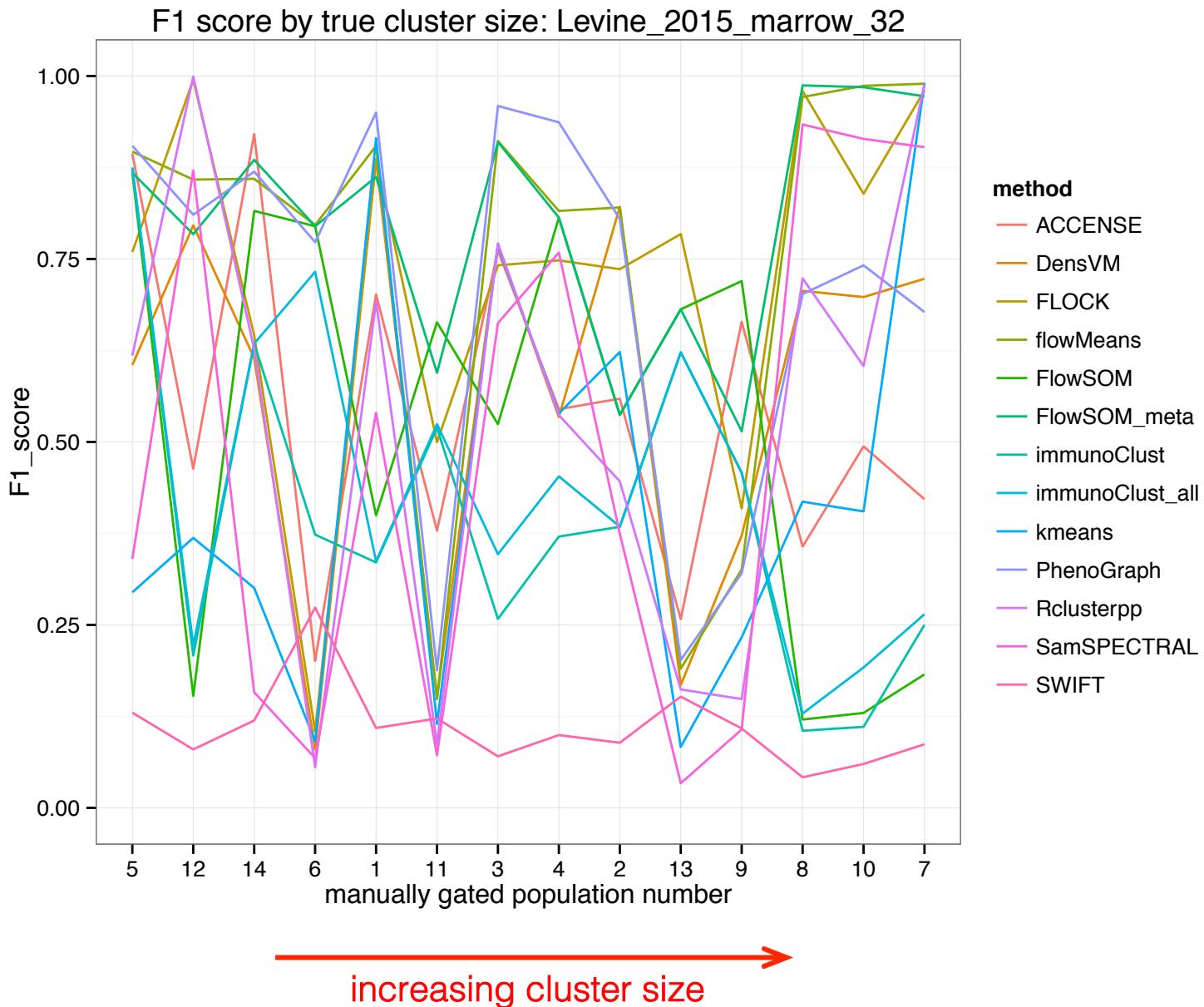
PhenoGraph



# Results

F1 score by  
individual  
true cluster

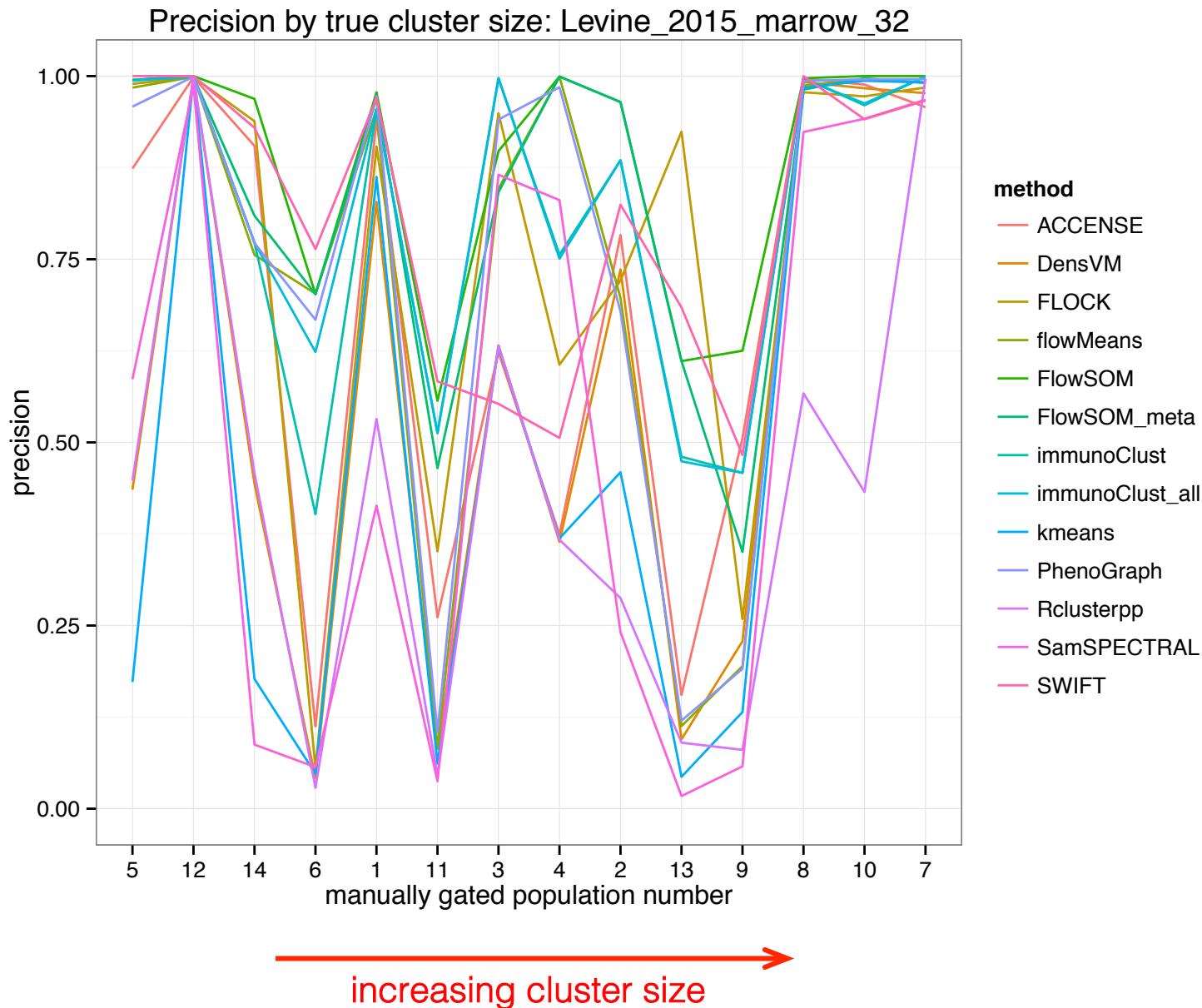
clusters  
ranked by  
size



# Results

precision by  
individual  
true cluster

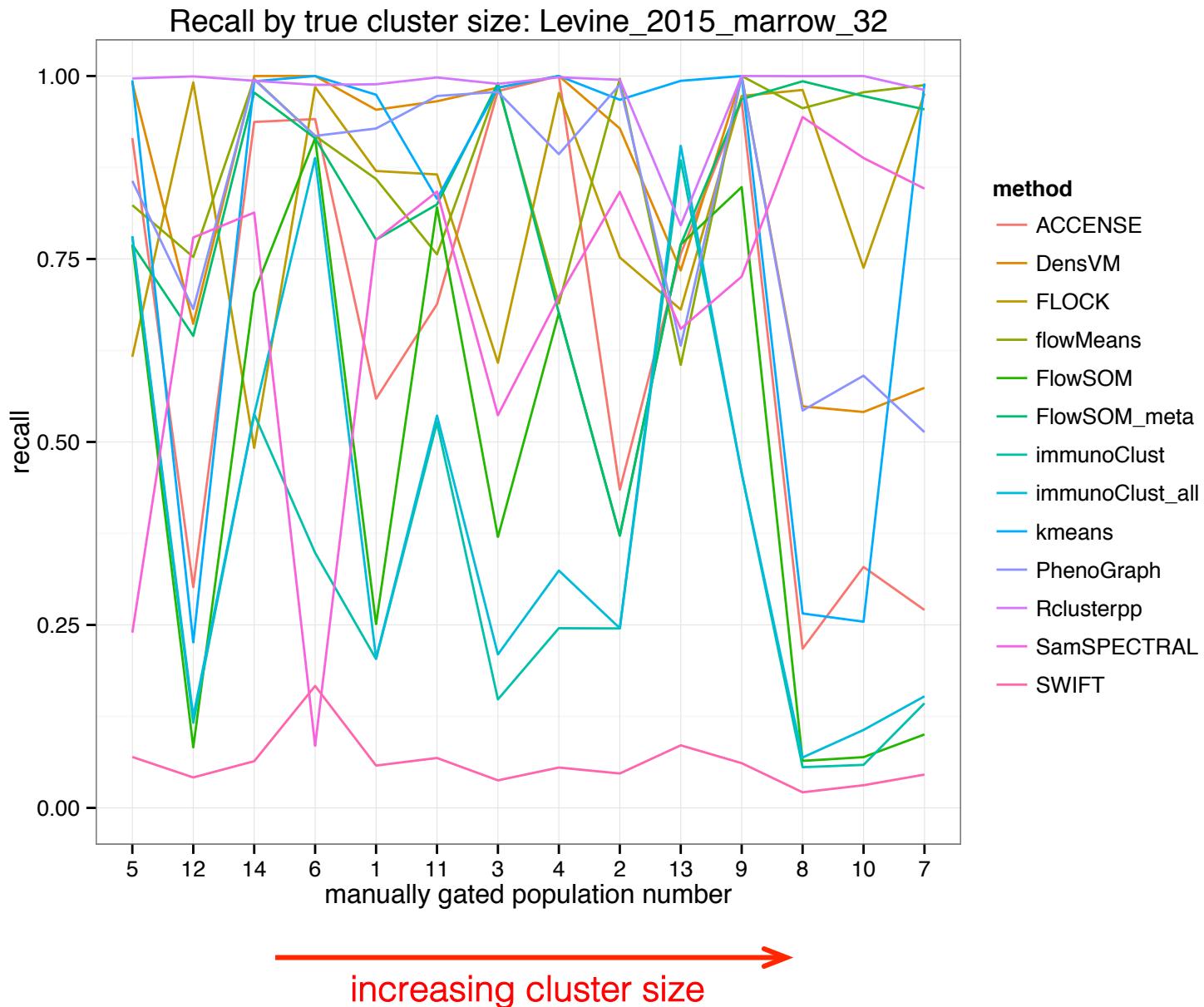
# clusters ranked by size



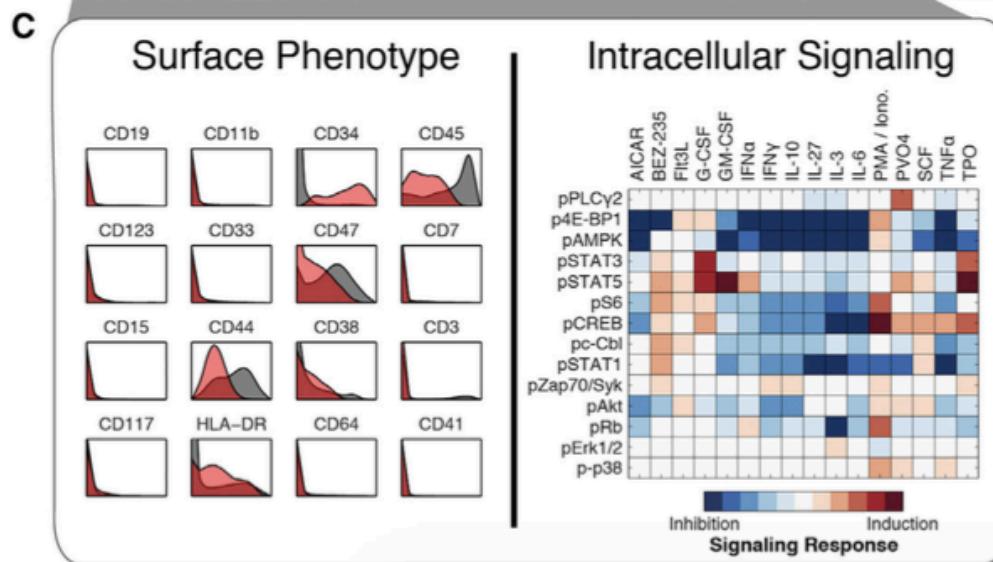
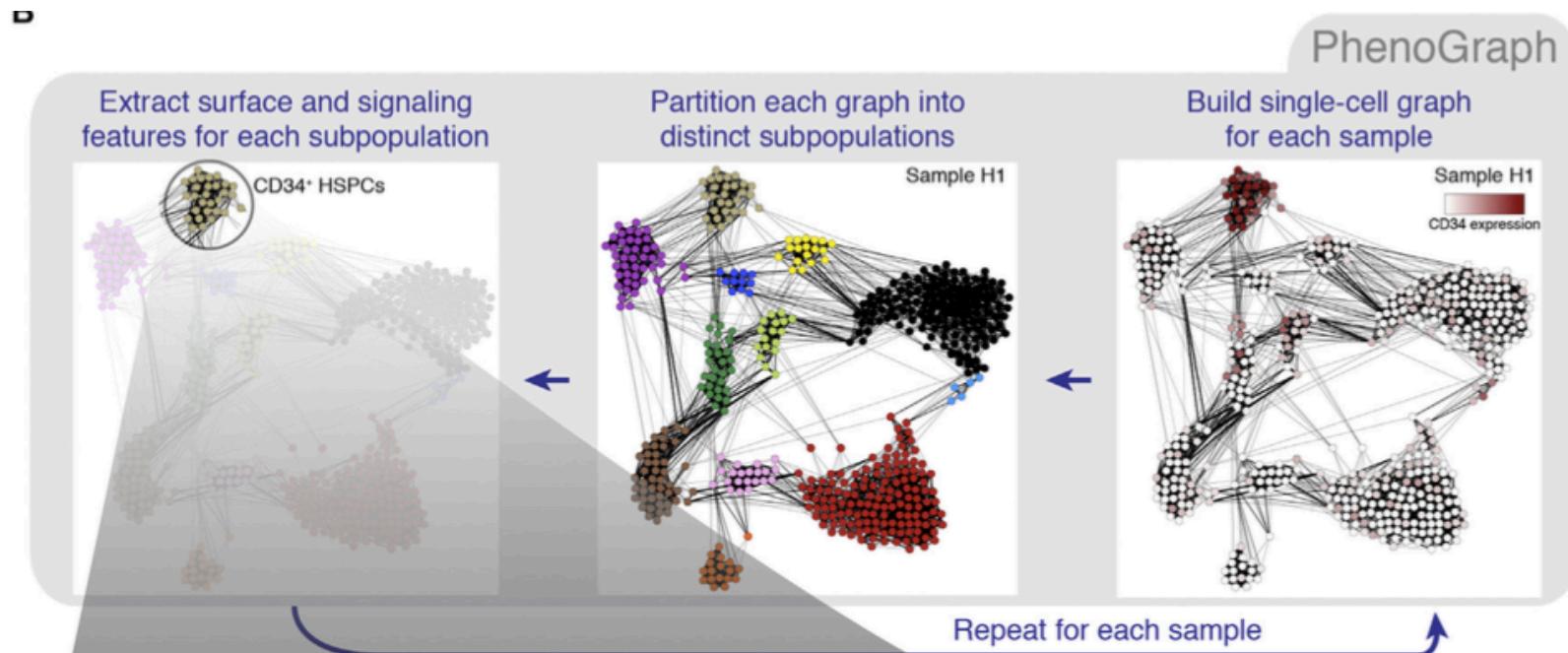
## Results

recall by  
individual  
true cluster

clusters  
ranked by  
size



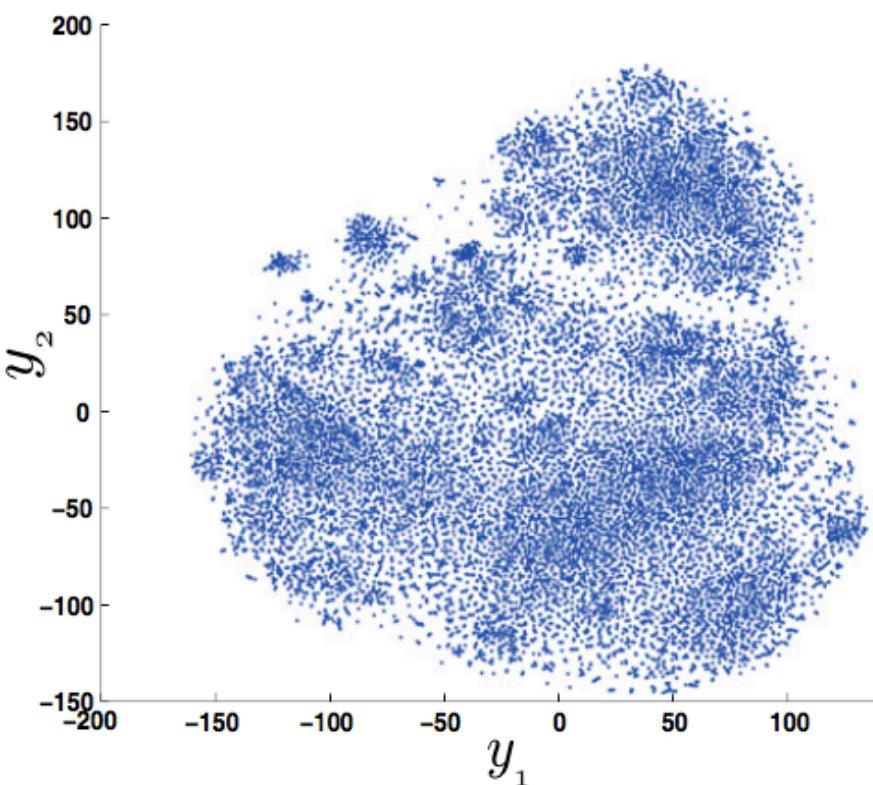
# PhenoGraph



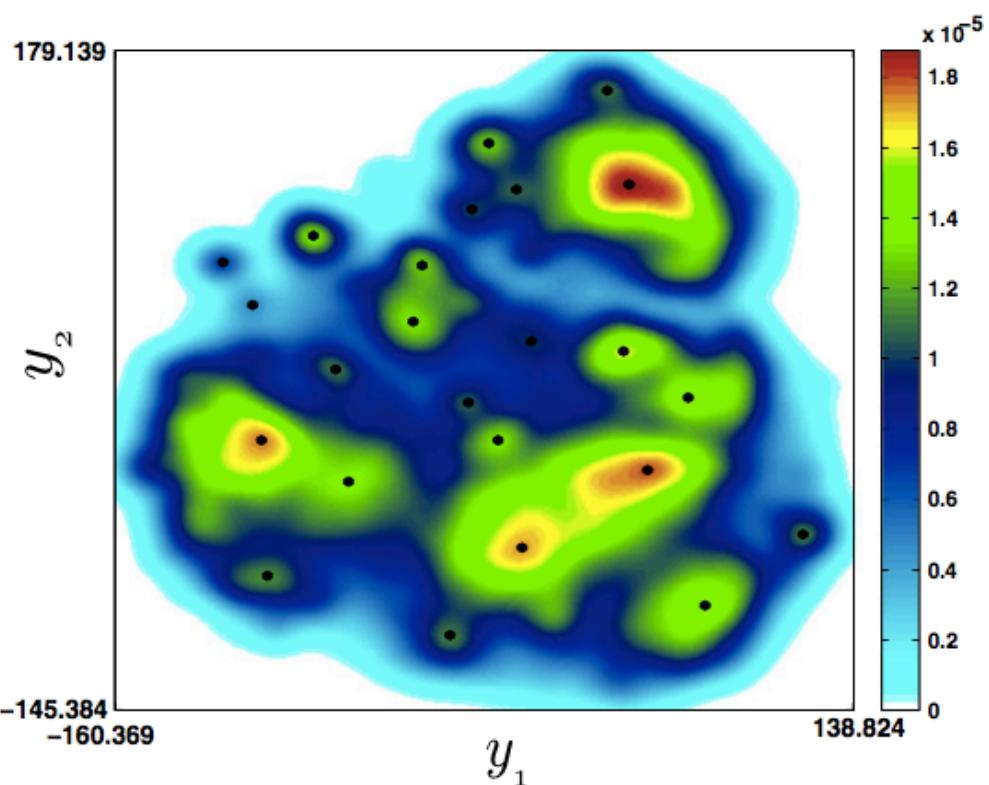
Levine et al. (2015), Fig. 1

# ACCENSE

C



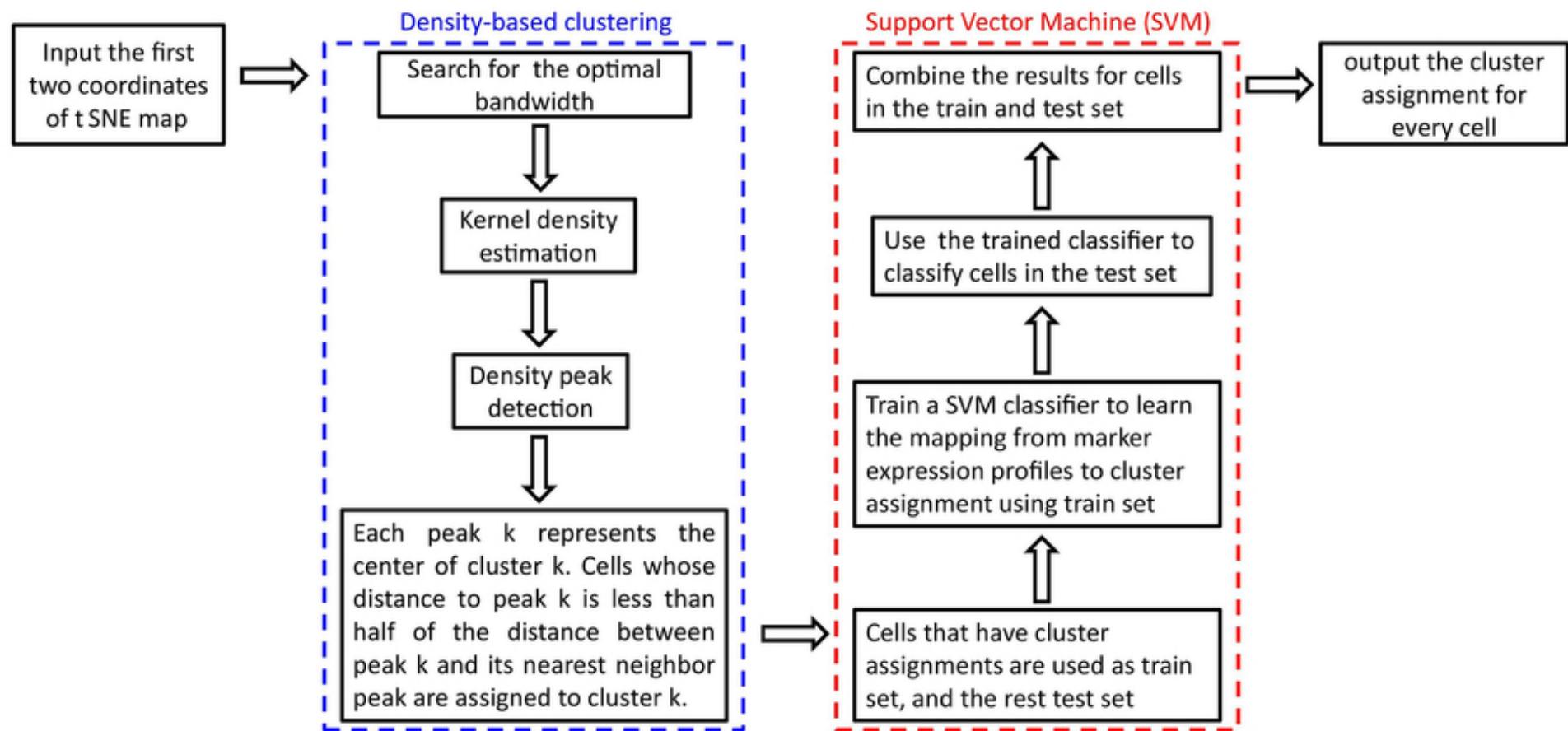
D



Shekhar et al. (2014), Fig. 1

# DensVM

DensVM flowchart



Becher et al. (2014), Supp. Fig. 2

# FlowSOM

van Gassen et al. (2015), Fig. 1

